

Tackling Documentation Debt: A Survey on Algorithmic Fairness Datasets

Anonymous Authors

ABSTRACT

Data-driven algorithms are studied and deployed in diverse domains to support critical decisions, directly impacting people's well-being. As a result, a growing community of researchers has been investigating the equity of existing algorithms and proposing novel ones, advancing the understanding of risks and opportunities of automated decision-making for historically disadvantaged populations. Progress in fair Machine Learning (ML) and equitable algorithm design hinges on data, which can be appropriately used only if adequately documented. Unfortunately, the research community, as a whole, suffers from a collective data documentation debt caused by a lack of information on specific resources (*opacity*) and scatteredness of available information (*sparsity*). In this work, we survey over two hundred datasets employed in algorithmic fairness research, producing standardized and searchable documentation for each of them. Moreover we rigorously identify the three most popular fairness datasets, namely Adult, COMPAS, and German Credit, for which we compile in-depth documentation. This unifying documentation effort targets *documentation sparsity* and supports multiple contributions. In the first part of this work, we summarize the merits and limitations of Adult, COMPAS, and German Credit, adding to and unifying recent scholarship, calling into question their suitability as general-purpose fairness benchmarks. To overcome this limitation, we document hundreds of available alternatives, annotating their domain and the algorithmic fairness tasks they support, along with additional properties of interest for fairness practitioners and researchers, including their format, cardinality, and the sensitive attributes they encode. In the second part, we summarize this information, zooming in on the tasks, domains, and roles of these resources. Overall, we assemble and summarize sparse information on hundreds of datasets into a single resource, which we make available to the community, with the aim of tackling the data documentation debt.

KEYWORDS

Algorithmic fairness, Data studies, Documentation debt.

ACM Reference Format:

Anonymous Authors. 2022. Tackling Documentation Debt: A Survey on Algorithmic Fairness Datasets. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 129 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, Dates, Venue

© 2022 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

Following the widespread study and application of data-driven algorithms in contexts that are central to people's well-being, a large community of researchers has coalesced around the growing field of algorithmic fairness and equity, investigating algorithms through the lens of justice, bias, power, harms, and equality. A line of work gaining traction in the field, intersecting with critical data studies, human-computer interaction, and computer-supported cooperative work, focuses on data transparency and standardized documentation processes to describe key characteristics of datasets [37, 182, 183, 227, 250, 365]. Most prominently, Gebru et al. [182] and Holland et al. [227] proposed two complementary documentation frameworks, called *Datasheets for Datasets* and *Dataset Nutrition Labels*, to improve data curation practices and favour more informed data selection and utilization for dataset users. Overall, this line of work has contributed to an unprecedented attention to dataset documentation in ML, including a novel track focused on datasets at the Conference on Neural Information Processing Systems (NeurIPS), an initiative to support dataset tracking in repositories for scholarly articles,¹ and dedicated works producing retrospective documentation for existing datasets [27, 178], auditing their properties [406] and tracing their usage [397].

In recent work, Bender et al. [38] propose the notion of *documentation debt*, in relation to training sets that are undocumented and too large to document retrospectively. We extend this definition to the collection of datasets employed in a given field of research. We see two components at work contributing to the documentation debt of a research community. On one hand, *opacity* is the result of poor documentation affecting single datasets, contributing to misunderstandings and misuse of specific resources. On the other hand, when relevant information exists but does not reach interested parties, there is a problem of documentation *sparsity*. One example that is particularly relevant for the algorithmic fairness community is represented by the German Credit dataset [501], a popular resource in this field. Many works of algorithmic fairness, including recent ones, carry out experiments on this dataset using sex as a protected attribute [23, 221, 328, 343, 398, 456, 528, 551], while existing yet overlooked documentation shows that this feature cannot be reliably retrieved [204]. Moreover, the mere fact that a dataset exists and is relevant to a given task or a given domain may be unknown. The BUPT Faces datasets, for instance, were presented as the second existing resource for face analysis with race annotations [529]. However several resources were already available at the time, including Labeled Faces in the Wild [211], UTK Face [578], Racial Faces in the Wild [530], and Diversity in Faces [361].²

¹<https://medium.com/paperswithcode/datasets-on-arxiv-1a5a8f7bd104>

²Hereafter, for brevity, we only report dataset names. The relevant references and additional information can be found in Appendix A.

To tackle the documentation debt of the algorithmic fairness community, we survey the datasets used in over 500 articles on fair ML and equitable algorithms, presented at seven major conferences, considering each edition in the period 2014–2021, and more than twenty domain-specific workshops in the same period. We find over 200 datasets employed in studies of algorithmic fairness, for which we produce compact and standardized documentation, called *data briefs*. Data briefs are intended as a lightweight format to document fundamental properties of data artifacts used in algorithmic fairness, including their purpose, their features, with particular attention to sensitive ones, the underlying labeling procedure, and the envisioned ML task, if any. To favor domain-based and task-based search from dataset users, data briefs also indicate the domain of the processes that produced the data (e.g., radiology) and list the fairness tasks studied on a given dataset (e.g. fair ranking). For this endeavour, we have contacted creators and knowledgeable practitioners identified as primary points of contact for the datasets. We received feedback (incorporated into the final version of the data briefs) from 77 curators and practitioners, whose contribution is acknowledged at the end of this article. Moreover, we identify and carefully analyze the three datasets most often utilized in the surveyed articles (Adult, COMPAS, and German Credit), retrospectively producing a datasheet and a nutrition label for each of them. From these documentation efforts, we extract a summary of the merits and limitations of popular algorithmic fairness benchmarks, and a categorization of alternative resources with respect to domains, tasks and roles they play in works of algorithmic fairness. Overall, we make the following contributions.

- **Unified analysis of popular fairness benchmarks.** We produce *datasheets* and *nutrition labels* for Adult, COMPAS, and German Credit, from which we extract a summary of their merits and limitations. We add to and unify recent scholarship on these datasets, calling into question their suitability as general-purpose fairness benchmarks due to contrived prediction tasks, noisy data, severe coding mistakes, limitations in encoding sensitive attributes, and age.
- **Survey of existing alternatives.** We compile standardized and compact documentation for over two hundred resources used in fair ML research, annotating their domain, the tasks they support, and the roles they play in works of algorithmic fairness. By assembling sparse information on hundreds of datasets into a single document, we aim to support multiple goals by researchers and practitioners, including domain-oriented and task-oriented search by dataset users. Contextually, we provide a novel taxonomy of tasks and domains investigated in algorithmic fairness research (summarized in Tables 2 and 3).

The rest of this work is organized as follows. Section 2 introduces related works. Section 3 presents the methodology and inclusion criteria of this survey. Section 4 analyzes the perks and limitations of the most popular datasets, namely Adult (§ 4.1), COMPAS (§ 4.2) and German Credit (§ 4.3), and provides an overall summary of their merits and limitations as fairness benchmarks (§ 4.4). Section 5 discusses alternative fairness resources from the perspective of the underlying domains (§ 5.1), the fair ML tasks they support (§ 5.2) and the roles they play (§ 5.3). Finally, Section 6 contains concluding

remarks and details the broader importance of this work for the research community. Interested readers may find the data briefs in Appendix A, followed by the detailed documentation produced for Adult (Appendix B), COMPAS (Appendix C) and German Credit (Appendix D).

2 RELATED WORK

2.1 Algorithmic fairness surveys

Multiple surveys about algorithmic fairness have been published in the literature [77, 358, 399]. These works typically focus on describing and classifying important measures of algorithmic fairness and methods to enhance it. Some articles also discuss sources of bias [358], software packages and projects which address fairness in ML [77], or describe selected sub-fields of algorithmic fairness [399]. Datasets are typically not emphasized in these works, which is also true of domain-specific surveys on algorithmic fairness, focused e.g. on ranking [404], Natural Language Processing (NLP) [483] and computational medicine [483]. As an exception, Pessach and Shmueli [399] and Zehlike et al. [565] list and briefly describe 12 popular algorithmic fairness datasets, and 19 datasets employed in fair ranking research, respectively.

2.2 Data studies

The work most closely related (and concurrently carried out) to ours is Le Quy et al. [302]. The authors perform a detailed analysis of 15 tabular datasets used in works of algorithmic fairness, listing important metadata (e.g. domain, protected attributes, collection period and location), and carrying out an exploratory analysis of the probabilistic relationship between features. Our work complements it by placing more emphasis on (1) a rigorous methodology for the selection of resources, (2) a wider selection of (over 200) datasets spanning different data types, including text, image, time-series, and tabular data, (3) a fine-grained evaluation of domains and tasks associated with each dataset. It will be interesting to see how different goals of the research community, such as selection of appropriate resources for experimentation and data studies, can benefit from the breadth and depth of both works.

Other works analyzing multiple datasets along specific lines have been carried out in recent years. Crawford and Paglen [120] focus on resources commonly used as training sets in computer vision, with attention to associated labels and underlying taxonomies. Fabrizzi et al. [158] also consider computer vision datasets, describing types of bias affecting them, along with methods for discovering and measuring bias. Peng et al. [397] analyze ethical concerns in three popular face and person recognition datasets, stemming from derivative datasets and models, lack of clarity of licenses, and dataset management practices. Geiger et al. [183] evaluate transparency in the documentation of labeling practices employed in over 100 datasets about Twitter. Leonelli and Tempini [308] study practices of collection, cleaning, visualization, sharing, and analysis across a variety of research domains. Romei and Ruggieri [436] survey techniques and data for discrimination analysis, focused on measuring, rather than enforcing, equity in human processes.

A different, yet related, family of articles provides deeper analyses of single datasets. Prabhu and Birhane [406] focus on Imagenet (ILSVRC 2012) which they analyze along the lines of consent,

problematic content, and individual re-identification. Kizhner et al. [278] study issues of representation in the Google Arts and Culture project across countries, cities and institutions. Some works provide datasheets for a given resource, such as CheXpert [178] and the BookCorpus [27]. Among popular fairness datasets, COMPAS has drawn scrutiny from multiple works, analysing its numerical idiosyncrasies [31] and sources of bias [28]. Ding et al. [141] study numerical idiosyncrasies in the Adult dataset, and propose a novel version, for which they provide a datasheet. Grömping [204] discuss issues resulting from coding mistakes in German Credit.

Our work combines the breadth and the depth of multi-dataset and single-dataset studies. On one hand, we survey numerous resources used in works of algorithmic fairness, analyzing them across multiple dimensions. On the other hand, we identify the most popular resources, compiling their datasheet and nutrition label, and summarize their perks and limitations. Moreover, by making our data briefs available, we hope to contribute a useful tool to the research community, favouring further data studies and analyses, as outlined in Section 6.

2.3 Documentation frameworks

Several data documentation frameworks have been proposed in the literature; three popular ones are described below. *Datasheets for Datasets* [182] are a general-purpose qualitative framework with over fifty questions covering key aspects of datasets, such as motivation, composition, collection, preprocessing, uses, distribution, and maintenance. Another qualitative framework is represented by *data statements* [37], which is tailored for NLP, requiring domain-specific information on language variety and speaker demographics. *Dataset Nutrition Labels* [227] describe a complementary, quantitative framework, focused on numerical aspects such as the marginal and joint distribution of variables.

Popular datasets require close scrutiny; for this reason we adopt these frameworks, producing three datasheets and nutrition labels for Adult, German Credit, and COMPAS. This approach, however, is not suited for a wider documentation effort with limited resources. For this reason, we propose and produce *data briefs*, a lightweight documentation format designed for algorithmic fairness datasets. Data briefs, described in Appendix A, include fields specific to fair ML, such sensitive attributes and tasks for which the dataset has been used in the algorithmic fairness literature.

3 METHODOLOGY

In this work, we consider (1) every article published in the proceedings of domain-specific conferences such as the ACM Conference on Fairness, Accountability, and Transparency (FAccT), and the AAAI/ACM Conference on Artificial Intelligence, Ethics and Society (AIES); (2) every article published in proceedings of well-known machine learning and data mining conferences, including the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), the Conference on Neural Information Processing Systems (NeurIPS), the International Conference on Machine Learning (ICML), the International Conference on Learning Representations (ICLR), the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD); (3) every article available from Past Network Events and Older Workshops and Events of the FAccT

network.³ We consider the period from 2014, the year of the first workshop on Fairness, Accountability, and Transparency in Machine Learning, to June 2021, thus including works presented at FAccT, ICLR, AIES, and CVPR in 2021.⁴

To target works of algorithmic fairness, we select a subsample of these articles whose titles contain either of the following strings, where the star symbol represents the wildcard character: **fair** (targeting e.g. fairness, unfair), **bias** (biased, debiasing), *discriminat** (discrimination, discriminatory), **equal** (equality, unequal), **equi t** (equity, equitable), *disparate* (disparate impact), **pari t** (parity, disparities). These selection criteria are centered around equity-based notions of fairness, typically operationalized by equalizing some algorithmic property across individuals or groups of individuals. Through manual inspection by two authors, we discard articles where these keywords are used with a different meaning. Discarded works, for instance, include articles on handling pose distribution bias [585], compensating selection bias to improve accuracy without attention to sensitive attributes [268], enhancing desirable discriminating properties of models [91], or generally focused on model performance [317, 587]. This leaves us with 558 articles.

From the articles that pass this initial screening, we select datasets treated as important data artifacts, either being used to train/test an algorithm or undergoing a data audit, i.e., an in-depth analysis of different properties. We produce a data brief for these datasets by (1) reading the information provided in the surveyed articles, (2) consulting the provided references, and (3) reviewing scholarly articles or official websites found by querying popular search engines with the dataset name. We discard the following:

- Word Embeddings (WEs). We only consider the corpora they are trained on, provided WEs are trained as part of a given work and not taken off the shelf;
- toy datasets, i.e., simulations with no connection to real-world processes, unless they are used in more than one article, which we take as a sign of importance in the field;
- auxiliary resources that are only used as a minor source of ancillary information, such as the percentage of US residents in each state;
- datasets for which the available information is insufficient. This happens very seldom when points (1), (2), and (3) outlined above result in little to no information about the curators, purposes, features, and format of a dataset. For popular datasets, this is never the case.

For each of the 226 datasets satisfying the above criteria, we produce a data brief, available in Appendix A with a description of the underlying coding procedure. From this effort, we rigorously identify the three most popular resources, whose perks and limitations are summarized in the next section.

4 MOST POPULAR DATASETS

Figure 1 depicts the number of articles using each dataset, showing that dataset utilization in surveyed scholarly works follows a long

³<https://faccconference.org/network/>

⁴We are working on a yearly update covering more recent work, including articles presented at the ACM conference on Equity and Access in Algorithms, Mechanisms, and Optimization.

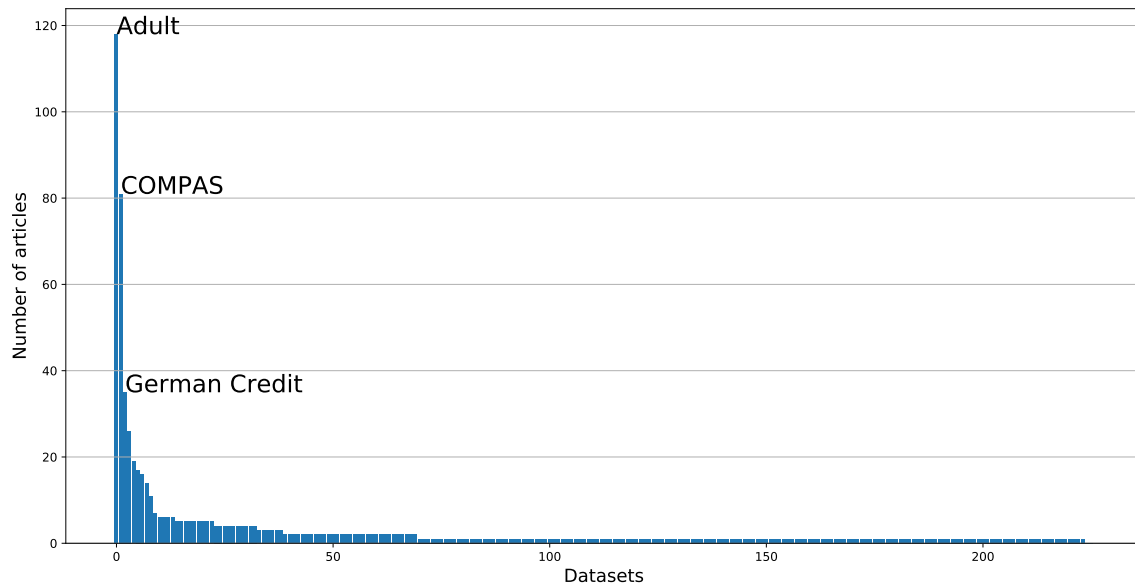


Figure 1: Utilization of datasets in fairness research follows a long tail distribution.

tail distribution. Over 100 datasets are only used once, also because some of these resources are not publicly available. Complementing this long tail is a short head of nine resources used in ten or more articles. These datasets are Adult (118 usages), COMPAS (81), German Credit (35), Communities and Crime (26), Bank Marketing (19), Law School (17), CelebA (16), MovieLens (14), and Credit Card Default (11). The tenth most used resource is the toy dataset from Zafar et al. [564], used in 7 articles. In this section, we summarize positive and negative aspects of the three most popular datasets, namely Adult, COMPAS, and German Credit, informed by extensive documentation in Appendices B, C and D.

4.1 Adult

The Adult dataset was created as a resource to benchmark the performance of machine learning algorithms on socially relevant data. Each instance is a person who responded to the March 1994 US Current Population Survey, represented along demographic and socio-economic dimensions, with features describing their profession, education, age, sex, race, personal, and financial condition. The dataset was extracted from the census database, preprocessed, and donated to UCI Machine Learning Repository in 1996 by Ronny Kohavi and Barry Becker. A binary variable encoding whether respondents’ income is above \$50,000 was chosen as the target of the prediction task associated with this resource.

Adult inherits some positive sides from the best practices employed by the US Census Bureau. Although later filtered somewhat arbitrarily, the original sample was designed to be representative of the US population. Trained and compensated interviewers collected the data. Attributes in the dataset are self-reported and provided by consensual respondents. Finally, the original data from the US Census Bureau is well documented, and its variables can be mapped to Adult by consulting the original documentation [505], except for a variable denominated `fnlwgt`, whose precise meaning is unclear.

A negative aspect of this dataset is the contrived prediction task associated with it. Income prediction from socio-economic factors is a task whose social utility appears rather limited. Even discounting this aspect, the arbitrary \$50,000 threshold for the binary prediction task is high, and model properties such as accuracy and fairness are very sensitive to it [141]. Furthermore, there are several sources of noise affecting the data. Roughly 7% of the data points have missing values, plausibly due to issues with data recording and coding, or respondents’ inability to recall information. Moreover, the tendency in household surveys for respondents to under-report their income is a common concern of the Census Bureau [370]. Another source of noise is top-coding of the variable “capital-gain” (saturation to \$99,999) to avoid the re-identification of certain individuals [505]. Finally, the dataset is rather old; sensitive attribute “race” contains the outdated “Asian Pacific Islander” class. It is worth noting that a set of similar resources was recently made available, allowing more current socio-economic studies of the US population [141].

4.2 COMPAS

This dataset was created for an external audit of racial biases in the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) risk assessment tool developed by Northpointe (now Equivant), which estimates the likelihood of a defendant becoming a recidivist. Instances represent defendants scored by COMPAS in Broward County, Florida, between 2013–2014, reporting their demographics, criminal record, custody and COMPAS scores. Defendants’ public criminal records were obtained from the Broward County Clerk’s Office website matching them based on date of birth, first and last names. The dataset was augmented with jail records and COMPAS scores provided by the Broward County Sheriff’s Office. Finally, public incarceration records were downloaded from the Florida Department of Corrections website. Instances are associated with two target variables (`is_recid` and

Table 1: Limitations of popular algorithmic fairness datasets.

	Adult	COMPAS	German Credit
Age	Old (1994)	Recent (2013–2016)	Very old (1973–1975)
Prediction task	Contrived (income > 50K\$)	Realistic (recidivism)	Realistic (creditworthiness)
Sensitive attributes	Outdated racial categories	Outdated racial categories	Sex cannot be retrieved
Sources of noise	Top-coding; tendency to under-report income	Data leakage; label bias; clerical errors	Incorrect code table
Sample representativeness	US working population	Convenience sample (Broward County)	Artificial sample (credit granted, negative class oversampled)
Preprocessing needed	Handling missing values (7%)	Handling missing values (80%); removing redundant features; ground truth on detainment	None
Additional concerns	Accuracy and fairness are sensitive to arbitrary 50K\$ threshold	Potential for misguided discussion on criminal justice	Interpretability and exploratory analyses are invalid

is_violent_recid), indicating whether defendants were booked in jail for a criminal offense (potentially violent) that occurred after their COMPAS screening but within two years.

On the upside, this dataset is recent and captures some relevant aspects of the COMPAS risk assessment tool and the criminal justice system in Broward County. On the downside, it was compiled from disparate sources, hence clerical errors and mismatches are present [301]. Moreover, in its official release [408], the COMPAS dataset features redundant variables and data leakage due to spuriously time-dependent recidivism rates [31]. For these reasons, researchers must perform further preprocessing in addition to the standard one by ProPublica. More subjective choices are required of researchers interested in counterfactual evaluation of risk-assessment tools, due to the absence of a clear indication of whether defendants were detained or released pre-trial [367]. The lack of a standard preprocessing protocol beyond the one by ProPublica [408], which is insufficient to handle these factors, may cause issues of reproducibility and difficulty in comparing methods. Moreover, according to Northpointe’s response to the ProPublica’s study, several risk factors considered by the COMPAS algorithm are absent from the dataset [140]. As an additional concern, race categories lack Native Hawaiian or Other Pacific Islander, while Hispanic is redefined as race instead of ethnicity [28]. Finally, defendants’ personal information (e.g. race and criminal history) is available in conjunction with obvious identifiers, making re-identification of defendants trivial.

COMPAS also represents a case of a broad phenomenon which can be termed *data bias*. With terminology from [175], when it comes to datasets encoding complex human phenomena, there is often a disconnect between the construct space (what we aim to measure) and the observed space (what we end up observing). This may be especially problematic if the difference between construct and observation is uneven across individuals or groups. COMPAS, for example, is a dataset about criminal offense. Offense is central to the prediction target Y , aimed at encoding recidivism, and to the available covariates X , summarizing criminal history. However, the COMPAS dataset (observed space) is an imperfect proxy for the criminal patterns it should summarize (construct space). The prediction labels Y actually encode re-arrest, instead of re-offense [301],

and are thus clearly influenced by spatially differentiated policing practices [173]. This is also true of criminal history encoded in COMPAS covariates, again mediated by arrest and policing practices which may be racially biased [28, 348]. As a result, the true fairness of an algorithm, just like its accuracy, may differ significantly from what is reported on biased data. For example, algorithms that achieve equality of true positive rates across sensitive groups on COMPAS are deemed fair under the *equal opportunity* measure [215]. However, if both the training set on which this objective is enforced and the test set on which it is measured are affected by race-dependent noise described above, those algorithms are only “fair” in an abstract observed space, but not in real construct space we ultimately care about [175].

Overall, these considerations paint a mixed picture for a dataset of high social relevance that was extremely useful to catalyze attention on algorithmic fairness issues, displaying at the same time several limitations in terms of its continued use as a flexible benchmark for fairness studies of all sorts. In this regard, Bao et al. [28] suggest avoiding the use of COMPAS to demonstrate novel approaches in algorithmic fairness, as considering the data without proper context may bring to misleading conclusions which could misguidedly enter the broader debate on criminal justice and risk assessment.

4.3 German Credit

The German Credit dataset was created to study the problem of automated credit decisions at a regional Bank in southern Germany. Instances represent loan applicants from 1973 to 1975, who were deemed creditworthy and were granted a loan, bringing about a natural selection bias. Within this sample, bad credits are oversampled to favour a balance in target classes [204]. The data summarizes their financial situation, credit history, and personal situation, including housing and number of liable people. A binary variable encoding whether each loan recipient punctually payed every installment is the target of a classification task. Among the covariates, marital status and sex are jointly encoded in a single variable. Many documentation mistakes are present in the UCI entry associated

with this resource [501]. A revised version with correct variable encodings, called South German Credit, was donated to UCI Machine Learning Repository [503] with an accompanying report [204].

The greatest upside of this dataset is the fact that it captures a real-world application of credit scoring at a bank. On the downside, the data is half a century old, significantly limiting the societally useful insights that can be gleaned from it. Most importantly, the popular release of this dataset [501] comes with highly inaccurate documentation which contains wrong variable codings. For example, the variable reporting whether loan recipients are foreign workers has its coding reversed, so that, apparently, fewer than 5% of the loan recipients in the dataset would be German. Luckily, this error has no impact on numerical results obtained from this dataset, as it is irrelevant at the level of abstraction afforded by raw features, with the exception of potentially counterintuitive explanations in works of interpretability and exploratory analysis [302]. This coding error, along with others discussed in Grömping [204] was corrected in a novel release of the dataset [503]. Unfortunately and most importantly for the fair ML community, retrieving the sex of loan applicants is simply not possible, unlike the original documentation suggested. This is due to the fact that one value of this feature was used to indicate both women who are divorced, separated or married and men who are single, while the original documentation reported each feature value to correspond to same-sex applicants (either male-only or female-only). This particular coding error ended up having a non-negligible impact on the fair ML community, where many works studying group fairness extract sex from the joint variable and use it as a sensitive attribute, even years after the redacted documentation was published [528]. These coding mistakes are part of a documentation debt whose influence continues to affect the algorithmic fairness community.

4.4 Summary

Adult, COMPAS and German Credit are the most used datasets in the surveyed algorithmic fairness literature, despite the limitations summarized in Table 1. Their status as de facto fairness benchmarks is probably due to their use in seminal works [69, 396] and influential articles [15] on algorithmic fairness. Once this fame was created, researchers had clear incentives to study novel problems and approaches on these datasets, which have become even more established benchmarks in the algorithmic fairness literature as a result [28]. On close scrutiny, the fundamental merit of these datasets lies in originating from human processes, encoding protected attributes, and having different base rates for the target variable across sensitive groups. Their use in recent works on algorithmic fairness can be interpreted as a signal that the authors have basic awareness of default data practices in the field and that the data was not made up to fit the algorithm. Overarching claims of significance in real-world scenarios stemming from experiments on these datasets should be met with skepticism. Experiments that claim extracting a sex variable from the German Credit dataset should be considered noisy at best. As for alternatives, Bao et al. [28] suggest employing well-designed simulations. A complementary avenue is to seek different datasets that are relevant for the problem at hand. We hope that the two hundred data briefs accompanying this work will prove useful in this regard, favouring both domain-oriented

and task-oriented searches, according to the classification discussed in the next section.

5 EXISTING ALTERNATIVES

In this section, we discuss existing fairness resources from three different perspectives. In section 5.1 we describe the different domains spanned by fairness datasets. In section 5.2 we provide a categorization of fairness tasks supported by the same resources. In section 5.3 we discuss the different roles played by these datasets in fairness research, such as supporting training and benchmarking.

5.1 Domain

Algorithmic fairness concerns arise in any domain where Automated Decision Making (ADM) systems may influence human well-being. Unsurprisingly, the datasets in our survey reflect a variety of areas where ADM systems are studied or deployed, including criminal justice, education, search engines, online marketplaces, emergency response, social media, medicine, and hiring. In Figure 2, we report a subdivision of the surveyed datasets in different macrodomains.⁵ We mostly follow the area-category taxonomy by Scimago,⁶ departing from it where appropriate. For example, we consider computer vision and linguistics macrodomains of their own for the purposes of algorithmic fairness, as much fair ML work has been published in both disciplines. Below we present a description of each macrodomain and its main subdomains, summarized in detail in Table 3 (Appendix A).

Computer Science. Datasets from this macrodomain are very well represented, comprising *information systems*, *social media*, *library and information sciences*, *computer networks*, and *signal processing*. *Information systems* heavily feature datasets on search engines for various items such as text, images, worker profiles, and real estate, retrieved in response to queries issued by users (Occupations in Google Images, Scientist+Painter, Zillow Searches, Barcelona Room Rental, Burst, TaskRabbit, Online Freelance Marketplaces, Bing US Queries, Symptoms in Queries). Other datasets represent problems of item recommendation, covering products, businesses, and movies (Amazon Recommendations, Amazon Reviews, Google Local, MovieLens, FilmTrust). The remaining datasets in this subdomain represent knowledge bases (Freebase15k-237, Wikidata) and automated screening systems (CVs from Singapore, Pymetrics Bias Group). Datasets from *social media* that are not focused on links and relationships between people are also considered part of computer science in this survey. These resources are often focused on text, powering tools, and analyses of hate speech and toxicity (Civil Comments, Twitter Abusive Behavior, Twitter Offensive Language, Twitter Hate Speech Detection, Twitter Online Harassment), dialect (TwitterAAE), and political leaning (Twitter Presidential Politics). Twitter is by far the most represented platform, while datasets from Facebook (German Political Posts), Steemit (Steemit), Instagram (Instagram Photos), Reddit (RtGender, Reddit Comments), Fitocracy (RtGender), and YouTube (YouTube Dialect Accuracy) are also present. Datasets from *library and information sciences* are mainly focused on academic collaboration networks (Cora Papers,

⁵The total exceeds 226 due to multiple domains being applicable to some dataset.

⁶See the “subject area” and “subject category” drop down menus from <https://www.scimagojr.com/journalrank.php>, accessed on March 15, 2022

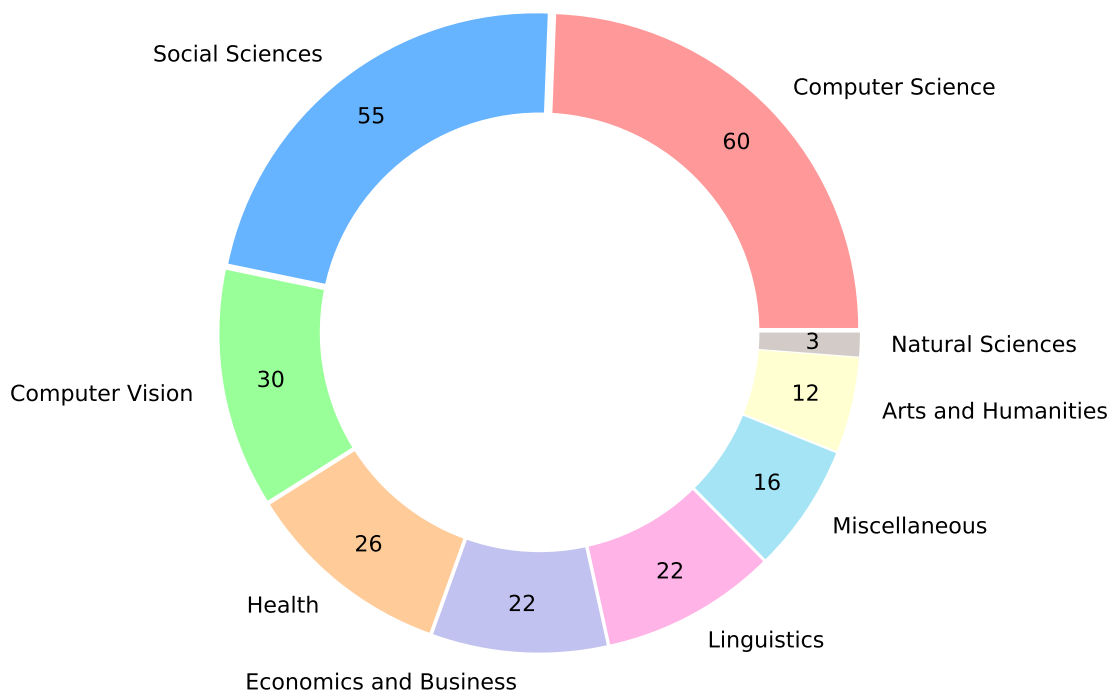


Figure 2: Datasets employed in fairness research span diverse domains. See Table 3 for a detailed breakdown.

CiteSeer Papers, PubMed Diabetes Papers, ArnetMiner Citation Network, 4area, Academic Collaboration Networks), except for a dataset about peer review of scholarly manuscripts (Paper-Reviewer Matching).

Social Sciences. Datasets from social sciences are also plentiful, spanning *law*, *education*, *social networks*, *demography*, *social work*, *political science*, *transportation*, *sociology* and *urban studies*. *Law* datasets are mostly focused on recidivism (Crowd Judgement, COMPAS, Recidivism of Felons on Probation, State Court Processing Statistics, Los Angeles City Attorney’s Office Records) and crime prediction (Strategic Subject List, Philadelphia Crime Incidents, Stop, Question and Frisk, Real-Time Crime Forecasting Challenge, Dallas Police Incidents, Communities and Crime), with a granularity spanning the range from individuals to communities. In the area of *education* we find datasets that encode application processes (Nursery, IIT-JEE), student performance (Student, Law School, UniGe, ILEA, US Student Performance, Indian Student Performance, EdGap, Berkeley Students), including attempts at automated grading (Automated Student Assessment Prize), and placement information after school (Campus Recruitment). Some datasets on student performance support studies of differences across schools and educational systems, for which they report useful features (Law School, ILEA, EdGap), while the remaining datasets are more focused on differences in the individual condition and outcome for students, typically within the same institution. Datasets about *social networks* mostly concern online social networks (Facebook Ego-networks, Facebook Large Network, Pokec Social Network, Rice Facebook Network,

Twitch Social Networks, University Facebook Networks), except for High School Contact and Friendship Network, also featuring offline relations. *Demography* datasets comprise census data from different countries (Dutch Census, Indian Census, National Longitudinal Survey of Youth, Section 203 determinations, US Census Data (1990)). Datasets from *social work* cover complex personal and social problems, including child maltreatment prevention (Allegheny Child Welfare), emergency response (Harvey Rescue) and drug abuse prevention (Homeless Youths’ Social Networks, DrugNet). Resources from *political science* describe registered voters (North Carolina Voters), electoral precincts (MGGG States), polling (2016 US Presidential Poll), and sortition (Climate Assembly UK). *Transportation* data summarizes trips and fares from taxis (NYC Taxi Trips, Shanghai Taxi Trajectories), ride-hailing (Chicago Ridesharing, Ride-hailing App), and bike sharing services (Seoul Bike Sharing), along with public transport coverage (Equitable School Access in Chicago). *Sociology* resources summarize online (Libimseti) and offline dating (Columbia University Speed Dating). Finally, we assign SafeGraph Research Release to *urban studies*.

Computer Vision. This is an area of early success for artificial intelligence, where fairness typically concerns learned representations and equality of performance across classes. The surveyed articles feature several popular datasets on image classification (ImageNet, MNIST, Fashion MNIST, CIFAR), visual question answering (Visual Question Answering), segmentation and captioning (MS-COCO, Open Images Dataset). We find over ten face analysis datasets (Labeled Faces in the Wild, UTK Face, Adience, FairFace,

IJB-A, CelebA, Pilot Parliaments Benchmark, MS-Celeb-1M, Diversity in Faces, Multi-task Facial Landmark, Racial Faces in the Wild, BUPT Faces), including one from experimental psychology (FACES), for which fairness is most often intended as the robustness of classifiers across different subpopulations, without much regard for downstream benefits or harms to these populations. Synthetic images are popular to study the relationship between fairness and disentangled representations (dSprites, Cars3D, shapes3D). Similar studies can be conducted on datasets with spurious correlations between subjects and backgrounds (Waterbirds, Benchmarking Attribution Methods) or gender and occupation (Athletes and health professionals). Finally, the Image Embedding Association Test dataset is a fairness benchmark to study biases in image embeddings across religion, gender, age, race, sexual orientation, disability, skin tone, and weight. It is worth noting that this significant proportion of computer vision datasets is not an artifact of including CVPR in the list of candidate conferences, which contributed just five additional datasets (Multi-task Facial Landmark, Office31, Racial Faces in the Wild, BUPT Faces, Visual Question Answering).

Health. This macrodomain, comprising medicine, psychology and pharmacology displays a notable diversity of subdomains interested by fairness concerns. Specialties represented in the surveyed datasets are mostly medical, including *public health* (Antelope Valley Networks, Willingness-to-Pay for Vaccine, Kidney Matching, Kidney Exchange Program), *cardiology* (Heart Disease, Arrhythmia, Framingham), *endocrinology* (Diabetes 130-US Hospitals, Pima Indians Diabetes Dataset), *health policy* (Heritage Health, MEPS-HC). Specialties such as *radiology* (National Lung Screening Trial, MIMIC-CXR-JPG, CheXpert) and *dermatology* (SIIM-ISIC Melanoma Classification, HAM10000) feature several image datasets for their strong connections with medical imaging. Other specialties include *critical care medicine* (MIMIC-III), *neurology* (Epileptic Seizures), *pediatrics* (Infant Health and Development Program), *sleep medicine* (Apnea), *nephrology* (Renal Failure), *pharmacology* (Warfarin) and *psychology* (Drug Consumption, FACES). These datasets are often extracted from care data of multiple medical centers to study problems of automated diagnosis. Resources derived from longitudinal studies, including Framingham and Infant Health and Development Program are also present. Works of algorithmic fairness in this domain are typically concerned with obtaining models with similar performance for patients across race and sex.

Linguistics. In addition to the textual resources we already described, such as the ones derived from social media, several datasets employed in algorithmic fairness literature can be assigned to the domain of linguistics and Natural Language Processing (NLP). There are many examples of resources curated to be fairness benchmarks for different tasks, including machine translation (Bias in Translation Templates), sentiment analysis (Equity Evaluation Corpus), coreference resolution (Winogender, Winobias, GAP Coreference), named entity recognition (In-Situ), language models (BOLD) and word embeddings (WEAT). Other datasets have been considered for their size and importance for pretraining text representations (Wikipedia dumps, One billion word benchmark, BookCorpus, WebText) or their utility as NLP benchmarks (GLUE, Business Entity Resolution). Speech recognition resources have also been considered (TIMIT).

Economics and Business. This macrodomain comprises datasets from *economics*, *finance*, *marketing*, and *management information systems*. *Economics* datasets mostly consist of census data focused on wealth (Adult, US Family Income, Poverty in Colombia, Costa Rica Household Survey) and other resources which summarize employment (ANPE), tariffs (U.S. Harmonized Tariff Schedules), insurance (Italian Car Insurance), and division of goods (Spliddit Divide Goods). *Finance* resources feature data on microcredit and peer-to-peer lending (Mobile Money Loans, Kiva, Prosper Loans Network), mortgages (HMDA), loans (German Credit, Credit Elasticities), credit scoring (FICO) and default prediction (Credit Card Default). *Marketing* datasets describe marketing campaigns (Bank Marketing), customer data (Wholesale) and advertising bids (Yahoo! A1 Search Marketing). Finally, datasets from *management information systems* summarize information about automated hiring (CVs from Singapore, Pymetrics Bias Group) and employee retention (IBM HR Analytics).

Miscellaneous. This macrodomain contains several datasets originating from the *news* domain (Yow news, Guardian Articles, Latin Newspapers, Adressa, Reuters 50 50, New York Times Annotated Corpus, TREC Robust04). Other resources include datasets on food (Sushi), sports (Fantasy Football, FIFA 20 Players, Olympic Athletes), and toy datasets (Toy Dataset 1–4).

Arts and Humanities. In this area we mostly find *literature* datasets, which contain text from literary works (Shakespeare, Curatr British Library Digital Corpus, Victorian Era Authorship Attribution, Nominees Corpus, Riddle of Literary Quality), which are typically studied with NLP tools. Other datasets in this domain include domain-specific information systems about books (Goodreads Reviews), *movies* (MovieLens) and *music* (Last.fm, Million Song Dataset, Million Playlist Dataset).

Natural Sciences. This domain is represented with three datasets from *biology* (iNaturalist), *biochemistry* (PP-Pathways) and *plant science*, with the classic Iris dataset.

As a whole, many of these datasets encode fundamental human activities where algorithms and ADM systems have been studied and deployed. Alertness and attention to equity seems especially important in specific domains, including social sciences, computer science, medicine, and economics. Here the potential for impact may result in large benefits, but also great harm, particularly for vulnerable populations and minorities, more likely to be neglected during the design, training, and testing of an ADM. After concentrating on domains, in the next section we analyze the variety of tasks studied in works of algorithmic fairness and supported by these datasets.

5.2 Task and setting

Researchers and practitioners are showing an increasing interest in algorithmic fairness, proposing solutions for many different *tasks*, including fair classification, regression and ranking. At the same time, the academic community is developing an improved understanding of important challenges that run across different tasks in the algorithmic fairness space [107], also thanks to practitioner surveys [230] and studies of specific legal challenges [13]. To exemplify, the presence of noise corrupting labels for sensitive attributes

represents a challenge that may apply across different tasks, including fair classification, regression and ranking. We refer to these challenges as *settings*, describing them in the second part of this section. While our work focuses on fair ML datasets, it is cognizant of the wide variety of tasks tackled in the algorithmic fairness literature, which are captured in a specific field of our data briefs. In this section we provide an overview of common tasks and settings studied on these datasets, showing their variety and diversity. Table 2 summarizes these tasks, listing the three most used datasets for each one. When describing a task, we explicitly highlight datasets that are particularly relevant to it, even when outside of the top three.

5.2.1 Task. Fair classification [70, 146] is the most common task by far. Typically, it involves equalizing some measure of interest across subpopulations, such as the recall, precision, or accuracy for different racial groups. On the other hand, individually fair classification focuses on the idea that similar individuals (low distance in the covariate space) should be treated similarly (low distance in the outcome space), often formalized as a Lipschitz condition. Unsurprisingly, the most common datasets for fair classification are the most popular ones overall (§ 4), i.e., Adult, COMPAS, and German Credit.

Fair regression [41] concentrates on models that predict a real-valued target, requiring the average loss to be balanced across groups. Individual fairness in this context may require losses to be as uniform as possible across all individuals. Fair regression is a less popular task, often studied on the Communities and Crime dataset, where the task is predicting the rate of violent crimes in different communities.

Fair ranking [553] requires ordering candidate items based on their relevance to a current need. Fairness in this context may concern both the people producing the items that are being ranked (e.g. artists) and those consuming the items (users of a music streaming platform). It is typically studied in applications of recommendation (MovieLens, Amazon Recommendations, Last.fm, Million Song Dataset, Adressa) and search engines (Yahoo! c14B Learning to Rank, Microsoft Learning to Rank, TREC Robust04).

Fair matching [283] is similar to ranking as they are both tasks defined on two-sided markets. This task, however, is focused on highlighting and matching pairs of items on both sides of the market, without emphasis on the ranking component. Datasets for this task are from diverse domains, including dating (Libimseti, Columbia University Speed Dating) transportation (NYC Taxi Trips, Ride-hailing App) and organ donation (Kidney Matching, Kidney Exchange Program).

Fair risk assessment [116] studies algorithms that score instances in a dataset according to a predefined type of risk. Relevant domains include healthcare and criminal justice. Key differences with respect to classification are an emphasis on real-valued scores rather than labels, and awareness that the risk assessment process can lead to interventions impacting the target variable. For this reason, fairness concerns are often defined in a counterfactual fashion. The most popular dataset for this task is COMPAS, followed by datasets from medicine (IHDP, Stanford Medicine Research Data Repository), social work (Allegheny Child Welfare), Economics (ANPE) and Education (EdGap).

Fair representation learning [122] concerns the study of features learnt by models as intermediate representations for inference tasks. A popular line of work in this space, called *disentanglement*, aims to learn representations where a single factor of import corresponds to a single feature. Ideally, this approach should select representations where sensitive attributes cannot be used as proxies for target variables. Cars3D and dSprites are popular datasets for this task, consisting of synthetic images depicting controlled shape types under a controlled set of rotations. Post-processing approaches are also applicable to obtain fair representations from biased ones via debiasing.

Fair clustering [99] is an unsupervised task concerned with the division of a sample into homogenous groups. Fairness may be intended as an equitable representation of protected subpopulations in each cluster, or in terms of average distance from the cluster center. While Adult is the most common dataset for problems of fair clustering, other resources often used for this task include Bank Marketing, Diabetes 130-US Hospitals, Credit Card Default and US Census Data (1990).

Fair anomaly detection [567], also called **outlier detection** [128], is aimed at identifying surprising or anomalous points in a dataset. Fairness requirements involve equalizing salient quantities (e.g. acceptance rate, recall, precision, distribution of anomaly scores) across populations of interest. This problem is particularly relevant for members of minority groups, who, in the absence of specific attention to dataset inclusivity, are less likely to fit the norm in the feature space.

Fair districting [450] is the division of a territory into electoral districts for political elections. Fairness notions brought forth in this space are either outcome-based, requiring that seats earned by a party roughly match their share of the popular vote, or procedure-based, ignoring outcomes and requiring that counties or municipalities are split as little as possible. MGGG States is a reference resource for this task, providing precinct-level aggregated information about demographics and political leaning of voters in US districts.

Fair task assignment and truth discovery [189, 316] are different subproblems in the same area, focused on the subdivision of work and the aggregation of answers in crowdsourcing. Here fairness may be intended concerning errors in the aggregated answer, requiring errors to be balanced across subpopulations of interest, or in terms of the work load imposed to workers. A dataset suitable for this task is Crowd Judgement, containing crowd-sourced recidivism predictions.

Fair spatio-temporal process learning [454] focuses on the estimation of models for processes which evolve in time and space. Surveyed applications include crime forecasting (Real-Time Crime Forecasting Challenge, Dallas Police Incidents) and disaster relief (Harvey Rescue), with fairness requirements focused on equalization of performance across different neighbourhoods and special attention to their racial composition.

Fair graph diffusion [160] models and optimizes the propagation of information and influence over networks, and its probability of reaching individuals of different sensitive groups. Applications include obesity prevention (Antelope Valley Networks) and drug-use prevention (Homeless Youths' Social Networks). **Fair graph augmentation** [423] is a similar task, defined on graphs which define

Table 2: Most used datasets by fairness task and setting.

Task	Datasets
Fair classification	Adult; COMPAS; German Credit
Fair regression	Communities and Crime; Law School; Student
Fair ranking	MovieLens; German Credit; Kiva
Fair matching	NYC Taxi Trips; Libimseti; Columbia University Speed Dating
Fair risk assessment	COMPAS; Allegheny Child Welfare; Infant Health and Development Program (IHDP)
Fair representation learning	Adult; COMPAS; dSprites
Fair clustering	Adult; Bank Marketing; Diabetes 130-US Hospitals
Fair anomaly detection	Adult; MNIST; Credit Card Default
Fair districting	MGGG States
Fair task assignment	Crowd Judgement; COMPAS
Fair spatio-temporal process learning	Real-Time Crime Forecasting Challenge; Dallas Police Incidents; Harvey Rescue
Fair graph diffusion/augmentation	University Facebook Networks; Antelope Valley Networks; Rice Facebook Network
Fair resource allocation/subset selection	ML Fairness Gym; US Federal Judges; Climate Assembly UK
Fair data summarization	Adult; Student; Credit Card Default
Fair data generation	CelebA; MovieLens; shapes3D
Fair graph mining	MovieLens; Freebase15k-237; PP-Pathways
Fair pricing	Willingness-to-Pay for Vaccine; Credit Elasticities; Italian Car Insurance
Fair advertising	Yahoo! A1 Search Marketing; North Carolina Voters; Instagram Photos
Fair routing	Shanghai Taxi Trajectories
Fair entity resolution	Winogender; Winobias; Business Entity Resolution
Fair sentiment analysis	Popular Baby Names; Equity Evaluation Corpus (EEC); TwitterAAE
Bias in word embeddings	Wikipedia dumps; Word Embedding Association Test (WEAT); Popular Baby Names
Bias in language models	TwitterAAE; BOLD; GLUE
Fair machine translation	Bias in Translation Templates
Fair speech recognition	YouTube Dialect Accuracy; TIMIT

Setting	Datasets
Rich-subgroup fairness	Adult; COMPAS; Communities and Crime
Fairness under unawareness	Adult; COMPAS; HMDA
Limited-label fairness	Adult; German Credit; COMPAS
Robust fairness	COMPAS; Adult; MEPS-HC
Dynamical fairness	FICO; ML Fairness Gym; COMPAS
Preference-based fairness	Adult; COMPAS; Toy Dataset 1
Multi-stage fairness	Adult; Heritage Health; Twitter Offensive Language
Fair few-shot learning	Communities and Crime; Toy Dataset 1; Mobile Money Loans
Fair private learning	UTK Face; CheXpert; FairFace
Fair federated learning	Vehicle; Sentiment140; Shakespeare
Fair incremental learning	ImageNet; CIFAR
Fair active learning	Adult; German Credit; Heart Disease
Fair selective classification	CheXpert; CelebA; Civil Comments

access to resources based on existing infrastructure (e.g. transportation), which can be augmented under a budget to increase equity. This task has been proposed to improve school access (Equitable School Access in Chicago) and information availability in social networks (Facebook100).

Fair resource allocation/subset selection [20, 238] can often be formalized as a classification problem with constraints on the number of positives. Fairness requirements are similar to those of classification. Subset selection may be employed to choose a group of people from a wider set for a given task (US Federal Judges, Climate Assembly UK). Resource allocation concerns the division

of goods (Spliddit Divide Goods) and resources (ML Fairness Gym, German Credit).

Fair data summarization [78] refers to presenting a summary of datasets that is equitable to subpopulations of interest. It may involve finding a small subset representative of a larger dataset (strongly linked to subset selection) or selecting the most important features (dimensionality reduction). Approaches for this task have been applied to select a subset of images (Scientist+Painter) or customers (Bank Marketing), that represent the underlying population across sensitive demographics.

Fair data generation [548] deals with generating “fair” data points and labels, which can be used as training or test sets. Approaches in this space may be used to ensure an equitable representation of protected categories in data generation processes learnt from biased datasets (CelebA, IBM HR Analytics), and to evaluate biases in existing classifiers (MS-Celeb-1M). Data generation may also be limited to synthesizing artificial sensitive attributes [64].

Fair graph mining [262] focuses on representations and prediction tasks on graph structures. Fairness may be defined either as a lack of bias in representations, or with respect to a final inference task defined on the graph. Fair graph mining approaches have been applied to knowledge bases (Freebase15k-237, Wikidata), collaboration networks (CiteSeer Paper, Academic Collaboration Networks) and social network datasets (Facebook Large Network, Twitch Social Networks).

Fair pricing [260] concerns learning and deploying an optimal pricing policy for revenue while maintaining equity of access to services and consumer welfare across sensitive groups. Datasets employed in fair pricing are from the economics (Credit Elasticities, Italian Car Insurance), transportation (Chicago Ridesharing), and public health domains (Willingness-to-Pay for Vaccine).

Fair advertising [79] is also concerned with access to goods and services. It comprises both bidding strategies and auction mechanisms which may be modified to reduce discrimination with respect to the gender or race composition of the audience that sees an ad. One publicly available dataset for this subtask is Yahoo! A1 Search Marketing.

Fair routing [411] is the task of suggesting an optimal path from a starting location to a destination. For this task, experimentation has been carried out on a semi-synthetic traffic dataset (Shanghai Taxi Trajectories). The proposed fairness measure requires equalizing the driving cost per customer across all drivers.

Fair entity resolution [119] is a task focused on deciding whether multiple records refer to the same entity, which is useful, for instance, for the construction and maintenance of knowledge bases. Business Entity Resolution is a proprietary dataset for fair entity resolution, where constraints of performance equality across chain and non-chain businesses can be tested. Winogender and Winobias are publicly available datasets developed to study gender biases in pronoun resolution.

Fair sentiment analysis [277] is a well-established instance of fair classification, where text snippets are typically classified as positive, negative, or neutral depending on the sentiment they express. Fairness is intended with respect to the entities mentioned in the text (e.g. men and women). The central idea is that the estimated sentiment for a sentence should not change if female entities (e.g. “her”, “woman”, “Mary”) are substituted with their male counterparts (“him”, “man”, “James”). The Equity Evaluation Corpus is a benchmark developed to assess gender and race bias in sentiment analysis models.

Bias in Word Embeddings (WEs) [49] is the study of undesired semantics and stereotypes captured by vectorial representations of words. WEs are typically trained on large text corpora (Wikipedia dumps) and audited for associations between gendered words (or other words connected to sensitive attributes) and stereotypical or harmful concepts, such as the ones encoded in WEAT.

Bias in Language Models (LMs) [51] is, quite similarly, the study of biases in LMs, which are flexible models of human language based on contextualized word representations, which can be employed in a variety of linguistics and NLP tasks. LMs are trained on large text corpora from which they may learn spurious correlations and stereotypes. The BOLD dataset is an evaluation benchmark for LMs, based on prompts that mention different socio-demographic groups. LMs complete these prompts into full sentences, which can be tested along different dimensions (sentiment, regard, toxicity, emotion and gender polarity).

Fair Machine Translation (MT) [478] concerns automatic translation of text from a source language into a target one. MT systems can exhibit gender biases, such as a tendency to translate gender-neutral pronouns from the source language into gendered pronouns of the target language in accordance with gender stereotypes. For example, a “nurse” mentioned in a gender-neutral context in the source sentence may be rendered with feminine grammar in the target language. Bias in Translation Templates is a set of short templates to test such biases.

Fair speech recognition [491] requires accurate annotation of spoken language into text across different demographics. YouTube Dialect Accuracy is a dataset developed to audit the accuracy of YouTube’s automatic captions across two genders and five dialects of English. Similarly, TIMIT is a classic speech recognition dataset annotated with American English dialect and gender of speaker.

5.2.2 Setting. As noted at the beginning of this section, there are different *settings* (or challenges) that run across many tasks described above. Some of these settings are specific to fair ML, such as ensuring fairness across an exponential number of groups, or in the presence of noisy labels for sensitive attributes. Other settings are connected with common ML challenges, including few-shot and privacy-preserving learning. Below we describe common settings encountered in the surveyed articles. Most of these settings are tested on fairness datasets which are popular overall, i.e. Adult, COMPAS and German Credit. We highlight situations where this is not the case, potentially due to a given challenge arising naturally in some other dataset.

Rich-subgroup fairness [269] is a setting where fairness properties are required to hold not only for a limited number of protected groups, but across an exponentially large number of subpopulations. This line of work represents an attempt to bridge the normative reasoning underlying individual and group fairness.

Fairness under unawareness is a general expression to indicate problems where sensitive attributes are missing [94], encrypted [272] or corrupted by noise [299]. These problems respond to real-world challenges related to the confidential nature of protected attributes, that individuals may wish to hide, encrypt, or obfuscate. This setting is most commonly studied on highly popular fairness dataset (Adult, COMPAS), moderately popular ones (Law School and Credit Card Default), and a dataset about home mortgage applications in the US (HMDA).

Limited-label fairness comprises settings with limited information on the target variable, including situations where labelled instances are few [248], noisy [528], or only available in aggregate form [443].

Robust fairness problems arise under perturbations to the training set [237], adversarial attacks [378] and dataset shift [464]. This line of research is often connected with work in robust machine learning, extending the stability requirements beyond accuracy-related metrics to fairness-related ones.

Dynamical fairness [124, 322] entails repeated decisions in changing environments, potentially affected by the very algorithm that is being studied. Works in this space study the co-evolution of algorithms and populations on which they act over time. For example, an algorithm that achieves equality of acceptance rates across protected groups in a static setting may generate further incentives for the next generation of individuals from historically disadvantaged groups. Popular resources for this setting are FICO and the ML Fairness GYM.

Preference-based fairness [563] denotes work informed, explicitly or implicitly, by the preferences of stakeholders. For people subjected to a decision this is related to notions of envy-freeness and loss aversion [10]; alternatively, policy-makers can express indications on how to trade-off different fairness measures [574], or experts can provide demonstrations of fair outcomes [177].

Multi-stage fairness [338] refers to settings where several decision makers coexist in a compound decision-making process. Decision makers, both humans and algorithmic, may act with different levels of coordination. A fundamental question in this setting is how to ensure fairness under composition of different decision mechanisms.

Fair few-shot learning [580] aims at developing fair ML solutions in the presence of a small amount of data samples. The problem is closely related to, and possibly solved by, **fair transfer learning** [117] where the goal is to exploit the knowledge gained on a problem to solve a different but related one. Datasets where this setting arises naturally are Communities and Crime, where one may restrict the training set to a subset of US states, and Mobile Money Loans, which consists of data from different African countries.

Fair private learning [22, 247] studies the interplay between privacy-preserving mechanisms and fairness constraints. Works in this space consider the equity of machine learning models designed to avoid leakage of information about individuals in the training set. Common domains for datasets employed in this setting are face analysis (UTK Face, FairFace, Diversity in Face) and medicine (CheXpert, SIIM-ISIC Melanoma Classification, MIMIC-CXR-JPG).

Additional settings that are less common include **fair federated learning** [314], where algorithms are trained across multiple decentralized devices, **fair incremental learning** [579], where novel classes may be added to the learning problem over time, **fair active learning** [383], allowing for the acquisition of novel information during inference and **fair selective classification** [253], where predictions are issued only if model confidence is above a certain threshold.

Overall, we found a variety of tasks defined on fairness datasets, ranging from generic, such as *fair classification*, to narrow and specifically defined on certain datasets, such as *fair districting* on MGGG States and *fair truth discovery* on Crowd Judgement. Orthogonally to this dimension, many settings or challenges may arise to complicate these tasks, including noisy labels, system dynamics, and privacy concerns. Quite clearly, algorithmic fairness research

has been expanding in both directions, by studying a variety of tasks under diverse and challenging settings. In the next section, we analyze the roles played in scholarly works by the surveyed datasets.

5.3 Role

The datasets used in algorithmic fairness research can play different roles. For example, some may be used to train novel algorithms, while others are suited to test existing algorithms from a specific point of view. Chapter 7 of Barocas et al. [33], describes six different roles of datasets in machine learning. We adopt their framework to analyse fair ML datasets, adding to the taxonomy two roles that are specific to fairness research.

A source of real data. While synthetic datasets and simulations may be suited to demonstrate specific properties of a novel method, the usefulness of an algorithm is typically established on data from the real world. More than a sign of immediate applicability to important challenges, good performance on real-world sources of data signals that the researchers did not make up the data to suit the algorithm. This is likely the most common role for fairness datasets, especially common for the ones hosted on the UCI ML repository, including Adult, German Credit, Communities and Crime, Diabetes 130-US Hospitals, Bank Marketing, Credit Card Default, US Census Data (1990). These resources owe their popularity in fair ML research to being a product of human processes and to encoding protected attributes. Quite simply, they are sources of real human data.

A catalyst of domain-specific progress. Datasets can spur algorithmic insight and bring about domain-specific progress. Civil Comments is a great example of this role, powering the Jigsaw Unintended Bias in Toxicity Classification challenge. The challenge responds to a specific need in the space of automated moderation against toxic comments in online discussion. Early attempts at toxicity detection resulted in models which associate mentions of frequently attacked identities (e.g. gay) with toxicity, due to spurious correlations in training sets. The dataset and associated challenge tackle this issue by providing toxicity ratings for comments, along with labels encoding whether members of a certain group are mentioned, favouring measurement of undesired bias. Many other datasets can play a similar role, including, Winogender, Winobias and the Equity Evaluation Corpus. In a broader sense, COMPAS and the accompanying study [15] have been an important catalyst, not for a specific task, but for fairness research overall.

A way to numerically track progress on a problem. This role is common for machine learning benchmarks that also provide human performance baselines. Algorithmic methods approaching or surpassing these baselines are often considered a sign that the task is “solved” and that harder benchmarks are required [33]. Algorithmic fairness is a complicated, context-dependent, contested construct whose correct measurement is continuously debated. Due to this reason, we are unaware of any dataset having a similar role in the algorithmic fairness literature.

A resource to compare models. Practitioners interested in solving a specific problem may take a large set of algorithms and test them on a group of datasets that are representative of their problem, in order to select the most promising ones. For well-established

ML challenges, there are often leaderboards providing a concise comparison between algorithms for a given task, which may be used for model selection. This setting is rare in the fairness literature, also due to inherent difficulties in establishing a single measure of interest in the field. One notable exception is represented by Friedler et al. [176], who employed a suite of four datasets (Adult, COMPAS, German Credit, Ricci) to compare the performance of four different approaches to fair classification.

A source of pre-training data. Flexible, general-purpose models are often pre-trained to encode useful representations, which are later fine-tuned for specific tasks in the same domain. For example, large text corpora are often employed to train language models and word embeddings which are later specialized to support a variety of downstream NLP applications. Wikipedia dumps, for instance, are often used to train word embeddings and investigate their biases [62, 318, 393]. Several algorithmic fairness works aim to study and mitigate undesirable biases in learnt representations. Corpora like Wikipedia dumps are used to obtain representations via realistic training procedures that mimic common machine learning practice as closely as possible.

A source of training data. Models for a specific task are typically learnt from training sets that encode relations between features and target variable in a representative fashion. One example from the fairness literature is Large Movie Review, used to train sentiment analysis models, later audited for fairness [318]. For fairness audits, one alternative would be resorting to publicly available models, but sometimes a close control on training corpus and procedure is necessary. Indeed, it is interesting to study issues of model fairness in relation to biases present in the respective training corpora, which can help explain the causes of bias [62]. Some works measure biases in representations before and after fine-tuning on a training set and regard the difference as a measure of bias in the training set. Babaeianjelodar et al. [19] employ this approach to measure biases in RtGender, Civil Comments and datasets from GLUE.

A representative summary of a service. Much important work in the fairness literature is focused on measuring fairness and harms in the real world. This line of work includes audits of products and services, which rely on datasets extracted from the application of interest. Datasets created for this purpose include Amazon Recommendations, Pymetrics Bias Group, Occupations in Google Images, Zillow Searches, Online Freelance Marketplaces, Bing US Queries, YouTube Dialect Accuracy. Several other datasets were originally created for this purpose and later repurposed in the fairness literature as sources of real data, including Stop Question and Frisk, HMDA, Law School, and COMPAS.

An important source of data. Some datasets acquire a pivotal role in research and industry, to the point of being considered a de-facto standard for a given purpose. This status warrants closer scrutiny of the dataset, through which researchers aim to uncover potential biases and problematic aspects that may impact models and insights derived from the dataset. ImageNet, for instance, is a dataset with millions of images across thousands of categories. Since its release in 2011, this resource has been used to train, benchmark and compare hundreds of computer vision models. Given its status in machine learning research, ImageNet has been the subject of two quantitative investigations analyzing its biases and other problematic aspects in the person subtree, uncovering issues of representation [552] and

non-consensuality [406]. A different data bias audit was carried out on SafeGraph Research Release. SafeGraph data captures mobility patterns in the US, with data from nearly 50 million mobile devices obtained and maintained by Safegraph, a private data company. Their recent academic release has become a fundamental resource for pandemic research, to the point of being used by the Centers for Disease Control and Prevention to measure the effectiveness of social distancing measures [371]. To evaluate its representativeness for the overall US population, Coston et al. [115] have studied selection biases in this dataset.

In algorithmic fairness research, datasets play similar roles to the ones they play in machine learning according to Barocas et al. [33], including training, catalyzing attention, and signalling awareness of common data practices. One notable exception is that fairness datasets are not used to track algorithmic progress on a problem over time, likely due to the fact that there is no consensus on a single measure to be reported. On the other hand, two roles peculiar to fairness research are summarizing a service or product that is being audited, and representing an important resource whose biases and ethical aspects are particularly worthy of attention. We note that these roles are not mutually exclusive and that datasets can play multiple roles. COMPAS, for example, was originally curated to perform an audit of pretrial risk assessment tools and was later used extensively in fair ML research as a source of real human data, becoming, overall, a catalyst for fairness research and debate.

In sum, existing fairness datasets originate from a variety of domains, support diverse tasks, and play different roles in the algorithmic fairness literature. We hope our work will contribute to establishing principled data practices in the field, to guide an optimal usage of these resources.

6 CONCLUSIONS

Algorithmic fairness is a young research area, undergoing a fast expansion, with diverse contributions in terms of methodology and applications. Progress in the field hinges on different resources, including, very prominently, datasets. In this work, we have surveyed hundreds of datasets used in the fair ML and algorithmic equity literature to help the research community reduce its documentation debt, identify gaps, and improve the utilization of existing resources.

We have rigorously identified the most popular datasets in the literature, and carried out a thorough documentation effort for the three most popular ones (Adult, COMPAS and German Credit). Our work unifies and adds to recent literature on data studies, calling into question their current status of general-purpose fairness benchmarks, due to contrived prediction tasks, noisy data, severe coding mistakes, limitations in encoding sensitive attributes, and age. In a practical demonstration of documentation debt and its consequences, we find several works of algorithmic fairness using German Credit with sex as a protected attribute, while careful analysis of recent documentation shows that this feature cannot be reliably retrieved from the data.

We have documented over two hundred datasets to provide viable alternatives, annotating their domain, the tasks they support and discussing the roles they play in works of algorithmic fairness. We have shown that the processes generating the data belong to

many different domains, including, for instance, criminal justice, education, search engines, online marketplaces, emergency response, social media, medicine, hiring, and finance. At the same time, we have described a variety of tasks studied on these resources, ranging from generic, such as *fair regression*, to narrow such as *fair districting* and *fair truth discovery*. Overall, such diversity of domains and tasks provides a glimpse into the variety of human activities and applications that can be impacted by automated decision making, and that can benefit from fair ML and algorithmic equity research.

Dataset tasks, domains, and the whole metadata are made available in our data briefs (Appendix A), which we plan to update on a yearly basis. We envision several benefits for the algorithmic equity and data studies research communities, including (1) informing the choice of datasets for experimental evaluations of fair algorithms, including domain-oriented and task-oriented search, (2) directing studies of data bias, and other quantitative and qualitative analyses, including retrospective documentation efforts, towards popular (or otherwise important) resources, (3) identifying areas and sub-problems that are understudied in the algorithmic fairness literature, and (4) supporting multi-dataset studies, focused on resources united by a common characteristic, such as encoding a given sensitive attribute [448], concerning computer vision [158], or being popular in the fairness literature [302].

In this work, we have targeted the collective documentation debt of the algorithmic fairness community, resulting from the opacity surrounding certain resources and the sparsity of existing documentation. We have mainly targeted sparsity in a centralized documentation effort. Similarly to other types of data interventions, useful documentation can be produced after release, but, as shown in this work, the documentation debt may propagate nonetheless. In a mature research community, curators, users and reviewers can all contribute to cultivating a data documentation culture and keep the overall documentation debt in check.

ACKNOWLEDGMENTS

The authors would like to thank the following researchers and dataset creators for the useful feedback on the data briefs: Alain Barrat, Luc Behaghel, Asia Biega, Marko Bohanec, Chris Burgess, Robin Burke, Alejandro Noriega Campero, Margarida Carvalho, Abhijnan Chakraborty, Robert Cheetham, Paulo Cortez, Thomas Davidson, Maria De-Arteaga, Lucas Dixon, Danijela Djordjević, Michele Donini, Marco Duarte, Natalie Ebner, Fehrman, H. Altay Guvenir, Moritz Hardt, Yu Hen Hu, Irina Higgins, Won Ik Cho, Rachel Huddart, Lalana Kagal, Dean Karlan, Vijay Keswani, Been Kim, Hyunjik Kim, Jiwon Kim, Svetlana Kiritchenko, Joseph A. Konstan, Varun Kumar, Jeremy Andrew Irvin, Jamie N. Larson, Jure Leskovec, Andrea Lodi, Oisín Mac Aodha, Loic Matthey, Julian McAuley, Brendan McMahan, Sergio Moro, Luca Oneto, Orestis Papakyriakopoulos, Stephen Robert Pfohl, Christopher G. Potts, Mike Redmond, Kit Rodolfa, Ben Roshan, Veronica Rotemberg, Rachel Rudinger, Sivan Sabato, Kate Saenko, Mark D. Shermis, Daniel Slunge, David Solans, Luca Soldaini, Efstathios Stamatatos, Ryan Steed, Rachael Tatman, Schrasing Tong, Alan Tsang, Sathishkumar V E, Andreas van Cranenburgh, Lucy Vasserman, Roland Vollgraf, Alex Wang, Zeerak Waseem, Kellie Webster, Pang Wei Koh, Bryan

Wilder, Nick Wilson, I-Cheng Yeh, Elad Yom-Tov, Neil Yorke-Smith, Michal Zabovskyy, Yukun Zhu.

REFERENCES

- [1] Mohsen Abbasi, Aditya Bhaskara, and Suresh Venkatasubramanian. 2021. Fair Clustering via Equitable Group Representations. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 504–514. <https://doi.org/10.1145/3442188.3445913>
- [2] Robert Adragna, Elliot Creager, David Madras, and Richard Zemel. 2020. Fairness and Robustness in Invariant Learning: A Case Study in Toxicity Classification. arXiv:2011.06485 [cs.LG] NeurIPS 2020 workshop: "Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)".
- [3] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. 2018. A Reductions Approach to Fair Classification. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholmsmässan, Stockholm Sweden, 60–69. <http://proceedings.mlr.press/v80/agarwal18a.html>
- [4] Alekh Agarwal, Miroslav Dudík, and Zhiwei Steven Wu. 2019. Fair Regression: Quantitative Definitions and Reduction-Based Algorithms. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 120–129. <http://proceedings.mlr.press/v97/agarwal19d.html>
- [5] Monica Agrawal, Marinka Zitnik, Jure Leskovec, et al. 2018. Large-scale analysis of disease pathways in the human interactome. In *PSB*. World Scientific, 111–122.
- [6] Sara Ahmadian, Alessandro Epasto, Marina Knittel, Ravi Kumar, Mohammad Mahdian, Benjamin Moseley, Philip Pham, Sergei Vassilvitskii, and Yuyan Wang. 2020. Fair Hierarchical Clustering. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/f10f2da9a238b746d2bac55759915f0d-Abstract.html>
- [7] Osman Aka, Ken Burke, Alex Bauerle, Christina Greer, and Margaret Mitchell. 2021. *Measuring Model Biases in the Absence of Ground Truth*. Association for Computing Machinery, New York, NY, USA, 327–335. <https://doi.org/10.1145/3461702.3462557>
- [8] Junaid Ali, Mahmoudreza Babaei, Abhijnan Chakraborty, Baharan Mirza-soleiman, Krishna P. Gummadi, and Adish Singla. 2019. On the Fairness of Time-Critical Influence Maximization in Social Networks. arXiv:1905.06618 [cs.SI] NeurIPS 2019 workshop: "Human-Centric Machine Learning".
- [9] Junaid Ali, Preethi Lahoti, and Krishna P. Gummadi. 2021. *Accounting for Model Uncertainty in Algorithmic Discrimination*. Association for Computing Machinery, New York, NY, USA, 336–345. <https://doi.org/10.1145/3461702.3462630>
- [10] Junaid Ali, Muhammad Bilal Zafar, Adish Singla, and Krishna P. Gummadi. 2019. Loss-Aversively Fair Classification. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society (Honolulu, HI, USA) (AI/ES '19)*. Association for Computing Machinery, New York, NY, USA, 211–218. <https://doi.org/10.1145/3306618.3314266>
- [11] Alexander Amini, Ava P. Soleimany, Wilko Schwarting, Sangeeta N. Bhatia, and Daniela Rus. 2019. Uncovering and Mitigating Algorithmic Bias through Learned Latent Structure. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society (Honolulu, HI, USA) (AI/ES '19)*. Association for Computing Machinery, New York, NY, USA, 289–295. <https://doi.org/10.1145/3306618.3314243>
- [12] Edgar Anderson. 1936. The species problem in Iris. *Annals of the Missouri Botanical Garden* 23, 3 (1936), 457–509.
- [13] McKane Andrus, Elena Spitzer, Jeffrey Brown, and Alice Xiang. 2021. What We Can't Measure, We Can't Understand: Challenges to Demographic Data Procurement in the Pursuit of Fairness. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 249–260. <https://doi.org/10.1145/3442188.3445888>
- [14] Ralph G Andrzejak, Klaus Lehnertz, Florian Mormann, Christoph Rieke, Peter David, and Christian E Elger. 2001. Indications of nonlinear deterministic and finite-dimensional structures in time series of brain electrical activity: Dependence on recording region and brain state. *Physical Review E* 64, 6 (2001), 061907.
- [15] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine Bias. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

- [16] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. 2020. Invariant Risk Minimization. arXiv:1907.02893 [stat.ML]
- [17] James Atwood, Hansa Srinivasan, Yoni Halpern, and D Sculley. 2019. Fair treatment allocations in social networks. arXiv:1911.05489 [cs.SI] NeurIPS 2019 workshop: "Fair ML for Health".
- [18] Pranjal Awasthi, Alex Beutel, Matthäus Kleindessner, Jamie Morgenstern, and Xuezhi Wang. 2021. Evaluating Fairness of Machine Learning Models Under Uncertain and Incomplete Information. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (FAccT '21). Association for Computing Machinery, New York, NY, USA, 206–214. <https://doi.org/10.1145/3442188.3445884>
- [19] Marzieh Bahaianjelodar, Stephen Lorenz, Josh Gordon, Jeanna Matthews, and Evan Freitag. 2020. Quantifying Gender Bias in Different Corpora. In *Companion Proceedings of the Web Conference 2020* (Taipei, Taiwan) (WWW '20). Association for Computing Machinery, New York, NY, USA, 752–759. <https://doi.org/10.1145/3366424.3383559>
- [20] Moshe Babaioff, Noam Nisan, and Inbal Talgam-Cohen. 2019. Fair Allocation through Competitive Equilibrium from Generic Incomes. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (Atlanta, GA, USA) (FAccT '19). Association for Computing Machinery, New York, NY, USA, 180. <https://doi.org/10.1145/3287560.3287582>
- [21] Arturs Backurs, Piotr Indyk, Krzysztof Onak, Baruch Schieber, Ali Vakilian, and Tal Wagner. 2019. Scalable Fair Clustering. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 405–413. <http://proceedings.mlr.press/v97/backurs19a.html>
- [22] Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. 2019. Differential Privacy Has Disparate Impact on Model Accuracy. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2019/file/fc0de4e0396fff257ea362983c2dda5a-Paper.pdf>
- [23] Sina Baharlouei, Maher Nouiehed, Ahmad Beirami, and Meisam Razaviyayn. 2020. Rényi Fair Inference. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=HkgsUjrtDB>
- [24] Michiel A. Bakker, Duy Patrick Tu, Krishna P. Gummadi, Alex Sandy Pentland, Kush R. Varshney, and Adrian Weller. 2021. *Beyond Reasonable Doubt: Improving Fairness in Budget-Constrained Decision Making Using Confidence Thresholds*. Association for Computing Machinery, New York, NY, USA, 346–356. <https://doi.org/10.1145/3461702.3462575>
- [25] Michiel A. Bakker, Duy Patrick Tu, Humberto Riverón Valdés, Krishna P. Gummadi, Kush R. Varshney, Adrian Weller, and Alex Pentland. 2019. DADI: Dynamic Discovery of Fair Information with Adversarial Reinforcement Learning. arXiv:1910.13983 [cs.LG] NeurIPS 2019 workshop: "Human-Centric Machine Learning".
- [26] Ari Ball-Burack, Michelle Seng Ah Lee, Jennifer Cobbe, and Jatinder Singh. 2021. Differential Tweetment: Mitigating Racial Dialect Bias in Harmful Tweet Detection. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (FAccT '21). Association for Computing Machinery, New York, NY, USA, 116–128. <https://doi.org/10.1145/3442188.3445875>
- [27] Jack Bandy and Nicholas Vincent. 2021. Addressing "Documentation Debt" in Machine Learning Research: A Retrospective Datasheet for BookCorpus. arXiv preprint arXiv:2105.05241 (2021).
- [28] Michelle Bao, Angela Zhou, Samantha Zottola, Brian Brubach, Sarah Desmarais, Aaron Horowitz, Kristian Lum, and Suresh Venkatasubramanian. 2021. It's COMPASlicated: The Messy Relationship between RAI Datasets and Algorithmic Fairness Benchmarks. arXiv preprint arXiv:2106.05498 (2021).
- [29] Chelsea Barabas, Karthik Dinakar, and Colin Doyle. 2019. The Problems With Risk Assessment Tools. <https://www.nytimes.com/2019/07/17/opinion/pretrial-ai.html>
- [30] Michael Barbaro. 2007. In Apparel, All Tariffs Aren't Created Equal. <https://www.nytimes.com/2007/04/28/business/28gender.html>
- [31] Matias Barenstein. 2019. ProPublica's COMPAS Data Revisited. arXiv preprint arXiv:1906.04711 (2019).
- [32] Anamika Barman-Adhikari, Stephanie Begun, Eric Rice, Amanda Yoshioka-Maxwell, and Andrea Perez-Portillo. 2016. Sociometric network structure and its association with methamphetamine use norms among homeless youth. *Social science research* 58 (2016), 292–308.
- [33] Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2019. *Fairness and Machine Learning*. fairmlbook.org. <http://www.fairmlbook.org>.
- [34] Jean-Patrick Baudry, Margarida Cardoso, Gilles Celeux, Maria José Amorim, and Ana Sousa Ferreira. 2015. Enhancing the selection of a model-based clustering with external categorical variables. *Advances in Data Analysis and Classification* 9, 2 (2015), 177–196.
- [35] Luc Behaghel, Bruno Crépon, and Marc Gurgand. 2014. Private and Public Provision of Counseling to Job Seekers: Evidence from a Large Controlled Experiment. *American Economic Journal: Applied Economics* 6, 4 (October 2014), 142–74. <https://doi.org/10.1257/app.6.4.142>
- [36] Clara Belitz, Lan Jiang, and Nigel Bosch. 2021. *Automating Procedurally Fair Feature Selection in Machine Learning*. Association for Computing Machinery, New York, NY, USA, 379–389. <https://doi.org/10.1145/3461702.3462585>
- [37] Emily M. Bender and Batya Friedman. 2018. Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science. *Transactions of the Association for Computational Linguistics* 6 (2018), 587–604. https://doi.org/10.1162/tacl_a_00041
- [38] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? (FAccT '21). Association for Computing Machinery, New York, NY, USA, 610–623. <https://doi.org/10.1145/3442188.3445922>
- [39] Suman Bera, Deeparnab Chakrabarty, Nicolas Flores, and Maryam Negahbani. 2019. Fair Algorithms for Clustering. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 4954–4965. <https://proceedings.neurips.cc/paper/2019/file/fc192b0c0d270dbf41870a63a8c76c2f-Paper.pdf>
- [40] Elena Beretta, Antonio Vetro, Bruno Lepri, and Juan Carlos De Martin. 2021. Detecting Discriminatory Risk through Data Annotation Based on Bayesian Inferences. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (FAccT '21). Association for Computing Machinery, New York, NY, USA, 794–804. <https://doi.org/10.1145/3442188.3445940>
- [41] Richard Berk, Hoda Heidari, Shahin Jabbari, Matthew Joseph, Michael Kearns, Jamie Morgenstern, Seth Neel, and Aaron Roth. 2017. A Convex Framework for Fair Regression. arXiv:1706.02409 [cs.LG] KDD 2017 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [42] Thierry Bertin-Mahieux, Daniel P. W. Ellis, Brian Whitman, and Paul Lamere. 2011. The Million Song Dataset. In *Proceedings of the 12th International Society for Music Information Retrieval Conference*. ISMIR, Miami, United States, 591–596. <https://doi.org/10.5281/zenodo.1415820>
- [43] Alex Beutel, Jilin Chen, Zhe Zhao, and Ed H. Chi. 2017. Data Decisions and Theoretical Implications when Adversarially Learning Fair Representations. arXiv:1707.00075 [cs.LG] KDD 2017 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [44] Arpita Biswas and Suvam Mukherjee. 2021. *Ensuring Fairness under Prior Probability Shifts*. Association for Computing Machinery, New York, NY, USA, 414–424. <https://doi.org/10.1145/3461702.3462596>
- [45] Emily Black and Matt Fredrikson. 2021. Leave-One-out Unfairness. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (FAccT '21). Association for Computing Machinery, New York, NY, USA, 285–295. <https://doi.org/10.1145/3442188.3445894>
- [46] Emily Black, Samuel Yeom, and Matt Fredrikson. 2020. FlipTest: Fairness Testing via Optimal Transport. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (FAccT '20). Association for Computing Machinery, New York, NY, USA, 111–121. <https://doi.org/10.1145/3351095.3372845>
- [47] Su Lin Blodgett and Brendan O'Connor. 2017. Racial Disparity in Natural Language Processing: A Case Study of Social Media African-American English. arXiv:1707.00061 [cs.CY] KDD 2017 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [48] Su Lin Blodgett, Lisa Green, and Brendan O'Connor. 2016. Demographic Dialectal Variation in Social Media: A Case Study of African-American English. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Austin, Texas, 1119–1130. <https://doi.org/10.18653/v1/D16-1120>
- [49] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. In *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Eds.), Vol. 29. Curran Associates, Inc., 4349–4357. <https://proceedings.neurips.cc/paper/2016/file/a486cd07e4ac3d270571622f4f316ce5-Paper.pdf>
- [50] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating Embeddings for Modeling Multi-relational Data. In *Advances in Neural Information Processing Systems*, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger (Eds.), Vol. 26. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2013/file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf>
- [51] Shikha Bordia and Samuel R. Bowman. 2019. Identifying and Reducing Gender Bias in Word-Level Language Models. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*. Association for Computational Linguistics, Minneapolis, Minnesota, 7–15. <https://doi.org/10.18653/v1/N19-3002>
- [52] Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2019. Nuanced Metrics for Measuring Unintended Bias with Real Data for

- Text Classification. In *Companion Proceedings of The 2019 World Wide Web Conference* (San Francisco, USA) (WWW '19). Association for Computing Machinery, New York, NY, USA, 491–500. <https://doi.org/10.1145/3308560.3317593>
- [53] Avishk Bose and William Hamilton. 2019. Compositional Fairness Constraints for Graph Embeddings. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 715–724. <http://proceedings.mlr.press/v97/bose19a.html>
- [54] Amanda Bower, Hamid Eftekhari, Mikhail Yurochkin, and Yuekai Sun. 2021. Individually Fair Rankings. In *International Conference on Learning Representations*. https://openreview.net/forum?id=71zCSP_HuBN
- [55] Amanda Bower, Laura Niss, Yuekai Sun, and Alexander Vargo. 2018. Debiasing representations by removing unwanted variation due to protected attributes. arXiv:1807.00461 [cs.CY] ICML 2018 workshop: “Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)”.
- [56] Tim Brennan, William Dieterich, and Beate Ehret. 2009. Evaluating the Predictive Validity of the Compas Risk and Needs Assessment System. *Criminal Justice and Behavior* 36, 1 (2009), 21–40. <https://doi.org/10.1177/0093854808326545> arXiv:https://doi.org/10.1177/0093854808326545
- [57] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. Openai gym. *arXiv preprint arXiv:1606.01540* (2016).
- [58] Jeanne Brooks-Gunn, Fong-ruey Liaw, and Pamela Kato Klebanov. 1992. Effects of early intervention on cognitive function of low birth weight preterm infants. *The Journal of pediatrics* 120, 3 (1992), 350–359.
- [59] Lukas Brozovsky and V. Petricek. 2007. Recommender System for Online Dating Service. *ArXiv abs/cs/0703042* (2007).
- [60] Lukáš Brožovský. 2006. *Recommender System for a Dating Service*. Master’s thesis, Charles University in Prague, Prague, Czech Republic. <http://colfi.wz.cz/colfi.pdf>
- [61] Brian Brubach, Darshan Chakrabarti, John Dickerson, Samir Khuller, Aravind Srinivasan, and Leonidas Tsipenak. 2020. A Pairwise Fair and Community-preserving Approach to k-Center Clustering. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 1178–1189. <http://proceedings.mlr.press/v119/brubach20a.html>
- [62] Marc-Etienne Brunet, Colleen Alkalay-Houlihan, Ashton Anderson, and Richard Zemel. 2019. Understanding the Origins of Bias in Word Embeddings. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 803–811. <http://proceedings.mlr.press/v97/brunet19a.html>
- [63] Joy Buolamwini and Timnit Gebru. 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A. Friedler and Christo Wilson (Eds.). PMLR, New York, NY, USA, 77–91. <http://proceedings.mlr.press/v81/buolamwini18a.html>
- [64] Robin Burke, Jackson Kontny, and Nasim Sonboli. 2018. Synthetic Attribute Data for Evaluating Consumer-side Fairness. arXiv:1809.04199 [cs.CY] RecSys 2018 workshop: “Workshop on Responsible Recommendation (FAT/Rec)”.
- [65] Robin Burke, Nasim Sonboli, and Aldo Ordonez-Gauger. 2018. Balanced Neighborhoods for Multi-sided Fairness in Recommendation. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A. Friedler and Christo Wilson (Eds.). PMLR, New York, NY, USA, 202–214. <http://proceedings.mlr.press/v81/burke18a.html>
- [66] Maarten Buyl and Tijl De Be. 2020. DeBayes: a Bayesian Method for Debiasing Network Embeddings. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 1220–1229. <http://proceedings.mlr.press/v119/buyl20a.html>
- [67] William Cai, Johann Gaebler, Nikhil Garg, and Sharad Goel. 2020. Fair Allocation through Selective Information Acquisition. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society* (New York, NY, USA) (AI/ES '20). Association for Computing Machinery, New York, NY, USA, 22–28. <https://doi.org/10.1145/3375627.3375823>
- [68] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. 2018. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097* (2018).
- [69] Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. 2009. Building Classifiers with Independence Constraints. In *2009 IEEE International Conference on Data Mining Workshops*. 13–18. <https://doi.org/10.1109/ICDMW.2009.83>
- [70] Toon Calders and Sicco Verwer. 2010. Three Naive Bayes Approaches for Discrimination-Free Classification. *Data Min. Knowl. Discov.* 21, 2 (Sept. 2010), 277–292. <https://doi.org/10.1007/s10618-010-0190-x>
- [71] Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356, 6334 (14 April 2017), 183–186. <https://doi.org/10.1126/science.aal4230>
- [72] Flavio Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R Varshney. 2017. Optimized Pre-Processing for Discrimination Prevention. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., 3992–4001. <https://proceedings.neurips.cc/paper/2017/file/9a49a25d845a483fae4be7e341368e36-Paper.pdf>
- [73] Ran Canetti, Aloni Cohen, Nishanth Dikkala, Govind Ramnarayan, Sarah Schefler, and Adam Smith. 2019. From Soft Classifiers to Hard Decisions: How Fair Can We Be?. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (Atlanta, GA, USA) (FAT* '19). Association for Computing Machinery, New York, NY, USA, 309–318. <https://doi.org/10.1145/3287560.3287561>
- [74] Ioannis Caragiannis, David Kurokawa, Hervé Moulin, Ariel D. Procaccia, Nisarg Shah, and Junxing Wang. 2016. The Unreasonable Fairness of Maximum Nash Welfare. In *Proceedings of the 2016 ACM Conference on Economics and Computation* (Maastricht, The Netherlands) (EC '16). Association for Computing Machinery, New York, NY, USA, 305–322. <https://doi.org/10.1145/2940716.2940726>
- [75] Rodrigo L. Cardoso, Wagner Meira Jr., Virgilio Almeida, and Mohammed J. Zaki. 2019. A Framework for Benchmarking Discrimination-Aware Models in Machine Learning. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society* (Honolulu, HI, USA) (AI/ES '19). Association for Computing Machinery, New York, NY, USA, 437–444. <https://doi.org/10.1145/3306618.3314262>
- [76] Margarida Carvalho and Andrea Lodi. 2019. Game theoretical analysis of kidney exchange programs. *arXiv preprint arXiv:1911.09207* (2019).
- [77] Simon Caton and Christian Haas. 2020. Fairness in machine learning: A survey. *arXiv preprint arXiv:2010.04053* (2020).
- [78] Elisa Celis, Vijay Keswani, Damian Straszak, Amit Deshpande, Tarun Kathuria, and Nisheeth Vishnoi. 2018. Fair and Diverse DPP-Based Data Summarization. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholm, Sweden, 716–725. <http://proceedings.mlr.press/v80/celis18a.html>
- [79] Elisa Celis, Anay Mehrotra, and Nisheeth Vishnoi. 2019. Toward Controlling Discrimination in Online Ad Auctions. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 4456–4465. <http://proceedings.mlr.press/v97/mehrotra19a.html>
- [80] L. Elisa Celis, Amit Deshpande, Tarun Kathuria, and Nisheeth K. Vishnoi. 2016. How to be Fair and Diverse? arXiv:1610.07183 [cs.LG] DTL 2016 workshop: “Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)”.
- [81] L. Elisa Celis, Lingxiao Huang, Vijay Keswani, and Nisheeth K. Vishnoi. 2019. Classification with Fairness Constraints: A Meta-Algorithm with Provable Guarantees. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (Atlanta, GA, USA) (FAT* '19). Association for Computing Machinery, New York, NY, USA, 319–328. <https://doi.org/10.1145/3287560.3287586>
- [82] L. Elisa Celis and Vijay Keswani. 2020. Implicit Diversity in Image Summarization. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW2 (Oct 2020), 1–28. <https://doi.org/10.1145/3415210>
- [83] L. Elisa Celis, Vijay Keswani, and Nisheeth Vishnoi. 2020. Data preprocessing to mitigate bias: A maximum entropy based approach. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 1349–1359. <http://proceedings.mlr.press/v119/celis20a.html>
- [84] L. Elisa Celis, Anay Mehrotra, and Nisheeth K. Vishnoi. 2020. Interventions for Ranking in the Presence of Implicit Bias. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (FAT* '20). Association for Computing Machinery, New York, NY, USA, 369–380. <https://doi.org/10.1145/3351095.3372858>
- [85] O. Celma. 2010. *Music Recommendation and Discovery in the Long Tail*. Springer.
- [86] Elias Chaibub Neto. 2020. A Causal Look at Statistical Definitions of Discrimination. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Virtual Event, CA, USA) (KDD '20). Association for Computing Machinery, New York, NY, USA, 873–881. <https://doi.org/10.1145/3394486.3403130>
- [87] Abhijnan Chakraborty, Gourab K. Patro, Niloy Ganguly, Krishna P. Gummadi, and Patrick Loiseau. 2019. Equality of Voice: Towards Fair Representation in Crowdsourced Top-K Recommendations. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (Atlanta, GA, USA) (FAT* '19). Association for Computing Machinery, New York, NY, USA, 129–138. <https://doi.org/10.1145/3287560.3287570>
- [88] Olivier Chapelle and Yi Chang. 2010. Yahoo! Learning to Rank Challenge Overview. In *Proceedings of the 2010 International Conference on Yahoo! Learning to Rank Challenge - Volume 14* (Haifa, Israel) (YLRC'10). JMLR.org, 1–24.
- [89] Harshal A. Chaudhari, Sangdi Lin, and Ondrej Linda. 2020. A General Framework for Fairness in Multistakeholder Recommendations. arXiv:2009.02423 [cs.AI] RecSys 2020 workshop: “3rd FAccTRec Workshop on Responsible Recommendation”.
- [90] Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Philipp Koehn, and Tony Robinson. 2014. One Billion Word Benchmark for Measuring

- Progress in Statistical Language Modeling. In *INTERSPEECH-2014*.
- [91] Binghui Chen, Weihong Deng, and Haifeng Shen. 2018. Virtual Class Enhanced Discriminative Embedding Learning. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2018/file/d79aac075930c83c2f1e369a511148fe-Paper.pdf>
- [92] Ching-Wei Chen, Paul Lamere, Markus Schedl, and Hamed Zamani. 2018. Recsys Challenge 2018: Automatic Music Playlist Continuation (*RecSys '18*). Association for Computing Machinery, New York, NY, USA, 527–528. <https://doi.org/10.1145/3240323.3240342>
- [93] Irene Chen, Fredrik D Johansson, and David Sontag. 2018. Why Is My Classifier Discriminatory?. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2018/file/1f1baa5b8edac74eb4ea329f14a0361-Paper.pdf>
- [94] Jiahao Chen, Nathan Kallus, Xiaojie Mao, Geoffrey Svacha, and Madeleine Udell. 2019. Fairness Under Unawareness: Assessing Disparity When Protected Class Is Unobserved. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (Atlanta, GA, USA) (*FAT* '19*). Association for Computing Machinery, New York, NY, USA, 339–348. <https://doi.org/10.1145/3287560.3287594>
- [95] Xingyu Chen, Brandon Fain, Liang Lyu, and Kamesh Munagala. 2019. Proportionally Fair Clustering. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 1032–1041. <http://proceedings.mlr.press/v97/chen19d.html>
- [96] Yan Chen, Christopher Mahoney, Isabella Grasso, Esmā Wali, Abigail Matthews, Thomas Middleton, Mariama Njie, and Jeanna Matthews. 2021. *Gender Bias and Under-Representation in Natural Language Processing Across Human Languages*. Association for Computing Machinery, New York, NY, USA, 24–34. <https://doi.org/10.1145/3461702.3462530>
- [97] Pengyu Cheng, Weituo Hao, Siyang Yuan, Shijing Si, and Lawrence Carin. 2021. FairFil: Contrastive Neural Debiasing Method for Pretrained Text Encoders. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=N6JECd-PI5w>
- [98] Victoria Cheng, Vinith M. Suriyakumar, Natalie Dullerud, Shalmali Joshi, and Marzyeh Ghassemi. 2021. Can You Fake It Until You Make It? Impacts of Differentially Private Synthetic Data on Downstream Classification Fairness. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (*FAccT '21*). Association for Computing Machinery, New York, NY, USA, 149–160. <https://doi.org/10.1145/3442188.3445879>
- [99] Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, and Sergei Vassilvitskii. 2017. Fair Clustering Through Fairlets. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., 5029–5037. <https://proceedings.neurips.cc/paper/2017/file/978fce5bcc4ecc88ad48ce3914124a2-Paper.pdf>
- [100] Ashish Chiplunkar, Sagar Kale, and Sivaramkrishnan Natarajan Ramamoorthy. 2020. How to Solve Fair k-Center in Massive Data Models. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 1877–1886. <http://proceedings.mlr.press/v119/chiplunkar20a.html>
- [101] Jaewoong Cho, Gyeongjo Hwang, and Changho Suh. 2020. A Fair Classifier Using Kernel Density Estimation. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 15088–15099. <https://proceedings.neurips.cc/paper/2020/file/ac3870fcad1cfc367825cda0101ee62-Paper.pdf>
- [102] Won Ik Cho, Jiwon Kim, Jaeyoung Yang, and Nam Soo Kim. 2021. Towards Cross-Lingual Generalization of Translation Gender Bias. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (*FAccT '21*). Association for Computing Machinery, New York, NY, USA, 449–457. <https://doi.org/10.1145/3442188.3445907>
- [103] Kristy Choi, Aditya Grover, Trisha Singh, Rui Shu, and Stefano Ermon. 2020. Fair Generative Modeling via Weak Supervision. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 1887–1898. <http://proceedings.mlr.press/v119/choi20a.html>
- [104] Yoojung Choi, Meihua Dang, and Guy Van den Broeck. 2020. Group Fairness by Probabilistic Modeling with Latent Fair Decisions. arXiv:2009.09031 [cs.LG] NeurIPS 2020 workshop: “Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)”.
- [105] A. Chouldechova. 2017. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data* 5 2 (2017), 153–163.
- [106] Alexandra Chouldechova and Max G'Sell. 2017. Fairer and more accurate, but for whom? arXiv:1707.00046 [stat.AP] KDD 2017 workshop: “Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)”.
- [107] Alexandra Chouldechova and Aaron Roth. 2020. A Snapshot of the Frontiers of Fairness in Machine Learning. *Commun. ACM* 63, 5 (April 2020), 82–89. <https://doi.org/10.1145/3376898>
- [108] Ching-Yao Chuang and Youssef Mroueh. 2021. Fair Mixup: Fairness via Interpolation. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=DN15s5BXeBn>
- [109] Evgenii Chzhen, Christophe Denis, Mohamed Hebiri, Luca Oneto, and Massimiliano Pontil. 2019. Leveraging Labeled and Unlabeled Data for Consistent Fair Binary Classification. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 12760–12770. <https://proceedings.neurips.cc/paper/2019/file/ba51e6158bca8f0d0d834950251e693-Paper.pdf>
- [110] Evgenii Chzhen, Christophe Denis, Mohamed Hebiri, Luca Oneto, and Massimiliano Pontil. 2020. Fair regression via plug-in estimator and recalibration with statistical guarantees. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/ddd808772c035aed16d42ad3559be5f-Abstract.html>
- [111] Evgenii Chzhen, Christophe Denis, Mohamed Hebiri, Luca Oneto, and Massimiliano Pontil. 2020. Fair regression with Wasserstein barycenters. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 7321–7331. <https://proceedings.neurips.cc/paper/2020/file/51cdbd2611e844ece5d80878eb770436-Paper.pdf>
- [112] Sharon R Cohany, Anne E Polivka, and Jennefer M Rothgeb. 1994. Revisions in the current population survey effective January 1994. *Emp. & Earnings* 41 (1994), 13.
- [113] Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. 2017. Algorithmic Decision Making and the Cost of Fairness. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (Halifax, NS, Canada) (*KDD '17*). Association for Computing Machinery, New York, NY, USA, 797–806. <https://doi.org/10.1145/3097983.3098095>
- [114] P. Cortez and A. M. G. Silva. 2008. Using data mining to predict secondary school student performance. In *Proceedings of 5th Future Business Technology Conference*.
- [115] Amanda Coston, Neel Guha, Derek Ouyang, Lisa Lu, Alexandra Chouldechova, and Daniel E. Ho. 2021. Leveraging Administrative Data for Bias Audits: Assessing Disparate Coverage with Mobility Data for COVID-19 Policy. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (*FAccT '21*). Association for Computing Machinery, New York, NY, USA, 173–184. <https://doi.org/10.1145/3442188.3445881>
- [116] Amanda Coston, Alan Mishler, Edward H. Kennedy, and Alexandra Chouldechova. 2020. Counterfactual Risk Assessments, Evaluation, and Fairness. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (*FAT* '20*). Association for Computing Machinery, New York, NY, USA, 582–593. <https://doi.org/10.1145/3351095.3372851>
- [117] Amanda Coston, Karthikeyan Natesan Ramamurthy, Dennis Wei, Kush R. Varshney, Skyler Speakman, Zairah Mustahsan, and Supriyo Chakraborty. 2019. Fair Transfer Learning with Missing Protected Attributes. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society* (Honolulu, HI, USA) (*AIES '19*). Association for Computing Machinery, New York, NY, USA, 91–98. <https://doi.org/10.1145/3306618.3314236>
- [118] Andrew Cotter, Maya Gupta, Heinrich Jiang, Nathan Srebro, Karthik Sridharan, Serena Wang, Blake Woodworth, and Seungil You. 2018. Training Fairness-Constrained Classifiers to Generalize. ICML 2018 workshop: “Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)”.
- [119] Andrew Cotter, Maya Gupta, Heinrich Jiang, Nathan Srebro, Karthik Sridharan, Serena Wang, Blake Woodworth, and Seungil You. 2019. Training Well-Generalizing Classifiers for Fairness Metrics and Other Data-Dependent Constraints. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 1397–1405. <http://proceedings.mlr.press/v97/cotter19b.html>
- [120] Kate Crawford and Trevor Paglen. 2021. Excavating AI: the Politics of Images in Machine Learning Training Sets. <https://excavating.ai/>
- [121] Elliot Creager, Joern-Henrik Jacobsen, and Richard Zemel. 2021. Exchanging Lessons Between Algorithmic Fairness and Domain Generalization. <https://openreview.net/forum?id=DC1Im3MkGG> NeurIPS 2020 workshop: “Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)”.
- [122] Elliot Creager, David Madras, Joern-Henrik Jacobsen, Marissa Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. 2019. Flexibly Fair Representation Learning by Disentanglement. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 1436–1445. <http://proceedings.mlr.press/v97/creager19a.html>

- [123] Elliot Creager, David Madras, Toniann Pitassi, and Richard Zemel. 2020. Causal Modeling for Fairness In Dynamical Systems. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Duménil III and Aarti Singh (Eds.). PMLR, Virtual, 2185–2195. <http://proceedings.mlr.press/v119/creager20a.html>
- [124] Alexander D'Amour, Hansa Srinivasan, James Atwood, Pallavi Baljekar, D. Sculley, and Yoni Halpern. 2020. Fairness is Not Static: Deeper Understanding of Long Term Fairness via Simulation Studies. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (FAT* '20). Association for Computing Machinery, New York, NY, USA, 525–534. <https://doi.org/10.1145/3351095.3372878>
- [125] Abhisek Dash, Abhijnan Chakraborty, Saptarshi Ghosh, Animesh Mukherjee, and Krishna P. Gummadi. 2021. When the Umpire is Also a Player: Bias in Private Label Product Recommendations on E-Commerce Marketplaces. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (FAccT '21). Association for Computing Machinery, New York, NY, USA, 873–884. <https://doi.org/10.1145/3442188.3445944>
- [126] Somalee Datta, Jose Posada, Garrick Olson, Wencheng Li, Ciaran O'Reilly, Deepa Balraj, Joseph Mesterhazy, Joseph Pallas, Priyamvada Desai, and Nigam Shah. 2020. A new paradigm for accelerating clinical data science at Stanford Medicine. *arXiv preprint arXiv:2003.10534* (2020).
- [127] Kurtis Evan David, Qiang Liu, and Ruth Fong. 2020. Debiasing Convolutional Neural Networks via Meta Orthogonalization. arXiv:2011.07453 [cs.LG] NeurIPS 2020 workshop: "Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)".
- [128] Ian Davidson and Selvan Suntiha Ravi. 2020. A framework for determining the fairness of outlier detection. In *ECAI 2020*. IOS Press, 2465–2472.
- [129] Thomas Davidson, Dana Warmusley, Michael W. Macy, and Ingmar Weber. 2017. Automated Hate Speech Detection and the Problem of Offensive Language. In *Proceedings of the Eleventh International Conference on Web and Social Media, ICWSM 2017, Montréal, Québec, Canada, May 15-18, 2017*. AAAI Press, 512–515. <https://aaai.org/ocs/index.php/ICWSM/ICWSM17/paper/view/15665>
- [130] Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. 2019. Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (Atlanta, GA, USA) (FAT* '19). Association for Computing Machinery, New York, NY, USA, 120–128. <https://doi.org/10.1145/3287560.3287572>
- [131] Pieter Delobelle, Paul Temple, Gilles Perrouin, Benoît Frénay, Patrick Heymans, and Bettina Berendt. 2020. Ethical Adversaries: Towards Mitigating Unfairness with Adversarial Machine Learning. arXiv:2005.06852 [cs.LG] ECMLPKDD 2020 workshop: "BIAS 2020: Bias and Fairness in AI".
- [132] J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. 248–255. <https://doi.org/10.1109/CVPR.2009.5206848>
- [133] Ketki V. Deshpande, Shimei Pan, and James R. Foulds. 2020. Mitigating Demographic Bias in AI-Based Resume Filtering. In *Adjunct Publication of the 28th ACM Conference on User Modeling, Adaptation and Personalization* (Genoa, Italy) (UMAP '20 Adjunct). Association for Computing Machinery, New York, NY, USA, 268–275. <https://doi.org/10.1145/3386392.3399569>
- [134] Robert Detrano, Andras Janosi, Walter Steinbrunn, Matthias Pfisterer, Johann-Jakob Schmid, Sarbjit Sandhu, Kern H. Guppy, Stella Lee, and Victor Froelicher. 1989. International application of a new probability algorithm for the diagnosis of coronary artery disease. *The American Journal of Cardiology* 64, 5 (1989), 304–310.
- [135] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186.
- [136] Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruk-sachatkun, Kai-Wei Chang, and Rahul Gupta. 2021. BOLD: Dataset and Metrics for Measuring Biases in Open-Ended Language Generation. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (FAccT '21). Association for Computing Machinery, New York, NY, USA, 862–872. <https://doi.org/10.1145/3442188.3445924>
- [137] Emily Diana, Wesley Gill, Michael Kearns, Krishnaram Kenthapadi, and Aaron Roth. 2021. *Minimax Group Fairness: Algorithms and Experiments*. Association for Computing Machinery, New York, NY, USA, 66–76. <https://doi.org/10.1145/3461702.3462523>
- [138] Cyrus DiCiccio, Sriram Vasudevan, Kinjal Basu, Krishnaram Kenthapadi, and Deepak Agarwal. 2020. *Evaluating Fairness Using Permutation Tests*. Association for Computing Machinery, New York, NY, USA, 1467–1477. <https://doi.org/10.1145/3394486.3403199>
- [139] Charles Dickens, Rishika Singh, and Lise Getoor. 2020. HyperFair: A Soft Approach to Integrating Fairness Criteria. arXiv:2009.08952 [cs.LR] RecSys 2020 workshop: "3rd FAccTRec Workshop on Responsible Recommendation".
- [140] William Dieterich, Christina Mendoza, and Tim Brennan. 2016. COMPAS risk scales: Demonstrating accuracy equity and predictive parity.
- [141] Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. 2021. Retiring adult: New datasets for fair machine learning. *Advances in Neural Information Processing Systems* 34 (2021).
- [142] Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and Mitigating Unintended Bias in Text Classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society* (New Orleans, LA, USA) (AI/ES '18). Association for Computing Machinery, New York, NY, USA, 67–73. <https://doi.org/10.1145/3278721.3278729>
- [143] Michele Donini, Luca Oneto, Shai Ben-David, John S Shawe-Taylor, and Massimiliano Pontil. 2018. Empirical Risk Minimization Under Fairness Constraints. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc., 2791–2801. <https://proceedings.neurips.cc/paper/2018/file/83cdce08fb90370fcf53bdd56604ff-Paper.pdf>
- [144] Julia Dressel and Hany Farid. 2018. The accuracy, fairness, and limits of predicting recidivism. *Science advances* 4, 1 (2018), ea05580.
- [145] Marco F Duarte and Yu Hen Hu. 2004. Vehicle classification in distributed sensor networks. *J. Parallel and Distrib. Comput.* 64, 7 (2004), 826 – 838. <https://doi.org/10.1016/j.jpdc.2004.03.020> Computing and Communication in Distributed Sensor Networks.
- [146] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through Awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference* (Cambridge, Massachusetts) (ITCS '12). Association for Computing Machinery, New York, NY, USA, 214–226. <https://doi.org/10.1145/2090236.2090255>
- [147] Cynthia Dwork, Nicole Immorlica, Adam Tauman Kalai, and Max Leiserson. 2017. Decoupled classifiers for fair and efficient machine learning. arXiv:1707.06613 [cs.LG] KDD 2017 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [148] Cynthia Dwork, Nicole Immorlica, Adam Tauman Kalai, and Max Leiserson. 2018. Decoupled Classifiers for Group-Fair and Efficient Machine Learning. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A. Friedler and Christo Wilson (Eds.). PMLR, New York, NY, USA, 119–133. <http://proceedings.mlr.press/v81/dwork18a.html>
- [149] Natalie C Ebner, Michaela Riediger, and Ulman Lindenberger. 2010. FACES—A database of facial expressions in young, middle-aged, and older women and men: Development and validation. *Behavior research methods* 42, 1 (2010), 351–362.
- [150] Eran Eiding, Roei Enbar, and Tal Hassner. 2014. Age and Gender Estimation of Unfiltered Faces. *IEEE Transactions on Information Forensics and Security* 9, 12 (2014), 2170–2179. <https://doi.org/10.1109/TIFS.2014.2359646>
- [151] Michael D. Ekstrand, Mucun Tian, Ion Madrazo Azpiaz, Jennifer D. Ekstrand, Oghenemaro Anuyah, David McNeill, and Maria Soledad Pera. 2018. All The Cool Kids, How Do They Fit In?: Popularity and Demographic Biases in Recommender Evaluation and Effectiveness. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A. Friedler and Christo Wilson (Eds.). PMLR, New York, NY, USA, 172–186. <http://proceedings.mlr.press/v81/ekstrand18b.html>
- [152] Khaled El Emam, Luk Arbuckle, Gunes Koru, Benjamin Eze, Lisa Gaudette, Emilio Neri, Sean Rose, Jeremy Howard, and Jonathan Gluck. 2012. De-identification methods for open health data: the case of the Heritage Health Prize claims dataset. *Journal of medical Internet research* 14, 1 (2012), e33.
- [153] Marwa El Halabi, Slobodan Mitrović, Ashkan Norouzi-Fard, Jakab Tardos, and Jakub M Tarnawski. 2020. Fairness in Streaming Submodular Maximization: Algorithms and Hardness. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 13609–13622. <https://proceedings.neurips.cc/paper/2020/file/9d752cb08ef466fc480fba981cfa44a1-Paper.pdf>
- [154] Hadi Elzayn, Shahin Jabbari, Christopher Jung, Michael Kearns, Seth Neel, Aaron Roth, and Zachary Schutzman. 2019. Fair Algorithms for Learning in Allocation Problems. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (Atlanta, GA, USA) (FAT* '19). Association for Computing Machinery, New York, NY, USA, 170–179. <https://doi.org/10.1145/3287560.3287571>
- [155] L. Epstein, W.M. Landes, and R.A. Posner. 2013. *The Behavior of Federal Judges: A Theoretical and Empirical Study of Rational Choice*. Harvard University Press. <https://books.google.it/books?id=RcQEBeic3ecC>
- [156] Equivalent. 2019. Practitioner's Guide to COMPAS Core. <https://www.equivalent.com/wp-content/uploads/Practitioners-Guide-to-COMPAS-Core-040419.pdf>
- [157] Seyed Esmaili, Brian Brubach, Leonidas Tsapenakas, and John Dickerson. 2020. Probabilistic Fair Clustering. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 12743–12755. <https://proceedings.neurips.cc/paper/2020/file/95f2b84de5660df45c8a34933a2e66f-Paper.pdf>

- [158] Simone Fabbrizzi, Symeon Papadopoulos, Eirini Ntoutsi, and Ioannis Kompatsaris. 2021. A survey on bias in visual datasets. *arXiv preprint arXiv:2107.07919* (2021).
- [159] Alessandro Fabris, Alan Mishler, Stefano Gottardi, Mattia Carletti, Matteo Daicampi, Gian Antonio Susto, and Gianmaria Silvello. 2021. *Algorithmic Audit of Italian Car Insurance: Evidence of Unfairness in Access and Pricing*. Association for Computing Machinery, New York, NY, USA, 458–468. <https://doi.org/10.1145/3461702.3462569>
- [160] Golnoosh Farnad, Behrouz Babaki, and Michel Gendreau. 2020. A Unifying Framework for Fairness-Aware Influence Maximization. In *Companion Proceedings of the Web Conference 2020 (Taipei, Taiwan) (WWW '20)*. Association for Computing Machinery, New York, NY, USA, 714–722. <https://doi.org/10.1145/3366424.3383555>
- [161] Golnoosh Farnadi, Behrouz Babaki, and Margarida Carvalho. 2019. Enhancing Fairness in Kidney Exchange Program by Ranking Solutions. arXiv:1911.05489 [cs.SI] NeurIPS 2019 workshop: “Fair ML for Health”.
- [162] Golnoosh Farnadi, Pigi Kouki, Spencer K. Thompson, Sriram Srinivasan, and Lise Getoor. 2018. A Fairness-aware Hybrid Recommender System. arXiv:1809.09030 [cs.IR] RecSys 2018 workshop: “Workshop on Responsible Recommendation (FAT/Rec)”.
- [163] Elaine Fehrman, Vincent Egan, Alexander N Gorban, Jeremy Levesley, Evgeny M Mirkes, and Awaz K Muhammad. 2019. *Personality Traits and Drug Consumption: A story told by data*. Springer.
- [164] Elaine Fehrman, Awaz K. Muhammad, Evgeny M. Mirkes, Vincent Egan, and Alexander N. Gorban. 2017. The Five Factor Model of Personality and Evaluation of Drug Consumption Risk. In *Data Science*, Francesco Palumbo, Angela Montanari, and Maurizio Vichi (Eds.). Springer International Publishing, Cham, 231–242.
- [165] Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and Removing Disparate Impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (Sydney, NSW, Australia) (KDD '15)*. Association for Computing Machinery, New York, NY, USA, 259–268. <https://doi.org/10.1145/2783258.2783311>
- [166] Andres Ferraro, Dmitry Bogdanov, Xavier Serra, and Jason Yoon. 2019. Artist and style exposure bias in collaborative filtering based music recommendations. arXiv:1911.04827 [cs.IR] ISMIR 2019 workshop: “Workshop on Designing Human-Centric MIR Systems”.
- [167] Benjamin Fish, Jeremy Kun, and Á. Lelkes. 2015. Fair Boosting : a Case Study. IJML 2015 workshop: “Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)”.
- [168] Joseph Fisher, Dave Palfrey, Christos Christodoulopoulos, and Arpit Mittal. 2020. Measuring Social Bias in Knowledge Graph Embeddings. arXiv:1912.02761 [cs.CL] AKBC 2020 workshop: “Bias in Automatic Knowledge Graph Construction”.
- [169] Ronald A Fisher. 1936. The use of multiple measurements in taxonomic problems. *Annals of eugenics* 7, 2 (1936), 179–188.
- [170] Raymond Fisman, Sheena Iyengar, Emir Kamenica, and Itamar Simonson. 2006. Gender Differences in Mate Selection: Evidence From a Speed Dating Experiment. *The Quarterly Journal of Economics* 121 (02 2006), 673–697. <https://doi.org/10.1162/qjec.2006.121.2.673>
- [171] Bailey Flanigan, Paul Gözl, Anupam Gupta, and Ariel D. Procaccia. 2020. Neutralizing Self-Selection Bias in Sampling for Sortition. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/48237d9f2dea8c74c2a72126cf63d933-Abstract.html>
- [172] Omar U. Florez. 2019. On the Unintended Social Bias of Training Language Generation Models with Data from Local Media. arXiv:1911.00461 [cs.CL] NeurIPS 2019 workshop: “Human-Centric Machine Learning”.
- [173] Riccardo Fogliato, Alice Xiang, Zachary Lipton, Daniel Nagin, and Alexandra Chouldechova. 2021. On the validity of arrest as a proxy for offense: Race and the likelihood of arrest for violent crimes. In *Proceedings of the 4th AAAI/ACM Conference on AI, Ethics, and Society (AIES 2021)*. Virtual Event, 100–111. <https://doi.org/10.1145/3461702.3462538>
- [174] Antigoni-Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large Scale Crowdsourcing and Characterization of Twitter Abusive Behavior. In *Proceedings of the Twelfth International Conference on Web and Social Media, ICWSM 2018, Stanford, California, USA, June 25–28, 2018*. AAAI Press, 491–500. <https://aaai.org/ocs/index.php/ICWSM/ICWSM18/paper/view/17909>
- [175] Sorelle A. Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. 2021. The (Im)Possibility of Fairness: Different Value Systems Require Different Mechanisms for Fair Decision Making. *Commun. ACM* 64, 4 (mar 2021), 136–143. <https://doi.org/10.1145/3433949>
- [176] Sorelle A. Friedler, Carlos Scheidegger, Suresh Venkatasubramanian, Sonam Choudhary, Evan P. Hamilton, and Derek Roth. 2019. A Comparative Study of Fairness-Enhancing Interventions in Machine Learning. In *Proceedings of the Conference on Fairness, Accountability, and Transparency (Atlanta, GA, USA) (FAT* '19)*. Association for Computing Machinery, New York, NY, USA, 329–338. <https://doi.org/10.1145/3287560.3287589>
- [177] Sainyam Galhotra, Sandhya Saisubramanian, and Shlomo Zilberstein. 2021. *Learning to Generate Fair Clusters from Demonstrations*. Association for Computing Machinery, New York, NY, USA, 491–501. <https://doi.org/10.1145/3461702.3462558>
- [178] Christian Garbin, Pranav Rajpurkar, Jeremy Irvin, Matthew P. Lungren, and Oge Marques. 2021. Structured dataset documentation: a datasheet for CheXpert. arXiv:2105.03020 [eess.IV]
- [179] Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H. Chi, and Alex Beutel. 2019. Counterfactual Fairness in Text Classification through Robustness. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society (Honolulu, HI, USA) (AIES '19)*. Association for Computing Machinery, New York, NY, USA, 219–226. <https://doi.org/10.1145/3306618.3317950>
- [180] Joseph L. Gastwirth and Weiwei Miao. 2009. Formal statistical analysis of the data in disparate impact cases provides sounder inferences than the US government’s ‘four-fifths’ rule: an examination of the statistical evidence in Ricci v. DeStefano. *Law, Probability & Risk* 8, 2 (2009), 171–191.
- [181] Hancheng Ge, James Caverlee, and Haokai Lu. 2016. TAPER: A Contextual Tensor-Based Approach for Personalized Expert Recommendation. In *Proceedings of the 10th ACM Conference on Recommender Systems (Boston, Massachusetts, USA) (RecSys '16)*. Association for Computing Machinery, New York, NY, USA, 261–268. <https://doi.org/10.1145/2959100.2959151>
- [182] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2018. Datasheets for datasets. *arXiv preprint arXiv:1803.09010* (2018).
- [183] R. Stuart Geiger, Kevin Yu, Yanlai Yang, Mindy Dai, Jie Qiu, Rebekah Tang, and Jenny Huang. 2020. Garbage in, Garbage out? Do Machine Learning Application Papers in Social Computing Report Where Human-Labeled Training Data Comes From?. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 325–336. <https://doi.org/10.1145/3351095.3372862>
- [184] Andrew Gelman, Jeffrey Fagan, and Alex Kiss. 2007. An analysis of the New York City police department’s “stop-and-frisk” policy in the context of claims of racial bias. *Journal of the American statistical association* 102, 479 (2007), 813–823.
- [185] Emma J. Gerritse and Arjen P. de Vries. 2020. Effect of Debiasing on Information Retrieval. In *Bias and Social Aspects in Search and Recommendation*, Ludovico Boratto, Stefano Faralli, Mirko Marras, and Giovanni Stilo (Eds.). Springer International Publishing, Cham, 35–42.
- [186] Mehrdad Ghadiri, Samira Samadi, and Santosh Vempala. 2021. Socially Fair K-Means Clustering. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FACt '21)*. Association for Computing Machinery, New York, NY, USA, 438–448. <https://doi.org/10.1145/3442188.3445906>
- [187] Alec Go, Richa Bhayani, and Lei Huang. 2009. Twitter sentiment classification using distant supervision. *Processing* 150 (01 2009).
- [188] Naman Goel, Alfonso Amayuelas, Amit Deshpande, and Amit Sharma. 2020. The Importance of Modeling Data Missingness in Algorithmic Fairness: A Causal Perspective. arXiv:2012.11448 [cs.LG] NeurIPS 2020 workshop: “Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)”.
- [189] Naman Goel and Boi Faltings. 2019. Crowdsourcing with Fairness, Diversity and Budget Constraints. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society (Honolulu, HI, USA) (AIES '19)*. Association for Computing Machinery, New York, NY, USA, 297–304. <https://doi.org/10.1145/3306618.3314282>
- [190] Naman Goel, Mohammad Yaghini, and Boi Faltings. 2018. Non-Discriminatory Machine Learning through Convex Fairness Criteria. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society (New Orleans, LA, USA) (AIES '18)*. Association for Computing Machinery, New York, NY, USA, 116. <https://doi.org/10.1145/3278721.3278722>
- [191] Sharad Goel, Maya Perelman, Ravi Shroff, and David Alan Sklansky. 2017. Combating police discrimination in the age of big data. *New Criminal Law Review* 20, 2 (1 March 2017), 181–232. <https://doi.org/10.1525/nclr.2017.20.2.181>
- [192] Sharad Goel, Justin M Rao, Ravi Shroff, et al. 2016. Precinct or prejudice? Understanding racial disparities in New York City’s stop-and-frisk policy. *Annals of Applied Statistics* 10, 1 (2016), 365–394.
- [193] Paul Goelz, Anson Kahng, and Ariel D Procaccia. 2019. Paradoxes in Fair Machine Learning. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 8342–8352. <https://proceedings.neurips.cc/paper/2019/file/bbc92a647199b832ec90dc7f57074e9e-Paper.pdf>
- [194] Jennifer Golbeck, Zahra Ashktorab, Rashad O. Banjo, Alexandra Berlinger, Sid-dharth Bhagwan, Cody Buntain, Paul Cheakalos, Alicia A. Geller, Quint Gergory,

- Rajesh Kumar Gnanasekaran, Raja Rajan Gunasekaran, Kelly M. Hoffman, Jenny Hottle, Vichita Jienjitlert, Shivika Khare, Ryan Lau, Marianna J. Martindale, Shalmali Naik, Heather L. Nixon, Piyush Ramachandran, Kristine M. Rogers, Lisa Rogers, Meghna Sardana Sarin, Gaurav Shahane, Jayanee Thanki, Priyanka Vengataraman, Zijian Wan, and Derek Michael Wu. 2017. A Large Labeled Corpus for Online Harassment Research. In *Proceedings of the 2017 ACM on Web Science Conference* (Troy, New York, USA) (*WebSci '17*). Association for Computing Machinery, New York, NY, USA, 229–233. <https://doi.org/10.1145/3091478.3091509>
- [195] Harvey Goldstein. 1991. Multilevel Modelling of Survey Data. *Journal of the Royal Statistical Society. Series D (The Statistician)* 40, 2 (1991), 235–244. <http://www.jstor.org/stable/2348496>
- [196] Sixue Gong, Xiaoming Liu, and Anil K. Jain. 2021. Mitigating Face Recognition Bias via Group Adaptive Classifier. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 3414–3424.
- [197] Paula Gordaliza, Eustasio Del Barrio, Gamboa Fabrice, and Jean-Michel Loubes. 2019. Obtaining Fairness using Optimal Transport Theory. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 2357–2365. <http://proceedings.mlr.press/v97/gordaliza19a.html>
- [198] Joshua Gordon, Marzieh Babaeianjelodar, and Jeanna Matthews. 2020. Studying Political Bias via Word Embeddings. In *Companion Proceedings of the Web Conference 2020* (Taipei, Taiwan) (*WWW '20*). Association for Computing Machinery, New York, NY, USA, 760–764. <https://doi.org/10.1145/3366424.3383560>
- [199] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6904–6913.
- [200] Joseph Graffam, Alison J. Shinkfield, and Lesley Hardcastle. 2008. The Perceived Employability of Ex-Prisoners and Offenders. *International Journal of Offender Therapy and Comparative Criminology* 52, 6 (2008), 673–685. <https://doi.org/10.1177/0306624X07307783> arXiv:<https://doi.org/10.1177/0306624X07307783>
- [201] Ben Green and Yiling Chen. 2019. Disparate Interactions: An Algorithm-in-the-Loop Analysis of Fairness in Risk Assessments. In *Proceedings of the Conference on Fairness, Accountability, and Transparency (Atlanta, GA, USA) (FAT* '19)*. Association for Computing Machinery, New York, NY, USA, 90–99. <https://doi.org/10.1145/3287560.3287563>
- [202] Anthony G Greenwald, Debbie E McGhee, and Jordan LK Schwartz. 1998. Measuring individual differences in implicit cognition: the implicit association test. *Journal of personality and social psychology* 74, 6 (1998), 1464.
- [203] Nina Grgic-Hlaca, M. Zafar, K. Gummadi, and Adrian Weller. 2016. The Case for Process Fairness in Learning: Feature Selection for Fair Decision Making. NeurIPS 2016 workshop: "Machine Learning and the Law".
- [204] Ulrike Grömping. 2019. *South German Credit Data: Correcting a Widely Used Data Set. Report*. Technical Report. Beuth University of Applied Sciences Berlin. http://www1.beuth-hochschule.de/FB_II/reports/Report-2019-004.pdf
- [205] Jon Atle Gulla, Lemei Zhang, Peng Liu, Özlem Özgöbek, and Xiaomeng Su. 2017. The Adressa Dataset for News Recommendation. In *Proceedings of the International Conference on Web Intelligence (Leipzig, Germany) (WI '17)*. Association for Computing Machinery, New York, NY, USA, 1042–1048. <https://doi.org/10.1145/3106426.3109436>
- [206] Abdulmecit Gungor. 2018. *Benchmarking Authorship Attribution Techniques Using Over A Thousand Books by Fifty Victorian Era Novelists*. Master's thesis. Purdue University.
- [207] Guibing Guo, Jie Zhang, and Neil Yorke-Smith. 2016. A Novel Evidence-Based Bayesian Similarity Measure for Recommender Systems. *ACM Trans. Web* 10, 2, Article 8 (May 2016), 30 pages. <https://doi.org/10.1145/2856037>
- [208] Wei Guo and Aylin Caliskan. 2021. *Detecting Emergent Intersectional Biases: Contextualized Word Embeddings Contain a Distribution of Human-like Biases*. Association for Computing Machinery, New York, NY, USA, 122–133. <https://doi.org/10.1145/3461702.3462536>
- [209] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. 2016. MS-Celeb-1M: A Dataset and Benchmark for Large-Scale Face Recognition. In *Computer Vision – ECCV 2016*, Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (Eds.). Springer International Publishing, Cham, 87–102.
- [210] H Altay Guvenir, Burak Acar, Gulsen Demiroz, and Ayhan Cekin. 1997. A supervised machine learning algorithm for arrhythmia analysis. In *Computers in Cardiology 1997*. IEEE, 433–436.
- [211] Hu Han and Anil K. Jain. 2014. Age, Gender and Race Estimation from Unconstrained Face Images. http://biometrics.cse.msu.edu/Publications/Face/HanJain_UnconstrainedAgeGenderRaceEstimation_MSUTechReport2014.pdf
- [212] Anikó Hannák, Claudia Wagner, David Garcia, Alan Mislove, Markus Strohmaier, and Christo Wilson. 2017. Bias in Online Freelance Marketplaces: Evidence from TaskRabbit and Fiverr. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing* (Portland, Oregon, USA) (*CSCW '17*). Association for Computing Machinery, New York, NY, USA, 1914–1933. <https://doi.org/10.1145/2998181.2998327>
- [213] Sarel Har-Peled and Sepideh Mahabadi. 2019. Near Neighbor: Who is the Fairest of Them All?. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 13176–13187. <https://proceedings.neurips.cc/paper/2019/file/742141ceda6b8f6786609d31c8ef129f-Paper.pdf>
- [214] Elfarouk Harb and Ho Shan Lam. 2020. KFC: A Scalable Approximation Algorithm for k-center Fair Clustering. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 14509–14519. <https://proceedings.neurips.cc/paper/2020/file/a6d259bfbfa2062843ef543e21d7ec8e-Paper.pdf>
- [215] Moritz Hardt, Eric Price, Eric Price, and Nati Srebro. 2016. Equality of Opportunity in Supervised Learning. In *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Eds.), Vol. 29. Curran Associates, Inc., 3315–3323. <https://proceedings.neurips.cc/paper/2016/file/9d2682367c3935defcb1f9e247a97c0d-Paper.pdf>
- [216] F. Maxwell Harper and Joseph A. Konstan. 2015. The MovieLens Datasets: History and Context. *ACM Trans. Interact. Intell. Syst.* 5, 4, Article 19 (Dec. 2015), 19 pages. <https://doi.org/10.1145/2827872>
- [217] Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. 2018. Fairness Without Demographics in Repeated Loss Minimization. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholm, Sweden, 1929–1938. <http://proceedings.mlr.press/v80/hashimoto18a.html>
- [218] Ruining He, Wang-Cheng Kang, and Julian McAuley. 2017. Translation-Based Recommendation. In *Proceedings of the Eleventh ACM Conference on Recommender Systems (Como, Italy) (RecSys '17)*. Association for Computing Machinery, New York, NY, USA, 161–169. <https://doi.org/10.1145/3109859.3109882>
- [219] Ruining He and Julian McAuley. 2016. Ups and Downs. *Proceedings of the 25th International Conference on World Wide Web* (Apr 2016). <https://doi.org/10.1145/2872427.2883037>
- [220] Yuzi He, Keith Burghardt, Siyi Guo, and Kristina Lerman. 2020. Inherent Tradeoffs in the Fair Allocation of Treatments. arXiv:2010.16409 [cs.LG] NeurIPS 2020 workshop: "Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)".
- [221] Yuzi He, Keith Burghardt, and Kristina Lerman. 2020. A Geometric Solution to Fair Representations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society* (New York, NY, USA) (*AIES '20*). Association for Computing Machinery, New York, NY, USA, 279–285. <https://doi.org/10.1145/3375627.3375864>
- [222] Hoda Heidari, Claudio Ferrari, Krishna Gummadi, and Andreas Krause. 2018. Fairness Behind a Veil of Ignorance: A Welfare Analysis for Automated Decision Making. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc., 1265–1276. <https://proceedings.neurips.cc/paper/2018/file/be3159ad04564bfb90db9e32851ebf9c-Paper.pdf>
- [223] Hoda Heidari, Michele Loi, Krishna P. Gummadi, and Andreas Krause. 2019. A Moral Framework for Understanding Fair ML through Economic Models of Equality of Opportunity. In *Proceedings of the Conference on Fairness, Accountability, and Transparency (Atlanta, GA, USA) (FAT* '19)*. Association for Computing Machinery, New York, NY, USA, 181–190. <https://doi.org/10.1145/3287560.3287584>
- [224] Hoda Heidari, Vedant Nanda, and Krishna Gummadi. 2019. On the Long-term Impact of Algorithmic Decision Policies: Effort Unfairness and Feature Segregation through Social Learning. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 2692–2701. <http://proceedings.mlr.press/v97/heidari19a.html>
- [225] Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. 2018. Women Also Snowboard: Overcoming Bias in Captioning Models. In *Computer Vision – ECCV 2018*, Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss (Eds.). Springer International Publishing, Cham, 793–811.
- [226] I. Higgins, Loic Matthey, A. Pal, C. Burgess, Xavier Glorot, M. Botvinick, S. Mohamed, and Alexander Lerchner. 2017. beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. In *ICLR*.
- [227] Sarah Holland, Ahmed Hosny, Sarah Newman, Joshua Joseph, and Kasia Chmielinski. 2018. The dataset nutrition label: A framework to drive higher data quality standards. arXiv preprint arXiv:1805.03677 (2018).
- [228] John Hollywood, Kenneth McKay, Dulani Woods, and Denis Agniel. 2019. Real Time Crime Centers in Chicago. https://www.rand.org/content/dam/rand/pubs/research_reports/RR3200/RR3242/RAND_RR3242.pdf
- [229] Malcolm D. Holmes, Brad W. Smith, Adrienne B. Freng, and Ed A. Muñoz. 2008. Minority threat, crime control, and police resource allocation in the southwestern united states. *Crime & Delinquency* 54, 1 (2008), 128–152. <https://doi.org/10.1177/0011128707309718> arXiv:<https://doi.org/10.1177/0011128707309718>

- [230] Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé III, Miro Dudik, and Hanna Wallach. 2019. Improving fairness in machine learning systems: What do industry practitioners need?. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI 2019)*. Glasgow, UK, 1–16.
- [231] John Houvardas and Efstathios Stamatatos. 2006. N-gram feature selection for authorship identification. In *International conference on artificial intelligence: Methodology, systems, and applications*. Springer, 77–86.
- [232] Lily Hu and Yiling Chen. 2020. Fair Classification and Social Welfare. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 535–545. <https://doi.org/10.1145/3351095.3372857>
- [233] Yaowei Hu, Yongkai Wu, Lu Zhang, and Xintao Wu. 2020. Fair Multi-ple Decision Making Through Soft Interventions. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/d0921d442ee91b89ad95059d13df618-Abstract.html>
- [234] Wen Huan, Yongkai Wu, Lu Zhang, and Xintao Wu. 2020. Fairness through Equality of Effort. In *Companion Proceedings of the Web Conference 2020 (Taipei, Taiwan) (WWW '20)*. Association for Computing Machinery, New York, NY, USA, 743–751. <https://doi.org/10.1145/3366424.3383558>
- [235] Gary B. Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. 2007. Labeled faces in the wild: A database for studying face recognition in unconstrained environments.
- [236] Lingxiao Huang, Shaofeng Jiang, and Nisheeth Vishnoi. 2019. Coresets for clustering with fairness constraints. In *Advances in Neural Information Processing Systems*. 7589–7600.
- [237] Lingxiao Huang and Nisheeth Vishnoi. 2019. Stable and Fair Classification. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 2879–2890. <http://proceedings.mlr.press/v97/huang19e.html>
- [238] Lingxiao Huang, Julia Wei, and Elisa Celis. 2020. Towards Just, Fair and Interpretable Methods for Judicial Subset Selection. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (New York, NY, USA) (AI/ES '20)*. Association for Computing Machinery, New York, NY, USA, 293–299. <https://doi.org/10.1145/3375627.3375848>
- [239] J.J. Hull. 1994. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 16, 5 (1994), 550–554. <https://doi.org/10.1109/34.291440>
- [240] Sadiq Hussain, Neama Abdulaziz Dahan, Fadl Mutaher Ba-Alwbi, and Najoua Ribata. 2018. Educational data mining and analysis of students' academic performance using WEKA. *Indonesian Journal of Electrical Engineering and Computer Science* 9, 2 (2018), 447–459.
- [241] Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. 2020. Unintended Machine Learning Biases as Social Barriers for Persons with Disabilities. *SIGACCESS Access. Comput.* 125, Article 9 (March 2020), 1 pages. <https://doi.org/10.1145/3386296.3386305>
- [242] M. Häußler, Walter. 1979. Empirische Ergebnisse zu Diskriminationsverfahren bei Kreditscoringsystemen. <https://link.springer.com/article/10.1007/BF01917956>
- [243] International Warfarin Pharmacogenetics Consortium. 2009. Estimation of the warfarin dose with clinical and pharmacogenetic data. *New England Journal of Medicine* 360, 8 (2009), 753–764.
- [244] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *33rd AAAI Conference on Artificial Intelligence, AAAI 2019, 31st Innovative Applications of Artificial Intelligence Conference, IAAI 2019 and the 9th AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019 (33rd AAAI Conference on Artificial Intelligence, AAAI 2019, 31st Innovative Applications of Artificial Intelligence Conference, IAAI 2019 and the 9th AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019)*. AAAI Press, 590–597. Publisher Copyright: © 2019, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.; 33rd AAAI Conference on Artificial Intelligence, AAAI 2019, 31st Annual Conference on Innovative Applications of Artificial Intelligence, IAAI 2019 and the 9th AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019; Conference date: 27-01-2019 Through 01-02-2019.
- [245] Rashidul Islam, Shimei Pan, and James R. Foulds. 2021. *Can We Obtain Fairness For Free?* Association for Computing Machinery, New York, NY, USA, 586–596. <https://doi.org/10.1145/3461702.3462614>
- [246] Shahin Jabbari, Han-Ching Ou, Himabindu Lakkaraju, and Milind Tambe. 2020. An Empirical Study of the Trade-Offs Between Interpretability and Fairness. In *ICML 2020 Workshop on Human Interpretability in Machine Learning, preliminary version*. ICML 2020 workshop: "Workshop on Human Interpretability in Machine Learning (WHI)".
- [247] Matthew Jagielski, Michael Kearns, Jieming Mao, Alina Oprea, Aaron Roth, Saeed Sharifi Malvajerdi, and Jonathan Ullman. 2019. Differentially Private Fair Learning. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 3000–3008. <http://proceedings.mlr.press/v97/jagielski19a.html>
- [248] Disi Ji, Padhraic Smyth, and Mark Steyvers. 2020. Can I Trust My Fairness Metric? Assessing Fairness with Unlabeled Data and Bayesian Inference. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/d83de59e10227072a9c034ce10029c39-Abstract.html>
- [249] Weijie Jiang and Zachary A. Pardos. 2021. *Towards Equity and Algorithmic Fairness in Student Grade Prediction*. Association for Computing Machinery, New York, NY, USA, 608–617. <https://doi.org/10.1145/3461702.3462623>
- [250] Eun Seo Jo and Timnit Gebru. 2020. Lessons from Archives: Strategies for Collecting Sociocultural Data in Machine Learning. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 306–316. <https://doi.org/10.1145/3351095.3372829>
- [251] Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. 2019. MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042* (2019).
- [252] Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. MIMIC-III, a freely accessible critical care database. *Scientific data* 3 (2016), 160035.
- [253] Erik Jones, Shiori Sagawa, Pang Wei Koh, Ananya Kumar, and Percy Liang. 2021. Selective Classification Can Magnify Disparities Across Groups. In *International Conference on Learning Representations*. https://openreview.net/forum?id=NOM_4BkQ05i
- [254] Matthew Jones, Huy Nguyen, and Thy Nguyen. 2020. Fair k-Centers via Maximum Matching. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 4940–4949. <http://proceedings.mlr.press/v119/jones20a.html>
- [255] Sangwon Jung, Donggyu Lee, Taeon Park, and Taesup Moon. 2021. Fair Feature Distillation for Visual Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 12115–12124.
- [256] Nathan Kallus, Xiaojie Mao, and Angela Zhou. 2020. Assessing Algorithmic Fairness with Unobserved Protected Class Using Data Combination. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 110. <https://doi.org/10.1145/3351095.3373154>
- [257] Nathan Kallus and Angela Zhou. 2018. Residual Unfairness in Fair Machine Learning from Prejudiced Data. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholmsmässan, Stockholm Sweden, 2439–2448. <http://proceedings.mlr.press/v80/kallus18a.html>
- [258] Nathan Kallus and Angela Zhou. 2019. Assessing Disparate Impact of Personalized Interventions: Identifiability and Bounds. In *Advances in Neural Information Processing Systems*. H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 3426–3437. <https://proceedings.neurips.cc/paper/2019/file/d54e99a6c03704e95e6965532dec148b-Paper.pdf>
- [259] Nathan Kallus and Angela Zhou. 2019. The Fairness of Risk Scores Beyond Classification: Bipartite Ranking and the XAUC Metric. In *Advances in Neural Information Processing Systems*. H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 3438–3448. <https://proceedings.neurips.cc/paper/2019/file/73e0f7487b8e5297182c5a711d20bf26-Paper.pdf>
- [260] Nathan Kallus and Angela Zhou. 2021. Fairness, Welfare, and Equity in Personalized Pricing. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAcT '21)*. Association for Computing Machinery, New York, NY, USA, 296–314. <https://doi.org/10.1145/3442188.3445895>
- [261] Toshihiro Kamishima. 2003. Nantonac Collaborative Filtering: Recommendation Based on Order Responses. In *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (Washington, D.C.) (KDD '03)*. Association for Computing Machinery, New York, NY, USA, 583–588. <https://doi.org/10.1145/956750.956823>
- [262] Jian Kang, Jingrui He, Ross Maciejewski, and Hanghang Tong. 2020. *InFoRM: Individual Fairness on Graph Mining*. Association for Computing Machinery, New York, NY, USA, 379–389. <https://doi.org/10.1145/3394486.3403080>

- [263] William B Kannel and Daniel L McGee. 1979. Diabetes and cardiovascular disease: the Framingham study. *Jama* 241, 19 (1979), 2035–2038.
- [264] Chen Karako and Putra Manggala. 2018. Using Image Fairness Representations in Diversity-Based Re-ranking for Recommendations. arXiv:1809.03577 [cs.LR] UMAP 2018 workshop: “Fairness in User Modeling, Adaptation and Personalization (FairUMAP)”.
- [265] Kimmo Karkkainen and Jungseock Joo. 2021. FairFace: Face Attribute Dataset for Balanced Race, Gender, and Age for Bias Measurement and Mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1548–1558.
- [266] Dean S Karlan and Jonathan Zinman. 2008. Credit elasticities in less-developed economies: Implications for microfinance. *American Economic Review* 98, 3 (2008), 1040–68.
- [267] Maximilian Kasy and Rediet Abebe. 2021. Fairness, Equality, and Power in Algorithmic Decision-Making. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 576–586. <https://doi.org/10.1145/3442188.3445919>
- [268] Masahiro Kato, Takeshi Teshima, and Junya Honda. 2019. Learning from Positive and Unlabeled Data with a Selection Bias. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rJzLciCqKm>
- [269] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. 2018. Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup Fairness. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholmsmässan, Stockholm Sweden, 2564–2572. <http://proceedings.mlr.press/v80/kearns18a.html>
- [270] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. 2019. An Empirical Study of Rich Subgroup Fairness for Machine Learning. In *Proceedings of the Conference on Fairness, Accountability, and Transparency (Atlanta, GA, USA) (FAT* '19)*. Association for Computing Machinery, New York, NY, USA, 100–109. <https://doi.org/10.1145/3287560.3287592>
- [271] Vijay Keswani, Matthew Lease, and Krishnamurthy Kenthapadi. 2021. *Towards Unbiased and Accurate Deferral to Multiple Experts*. Association for Computing Machinery, New York, NY, USA, 154–165. <https://doi.org/10.1145/3461702.3462516>
- [272] Niki Kilbertus, Adria Gascon, Matt Kusner, Michael Veale, Krishna Gummadi, and Adrian Weller. 2018. Blind Justice: Fairness with Encrypted Sensitive Attributes. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholmsmässan, Stockholm Sweden, 2630–2639. <http://proceedings.mlr.press/v80/kilbertus18a.html>
- [273] Eugenia Kim, De'Aira Bryant, Deepak Srikanth, and Ayanna Howard. 2021. *Age Bias in Emotion Detection: An Analysis of Facial Emotion Recognition Performance on Young, Middle-Aged, and Older Adults*. Association for Computing Machinery, New York, NY, USA, 638–644. <https://doi.org/10.1145/3461702.3462609>
- [274] Hyunjik Kim and Andriy Mnih. 2018. Disentangling by Factorising. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, 2649–2658. <http://proceedings.mlr.press/v80/kim18b.html>
- [275] Joon Sik Kim, Jiahao Chen, and Ameet Talwalkar. 2020. FACT: A Diagnostic for Group Fairness Trade-offs. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 5264–5274. <http://proceedings.mlr.press/v119/kim20a.html>
- [276] Michael P. Kim, Amirata Ghorbani, and James Zou. 2019. Multiaccuracy: Black-Box Post-Processing for Fairness in Classification. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society (Honolulu, HI, USA) (AI/ES '19)*. Association for Computing Machinery, New York, NY, USA, 247–254. <https://doi.org/10.1145/3306618.3314287>
- [277] Svetlana Kiritchenko and Saif Mohammad. 2018. Examining Gender and Race Bias in Two Hundred Sentiment Analysis Systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics. Association for Computational Linguistics, New Orleans, Louisiana, 43–53*. <https://doi.org/10.18653/v1/S18-2005>
- [278] Inna Kizhner, Melissa Terras, Maxim Rumyantsev, Valentina Khokhlova, Elisaveta Demeshkova, Ivan Rudov, and Julia Afanasieva. 2020. Digital cultural colonialism: measuring bias in aggregated digitized content held in Google Arts and Culture. *Digital Scholarship in the Humanities* 36, 3 (12 2020), 607–640. <https://doi.org/10.1093/lc/fqaa055> arXiv:https://academic.oup.com/dsh/article-pdf/36/3/607/40873280/fqaa055.pdf
- [279] Brendan F. Klare, Ben Klein, Emma Taborsky, Austin Blanton, Jordan Cheney, Kristen Allen, Patrick Grother, Alan Mah, and Anil K. Jain. 2015. Pushing the Frontiers of Unconstrained Face Detection and Recognition: IARPA Janus Benchmark A. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [280] Matthäus Kleindessner, Pranjal Awasthi, and Jamie Morgenstern. 2019. Fair k-Center Clustering for Data Summarization. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 3448–3457. <http://proceedings.mlr.press/v97/kleindessner19a.html>
- [281] Matthäus Kleindessner, Samira Samadi, Pranjal Awasthi, and Jamie Morgenstern. 2019. Guarantees for Spectral Clustering with Fairness Constraints. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 3458–3467. <http://proceedings.mlr.press/v97/kleindessner19b.html>
- [282] Peter Knees and Moritz Hübler. 2019. Towards Uncovering Dataset Biases: Investigating Record Label Diversity in Music Playlists. ISMIR 2019 workshop: “Workshop on Designing Human-Centric MIR Systems”.
- [283] Ari Kobren, Barna Saha, and Andrew McCallum. 2019. Paper Matching with Local Fairness Constraints. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (Anchorage, AK, USA) (KDD '19)*. Association for Computing Machinery, New York, NY, USA, 1247–1257. <https://doi.org/10.1145/3292500.3330899>
- [284] Vid Kocijan, Oana-Maria Camburu, and Thomas Lukasiewicz. 2020. The Gap on GAP: Tackling the Problem of Differing Data Distributions in Bias-Measuring Datasets. arXiv:2011.01837 [cs.CL] NeurIPS 2020 workshop: “Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)”.
- [285] Ron Kohavi. 1996. Scaling up the Accuracy of Naive-Bayes Classifiers: A Decision-Tree Hybrid. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (Portland, Oregon) (KDD '96)*. AAAI Press, 202–207.
- [286] Junpei Komiyama, Akiko Takeda, Junya Honda, and Hajime Shima. 2018. Non-convex Optimization for Regression with Fairness Constraints. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholmsmässan, Stockholm Sweden, 2737–2746. <http://proceedings.mlr.press/v80/komiyama18a.html>
- [287] Giannis Konstantakis, Gianinis Promponas, Manthos Dretakis, and Panagiotis Papadakos. 2020. Bias Goggles: Exploring the Bias of Web Domains Through the Eyes of Users. In *Bias and Social Aspects in Search and Recommendation*, Ludovico Boratto, Stefano Faralli, Mirko Marras, and Giovanni Stilo (Eds.). Springer International Publishing, Cham, 66–71.
- [288] Corina Koolen. 2018. *Reading beyond the female: The relationship between perception of author gender and literary quality*. Ph.D. Dissertation. University of Amsterdam.
- [289] Corina Koolen and Andreas van Cranenburgh. 2017. These are not the Stereotypes You are Looking For: Bias and Fairness in Authorial Gender Attribution. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*. Association for Computational Linguistics, Valencia, Spain, 12–22. <https://doi.org/10.18653/v1/W17-1602>
- [290] A. Krizhevsky. 2009. Learning Multiple Layers of Features from Tiny Images.
- [291] Caitlin Kuhlman, Walter Gerych, and Elke Rundensteiner. 2021. Measuring Group Advantage: A Comparative Study of Fair Ranking Metrics. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society (Virtual Event, USA) (AI/ES '21)*. Association for Computing Machinery, New York, NY, USA, 674–682. <https://doi.org/10.1145/3461702.3462588>
- [292] Caitlin Kuhlman and Elke Rundensteiner. 2020. Rank Aggregation Algorithms for Fair Consensus. *Proc. VLDB Endow.* 13, 12 (jul 2020), 2706–2719. <https://doi.org/10.14778/3407790.3407855>
- [293] Juhı Kulshrestha, Motahhare Eslami, Johnnatan Messias, Muhammad Bilal Zafar, Saptarshi Ghosh, Krishna P. Gummadi, and Karrie Karahalios. 2017. Quantifying Search Bias: Investigating Sources of Bias for Political Searches in Social Media. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing (Portland, Oregon, USA) (CSCW '17)*. Association for Computing Machinery, New York, NY, USA, 417–432. <https://doi.org/10.1145/2998181.2998321>
- [294] Nicholas Kushmerick. 1999. Learning to Remove Internet Advertisements. In *Proceedings of the Third Annual Conference on Autonomous Agents (Seattle, Washington, USA) (AGENTS '99)*. Association for Computing Machinery, New York, NY, USA, 175–181. <https://doi.org/10.1145/301136.301186>
- [295] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual Fairness. In *Advances in Neural Information Processing Systems*, I. Guyton, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., 4066–4076. <https://proceedings.neurips.cc/paper/2017/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf>
- [296] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. 2020. The open images dataset v4. *International Journal of Computer Vision* 128, 7 (2020), 1956–1981.
- [297] Preethi Lahoti, Alex Beutel, Jilin Chen, Kang Lee, Flavien Prost, Nithum Thain, Xuezhi Wang, and Ed Chi. 2020. Fairness without Demographics through

- Adversarially Reweighted Learning. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 728–740. <https://proceedings.neurips.cc/paper/2020/file/07fc15c9d169ee48573edd749d25945d-Paper.pdf>
- [298] Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. 2015. Human-level concept learning through probabilistic program induction. *Science* 350, 6266 (2015), 1332–1338. <https://doi.org/10.1126/science.aab3050> arXiv:<https://science.sciencemag.org/content/350/6266/1332.full.pdf>
- [299] Alex Lamy, Ziyuan Zhong, Aditya K Menon, and Nakul Verma. 2019. Noise-tolerant fair classification. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 294–306. <https://proceedings.neurips.cc/paper/2019/file/8d5e957f297893487bd98fa830fa6413-Paper.pdf>
- [300] Chao Lan and Jun Huan. 2017. Discriminatory Transfer. arXiv:1707.00780 [cs.CY] KDD 2017 workshop: “Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)”.
- [301] Jeff Larson, Surya Mattu, Lauren Kirchner, and Julia Angwin. 2016. How We Analyzed the COMPAS Recidivism Algorithm. <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>
- [302] Tai Le Quy, Arjun Roy, Vasileios Iosifidis, Wenbin Zhang, and Eirini Ntoutsi. 2022. A survey on datasets for fairness-aware machine learning. *WIRES Data Mining and Knowledge Discovery* n/a, n/a (2022), e1452. <https://doi.org/10.1002/widm.1452> arXiv:<https://wires.onlinelibrary.wiley.com/doi/pdf/10.1002/widm.1452>
- [303] Susan Leavy, Gardine Meaney, Karen Wade, and Derek Greene. 2019. *Curatr: A Platform for Semantic Analysis and Curation of Historical Literary Texts*. 354–366. https://doi.org/10.1007/978-3-030-36599-8_31
- [304] Susan Leavy, Gardine Meaney, Karen Wade, and Derek Greene. 2020. Mitigating Gender Bias in Machine Learning Data Sets. In *Bias and Social Aspects in Search and Recommendation*, Ludovico Boratto, Stefano Faralli, Mirko Marras, and Giovanni Stilo (Eds.). Springer International Publishing, Cham, 12–26.
- [305] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324. <https://doi.org/10.1109/5.726791>
- [306] Y. LeCun, Fu Jie Huang, and L. Bottou. 2004. Learning methods for generic object recognition with invariance to pose and lighting. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, Vol. 2. II–104 Vol.2. <https://doi.org/10.1109/CVPR.2004.1315150>
- [307] Hansol Lee and René F. Kizilcec. 2020. Evaluation of Fairness Trade-offs in Predicting Student Success. arXiv:2007.00088 [cs.CY] International Conference on Educational Data Mining workshop: “Fairness, Accountability, and Transparency, in Educational Data (Mining)”.
- [308] Sabina Leonelli and Nicolò Tempini. 2020. *Data journeys in the sciences*. Springer Nature.
- [309] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. 2007. Graph evolution: Densification and shrinking diameters. *ACM transactions on Knowledge Discovery from Data (TKDD)* 1, 1 (2007), 2–es.
- [310] Jure Leskovec and Julian McAuley. 2012. Learning to Discover Social Circles in Ego Networks. In *Advances in Neural Information Processing Systems*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger (Eds.), Vol. 25. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2012/file/7a614fd06c325499f1680b9896beedeb-Paper.pdf>
- [311] Nixie S Lesmana, Xuan Zhang, and Xiaohui Bei. 2019. Balancing Efficiency and Fairness in On-Demand Ridesourcing. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 5309–5319. <https://proceedings.neurips.cc/paper/2019/file/3070e6addcd702cb58de5d7897bfdae1-Paper.pdf>
- [312] Peizhao Li, Yifei Wang, Han Zhao, Pengyu Hong, and Hongfu Liu. 2021. On Dyadic Fairness: Exploring and Mitigating Bias in Graph Connections. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=xgGS6PmzNq6>
- [313] Peizhao Li, Han Zhao, and Hongfu Liu. 2020. Deep Fair Clustering for Visual Learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [314] Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. 2020. Fair Resource Allocation in Federated Learning. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=ByexEISYDr>
- [315] Yanying Li, Yue Ning, Rong Liu, Ying Wu, and Wendy Hui Wang. 2020. Fairness of Classification Using Users' Social Relationships in Online Peer-To-Peer Lending. In *Companion Proceedings of the Web Conference 2020 (Taipei, Taiwan) (WWW '20)*. Association for Computing Machinery, New York, NY, USA, 733–742. <https://doi.org/10.1145/3366424.3383557>
- [316] Yanying Li, Haipei Sun, and Wendy Hui Wang. 2020. *Towards Fair Truth Discovery from Biased Crowdsourced Answers*. Association for Computing Machinery, New York, NY, USA, 599–607. <https://doi.org/10.1145/3394486.3403102>
- [317] Zhi Li, Hongke Zhao, Qi Liu, Zhenya Huang, Tao Mei, and Enhong Chen. 2018. Learning from History and Present: Next-Item Recommendation via Discriminatively Exploiting User Behaviors. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (London, United Kingdom) (KDD '18)*. Association for Computing Machinery, New York, NY, USA, 1734–1743. <https://doi.org/10.1145/3219819.3220014>
- [318] Lizhen Liang and Daniel E. Acuna. 2020. Artificial Mental Phenomena: Psychophysics as a Framework to Detect Perception Biases in AI Models. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 403–412. <https://doi.org/10.1145/3351095.3375623>
- [319] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: Common Objects in Context. In *Computer Vision – ECCV 2014*, David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars (Eds.). Springer International Publishing, Cham, 740–755.
- [320] Zachary Lipton, Julian McAuley, and Alexandra Chouldechova. 2018. Does mitigating ML's impact disparity require treatment disparity?. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2018/file/8e0384779e58ce2af40eb365b318cc32-Paper.pdf>
- [321] David Liu, Zohair Shafi, William Fleisher, Tina Eliassi-Rad, and Scott Alfeld. 2021. RAWLSNET: Altering Bayesian Networks to Encode Rawlsian Fair Equality of Opportunity. Association for Computing Machinery, New York, NY, USA, 745–755. <https://doi.org/10.1145/3461702.3462618>
- [322] Lydia T. Liu, Sarah Dean, Esther Rolf, Max Simchowitz, and Moritz Hardt. 2018. Delayed Impact of Fair Machine Learning. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholmsmässan, Stockholm Sweden, 3150–3158. <http://proceedings.mlr.press/v80/liu18c.html>
- [323] Lydia T. Liu, Max Simchowitz, and Moritz Hardt. 2019. The Implicit Fairness Criterion of Unconstrained Learning. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 4051–4060. <http://proceedings.mlr.press/v97/liu19f.html>
- [324] Lydia T. Liu, Ashia Wilson, Nika Haghtalab, Adam Tauman Kalai, Christian Borgs, and Jennifer Chayes. 2020. The Disparate Equilibria of Algorithmic Decision Making When Individuals Invest Rationally. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 381–391. <https://doi.org/10.1145/3351095.3372861>
- [325] Weiwen Liu and Robin Burke. 2018. Personalizing Fairness-aware Re-ranking. arXiv:1809.02921 [cs.IR] RecSys 2018 workshop: “Workshop on Responsible Recommendation (FAT/Rec)”.
- [326] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. arXiv:1411.7766 [cs.CV]
- [327] Francesco Locatello, Gabriele Abbati, Thomas Rainforth, Stefan Bauer, Bernhard Schölkopf, and Olivier Bachem. 2019. On the Fairness of Disentangled Representations. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 14611–14624. <https://proceedings.neurips.cc/paper/2019/file/1b486d7a5189eb8d8c46af64b0d1b4-Paper.pdf>
- [328] Michael Lohaus, Michael Perrot, and Ulrike Von Luxburg. 2020. Too Relaxed to Be Fair. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 6360–6369. <http://proceedings.mlr.press/v119/lohaus20a.html>
- [329] Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard S. Zemel. 2016. The Variational Fair Autoencoder. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1511.00830>
- [330] H. Lowe, Todd A. Ferris, P. Hernandez, and S. Weber. 2009. STRIDE - An Integrated Standards-Based Translational Research Informatics Platform. *AMIA ... Annual Symposium proceedings. AMIA Symposium 2009 (2009)*, 391–5.
- [331] Qing Lu and Lise Getoor. 2003. Link-Based Classification. In *Proceedings of the Twentieth International Conference on Artificial Intelligence on Machine Learning (Washington, DC, USA) (ICML '03)*. AAAI Press, 496–503.
- [332] Kristian Lum and James Johndrow. 2016. A statistical framework for fair predictive algorithms. arXiv:1610.08077 [stat.ML] DTL 2016 workshop: “Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)”.
- [333] B. T. Luong, S. Ruggieri, and F. Turini. 2016. Classification Rule Mining Supported by Ontology for Discrimination Discovery. In *2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW)*. 868–875. <https://doi.org/10.1109/ICDMW.2016.0128>
- [334] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning Word Vectors for Sentiment Analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Portland, Oregon, USA, 142–150. <https://www.aclweb.org/anthology/P11-1>

1015

- [335] Nitin Madnani, Anastassia Loukina, Alina von Davier, Jill Burstein, and Aoife Cahill. 2017. Building Better Open-Source Tools to Support Fairness in Automated Scoring. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*. Association for Computational Linguistics, Valencia, Spain, 41–52. <https://doi.org/10.18653/v1/W17-1605>
- [336] David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. 2018. Learning Adversarially Fair and Transferable Representations. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholm, Sweden, 3384–3393. <http://proceedings.mlr.press/v80/madras18a.html>
- [337] David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. 2019. Fairness through Causal Awareness: Learning Causal Latent-Variable Models for Biased Data. In *Proceedings of the Conference on Fairness, Accountability, and Transparency (Atlanta, GA, USA) (FAT* '19)*. Association for Computing Machinery, New York, NY, USA, 349–358. <https://doi.org/10.1145/3287560.3287564>
- [338] David Madras, Toni Pitassi, and Richard Zemel. 2018. Predict Responsibly: Improving Fairness and Accuracy by Learning to Defer. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc., 6147–6157. <https://proceedings.neurips.cc/paper/2018/file/09d37c08f7b129e96277388757530c72-Paper.pdf>
- [339] Sepideh Mahabadi and Ali Vakilian. 2020. Individual Fairness for k-Clustering. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 6586–6596. <http://proceedings.mlr.press/v119/mahabadi20a.html>
- [340] Subha Maity, Songkai Xue, Mikhail Yurochkin, and Yuekai Sun. 2021. Statistical inference for individual fairness. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=z9k8BWL-2u>
- [341] Debmalya Mandal, Samuel Deng, Suman Jana, Jeannette M. Wing, and Daniel J. Hsu. 2020. Ensuring Fairness Beyond the Training Data. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/d6539d3b57159bab6a72e106beb45bd-Abstract.html>
- [342] Varun Manjunatha, Nirat Saini, and Larry S. Davis. 2019. Explicit Bias Discovery in Visual Question Answering Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [343] Natalia Martinez, Martin Bertran, and Guillermo Sapiro. 2020. Minimax Pareto Fairness: A Multi Objective Perspective. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 6755–6764. <http://proceedings.mlr.press/v119/martinez20a.html>
- [344] Jeremie Mary, Clément Calauzènes, and Noureddine El Karoui. 2019. Fairness-Aware Learning for Continuous Attributes and Treatments. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 4382–4391. <http://proceedings.mlr.press/v97/mary19a.html>
- [345] Rossana Mastrandrea, Julie Fournet, and Alain Barrat. 2015. Contact Patterns in a High School: A Comparison between Data Collected Using Wearable Sensors, Contact Diaries and Friendship Surveys. *PLOS ONE* 10, 9 (Sep 2015), e0136497. <https://doi.org/10.1371/journal.pone.0136497>
- [346] Nicholas Mattei, Abdallah Saffidine, and Toby Walsh. 2018. An Axiomatic and Empirical Analysis of Mechanisms for Online Organ Matching. In *Proceedings of the 7th International Workshop on Computational Social Choice (COMSOC)*.
- [347] Nicholas Mattei, Abdallah Saffidine, and Toby Walsh. 2018. Fairness in Deceased Organ Matching. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society (New Orleans, LA, USA) (AI/ES '18)*. Association for Computing Machinery, New York, NY, USA, 236–242. <https://doi.org/10.1145/3278721.3278749>
- [348] Sandra G Mayson. 2018. Bias in, bias out. *Yale IJ* 128 (2018), 2218.
- [349] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton van den Hengel. 2015. Image-Based Recommendations on Styles and Substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval (Santiago, Chile) (SIGIR '15)*. Association for Computing Machinery, New York, NY, USA, 43–52. <https://doi.org/10.1145/2766462.2767755>
- [350] Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. 2000. Automating the construction of internet portals with machine learning. *Information Retrieval* 3, 2 (2000), 127–163.
- [351] Daniel McDuff, Shuang Ma, Yale Song, and Ashish Kapoor. 2019. Characterizing Bias in Classifiers using Generative Models. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 5403–5414. <https://proceedings.neurips.cc/paper/2019/file/7f018eb7b301a66658931cb8a93fd6e8-Paper.pdf>
- [352] Brian McFee, Thierry Bertin-Mahieux, Daniel P.W. Ellis, and Gert R.G. Lanckriet. 2012. The Million Song Dataset Challenge. In *Proceedings of the 21st International Conference on World Wide Web (Lyon, France) (WWW '12 Companion)*. Association for Computing Machinery, New York, NY, USA, 909–916. <https://doi.org/10.1145/2187980.2188222>
- [353] Laura McKenna. 2019. A History of the Current Population Survey and Disclosure Avoidance. <https://www2.census.gov/adrm/CED/Papers/FY20/2019-04-McKenna-cps%20and%20da.pdf>
- [354] Laura McKenna. 2019. A History of the U.S. Census Bureau's Disclosure Review Board. <https://www2.census.gov/adrm/CED/Papers/FY20/2019-04-McKenna-DRB.pdf>
- [355] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 2017. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (Proceedings of Machine Learning Research, Vol. 54)*, Aarti Singh and Jerry Zhu (Eds.). PMLR, Fort Lauderdale, FL, USA, 1273–1282. <http://proceedings.mlr.press/v54/mcmahan17a.html>
- [356] Daniel McNamara. 2019. Equalized Odds Implies Partially Equalized Outcomes Under Realistic Assumptions. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society (Honolulu, HI, USA) (AI/ES '19)*. Association for Computing Machinery, New York, NY, USA, 313–320. <https://doi.org/10.1145/3306618.3314290>
- [357] Christopher Meek, Bo Thiesson, and David Heckerman. 2002. The learning-curve sampling method applied to model-based clustering. *Journal of Machine Learning Research* 2, Feb (2002), 397–418.
- [358] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A Survey on Bias and Fairness in Machine Learning. *ACM Comput. Surv.* 54, 6, Article 115 (jul 2021), 35 pages. <https://doi.org/10.1145/3457607>
- [359] Anay Mehrotra and L. Elisa Celis. 2021. Mitigating Bias in Set Selection with Noisy Protected Attributes. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 237–248. <https://doi.org/10.1145/3442188.3445887>
- [360] Rishabh Mehrotra, Ashton Anderson, Fernando Diaz, Amit Sharma, Hanna Wallach, and Emine Yilmaz. 2017. Auditing Search Engines for Differential Satisfaction Across Demographics. In *Proceedings of the 26th International Conference on World Wide Web Companion (Perth, Australia) (WWW '17 Companion)*. International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 626–633. <https://doi.org/10.1145/3041021.3054197>
- [361] Michele Merler, Nalini Ratha, Rogerio S. Feris, and John R. Smith. 2019. Diversity in Faces. arXiv:1901.10436 [cs.CV]
- [362] Blossom Metevier, Stephen Giguere, Sarah Brockman, Ari Kobren, Yuriy Brun, Emma Brunskill, and Philip S. Thomas. 2019. Offline Contextual Bandits with High Probability Fairness Guarantees. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 14922–14933. <https://proceedings.neurips.cc/paper/2019/file/d69768b3da745b77e82cbbd88bac98-Paper.pdf>
- [363] Vishwali Mhasawade and Rumi Chunara. 2021. *Causal Multi-Level Fairness*. Association for Computing Machinery, New York, NY, USA, 784–794. <https://doi.org/10.1145/3461702.3462587>
- [364] Weiwen Miao. 2010. Did the results of promotion exams have a disparate impact on minorities? Using statistical evidence in Ricci v. DeStefano. *Journal of Statistics Education* 18, 3 (2010).
- [365] Milagros Miceli, Tianling Yang, Laurens Naudts, Martin Schuessler, Diana Serbanescu, and Alex Hanna. 2021. Documenting Computer Vision Datasets: An Invitation to Reflexive Data Practices. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 161–172.
- [366] Shachar Mirkin, Scott Nowson, Caroline Brun, and Julien Perez. 2015. Motivating Personality-aware Machine Translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Lisbon, Portugal, 1102–1108. <https://doi.org/10.18653/v1/D15-1130>
- [367] Alan Mishler, Edward H. Kennedy, and Alexandra Chouldechova. 2021. Fairness in Risk Assessment Instruments: Post-Processing to Achieve Counterfactual Equalized Odds. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 386–400. <https://doi.org/10.1145/3442188.3445902>
- [368] Shubhanshu Mishra, Sijun He, and Luca Belli. 2020. Assessing Demographic Bias in Named Entity Recognition. arXiv:2008.03415 [cs.CL] AKBC 2020 workshop: "Bias in Automatic Knowledge Graph Construction".
- [369] Alan Mislove, Bimal Viswanath, Krishna P. Gummadi, and Peter Druschel. 2010. You Are Who You Know: Inferring User Profiles in Online Social Networks. In *Proceedings of the Third ACM International Conference on Web Search and Data*

- Mining* (New York, New York, USA) (*WSDM '10*). Association for Computing Machinery, New York, NY, USA, 251–260. <https://doi.org/10.1145/1718487.1718519>
- [370] Jeffrey C Moore, Linda L Stinson, and Edward J Welniak. 2000. Income measurement error in surveys: A review. *Journal of Official Statistics-Stockholm-16, 4* (2000), 331–362.
- [371] Amanda Moreland, Christine Herlihy, Michael A Tynan, Gregory Sunshine, Russell F McCord, Charity Hilton, Jason Poovey, Angela K Werner, Christopher D Jones, Erika B Fulmer, et al. 2020. Timing of state and territorial COVID-19 stay-at-home orders and changes in population movement—United States, March 1–May 31, 2020. *Morbidity and Mortality Weekly Report* 69, 35 (2020), 1198.
- [372] S. Moro, P. Cortez, and P. Rita. 2014. A data-driven approach to predict the success of bank telemarketing. *Decis. Support Syst.* 62 (2014), 22–31.
- [373] Hussein Mozannar, Mesrob Ohannessian, and Nathan Srebro. 2020. Fair Learning with Private Demographic Data. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 7066–7075. <http://proceedings.mlr.press/v119/mozannar20a.html>
- [374] Debarghya Mukherjee, Mikhail Yurochkin, Moulinath Banerjee, and Yuekai Sun. 2020. Two Simple Ways to Learn Individual Fairness Metrics from Data. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 7097–7107. <http://proceedings.mlr.press/v119/mukherjee20a.html>
- [375] Madhumita Murgia. 2019. Microsoft quietly deletes largest public face recognition data set. <https://www.ft.com/content/7d3e0d6a-87a0-11e9-a028-86cea8523dc2>
- [376] Razieh Nabi, Daniel Malinsky, and Ilya Shpitser. 2019. Learning Optimal Fair Policies. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 4674–4682. <http://proceedings.mlr.press/v97/nabi19a.html>
- [377] Galileo Namata, Ben London, Lise Getoor, Bert Huang, and UMD EDU. 2012. Query-driven active surveying for collective classification. In *10th International Workshop on Mining and Learning with Graphs, Vol. 8*.
- [378] Vedant Nanda, Samuel Dooley, Sahil Singla, Soheil Feizi, and John P. Dickerson. 2021. Fairness Through Robustness: Investigating Robustness Disparity in Deep Learning. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (*FAccT '21*). Association for Computing Machinery, New York, NY, USA, 466–477. <https://doi.org/10.1145/3442188.3445910>
- [379] Vedant Nanda, Pan Xu, Karthik Abinav Sankararaman, John P. Dickerson, and Aravind Srinivasan. 2020. Balancing the Tradeoff between Profit and Fairness in Rideshare Platforms during High-Demand Hours. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society* (New York, NY, USA) (*AIES '20*). Association for Computing Machinery, New York, NY, USA, 131. <https://doi.org/10.1145/3375627.3375818>
- [380] Milad Nasr and Michael Carl Tschantz. 2020. Bidding Strategies with Gender Nondiscrimination Constraints for Online Ad Auctions. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (*FAT* '20*). Association for Computing Machinery, New York, NY, USA, 337–347. <https://doi.org/10.1145/3351095.3375783>
- [381] Ivoline C. Ngong, Krystal Maughan, and Joseph P. Near. 2020. Towards Auditability for Fairness in Deep Learning. arXiv:2012.00106 [cs.LG] NeurIPS 2020 workshop: “Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)”.
- [382] NLST Trial Research Team. 2011. The national lung screening trial: overview and study design. *Radiology* 258, 1 (2011), 243–253.
- [383] Alejandro Noriega-Campero, Michiel A. Bakker, Bernardo Garcia-Bulle, and Alex ‘Sandy’ Pentland. 2019. Active Fairness in Algorithmic Decision Making. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society* (Honolulu, HI, USA) (*AIES '19*). Association for Computing Machinery, New York, NY, USA, 77–83. <https://doi.org/10.1145/3306618.3314277>
- [384] Alejandro Noriega-Campero, Bernardo Garcia-Bulle, Luis Fernando Cantu, Michiel A. Bakker, Luis Tejerina, and Alex Pentland. 2020. Algorithmic Targeting of Social Policies: Fairness, Accuracy, and Distributed Governance. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (*FAT* '20*). Association for Computing Machinery, New York, NY, USA, 241–251. <https://doi.org/10.1145/3351095.3375784>
- [385] Desmond L. Nuttall, Harvey Goldstein, Robert Prosser, and Jon Rasbash. 1989. Differential school effectiveness. *International Journal of Educational Research* 13, 7 (1989), 769–776. [https://doi.org/10.1016/0883-0355\(89\)90027-X](https://doi.org/10.1016/0883-0355(89)90027-X)
- [386] Hikaru Ogura and Akiko Takeda. 2020. Convex Fairness Constrained Model Using Causal Effect Estimators. In *Companion Proceedings of the Web Conference 2020* (Taipei, Taiwan) (*WWW '20*). Association for Computing Machinery, New York, NY, USA, 723–732. <https://doi.org/10.1145/3366424.3383556>
- [387] Manuel Olave, Vladislav Rajkovic, and Marko Bohanec. 1989. An application for admission in public school systems. *Expert Systems in Public Administration* 1 (1989), 145–160.
- [388] Luca Oneto, Michele Donini, Amon Elders, and Massimiliano Pontil. 2019. Taking Advantage of Multitask Learning for Fair Classification. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society* (Honolulu, HI, USA) (*AIES '19*). Association for Computing Machinery, New York, NY, USA, 227–237. <https://doi.org/10.1145/3306618.3314255>
- [389] Luca Oneto, Michele Donini, Giulia Luise, Carlo Ciliberto, Andreas Maurer, and Massimiliano Pontil. 2020. Exploiting MMD and Sinkhorn Divergences for Fair and Transferable Representation Learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/af9c0e0c1dee63e5cad8b7ed1a5be96-Abstract.html>
- [390] Luca Oneto, Michele Donini, Andreas Maurer, and Massimiliano Pontil. 2019. Learning Fair and Transferable Representations. NeurIPS 2019 workshop: “Human-Centric Machine Learning”.
- [391] Luca Oneto, Anna Siri, Gianvittorio Luria, and Davide Anguita. 2017. Dropout Prediction at University of Genoa: a Privacy Preserving Data Driven Approach. In *ESANN*.
- [392] Akshat Pandey and Aylin Caliskan. 2021. *Disparate Impact of Artificial Intelligence Bias in Ridehailing Economy’s Price Discrimination Algorithms*. Association for Computing Machinery, New York, NY, USA, 822–833. <https://doi.org/10.1145/3461702.3462561>
- [393] Orestis Papakyriakopoulos, Simon Hegelich, Juan Carlos Medina Serrano, and Fabienne Marco. 2020. Bias in Word Embeddings. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (*FAT* '20*). Association for Computing Machinery, New York, NY, USA, 446–457. <https://doi.org/10.1145/3351095.3372843>
- [394] Dimitris Paraschakis and Bengt Nilsson. 2020. Matchmaking Under Fairness Constraints: a Speed Dating Case Study. ECIR 2020 workshop: “International Workshop on Algorithmic Bias in Search and Recommendation (BIAS 2020)”.
- [395] Gourab K Patro, Abhijnan Chakraborty, Niloy Ganguly, and Krishna P. Gumadi. 2019. Incremental Fairness in Two-Sided Market Platforms: On Smoothly Updating Recommendations. arXiv:1909.10005 [cs.SI] NeurIPS 2019 workshop: “Human-Centric Machine Learning”.
- [396] Dino Pedreshi, Salvatore Ruggieri, and Franco Turini. 2008. Discrimination-Aware Data Mining. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (Las Vegas, Nevada, USA) (*KDD '08*). Association for Computing Machinery, New York, NY, USA, 560–568. <https://doi.org/10.1145/1401890.1401959>
- [397] Kenny Peng, Arunesh Mathur, and Arvind Narayanan. 2021. Mitigating dataset harms requires stewardship: Lessons from 1000 papers. *arXiv preprint arXiv:2108.02922* (2021).
- [398] Valerio Perrone, Michele Donini, Muhammad Bilal Zafar, Robin Schmucker, Krishnaram Kenthapadi, and Cédric Archambeau. 2021. *Fair Bayesian Optimization*. Association for Computing Machinery, New York, NY, USA, 854–863. <https://doi.org/10.1145/3461702.3462629>
- [399] Dana Pessach and Erez Shmueli. 2020. Algorithmic fairness. *arXiv preprint arXiv:2001.09784* (2020).
- [400] Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep Contextualized Word Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 2227–2237.
- [401] Matthew E Peters and Dan Lecocq. 2013. Content extraction using diverse feature sets. In *Proceedings of the 22Nd International Conference on World Wide Web*. 89–90.
- [402] Stephen Pfohl, Ben Marafino, Adrien Coulet, Fatima Rodriguez, Latha Palaniappan, and Nigam H. Shah. 2019. Creating Fair Models of Atherosclerotic Cardiovascular Disease Risk. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society* (Honolulu, HI, USA) (*AIES '19*). Association for Computing Machinery, New York, NY, USA, 271–278. <https://doi.org/10.1145/3306618.3314278>
- [403] Michael Pinard. 2010. Collateral consequences of criminal convictions: Confronting issues of race and dignity. *NYUL Rev.* 85 (2010), 457.
- [404] Evaggelia Pitoura, Kostas Stefanidis, and Georgia Koutrika. 2021. Fairness in rankings and recommendations: An overview. *The VLDB Journal* (2021), 1–28.
- [405] Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. 2017. On Fairness and Calibration. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., 5680–5689. <https://proceedings.neurips.cc/paper/2017/file/b8b9c74ac526fffb2d39ab038d1cd7-Paper.pdf>
- [406] Vinay Uday Prabhu and Abeba Birhane. 2020. Large image datasets: A pyrrhic win for computer vision? *arXiv preprint arXiv:2006.16923* (2020).

- [407] Daniel Preotiu-Pietro and Lyle Ungar. 2018. User-Level Race and Ethnicity Predictors from Twitter Text. In *Proceedings of the 27th International Conference on Computational Linguistics*. Association for Computational Linguistics, Santa Fe, New Mexico, USA, 1534–1545. <https://www.aclweb.org/anthology/C18-1130>
- [408] ProPublica. 2016. COMPAS analysis github repository. <https://github.com/propublica/compas-analysis>
- [409] ProPublica. 2021. ProPublica Data Store Terms. <https://www.propublica.org/datastore/terms>
- [410] David Pujol, Ryan McKenna, Satya Kuppam, Michael Hay, Ashwin Machanavajjhala, and Gerome Miklau. 2020. Fair Decision Making Using Privacy-Protected Data. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 189–199. <https://doi.org/10.1145/3351095.3372872>
- [411] Shiyu Qian, Jian Cao, Frédéric Le Mouél, Issam Sahel, and Minglu Li. 2015. SCRAM: A Sharing Considered Route Assignment Mechanism for Fair Taxi Route Recommendations. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (Sydney, NSW, Australia) (KDD '15)*. Association for Computing Machinery, New York, NY, USA, 955–964. <https://doi.org/10.1145/2783258.2783261>
- [412] Tao Qin and Tie-Yan Liu. 2013. Introducing LETOR 4.0 Datasets. arXiv:1306.2597 [cs.IR]
- [413] Novi Quadrianto and Viktoriia Sharmanska. 2017. Recycling Privileged Learning and Distribution Matching for Fairness. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., 677–688. <https://proceedings.neurips.cc/paper/2017/file/250ef8b51c773f8dc8b4e867a9a02-Paper.pdf>
- [414] Novi Quadrianto, Viktoriia Sharmanska, and Oliver Thomas. 2019. Discovering Fair Representations in the Data Domain. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [415] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.
- [416] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- [417] Joanna Radin. 2017. "Digital natives": How medical and indigenous histories matter for big data. *Osiris* 32, 1 (2017), 43–64.
- [418] Edward Raff and Jared Sylvester. 2018. Gradient Reversal Against Discrimination. arXiv:1807.00392 [stat.ML] ICML 2018 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [419] Edward Raff, Jared Sylvester, and Steven Mills. 2018. Fair Forests: Regularized Tree Induction to Minimize Model Bias. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society (New Orleans, LA, USA) (AI/ES '18)*. Association for Computing Machinery, New York, NY, USA, 243–250. <https://doi.org/10.1145/3278721.3278742>
- [420] Aida Rahmattalabi, Phebe Vayanos, Anthony Fulginiti, Eric Rice, Bryan Wilder, Amulya Yadav, and Milind Tambe. 2019. Exploring Algorithmic Fairness in Robust Graph Covering Problems. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 15776–15787. <https://proceedings.neurips.cc/paper/2019/file/1d7e2aae840867027b7edd17b6aaa0e9-Paper.pdf>
- [421] Amifa Raj, Connor Wood, Ananda Montoly, and Michael D. Ekstrand. 2020. Comparing Fair Ranking Metrics. arXiv:2009.01311 [cs.IR] RecSys 2020 workshop: "3rd FAccTRec Workshop on Responsible Recommendation".
- [422] Inioluwa Deborah Raji and Joy Buolamwini. 2019. Actionable Auditing: Investigating the Impact of Publicly Naming Biased Performance Results of Commercial AI Products. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society (Honolulu, HI, USA) (AI/ES '19)*. Association for Computing Machinery, New York, NY, USA, 429–435. <https://doi.org/10.1145/3306618.3314244>
- [423] Govardana Sachithanandam Ramachandran, Ivan Brugere, Lav R. Varshney, and Caiming Xiong. 2021. GAEA: Graph Augmentation for Equitable Access via Reinforcement Learning. Association for Computing Machinery, New York, NY, USA, 884–894. <https://doi.org/10.1145/3461702.3462615>
- [424] Vikram V. Ramaswamy, Sunnie S. Y. Kim, and Olga Russakovsky. 2021. Fair Attribute Classification Through Latent Space De-Biasing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 9301–9310.
- [425] Veronica Red, Eric D Kelsic, Peter J Mucha, and Mason A Porter. 2011. Comparing community structure to characteristics in online collegiate social networks. *SIAM review* 53, 3 (2011), 526–543.
- [426] Michael Redmond and Alok Baveja. 2002. A data-driven software tool for enabling cooperative information sharing among police departments. *European Journal of Operational Research* 141 (09 2002), 660–678. [https://doi.org/10.1016/S0377-2217\(01\)00264-8](https://doi.org/10.1016/S0377-2217(01)00264-8)
- [427] U. Redmond and P. Cunningham. 2013. A temporal network analysis reveals the unprofitability of arbitrage in The Prosper Marketplace. *Expert Systems with Applications* 40, 9 (2013), 3715–3721. <https://doi.org/10.1016/j.eswa.2012.12.077>
- [428] Scott E Reed, Yi Zhang, Yuting Zhang, and Honglak Lee. 2015. Deep Visual Analogy-Making. In *Advances in Neural Information Processing Systems*, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (Eds.), Vol. 28. Curran Associates, Inc., 1252–1260. <https://proceedings.neurips.cc/paper/2015/file/e07413354875be01a996dc560274708e-Paper.pdf>
- [429] Ashkan Rezaei, Anqi Liu, Omid Memarrast, and Brian Ziebart. 2021. Robust Fairness under Covariate Shift. arXiv:2010.05166 [cs.LG] NeurIPS 2020 workshop: "Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)".
- [430] Chris Riederer and Augustin Chaintreau. 2017. The Price of Fairness in Location Based Advertising. <https://doi.org/10.18122/B2MD8C> RecSys 2017 workshop: "Workshop on Responsible Recommendation (FAT/Rec)".
- [431] Luc Rocher, Julien M Hendrickx, and Yves-Alexandre De Montjoye. 2019. Estimating the success of re-identifications in incomplete datasets using generative models. *Nature communications* 10, 1 (2019), 1–9.
- [432] Kit T. Rodolfa, Erika Salomon, Lauren Haynes, Iván Higuera Mendieta, Jamie Larson, and Rayid Ghani. 2020. Case Study: Predictive Fairness to Reduce Misdemeanor Recidivism through Social Service Interventions. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 142–153. <https://doi.org/10.1145/3351095.3372863>
- [433] Yuji Roh, Kangwook Lee, Steven Whang, and Changho Suh. 2020. FR-Train: A Mutual Information-Based Approach to Fair and Robust Training. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 8147–8157. <http://proceedings.mlr.press/v119/roh20a.html>
- [434] Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. 2021. FairBatch: Batch Selection for Model Fairness. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=YNnpaAKeCfx>
- [435] Yaniv Romano, Stephen Bates, and Emmanuel Candes. 2020. Achieving Equalized Odds by Resampling Sensitive Attributes. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 361–371. <https://proceedings.neurips.cc/paper/2020/file/03593ce517feac573fdaafa6dcedef61-Paper.pdf>
- [436] Andrea Romei and Salvatore Ruggieri. 2014. A multidisciplinary survey on discrimination analysis. *The Knowledge Engineering Review* 29, 5 (2014), 582–638. <https://doi.org/10.1017/S0269888913000039>
- [437] Veronica Rotemberg, Nicholas Kurtansky, Brigid Betz-Stablein, Liam Caffery, Emmanouil Chousakos, Noel Codella, Marc Combalia, Stephen Dusza, Pascale Guitera, David Gutman, et al. 2021. A patient-centric dataset of images and metadata for identifying melanomas using clinical context. *Scientific data* 8, 1 (2021), 1–8.
- [438] Benedek Rozemberczki, Carl Allen, and Rik Sarkar. 2021. Multi-scale attributed node embedding. *Journal of Complex Networks* 9, 2 (2021), cnab014.
- [439] Rachel Rudinger, Chandler May, and Benjamin Van Durme. 2017. Social Bias in Elicited Natural Language Inferences. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*. Association for Computational Linguistics, Valencia, Spain, 74–79. <https://doi.org/10.18653/v1/W17-1609>
- [440] Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender Bias in Coreference Resolution. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 8–14. <https://doi.org/10.18653/v1/N18-2002>
- [441] Anian Ruoss, Mislav Balunovic, Marc Fischer, and Martin Vechev. 2020. Learning Certified Individually Fair Representations. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 7584–7596. <https://proceedings.neurips.cc/paper/2020/file/55d491cf951b1b920900684d71419282-Paper.pdf>
- [442] Chris Russell, Matt J Kusner, Joshua Loftus, and Ricardo Silva. 2017. When Worlds Collide: Integrating Different Counterfactual Assumptions in Fairness. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., 6414–6423. <https://proceedings.neurips.cc/paper/2017/file/1271a7029c9df08643b631b02cf9e116-Paper.pdf>
- [443] Sivan Sabato and Elad Yom-Tov. 2020. Bounding the fairness and accuracy of classifiers from population statistics. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, Virtual, 8316–8325. <http://proceedings.mlr.press/v119/sabato20a.html>
- [444] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. 2010. Adapting Visual Category Models to New Domains. In *Proceedings of the 11th European Conference on Computer Vision: Part IV (Heraklion, Crete, Greece) (ECCV'10)*. Springer-Verlag, Berlin, Heidelberg, 213–226.
- [445] Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. 2020. Distributionally Robust Neural Networks. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=ryxGuJrFvS>

- [446] Samira Samadi, Uthaipon Tantipongpipat, Jamie H Morgenstern, Mohit Singh, and Santosh Vempala. 2018. The Price of Fair PCA: One Extra dimension. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc., 10976–10987. <https://proceedings.neurips.cc/paper/2018/file/cc4af25fa9d2d5c953496579b75f6f6c-Paper.pdf>
- [447] Yash Savani, Colin White, and Naveen Sundar Govindarajulu. 2020. Intra-Processing Methods for Debiasing Neural Networks. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/1d8d70dddf147d2d92a634817f01b239-Abstract.html>
- [448] Morgan Klaus Scheuerman, Kandrea Wade, Caitlin Lustig, and Jed R. Brubaker. 2020. How We've Taught Algorithms to See Identity: Constructing Race and Gender in Image Databases for Facial Analysis. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW1, Article 058 (May 2020), 35 pages. <https://doi.org/10.1145/3392866>
- [449] Candice Schumann, Susanna Ricco, Utsav Prabhu, Vittorio Ferrari, and Caroline Pantofaru. 2021. *A Step Toward More Inclusive People Annotations for Fairness*. Association for Computing Machinery, New York, NY, USA, 916–925. <https://doi.org/10.1145/3461702.3462594>
- [450] Zachary Schutzman. 2020. Trade-Offs in Fair Redistricting. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (New York, NY, USA) (AI/ES '20)*. Association for Computing Machinery, New York, NY, USA, 159–165. <https://doi.org/10.1145/3375627.3375802>
- [451] Shahar Segal, Yossi Adi, Benny Pinkas, Carsten Baum, Chaya Ganesh, and Joseph Keshet. 2021. *Fairness in the Eyes of the Data: Certifying Machine-Learning Models*. Association for Computing Machinery, New York, NY, USA, 926–935. <https://doi.org/10.1145/3461702.3462554>
- [452] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. 2008. Collective classification in network data. *AI magazine* 29, 3 (2008), 93–93.
- [453] Kulin Shah, Pooja Gupta, Amit Deshpande, and Chiranjib Bhattacharyya. 2021. *Rawlsian Fair Adaptation of Deep Learning Classifiers*. Association for Computing Machinery, New York, NY, USA, 936–945. <https://doi.org/10.1145/3461702.3462592>
- [454] Jin Shang, Mingxuan Sun, and Nina S.N. Lam. 2020. *List-Wise Fairness Criterion for Point Processes*. Association for Computing Machinery, New York, NY, USA, 1948–1958. <https://doi.org/10.1145/3394486.3403246>
- [455] Saeed Sharifi-Malvajardi, Michael Kearns, and Aaron Roth. 2019. Average Individual Fairness: Algorithms, Generalization and Experiments. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 8242–8251. <https://proceedings.neurips.cc/paper/2019/file/0e1feae55e360ff05fef58199b3fa521-Paper.pdf>
- [456] Shubham Sharma, Alan H. Gee, David Paydarfar, and Joydeep Ghosh. 2021. *FaiR-N: Fair and Robust Neural Networks for Structured Data*. Association for Computing Machinery, New York, NY, USA, 946–955. <https://doi.org/10.1145/3461702.3462559>
- [457] Shubham Sharma, Jette Henderson, and Joydeep Ghosh. 2020. CERTIFAI: A Common Framework to Provide Explanations and Analyse the Fairness and Robustness of Black-Box Models. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (New York, NY, USA) (AI/ES '20)*. Association for Computing Machinery, New York, NY, USA, 166–172. <https://doi.org/10.1145/3375627.3375812>
- [458] Shubham Sharma, Yunfeng Zhang, Jesús M. Ríos Aliaga, Djallel Bouneffouf, Vinod Muthusamy, and Kush R. Varshney. 2020. Data Augmentation for Discrimination Prevention and Bias Disambiguation. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (New York, NY, USA) (AI/ES '20)*. Association for Computing Machinery, New York, NY, USA, 358–364. <https://doi.org/10.1145/3375627.3375865>
- [459] Shubhranshu Shekhar, Neil Shah, and Leman Akoglu. 2021. FairOD: Fairness-Aware Outlier Detection. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society (Virtual Event, USA) (AI/ES '21)*. Association for Computing Machinery, New York, NY, USA, 210–220. <https://doi.org/10.1145/3461702.3462517>
- [460] Judy Hanwen Shen, Lauren Fratamico, Iyad Rahwan, and Alexander M. Rush. 2018. Darling or Babygirl? Investigating Stylistic Bias in Sentiment Analysis. KDD 2018 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [461] Mark D. Shermis. 2014. State-of-the-art automated essay scoring: Competition, results, and future directions from a United States demonstration. *Assessing Writing* 20 (2014), 53–76. <https://doi.org/10.1016/j.asw.2013.04.001>
- [462] Ashudeep Singh and Thorsten Joachims. 2018. Fairness of Exposure in Rankings. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (London, United Kingdom) (KDD '18)*. Association for Computing Machinery, New York, NY, USA, 2219–2228. <https://doi.org/10.1145/3219819.3220088>
- [463] Ashudeep Singh and Thorsten Joachims. 2019. Policy Learning for Fairness in Ranking. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 5426–5436. <https://proceedings.neurips.cc/paper/2019/file/9e82757e9a1c12cb710ad680db11f6f1-Paper.pdf>
- [464] Harvneet Singh, Rina Singh, Vishwali Mhasawade, and Rumi Chunara. 2021. Fairness Violations and Mitigation under Covariate Shift. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 3–13. <https://doi.org/10.1145/3442188.3445865>
- [465] Moninder Singh and Karthikeyan Natesan Ramamurthy. 2019. Understanding racial bias in health using the Medical Expenditure Panel Survey data. arXiv:1911.01509 [cs.LG] NeurIPS 2019 workshop: "Fair ML for Health".
- [466] Dylan Slack, Sorelle Friedler, and Emile Givental. 2019. Fair Meta-Learning: Learning How to Learn Fairly. <https://drive.google.com/file/d/1F5YF1Ar1hJ7l2H7zlsC35SzXOWqUyVW/view> NeurIPS 2019 workshop: "Human-Centric Machine Learning".
- [467] Dylan Slack, Sorelle Friedler, and Emile Givental. 2019. Fairness warnings. <https://drive.google.com/file/d/1eeu703ulWkhek0WEepYDwXg2KXSwOzc2/view> NeurIPS 2019 workshop: "Human-Centric Machine Learning".
- [468] Dylan Slack, Sorelle A. Friedler, and Emile Givental. 2020. Fairness Warnings and Fair-MAML: Learning Fairly with Minimal Data. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 200–209. <https://doi.org/10.1145/3351095.3372839>
- [469] Daniel Slunge. 2015. The willingness to pay for vaccination against tick-borne encephalitis and implications for public health policy: evidence from Sweden. *PLoS one* 10, 12 (2015), e0143875.
- [470] Jack W. Smith, J.E. Everhart, W.C. Dickson, W.C. Knowler, and R.S. Johannes. 1988. Using the ADAP Learning Algorithm to Forecast the Onset of Diabetes Mellitus. *Proceedings. Symposium on Computer Applications in Medical Care (November 1988)*, 261–265. <https://europepmc.org/articles/PMC2245318>
- [471] David Solans, Francesco Fabbri, Caterina Calsamiglia, Carlos Castillo, and Francesco Bonchi. 2021. *Comparing Equity and Effectiveness of Different Algorithms in an Application for the Room Rental Market*. Association for Computing Machinery, New York, NY, USA, 978–988. <https://doi.org/10.1145/3461702.3462600>
- [472] Nasim Sonboli and Robin Burke. 2019. Localized Fairness in Recommender Systems. In *Adjunct Publication of the 27th Conference on User Modeling, Adaptation and Personalization (Larnaca, Cyprus) (UMAP'19 Adjunct)*. Association for Computing Machinery, New York, NY, USA, 295–300. <https://doi.org/10.1145/3314183.3323845>
- [473] Nasim Sonboli, Robin Burke, Nicholas Mattei, Farzad Eskandarian, and Tian Gao. 2020. "And the Winner Is...": Dynamic Lotteries for Multi-group Fairness-Aware Recommendation. arXiv:2009.02590 [cs.LR] RecSys 2020 workshop: "3rd FAccTRec Workshop on Responsible Recommendation".
- [474] Skyler Speakman, Srihari Sridharan, and Isaac Markus. 2018. Three Population Covariate Shift for Mobile Phone-Based Credit Scoring. In *Proceedings of the 1st ACM SIGCAS Conference on Computing and Sustainable Societies (Menlo Park and San Jose, CA, USA) (COMPASS '18)*. Association for Computing Machinery, New York, NY, USA, Article 20, 7 pages.
- [475] Till Speicher, Muhammad Ali, Giridhari Venkatadri, Filipe Nunes Ribeiro, George Arvanitakis, Fabrício Benevenuto, Krishna P. Gummadi, Patrick Loiseau, and Alan Mislove. 2018. Potential for Discrimination in Online Targeted Advertising. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A. Friedler and Christo Wilson (Eds.). PMLR, New York, NY, USA, 5–19. <http://proceedings.mlr.press/v81/speicher18a.html>
- [476] Till Speicher, Hoda Heidari, Nina Grgic-Hlaca, Krishna P. Gummadi, Adish Singla, Adrian Weller, and Muhammad Bilal Zafar. 2018. A Unified Approach to Quantifying Algorithmic Unfairness: Measuring Individual & Group Unfairness via Inequality Indices. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (London, United Kingdom) (KDD '18)*. Association for Computing Machinery, New York, NY, USA, 2239–2248. <https://doi.org/10.1145/3219819.3220046>
- [477] Ryan Fox Squire. 2019. Measuring and Correcting Sampling Bias in SafeGraph Patterns for More Accurate Demographic Analysis. https://www.safegraph.com/blog/measuring-and-correcting-sampling-bias-for-accurate-demographic-analysis/?utm_source=content&utm_medium=referral&utm_campaign=colabnotebook&utm_content=panel_bias
- [478] Gabriel Stanovsky, Noah A. Smith, and Luke Zettlemoyer. 2019. Evaluating Gender Bias in Machine Translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 1679–1684. <https://doi.org/10.18653/v1/P19-1164>
- [479] Ryan Steed and Aylin Caliskan. 2021. Image Representations Learned With Unsupervised Pre-Training Contain Human-like Biases. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual*

- Event, Canada) (FAccT '21). Association for Computing Machinery, New York, NY, USA, 701–713. <https://doi.org/10.1145/3442188.3445932>
- [480] Beata Strack, Jonathan Deshazo, Chris Gennings, Juan Luis Olmo Ortiz, Sebastian Ventura, Krzysztof Cios, and John Clore. 2014. Impact of HbA1c Measurement on Hospital Readmission Rates: Analysis of 70,000 Clinical Database Patient Records. *BioMed research international* 2014 (04 2014), 781670. <https://doi.org/10.1155/2014/781670>
- [481] Tom Sühr, Asia J. Biega, Meike Zehlike, Krishna P. Gummadi, and Abhijnan Chakraborty. 2019. Two-Sided Fairness for Repeated Matchings in Two-Sided Markets: A Case Study of a Ride-Hailing Platform. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Anchorage, AK, USA) (KDD '19). Association for Computing Machinery, New York, NY, USA, 3082–3092. <https://doi.org/10.1145/3292500.3330793>
- [482] Tom Sühr, Sophie Hilgard, and Himabindu Lakkaraju. 2021. *Does Fair Ranking Improve Minority Outcomes? Understanding the Interplay of Human and Algorithmic Biases in Online Hiring*. Association for Computing Machinery, New York, NY, USA, 989–999. <https://doi.org/10.1145/3461702.3462602>
- [483] Tony Sun, Andrew Gaut, Shirlyng Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. Mitigating Gender Bias in Natural Language Processing: Literature Review. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 1630–1640. <https://doi.org/10.18653/v1/P19-1159>
- [484] Yizhou Sun, Jiawei Han, Jing Gao, and Yintao Yu. 2009. iTopicModel: Information Network-Integrated Topic Modeling. In *2009 Ninth IEEE International Conference on Data Mining*. 493–502. <https://doi.org/10.1109/ICDM.2009.43>
- [485] Nathaniel Swinger, Maria De-Arteaga, Neil Thomas Heffernan IV, Mark DM Leiserson, and Adam Tauman Kalai. 2019. What Are the Biases in My Word Embedding?. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society* (Honolulu, HI, USA) (AI/ES '19). Association for Computing Machinery, New York, NY, USA, 305–311. <https://doi.org/10.1145/3306618.3314270>
- [486] Lubos Takac and Michal Zabovsky. 2012. Data analysis in public social networks. In *International scientific conference and international workshop present day trends of innovations*, Vol. 1.
- [487] Yi Chern Tan and L. Elisa Celis. 2019. Assessing Social and Intersectional Biases in Contextualized Word Representations. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 13230–13241. <https://proceedings.neurips.cc/paper/2019/file/201d546992726352471cfe6b0df0a48-Paper.pdf>
- [488] Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. 2008. ArnetMiner: Extraction and Mining of Academic Social Networks. *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 990–998. <https://doi.org/10.1145/1401890.1402008>
- [489] Uthaiapon Tantipongpipat, Samira Samadi, Mohit Singh, Jamie H Morgenstern, and Santosh Vempala. 2019. Multi-Criteria Dimensionality Reduction with Applications to Fairness. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 15161–15171. <https://proceedings.neurips.cc/paper/2019/file/2201611d7a08ffda97e3e8c6b667a1bc-Paper.pdf>
- [490] Bahar Taskesen, Jose Blanchet, Daniel Kuhn, and Viet Anh Nguyen. 2021. A Statistical Test for Probabilistic Fairness. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (FAccT '21). Association for Computing Machinery, New York, NY, USA, 648–665. <https://doi.org/10.1145/3442188.3445927>
- [491] Rachael Tatman. 2017. Gender and Dialect Bias in YouTube's Automatic Captions. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*. Association for Computational Linguistics, Valencia, Spain, 53–59. <https://doi.org/10.18653/v1/W17-1606>
- [492] Mahbod Tavallaee, Ebrahim Bagheri, Wei Lu, and Ali A. Ghorbani. 2009. A detailed analysis of the KDD CUP 99 data set. In *2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications*. 1–6. <https://doi.org/10.1109/CISDA.2009.5356528>
- [493] Team Conduent Public Safety Solutions. 2018. Real Time Crime Forecasting Challenge: Post-Mortem Analysis Challenge Performance.
- [494] Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*. 142–147. <https://www.aclweb.org/anthology/W03-0419>
- [495] Schrasing Tong and Lalana Kagal. 2020. Investigating Bias in Image Classification using Model Explanations. arXiv:2012.05463 [cs.CV] ICML 2020 workshop: "Workshop on Human Interpretability in Machine Learning (WHI)".
- [496] Kristina Toutanova, Danqi Chen, Patrick Pantel, Hoifung Poon, Pallavi Choudhury, and Michael Gamon. 2015. Representing Text for Joint Embedding of Text and Knowledge Bases. <https://doi.org/10.18653/v1/D15-1174>
- [497] Alan Tsang, Bryan Wilder, Eric Rice, Milind Tambe, and Yair Zick. 2019. Group-Fairness in Influence Maximization. In *International Joint Conference on Artificial Intelligence*.
- [498] Connie W Tsao and Ramachandran S Vasan. 2015. Cohort Profile: The Framingham Heart Study (FHS): overview of milestones in cardiovascular epidemiology. *International journal of epidemiology* 44, 6 (2015), 1800–1813.
- [499] Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. 2018. The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data* 5, 1 (2018), 1–9.
- [500] Nikolaos Tziavelis, Ioannis Giannakopoulos, Katerina Doka, Nectarios Koziris, and Panagiotis Karras. 2019. Equitable Stable Matchings in Quadratic Time. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 457–467. <https://proceedings.neurips.cc/paper/2019/file/cb70ab375662576bd1ac5aaf16b3fca4-Paper.pdf>
- [501] UCI Machine Learning Repository. 1994. Statlog (German Credit Data) Data Set. [https://archive.ics.uci.edu/ml/datasets/statlog+\(german+credit+data\)](https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data))
- [502] UCI Machine Learning Repository. 1996. Adult Data Set. <https://archive.ics.uci.edu/ml/datasets/adult>
- [503] UCI Machine Learning Repository. 2019. South German Credit Data Set. <https://archive.ics.uci.edu/ml/datasets/South+German+Credit>
- [504] U.S. Dept. of Commerce Bureau of the Census. 1978. The Current Population Survey: Design and Methodology.
- [505] U.S. Dept. of Commerce Bureau of the Census. 1995. Current Population Survey: Annual Demographic File, 1994.
- [506] US Federal Reserve. 2007. Report to the congress on credit scoring and its effects on the availability and affordability of credit.
- [507] Berk Ustun, Yang Liu, and David Parkes. 2019. Fairness without Harm: Decoupled Classifiers with Preference Guarantees. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 6373–6382. <http://proceedings.mlr.press/v97/ustun19a.html>
- [508] Berk Ustun, M. Brandon Westover, Cynthia Rudin, and Matt T. Bianchi. 2016. Clinical Prediction Models for Sleep Apnea: The Importance of Medical History over Symptoms. *Journal of Clinical Sleep Medicine* 12, 02 (2016), 161–168. <https://doi.org/10.5664/jcsm.5476> arXiv:https://jcsa.aasm.org/doi/pdf/10.5664/jcsm.5476
- [509] Sathishkumar V E and Yongyun Cho. 2020. A rule-based model for Seoul Bike sharing demand prediction using weather data. *European Journal of Remote Sensing* 53, sup1 (2020), 166–183. <https://doi.org/10.1080/22797254.2020.1725789> arXiv:https://doi.org/10.1080/22797254.2020.1725789
- [510] Sathishkumar V E, Jangwoo Park, and Yongyun Cho. 2020. Using data mining techniques for bike sharing demand prediction in metropolitan city. *Computer Communications* 153 (2020), 353–366. <https://doi.org/10.1016/j.comcom.2020.02.007>
- [511] Rhema Vaithianathan, Emily Putnam-Hornstein, Nan Jiang, Parma Nand, and Tim Maloney. 2017. Developing Predictive Models to Support Child Maltreatment Hotline Screening Decisions: Allegheny County Methodology and Implementation. https://www.alleghenycountyanalytics.us/wp-content/uploads/2019/05/16-ACDHS-26_PredictiveRisk_Package_050119_FINAL-2.pdf
- [512] Isabel Valera, Adish Singla, and Manuel Gomez Rodriguez. 2018. Enhancing the Accuracy and Fairness of Human Decision Making. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc., 1769–1778. <https://proceedings.neurips.cc/paper/2018/file/0a113ef6b61820daa5611c870ed8d5ec-Paper.pdf>
- [513] Grant Van Horn, Elijah Cole, Sara Beery, Kimberly Wilber, Serge Belongie, and Oisín Mac Aodha. 2021. Benchmarking Representation Learning for Natural World Image Collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12884–12893.
- [514] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. 2018. The iNaturalist Species Classification and Detection Dataset. arXiv:1707.06642 [cs.CV]
- [515] Alexander Vargo, Fan Zhang, Mikhail Yurochkin, and Yuekai Sun. 2021. Individually Fair Gradient Boosting. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=JBAA9we1AL>
- [516] Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart M. Shieber. 2020. Investigating Gender Bias in Language Models Using Causal Mediation Analysis. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/92650b2e92217715fe312e6fa7b90d82-Abstract.html>
- [517] Prashanth Vijayaraghavan, Soroush Vosoughi, and Deb Roy. 2017. Twitter Demographic Classification Using Deep Multi-modal Multi-task Learning. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, Vancouver, Canada, 478–483. <https://doi.org/10.18653/v1/P17-2076>

- [518] Rob Voigt, David Jurgens, Vinodkumar Prabhakaran, Dan Jurafsky, and Yulia Tsvetkov. 2018. RTGender: A Corpus for Studying Differential Responses to Gender. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*. European Language Resources Association (ELRA), Miyazaki, Japan. <https://www.aclweb.org/anthology/L18-1445>
- [519] Julius von Kügelgen, Amir-Hossein Karimi, Umang Bhatt, Isabel Valera, Adrian Weller, and Bernhard Schölkopf. 2021. On the Fairness of Causal Algorithmic Recourse. arXiv:2010.06529 [cs.LG] NeurIPS 2020 workshop: "Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)".
- [520] Ellen Voorhees. 2005. Overview of the TREC 2005 Robust Retrieval Track. <https://trec.nist.gov/pubs/trec13/papers/ROBUST.OVERVIEW.pdf>
- [521] Christina Wadsworth, Francesca Vera, and Chris Piech. 2018. Achieving Fairness through Adversarial Learning: an Application to Recidivism Prediction. arXiv:1807.00199 [cs.LG] ICML 2018 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [522] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. 2011. The Caltech-UCSD Birds200-2011 Dataset. *Advances in Water Resources - ADV WATER RESOUR* (07 2011).
- [523] Mengting Wan and Julian McAuley. 2018. Item Recommendation on Monotonic Behavior Chains. In *Proceedings of the 12th ACM Conference on Recommender Systems (Vancouver, British Columbia, Canada) (RecSys '18)*. Association for Computing Machinery, New York, NY, USA, 86–94. <https://doi.org/10.1145/3240323.3240369>
- [524] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2019/file/4496bf24fe7fab6f046bf4923da8de6-Paper.pdf>
- [525] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rJ4km2R5t7>
- [526] Hanchen Wang, Nina Grgic-Hlaca, Preethi Lahoti, Krishna P. Gummadi, and Adrian Weller. 2019. An Empirical Study on Learning Fairness Metrics for COMPAS Data with Human Supervision. arXiv:1910.10255 [cs.CY] NeurIPS 2019 workshop: "Human-Centric Machine Learning".
- [527] Hao Wang, Berk Ustun, and Flavio Calmon. 2019. Repairing without Retraining: Avoiding Disparate Impact with Counterfactual Distributions. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 6618–6627. <http://proceedings.mlr.press/v97/wang19l.html>
- [528] Jialu Wang, Yang Liu, and Caleb Levy. 2021. Fair Classification with Group-Dependent Label Noise. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 526–536. <https://doi.org/10.1145/3442188.3445915>
- [529] Mei Wang and Weihong Deng. 2020. Mitigating Bias in Face Recognition Using Skewness-Aware Reinforcement Learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [530] Mei Wang, Weihong Deng, Jiani Hu, Xunqiang Tao, and Yaohai Huang. 2019. Racial faces in the wild: Reducing racial bias by information maximization adaptation network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 692–702.
- [531] Serena Wang, Wenshuo Guo, Harikrishna Narasimhan, Andrew Cotter, Maya Gupta, and Michael Jordan. 2020. Robust Optimization for Fairness with Noisy Protected Groups. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 5190–5203. <https://proceedings.neurips.cc/paper/2020/file/37d097caf1299d9aa79c2c2b843d2d78-Paper.pdf>
- [532] Tong Wang and Maytal Saar-Tsechansky. 2020. Augmented Fairness: An Interpretable Model Augmenting Decision-Makers' Fairness. arXiv:2011.08398 [cs.LG] NeurIPS 2020 workshop: "Algorithmic Fairness through the Lens of Causality and Interpretability (AFCI)".
- [533] Zeyu Wang, Klint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and Olga Russakovsky. 2020. Towards Fairness in Visual Recognition: Effective Strategies for Bias Mitigation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [534] Zeerak Waseem and Dirk Hovy. 2016. Hateful Symbols or Hateful People? Predictive Features for Hate Speech Detection on Twitter. In *Proceedings of the NAACL Student Research Workshop*. Association for Computational Linguistics, San Diego, California, 88–93. <https://doi.org/10.18653/v1/N16-2013>
- [535] Kellie Webster, Marta Recasens, Vera Axelrod, and Jason Baldridge. 2018. Mind the GAP: A Balanced Corpus of Gendered Ambiguous Pronouns. arXiv:1810.05201 [cs.CL]
- [536] Margaret Weeks, Scott Clair, Stephen Borgatti, Kim Radda, and Jean Schensul. 2002. Social Networks of Drug Users in High-Risk Sites: Finding the Connections. *AIDS and Behavior* 6 (06 2002), 193–206. <https://doi.org/10.1023/A:1015457400897>
- [537] Michael Wick, swetasudha panda, and Jean-Baptiste Tristan. 2019. Unlocking Fairness: a Trade-off Revisited. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 8783–8792. <https://proceedings.neurips.cc/paper/2019/file/373e4c5d8edfa8b74fd4b6791d0cf6dc-Paper.pdf>
- [538] L.F. Wightman, H. Ramsey, and Law School Admission Council. 1998. *LSAC National Longitudinal Bar Passage Study*. Law School Admission Council. <https://books.google.it/books?id=WdA7AQAIAAJ>
- [539] Bryan Wilder, Han Ching Ou, Kayla de la Haye, and Milind Tambe. 2018. Optimizing Network Structure for Preventative Health. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems (Stockholm, Sweden) (AAMAS '18)*. International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 841–849.
- [540] Josie V. Williams and Narges Razavian. 2019. Quantification of Bias in Machine Learning for Healthcare: A Case Study of Renal Failure Prediction. <https://drive.google.com/file/d/1dvjfvLIQVeeKaLRmIXfX6cVtzhkDQ0/view> NeurIPS 2019 workshop: "Fair ML for Health".
- [541] Robert Williamson and Aditya Menon. 2019. Fairness risk measures. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 6786–6797. <http://proceedings.mlr.press/v97/williamson19a.html>
- [542] Christo Wilson, Avijit Ghosh, Shan Jiang, Alan Mislove, Lewis Baker, Janelle Szary, Kelly Trindel, and Frida Polli. 2021. tBuilding and Auditing Fair Algorithms: A Case Study in Candidate Screening. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 666–677. <https://doi.org/10.1145/3442188.3445928>
- [543] Yongkai Wu, Lu Zhang, and Xintao Wu. 2018. On Discrimination Discovery and Removal in Ranked Data Using Causal Graph. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (London, United Kingdom) (KDD '18)*. Association for Computing Machinery, New York, NY, USA, 2536–2544. <https://doi.org/10.1145/3219819.3220087>
- [544] Yongkai Wu, Lu Zhang, Xintao Wu, and Hanghang Tong. 2019. PC-Fairness: A Unified Framework for Measuring Causality-based Fairness. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 3404–3414. <https://proceedings.neurips.cc/paper/2019/file/44a2e0804995faf8d2e3b084a1e2db1d-Paper.pdf>
- [545] Han Xiao, Kashif Rasul, and Roland Vollgraf. 2017. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. arXiv:1708.07747 [cs.LG]
- [546] Wenyi Xiao, Huan Zhao, Haojie Pan, Yangqiu Song, Vincent W. Zheng, and Qiang Yang. 2019. Beyond Personalization: Social Content Recommendation for Creator Equality and Consumer Satisfaction. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (Anchorage, AK, USA) (KDD '19)*. Association for Computing Machinery, New York, NY, USA, 235–245. <https://doi.org/10.1145/3292500.3330965>
- [547] Min Xie and Janet L Lauritsen. 2012. Racial context and crime reporting: A test of Black's stratification hypothesis. *Journal of Quantitative Criminology* 28, 2 (2012), 265–293.
- [548] Depeng Xu, Shuhan Yuan, Lu Zhang, and Xintao Wu. 2018. Fairgan: Fairness-aware generative adversarial networks. In *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 570–575.
- [549] Renzhe Xu, Peng Cui, Kun Kuang, Bo Li, Linjun Zhou, Zheyang Shen, and Wei Cui. 2020. *Algorithmic Decision Making with Conditional Fairness*. Association for Computing Machinery, New York, NY, USA, 2125–2135. <https://doi.org/10.1145/3394486.3403263>
- [550] Xingkun Xu, Yuge Huang, Pengcheng Shen, Shaoxin Li, Jilin Li, Feiyue Huang, Yong Li, and Zhen Cui. 2021. Consistent Instance False Positive Improves Fairness in Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 578–586.
- [551] Forest Yang, Mouhamadou Cisse, and Sanmi Koyejo. 2020. Fairness with Overlapping Groups; a Probabilistic Perspective. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 4067–4078. <https://proceedings.neurips.cc/paper/2020/file/29c0605a3bab4229e46723f89cf59d83-Paper.pdf>
- [552] Kaiyu Yang, Klint Qinami, Li Fei-Fei, Jia Deng, and Olga Russakovsky. 2020. Towards Fairer Datasets: Filtering and Balancing the Distribution of the People Subtree in the ImageNet Hierarchy. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (Barcelona, Spain) (FAccT '20)*. Association for Computing Machinery, New York, NY, USA, 547–558. <https://doi.org/10.1145/3351095.3375709>

- [553] Ke Yang and Julia Stoyanovich. 2017. Measuring Fairness in Ranked Outputs. In *Proceedings of the 29th International Conference on Scientific and Statistical Database Management* (Chicago, IL, USA) (SSDBM '17). Association for Computing Machinery, New York, NY, USA, Article 22, 6 pages. <https://doi.org/10.1145/3085504.3085526>
- [554] Mengjiao Yang and Been Kim. 2019. Benchmarking Attribution Methods with Relative Feature Importance. arXiv:1907.09701 [cs.LG]
- [555] Sirui Yao and Bert Huang. 2017. Beyond Parity: Fairness Objectives for Collaborative Filtering. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., 2921–2930. <https://proceedings.neurips.cc/paper/2017/file/e6384711491713d29bc63fc5eeb5ba4f-Paper.pdf>
- [556] Sirui Yao and Bert Huang. 2017. New Fairness Metrics for Recommendation that Embrace Differences. arXiv:1706.09838 [cs.CY] KDD 2017 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [557] I-Cheng Yeh and Che hui Lien. 2009. The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. *Expert Systems with Applications* 36, 2, Part 1 (2009), 2473 – 2480. <https://doi.org/10.1016/j.eswa.2007.12.020>
- [558] Seungeun Yi, Shirly Wang, Shalmali Joshi, and Marzyeh Ghassemi. 2019. Fair and Robust Treatment Effect Estimates: Estimation Under Treatment and Outcome Disparity with Deep Neural Models. <https://drive.google.com/file/d/1hUHRovnfzxnPaseITczzuQfvGU9jbT1I/view> NeurIPS 2019 workshop: "Fair ML for Health".
- [559] Sun Yi, Wang Xiaogang, and Tang Xiaoou. 2013. Deep Convolutional Network Cascade for Facial Point Detection. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*. 3476–3483. <https://doi.org/10.1109/CVPR.2013.446>
- [560] Mikhail Yurochkin, Amanda Bower, and Yuekai Sun. 2020. Training individually fair ML models with sensitive subspace robustness. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=B1gdxHFDH>
- [561] Mikhail Yurochkin and Yuekai Sun. 2021. SenSel: Sensitive Set Invariance for Enforcing Individual Fairness. In *International Conference on Learning Representations*. https://openreview.net/forum?id=DktZb97_Fx
- [562] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P. Gummadi. 2017. Fairness Beyond Disparate Treatment & Disparate Impact: Learning Classification without Disparate Mistreatment. In *Proceedings of the 26th International Conference on World Wide Web* (Perth, Australia) (WWW '17). International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 1171–1180. <https://doi.org/10.1145/3038912.3052660>
- [563] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, Krishna Gummadi, and Adrian Weller. 2017. From Parity to Preference-based Notions of Fairness in Classification. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., 229–239. <https://proceedings.neurips.cc/paper/2017/file/82161242827b703e6ac9c726942a1e4-Paper.pdf>
- [564] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi. 2017. Fairness Constraints: Mechanisms for Fair classification. In *Artificial Intelligence and Statistics*. PMLR, 962–970.
- [565] Meike Zehlike, Ke Yang, and Julia Stoyanovich. 2021. Fairness in Ranking: A Survey. arXiv preprint arXiv:2103.14000 (2021).
- [566] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. 2018. Mitigating Unwanted Biases with Adversarial Learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society* (New Orleans, LA, USA) (AI/ES '18). Association for Computing Machinery, New York, NY, USA, 335–340. <https://doi.org/10.1145/3278721.3278779>
- [567] Hongjing Zhang and Ian Davidson. 2021. Towards Fair Deep Anomaly Detection. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) (FAccT '21). Association for Computing Machinery, New York, NY, USA, 138–148. <https://doi.org/10.1145/3442188.3445878>
- [568] Haoran Zhang, Amy X. Lu, Mohamed Abdalla, Matthew McDermott, and Marzyeh Ghassemi. 2020. Hurtful Words: Quantifying Biases in Clinical Contextual Word Embeddings. In *Proceedings of the ACM Conference on Health, Inference, and Learning* (Toronto, Ontario, Canada) (CHIL '20). Association for Computing Machinery, New York, NY, USA, 110–120. <https://doi.org/10.1145/3368555.3384448>
- [569] Junzhe Zhang and Elias Bareinboim. 2018. Equality of Opportunity in Classification: A Causal Approach. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc., 3671–3681. <https://proceedings.neurips.cc/paper/2018/file/ff1418e8cc993fe8abcf3ce2003e5c5-Paper.pdf>
- [570] Lu Zhang, Yongkai Wu, and Xintao Wu. 2017. Achieving Non-Discrimination in Data Release. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (Halifax, NS, Canada) (KDD '17). Association for Computing Machinery, New York, NY, USA, 1335–1344. <https://doi.org/10.1145/3097983.3098167>
- [571] Xueru Zhang, Mohammadmahdi Khaliligarekani, Cem Tekin, and mingyan liu. 2019. Group Retention when Using Machine Learning in Sequential Decision Making: the Interplay between User Dynamics and Fairness. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 15269–15278. <https://proceedings.neurips.cc/paper/2019/file/7690dd4db7a92524c684e3191919eb6b-Paper.pdf>
- [572] Xueru Zhang, Ruibo Tu, Yang Liu, Mingyan Liu, Hedvig Kjellström, Kun Zhang, and Cheng Zhang. 2020. How do fair decisions fare in long-term qualification?. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/d6d231705f96d5a3aeb3a76402e49a3-Abstract.html>
- [573] Yi Zhang. 2005. *Bayesian Graphical Model for Adaptive Information Filtering*. Ph.D. Dissertation. Carnegie Mellon University.
- [574] Yunfeng Zhang, Rachel Bellamy, and Kush Varshney. 2020. Joint Optimization of AI Fairness and Utility: A Human-Centered Approach. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society* (New York, NY, USA) (AI/ES '20). Association for Computing Machinery, New York, NY, USA, 400–406. <https://doi.org/10.1145/3375627.3375862>
- [575] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. 2014. Facial Landmark Detection by Deep Multi-task Learning. https://doi.org/10.1007/978-3-319-10599-4_7
- [576] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. 2015. Learning deep representation for face alignment with auxiliary attributes. *IEEE transactions on pattern analysis and machine intelligence* 38, 5 (2015), 918–930.
- [577] Zhe Zhang and Daniel B. Neill. 2017. Identifying Significant Predictive Bias in Classifiers. arXiv:1611.08292 [stat.ML] KDD 2017 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [578] Zhifei Zhang, Yang Song, and Hairong Qi. 2017. Age Progression/Regression by Conditional Adversarial Autoencoder. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [579] Bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shu-Tao Xia. 2020. Maintaining Discrimination and Fairness in Class Incremental Learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [580] Chen Zhao, Changbin Li, Jincheng Li, and Feng Chen. 2020. Fair Meta-Learning For Few-Shot Classification. In *2020 IEEE International Conference on Knowledge Graph (ICKG)*. 275–282. <https://doi.org/10.1109/ICKG50248.2020.00047>
- [581] Han Zhao, Amanda Coston, Tameem Adel, and Geoffrey J. Gordon. 2020. Conditional Learning of Fair Representations. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=Hkekl0NFPr>
- [582] Han Zhao and Geoff Gordon. 2019. Inherent Tradeoffs in Learning Fair Representations. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., 15675–15685. <https://proceedings.neurips.cc/paper/2019/file/b4189d9de0fb2b9cce090bd1a15e3420-Paper.pdf>
- [583] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Copenhagen, Denmark, 2979–2989. <https://doi.org/10.18653/v1/D17-1323>
- [584] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 15–20. <https://doi.org/10.18653/v1/N18-2003>
- [585] Yunhan Zhao, Shu Kong, and Charless Fowlkes. 2021. Camera Pose Matters: Improving Depth Prediction by Mitigating Pose Distribution Bias. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 15759–15768.
- [586] Yong Zheng, Tanaya Dave, Neha Mishra, and Harshit Kumar. 2018. Fairness In Reciprocal Recommendations: A Speed-Dating Study. In *Adjunct Publication of the 26th Conference on User Modeling, Adaptation and Personalization* (Singapore, Singapore) (UMAP '18). Association for Computing Machinery, New York, NY, USA, 29–34. <https://doi.org/10.1145/3213586.3226207>
- [587] Yaoyao Zhong, Weihong Deng, Mei Wang, Jiani Hu, Jianteng Peng, Xunqiang Tao, and Yaohai Huang. 2019. Unequal-Training for Deep Face Recognition With Long-Tailed Noisy Data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [588] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2018. Places: A 10 Million Image Database for Scene Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40, 6 (2018), 1452–1464. <https://doi.org/10.1109/TPAMI.2017.2723009>
- [589] Yukun Zhu, Ryan Kiros, Richard Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning Books and Movies: Towards Story-like Visual Explanations by Watching Movies and Reading Books.

- arXiv:1506.06724 [cs.CV]
- [590] Ziwei Zhu, Jianling Wang, Yin Zhang, and James Caverlee. 2018. Fairness-Aware Recommendation of Information Curators. arXiv:1809.03040 [cs.IR] RecSys 2018 workshop: "Workshop on Responsible Recommendation (FAT/Rec)".
- [591] Indre Zliobaite. 2015. On the relation between accuracy and fairness in binary classification. arXiv:1505.05723 [cs.LG] ICML 2015 workshop: "Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)".
- [592] I. Žliobaite, F. Kamiran, and T. Calders. 2011. Handling Conditional Discrimination. In *2011 IEEE 11th International Conference on Data Mining*. 992–1001. <https://doi.org/10.1109/ICDM.2011.72>

A DATA BRIEFS

Data briefs were drafted by the first author and reviewed by the remaining authors. For over 95% of the surveyed datasets, we identified at least one contact involved in the data curation process or familiar with the dataset, who received a preliminary version of the respective data brief and a request for corrections and additions. Data briefs are meant as a short documentation format to provide key information on datasets used in fairness research. With reference to the Data Catalog Vocabulary (DCAT),⁷ data briefs refer to a Cataloged Resource, typically corresponding to a Dataset. Data briefs comprise the following fields:

Description. This is a free-text field reporting (1) the aim/purpose of a data artifact (i.e., why it was developed/collected), as stated by curators or inferred from context; (2) a high-level description of the available features; (3) the labeling procedure for annotated attributes, with special attention to sensitive ones, if any; (4) the envisioned ML task, if any. Corresponds to `dct:description` in DCAT.

Affiliation of creators. Typically derived from reports, articles, or official web pages presenting a dataset. Datasets can be derivatives of other datasets (e.g., Adult). We typically refer to the final resource while providing the prior context where appropriate. In DCAT vocabulary, it is the affiliation of a `dct:publisher` (for published resources) or a `dct:creator`.

Domain. The main field where the data is used (e.g., computer vision for ImageNet) or the field studying the processes and phenomena that produced the dataset (e.g., radiology for CheXpert).

Tasks in fairness literature. An indication of the task performed on the dataset in each surveyed article that uses the current resource.

Data spec. The main format of the data. The envisioned categories are text, image, time-series, tabular data, and pairs. The latter denotes a special type of tabular data where rows and columns correspond to entities and cells to a relation between them, such as relevance for query-document pairs, ratings for user-item pairs, co-authorship relation for author-author pairs. A “mixture” category was added for resources with multimodal data. Equivalent to `dct:type` in DCAT.

Sample size. Dataset cardinality.

Year. Last known update to the dataset. For resources whose collection and curation are ongoing (e.g., Framingham) we write “present”. Equivalent to `dct:modified`.

Sensitive features. Sensitive attributes in the dataset. These are typically explicitly annotated, but may include implicit ones, such as textual references to people and their demographics in text datasets. References to gender, for instance, can easily be retrieved from English-language text corpora based on intrinsically gendered words, such as she, man, aunt.

Link. A link to the website where the resource can be downloaded or requested. Equivalent to either `dcat:landingPage`.

Further information. Reference to works and web pages describing the dataset.

Following the algorithmic fairness literature, we define sensitive features as encoding membership to groups that are salient for society and have some special protection based on the law, including race, ethnicity, sex, gender, and age. We may occasionally stretch this definition and report features considered sensitive in some works, such as political leaning or education, so long as they reflect essential divisions in society. We also report domain-specific attributes considered sensitive in a given context, such as language for Section 203 determinations or brand ownership for Amazon Recommendations. We follow the language of the available documentation for the names and values of sensitive features, including distinctions between race and ethnicity. For datasets that report geographical information at any granularity (GPS coordinates, neighbourhoods, countries) we report “geography” among the sensitive attributes. If an article considers features to be sensitive in an arbitrary fashion (e.g., sepal width in the Iris dataset), we do not report it in the respective field.

For the dataset domain, we follow the area-category taxonomy defined by Scimago,⁸ with the addition of “news”, “social media”, “social networks”, “sports” and “food”. Table 3 contains a summary of the surveyed datasets through this domain-based taxonomy. Tasks in the fairness literature were labeled via open coding. The final taxonomy is detailed in Section 5.2. We distinguish between works that are more focused on evaluation rather than a proposal of novel solutions by writing, e.g. “fair ranking evaluation” instead of “fair ranking”. We use “evaluation” as a broad term for works focusing on analyses of algorithms, products, platforms, or datasets and their properties from multiple fairness and accuracy perspectives. With some abuse of nomenclature, we also use this label for works that focus on properties of fairness metrics [405]. Unless otherwise specified, “fairness evaluation” is about fair classification, which is the most common task. Exploratory approaches focused on discovering biases that are not fully specified ex-ante are indicated with the label “bias discovery”.

⁷<https://www.w3.org/TR/vocab-dcat-2/>

⁸<https://www.scimagojr.com/journalrank.php>

Domain	Sample datasets
Computer Science	
computer networks	KDD Cup 99
pattern recognition	Internet Ads
signal processing	Vehicle
social media	
toxicity and hate speech	Civil Comments, Wikipedia Toxic Comments, Twitter offensive language
political leaning	Twitter Presidential Politics
dialect	TwitterAAE
library and information sciences	
collaboration networks	Paper-Reviewer Matching, 4area, ArnetMiner Citation Network
peer review	Paper-Reviewer Matching
information systems	
search engines	Online Freelance Marketplaces, Bing US Queries, Symptoms in Queries
recommender systems	Amazon Recommendations, Amazon Reviews, MovieLens
knowledge bases	Freebase15k-237, Wikidata
Social Sciences	
urban studies	SafeGraph Research Release
social networks	University Facebook Networks, Pokec Social Network, Rice Facebook Network
demography	US Census Data (1990), Dutch Census, National Longitudinal Survey of Youth
sociology	Columbia University Speed Dating, Libimseti
law	
recidivism prediction	COMPAS, Recidivism of Felons on Probation, State Court Processing Statistics
crime prediction	Communities and Crime, Stop, Question and Frisk, Strategic Subject List
political science	
registered voters	North Carolina Voters
electoral precincts	MGGG States
polling	2016 US Presidential Poll
sortition	Climate Assembly UK
education	
application processes	Nursery, IIT-JEE
student performance	Student, Law School, UniGe
post-education placement	Campus Recruitment
social work	
child maltreatment prevention	Allegheny Child Welfare
emergency response	Harvey Rescue
drug abuse prevention	Homeless Youths' Social Networks, DrugNet
transportation	
taxi trips	NYC Taxi Trips, Shanghai Taxi Trajectories
ride hailing	Chicago Ridesharing, Ride-hailing App
bike sharing	Seoul Bike Sharing
public transport	Equitable School Access in Chicago
Computer Vision	
general purpose	ImageNet, MNIST, CIFAR
face analysis	CelebA, Pilot Parliaments Benchmar, FairFace
synthetic	dSprites, Cars3D, shapes3D
Health	
sleep medicine	Apnea
critical care medicine	MIMIC-III
public health	Kidney Exchange Program, Willingness-to-Pay for Vaccine, Kidney Matching
cardiology	Arrhythmia, Heart Disease, Framingham
neurology	Epileptic Seizures
pediatrics	Infant Health and Development Program (IHDP)
dermatology	HAM10000, SIIM-ISIC Melanoma Classification
medicine	Stanford Medicine Research Data Repository

pharmacology	Warfarin
endocrinology	Diabetes 130-US Hospitals, Pima Indians Diabetes Dataset (PIDD)
nephrology	Renal Failure
radiology	CheXpert, MIMIC-CXR-JPG, National Lung Screening Trial (NLST)
health policy	Heritage Health, MEPS-HC
applied psychology	Drug Consumption
experimental psychology	FACES
Economics and Business	
economics	
census	Adult, US Family Income, Poverty in Colombia
employment	ANPE
tariffs	U.S. Harmonized Tariff Schedule
insurance	Italian Car Insurance
division of goods	Spliddit Divide Goods
finance	
peer-to-peer lending	Mobile Money Loans, Kiva, Prosper Loans Network
mortgages	HMDA
other loans	German Credit, Credit Elasticities
credit scoring	FICO
default prediction	Credit Card Default
marketing	
marketing campaigns	Bank Marketing
advertising bids	Yahoo! A1 Search Marketing, Wholesale
management information systems	
automated hiring	Pymetrics Bias Group, CVs from Singapore
employee retention	IBM HR Analytics
Linguistics	
general purpose	Wikipedia dumps, One billion word benchmark, BookCorpus
fairness benchmarks	Bias in Translation Templates, Equity Evaluation Corpus, Winogender
Arts and Humanities	
music	Million Playlist Dataset (MPD), Million Song Dataset (MSD), Last.fm
literature	Goodreads Reviews, Riddle of Literary Quality, Nominees Corpus
movies	MovieLens, FilmTrust
Natural Sciences	
biology	iNaturalist Datasets
biochemistry	PP-Pathways
plant science	Iris
Miscellaneous	
news	TREC Robust04, New York Times Annotated Corpus, Reuters 50 50
sports	Fantasy Football, FIFA 20 Players, Olympic Athletes
food	Sushi

Table 3: A selection of datasets through the lens of the domain taxonomy.

A.1 2010 Frequently Occurring Surnames

- **Description:** this dataset reports all surnames occurring 100 or more times in the 2010 US Census, broken down by race (White, Black, Asian and Pacific Islander (API), American Indian and Alaskan Native only (AIAN), multiracial, or Hispanic).
- **Affiliation of creators:** US Census Bureau.
- **Domain:** linguistics.
- **Tasks in fairness literature:** fair subset selection under unawareness [359].
- **Data spec:** tabular data.
- **Sample size:** ~ 200K surnames.
- **Year:** 2016.
- **Sensitive features:** race.
- **Link:** https://www.census.gov/topics/population/genealogy/data/2010_surnames.html
- **Further info:** <https://www2.census.gov/topics/genealogy/2010surnames/surnames.pdf>

A.2 2016 US Presidential Poll

- **Description:** this dataset was collected and maintained by FiveThirtyEight, a website specialized in opinion poll analysis. This resource was developed with the goal of providing an aggregated estimate based on multiple polls, weighting each input according to sample size, recency, and historical accuracy of the polling organization. For each poll, the dataset provides the period of data collection, its sample size, the pollster conducting it, their rating, and a url linking to the source data.
- **Affiliation of creators:** FiveThirtyEight.
- **Domain:** political science.
- **Tasks in fairness literature:** limited-label fairness evaluation [443].
- **Data spec:** tabular data.
- **Sample size:** ~ 13K poll results.
- **Year:** 2016.
- **Sensitive features:** geography.
- **Link:** http://projects.fivethirtyeight.com/general-model/president_general_polls_2016.csv
- **Further info:** <https://projects.fivethirtyeight.com/2016-election-forecast/>

A.3 4area

- **Description:** this dataset was extracted from DBLP to study the problem of topic modeling on documents connected by links in a graph structure. The creators extracted from DBLP articles published at 20 major conferences from four related areas, i.e., database, data mining, machine learning, and information retrieval. Each author is associated with four continuous variables based on the fraction of research papers published in these areas. The associated task is the prediction of these attributes.
- **Affiliation of creators:** University of Illinois at Urbana-Champaign.
- **Domain:** library and information sciences.
- **Tasks in fairness literature:** fair clustering [214].
- **Data spec:** author-author pairs.
- **Sample size:** ~ 30K nodes (authors) connected by ~ 200K edges (co-author relations).
- **Year:** 2009.
- **Sensitive features:** author.
- **Link:** not available
- **Further info:** Sun et al. [484]

A.4 Academic Collaboration Networks

- **Description:** these dataset represent two collaboration networks from the preprint server arXiv, covering scientific papers submitted to the astrophysics (AstroPh) and condensed matter (CondMat) physics categories. Each node in the network is an author, with links indicating co-authorship of one or more articles. Nodes are indicated with ids, hence information about the researchers in the graph is not immediately available. These datasets were developed to study the evolution of graphs over time.
- **Affiliation of creators:** Carnegie Mellon University; Cornell University.
- **Domain:** library and information sciences.
- **Tasks in fairness literature:** fair graph mining [262].
- **Data spec:** author-author pairs.
- **Sample size:** ~19K nodes (authors) connected by ~ 200K edges (indications of co-authorship) (AstroPh). ~23K nodes connected by ~ 93K edges (CondMat).
- **Year:** 2009.
- **Sensitive features:** none.

- **Link:** <http://snap.stanford.edu/data/ca-AstroPh.html> (AstroPh) and <http://snap.stanford.edu/data/ca-CondMat.html> (CondMat)
- **Further info:** Leskovec et al. [309]

A.5 Adience

- **Description:** this resource was developed to favour the study of automated age and gender identification from images of faces. Photos were sourced from Flickr albums, among the ones automatically uploaded from iPhone and made available under Creative Commons license. All images were manually labeled for age, gender and identity “using both the images themselves and any available contextual information”. These annotations are fundamental for the tasks associated with this dataset, i.e. age and gender estimation. One author of Buolamwini and Gebru [63] labeled each image in Adience with Fitzpatrick skin type.
- **Affiliation of creators:** Adience; Open University of Israel.
- **Domain:** computer vision.
- **Tasks in fairness literature:** data bias evaluation [63], robust fairness evaluation [378].
- **Data spec:** image.
- **Sample size:** ~ 30K images of ~ 2K subjects.
- **Year:** 2014.
- **Sensitive features:** age, gender, skin type.
- **Link:** <https://talhassner.github.io/home/projects/Adience/Adience-data.html>
- **Further info:** Buolamwini and Gebru [63], Eidingen et al. [150]

A.6 Adressa

- **Description:** this dataset was curated as part of the RecTech project on recommendation technology owned by Adresseavisen (shortened to Adressa) a large Norwegian newspaper. It summarizes one week of traffic to the newspaper website by both subscribers and non-subscribers, during February 2017. The dataset describes reading events, i.e. a reader accessing an article, providing access timestamps and user information inferred from their IP. Specific information about the articles is also available, including author, keywords, body, and mentioned entities. The dataset curators also worked on an extended version of the dataset (Adressa 20M), ten times larger than the one described here.
- **Affiliation of creators:** Norwegian University of Science and Technology; Adresseavisen.
- **Domain:** news, information systems.
- **Tasks in fairness literature:** fair ranking [87].
- **Data spec:** user-article pairs.
- **Sample size:** ~ 3M ratings by ~ 15M readers over ~ 1K articles.
- **Year:** 2018.
- **Sensitive features:** geography.
- **Link:** <http://reclab.idi.ntnu.no/dataset/>
- **Further info:** [205]

A.7 Adult

- **Description:** this dataset was created as a resource to benchmark the performance of machine learning algorithms on socially relevant data. Each instance is a person who responded to the March 1994 US Current Population Survey, represented along demographic and socio-economic dimensions, with features describing their profession, education, age, sex, race, personal and financial condition. The dataset was extracted from the census database, preprocessed, and donated to UCI Machine Learning Repository in 1996 by Ronny Kohavi and Barry Becker. A binary variable encoding whether respondents’ income is above \$50,000 was chosen as the target of the prediction task associated with this resource. See Appendix B for extensive documentation.
- **Affiliation of creators:** Silicon Graphics Inc.
- **Domain:** economics.
- **Tasks in fairness literature:** fairness evaluation [75, 93, 138, 165, 176, 234, 245, 246, 275, 320, 323, 340, 381, 388, 405, 451, 457, 476, 519, 541, 591], fair classification [3, 23, 72, 81, 83, 101, 108, 119, 121, 131, 143, 165, 167, 190, 197, 221, 232, 328, 343, 363, 386, 398, 413, 418, 419, 433, 434, 447, 453, 456, 458, 515, 527, 544, 549, 551, 560, 561, 564, 566, 570], fair clustering [1, 6, 21, 39, 41, 43, 61, 99, 186, 214, 236, 339, 344, 532], fair clustering under unawareness [157], fair active classification [24, 25, 383], fair preference-based classification [10, 374, 507], fair classification under unawareness [272, 297, 373, 531], fair anomaly detection [459, 567], fairness evaluation under unawareness [18], robust fairness evaluation [45], data bias evaluation [40], rich-subgroup fairness evaluation [106, 270], fair representation learning [329, 336, 414, 441, 581, 582], fair multi-stage classification [188, 233], robust fair classification [237, 341, 429], dynamical fair classification [571], fair ranking evaluation [259], fair data summarization [36, 78, 100, 153, 254, 280], fair regression [4], limited-label fair classification [104, 109, 528], limited-label fairness evaluation [248], preference-based fair clustering [177].
- **Data spec:** tabular data.

- **Sample size:** ~ 50K instances.
- **Year:** 1996.
- **Sensitive features:** age, sex, race.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/adult>
- **Further info:** Ding et al. [141], Kohavi [285], McKenna [353, 354], UCI Machine Learning Repository [502], U.S. Dept. of Commerce Bureau of the Census [505]

A.8 Allegheny Child Welfare

- **Description:** this dataset stems from an initiative by the Allegheny County's Department of Human Services to develop assistive tools to support child maltreatment hotline screening decisions. Referrals received by Allegheny County via a hotline between September 2008 and April 2016 were assembled into a dataset. To obtain a relevant history and follow-up time for each referral, a subset of samples spanning the period from April 2010 to April 2014 is considered. Each data point pertains to a referral for suspected child abuse or neglect and contains a wealth of information from the integrated data management systems of Allegheny County. This data includes cross-sector administrative information for individuals associated with a report of child abuse or neglect, including data from child protective services, mental health services, drug, and alcohol services. The target to be estimated by risk models is future child harm, as measured e.g. by re-referrals, which complements the role of the screening staff who are focused on the information currently available about the referral.
- **Affiliation of creators:** Allegheny County Department of Human Services; Auckland University of Technology; University of Southern California; University of Auckland; University of California.
- **Domain:** social work.
- **Tasks in fairness literature:** fairness evaluation of risk assessment [116], fair risk assessment [367].
- **Data spec:** tabular data.
- **Sample size:** ~ 80K calls.
- **Year:** 2019.
- **Sensitive features:** age, race, gender of child.
- **Link:** not available
- **Further info:** Vaithianathan et al. [511]

A.9 Amazon Recommendations

- **Description:** this dataset was crawled to study anti-competitive behaviour on Amazon, and the extent to which Amazon's private label products are recommended on the platform. Considering the categories *backpack* and *battery*, where Amazon is known to have a strong private label presence, the creators gathered a set of organic and sponsored recommendations from Amazon.in, exploiting snowball sampling. Metadata for each product was also collected, including user rating, number of reviews, brand, seller.
- **Affiliation of creators:** Indian Institute of Technology; Max Planck Institute for Software Systems.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair ranking evaluation [125].
- **Data spec:** item-recommendation pairs.
- **Sample size:** ~ 1M recommendations associated with ~ 20K items.
- **Year:** 2021.
- **Sensitive features:** brand ownership.
- **Link:** not available
- **Further info:** Dash et al. [125]

A.10 Amazon Reviews

- **Description:** this is large-scale dataset of over ten million products and respective reviews on Amazon, spanning more than two decades. It was created to study the problem of image-based recommendation and its dynamics. Rich metadata are available for both products and reviews. Reviews consist of ratings, text, reviewer name, and review ID, while products include title, price, image, and sales rank of product.
- **Affiliation of creators:** University of California, San Diego.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair ranking [395].
- **Data spec:** user-product pairs (reviews).
- **Sample size:** ~ 200M reviews of products.
- **Year:** 2018.
- **Sensitive features:** none.
- **Link:** <https://nijianmo.github.io/amazon/index.html>

- **Further info:** He and McAuley [219], McAuley et al. [349]

A.11 ANPE

- **Description:** this dataset represents a large randomized controlled trial, assigning job seekers in France to a program run by the Public employment agency (ANPE), or to a program outsourced to private providers by the Unemployment insurance organization (Unédic). The data involves 400 public employment branches and over 200,000 job-seekers. Data about job seekers includes their demographics, their placement program and the subsequent duration of unemployment spells.
- **Affiliation of creators:** Paris School of Economics; Institute of Labor Economics; CREST; ANPE; Unédic; Direction de l'Animation de la Recherche et des Études Statistiques.
- **Domain:** economics.
- **Tasks in fairness literature:** fairness evaluation of risk assessment [258].
- **Data spec:** tabular data.
- **Sample size:** ~ 200K job seekers.
- **Year:** 2012.
- **Sensitive features:** age, gender, nationality.
- **Link:** <https://www.openicpsr.org/openicpsr/project/113904/version/V1/view?path=/openicpsr/113904/fcr:versions/V1/Archive&type=folder>
- **Further info:** Behaghel et al. [35]

A.12 Antelope Valley Networks

- **Description:** this a set of synthetic datasets generated to study the problem of influence maximization for obesity prevention. Samples of agents are generated to emulate the demographic and obesity distribution across regions in the Antelope Valley in California, exploiting data from the U.S. Census, the Los Angeles County Department of Public Health, and Los Angeles Times Mapping L.A. project. Each agent in the network has a geographic region, gender, ethnicity, age, and connections to other agents, which are more frequent for agents with similar attributes. Agents are also assigned a weight status, which may change based on interactions with other agents in their ego-network, emulating social learning.
- **Affiliation of creators:** National University of Singapore; National University of Southern California.
- **Domain:** public health.
- **Tasks in fairness literature:** fair graph diffusion [160].
- **Data spec:** agent-agent pairs.
- **Sample size:** ~ 20 synthetic networks, containing ~ 500 individuals each.
- **Year:** 2019.
- **Sensitive features:** ethnicity, gender, age, geography.
- **Link:** https://github.com/bwilder0/fair_influmax_code_release
- **Further info:** Tsang et al. [497], Wilder et al. [539]

A.13 Apnea

- **Description:** this dataset results from a sleep medicine study focused on establishing important factors for the automated diagnosis of Obstructive Sleep Apnea (OSA). The task associated with this dataset is the prediction of medical condition (OSA/no OSA) from available patient features, which include demographics, medical history, and symptoms.
- **Affiliation of creators:** Massachusetts Institute of Technology; Massachusetts General Hospital; Harvard Medical School.
- **Domain:** sleep medicine.
- **Tasks in fairness literature:** fair preference-based classification [507].
- **Data spec:** mixture (time series and tabular data).
- **Sample size:** ~ 2K patients.
- **Year:** 2016.
- **Sensitive features:** age, sex.
- **Link:** not available
- **Further info:** Ustun et al. [508]

A.14 ArnetMiner Citation Network

- **Description:** this dataset is one of the many resources made available by the ArnetMiner online service. The ArnetMiner system was developed for the extraction and mining of data from academic social networks, with a focus on profiling of researchers. The DBLP Citation Network is extracted from academic resources, such as DBLP, ACM and MAG (Microsoft Academic Graph). The dataset captures the relationships between scientific articles and their authors in a connected graph structure. It can be used for tasks such as community discovery, topic modeling, centrality and influence analysis. In its latest versions, the dataset comprises over 20 fields,

including paper title, keywords, abstract, venue, year, along with authors, and their affiliations. The ArnetMiner project was partially funded by the Chinese National High-tech R&D Program, the National Science Foundation of China, IBM China Research Lab, the Chinese Young Faculty Research Funding program and Minnesota China Collaborative Research Program.

- **Affiliation of creators:** Tsinghua University; IBM.
- **Domain:** library and information sciences.
- **Tasks in fairness literature:** fair graph mining [66].
- **Data spec:** article-article pairs.
- **Sample size:** ~ 5M papers connected by ~ 50M citations.
- **Year:** 2021.
- **Sensitive features:** author.
- **Link:** <http://www.arnetminer.org/citation>
- **Further info:** Tang et al. [488]; <https://www.aminer.org/>

A.15 Arrhythmia

- **Description:** data provenance for this set of patient records seems uncertain. The first work referencing this dataset dates to 1997 and details a machine learning approach for the diagnosis of arrhythmia, which presumably motivated its collection. Each data point describes a different patient; features include demographics, weight and height and clinical measurements from ECG signals, along with the diagnosis of a cardiologist into 16 different classes of arrhythmia (including none), which represents the target variable.
- **Affiliation of creators:** Bilkent University; Baskent University.
- **Domain:** cardiology.
- **Tasks in fairness literature:** fair classification [143, 344], robust fair classification [429], limited-label fair classification [109].
- **Data spec:** tabular data.
- **Sample size:** ~ 500 patients.
- **Year:** 1997.
- **Sensitive features:** age, sex.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/arrhythmia>
- **Further info:** Guvenir et al. [210]

A.16 Athletes and health professionals

- **Description:** the datasets were developed to study the effects of bias in image classification. The health professional dataset (doctors and nurses) contains race and gender as sensitive features and the athlete dataset (basketball and volleyball players) contains gender and jersey color as sensitive features. Each subgroup, separated by combinations of sensitive features, is roughly balanced at 200 images. The collected data was manually examined by the curators to remove stylized images and images containing both females and males.
- **Affiliation of creators:** Massachusetts Institute of Technology.
- **Domain:** computer vision.
- **Tasks in fairness literature:** bias discovery [495].
- **Data spec:** image.
- **Sample size:** ~ 800 images of athletes and ~ 500 images of health professionals.
- **Year:** 2020.
- **Sensitive features:** Gender (both), race (health professionals), jersey color (athletes).
- **Link:** <https://github.com/ghayat2/Datasets>
- **Further info:** Tong and Kagal [495]

A.17 Automated Student Assessment Prize (ASAP)

- **Description:** this dataset was collected to evaluate the feasibility of automated essay scoring. It consists of a collection of essays by US students in grade levels 7–10, rated by at least two human raters. The dataset comes with a predefined training/validation/test split and powers the Hewlett Foundation Automated Essay Scoring competition on Kaggle. The curators tried to remove personally identifying information from the essays using Named Entity Recognizer (NER) and several heuristics.
- **Affiliation of creators:** University of Akron; The Common Pool; OpenEd Solutions.
- **Domain:** education.
- **Tasks in fairness literature:** fair regression evaluation [335].
- **Data spec:** text.
- **Sample size:** ~ 20K student essays.
- **Year:** 2012.
- **Sensitive features:** none.

- **Link:** <https://www.kaggle.com/c/asap-aes/data/>
- **Further info:** Shermis [461]

A.18 Bank Marketing

- **Description:** often simply called *Bank* dataset in the fairness literature, this resource was produced to support a study of success factors in telemarketing of long-term deposits within a Portuguese bank, with data collected over the period 2008–2010. Each data point represents a telemarketing phone call and includes client-specific features (e.g. job, education), features about the marketing phone call (e.g. day of the week and duration) and meaningful environmental features (e.g. euribor). The classification target is a binary variable indicating client subscription to a term deposit.
- **Affiliation of creators:** Instituto Universitário de Lisboa (ISCTE-IUL), ISTAR, Lisboa; University of Minho.
- **Domain:** marketing.
- **Tasks in fairness literature:** fair classification [23, 137, 447, 453, 564], fair clustering [1, 6, 21, 39, 99, 214, 236, 339], fair data summarization [153], fair classification under unawareness [272], fairness evaluation [245, 320], limited-label fairness evaluation [248], preference-based fair clustering [177].
- **Data spec:** tabular data.
- **Sample size:** ~ 40K phone contacts.
- **Year:** 2012.
- **Sensitive features:** age.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/Bank+Marketing>
- **Further info:** Moro et al. [372]

A.19 Barcelona Room Rental

- **Description:** this dataset summarizes the operations of a room rental platform in Barcelona over 30 months, from January 2017 through June 2019. It contains information about over 60,000 users, divided into those seeking (seeker) and those listing (lister) a room. The data consists of lister-seeker pairs, such that a seeker is recommended for a room and lister. Recommendations are provided by a set of different recommender systems (recsys). For each pair, the data reports the rank in which each seeker was listed, the recsys providing the recommendation, and the post-recommendation interaction, if any, along with demographic information on both users. Textual indications of “gay-friendliness” in user profiles is treated as a sensitive feature (among others), as sexual orientation was previously found to be a discriminating factor in access to housing.
- **Affiliation of creators:** Universitat Pompeu Fabra; Eurecat; Institute for Political Economy and Governance; ISI Foundation.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair ranking evaluation [471].
- **Data spec:** lister-seeker pairs.
- **Sample size:** ~ 4M pairs.
- **Year:** 2021.
- **Sensitive features:** gender, age, spoken language, “gay-friendliness”.
- **Link:** not available
- **Further info:** Solans et al. [471]

A.20 Benchmarking Attribution Methods (BAM)

- **Description:** this dataset was developed to evaluate different explainability methods in computer vision. It was constructed by pasting object pixels from MS-COCO [319] into scene images from MiniPlaces [588]. Objects are rescaled to a variable proportion between one third and one half of the scene images onto which they are pasted. Both scene images and object images belong to ten different classes, for a total of 100 possible combinations. Scene images were chosen between the ones that do not contain the objects from the ten MS-COCO classes. This dataset enables users to freely control how each object is correlated with scenes, from which ground truth explanations can be formed. The creators also propose a few quantitative metrics to evaluate interpretability methods by either contrasting different inputs in the same dataset or contrasting two models with the same input.
- **Affiliation of creators:** Google.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair representation learning [127].
- **Data spec:** image.
- **Sample size:** ~ 100K images over 10 object classes and 10 image classes.
- **Year:** 2020.
- **Sensitive features:** none.
- **Link:** <https://github.com/google-research-datasets/bam>
- **Further info:** Yang and Kim [554]

A.21 Berkeley Students

- **Description:** this dataset holds anonymized student records at UC Berkeley from Spring 2012 through Fall 2019. It consists of enrollment information on a per-semester basis for tens of thousands of students. For each enrollment, student course scores are provided, along with student demographic information, including gender, race, entry status and parental income. The dataset supports evaluations of equity in educational outcome as well as grade predictions for academic support interventions. It is maintained by the University's Enterprise Data and Analytics unit.
- **Affiliation of creators:** University of California, Berkeley.
- **Domain:** education.
- **Tasks in fairness literature:** fair classification [249].
- **Data spec:** tabular data.
- **Sample size:** ~ 2M enrollments across ~ 80K students.
- **Year:** 2021.
- **Sensitive features:** gender, race.
- **Link:** not available
- **Further info:** Jiang and Pardos [249]

A.22 Bias in Bios

- **Description:** this dataset was developed as a large-scale study of gender bias in occupation classification. It consists of online biographies of professionals scraped from the Common Crawl. Biographies are detected in crawls when they match the regular expression "`<name> is a(n) <title>`", with `<title>` being one of twenty-eight common occupations. The gender of each person in the dataset is identified via the third person gendered pronoun, typically used in professional biographies. The envisioned task mirrors that of a job search automated system in a two-sided labor marketplace, i.e. automated occupation classification. The dataset curators provide python code to recreate the dataset from old Common Crawls.
- **Affiliation of creators:** Carnegie Mellon University; University of Massachusetts Lowell; Microsoft; LinkedIn.
- **Domain:** linguistics, information systems.
- **Tasks in fairness literature:** fairness evaluation [130], fair classification [561].
- **Data spec:** text.
- **Sample size:** ~ 400K biographies.
- **Year:** 2018.
- **Sensitive features:** gender.
- **Link:** <https://github.com/Microsoft/biosbias>
- **Further info:** De-Arteaga et al. [130]

A.23 Bias in Translation Templates

- **Description:** this resource was developed to study the problem of gender biases in machine translation. It consists of a set of short templates of the form `One thing about the man/woman, [he/she] is [a ##],` where `[he/she]` can be a gender-neutral or gender-specific pronoun, and `[a ##]` refers to a profession or conveys sentiment. Templates are built so that the part before the comma acts as a gender-specific clue, and the part after the comma contains information about gender and sentiment/profession. Accurate translations should correctly match the grammatical gender before and after the comma, in every word where it is required by the target language. The curators identify a set of languages to which this template is easily applicable, namely German, Korean, Portuguese, and Tagalog, which are chosen for their different properties with respect to grammatical gender. Depending on which language pair is being considered for translation, the curators identify a set of criteria for the evaluation of translation quality, with special emphasis on the correctness of grammatical gender.
- **Affiliation of creators:** Seoul National University.
- **Domain:** linguistics.
- **Tasks in fairness literature:** bias evaluation of machine translation [102].
- **Data spec:** text.
- **Sample size:** ~ 1K templates.
- **Year:** 2021.
- **Sensitive features:** gender.
- **Link:** <https://github.com/nolongerprejudice/tgbi-x>
- **Further info:** Cho et al. [102]

A.24 Bing US Queries

- **Description:** this dataset was created to investigate differential user satisfaction with the Bing search engine across different demographic groups. The authors selected log data of a random subset of Bing's desktop and laptop users from the English-speaking

US market over a two week period. The data was preprocessed by cleaning spam and bot queries, and it was enriched with user demographics, namely age (bucketed) and gender (binary), which were self-reported by users during account registration and automatically validated by the dataset curators. Moreover, queries were labeled with topic information. Finally, four different signals were extracted from search logs, namely graded utility, reformulation rate, page click count, and successful click count.

- **Affiliation of creators:** Microsoft.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair ranking evaluation [360].
- **Data spec:** query-result pairs.
- **Sample size:** ~ 30M (non-unique) queries issued by ~ 4M distinct users.
- **Year:** 2017.
- **Sensitive features:** age, gender.
- **Link:** not available
- **Further info:** Mehrotra et al. [360]

A.25 BOLD

- **Description:** this resource is a benchmark to measure biases of language models with respect to sensitive demographic attributes. The creators identified six attributes (e.g. race, profession) and values of said attribute (e.g. African American, flight nurse) for which they gather prompts from English Language Wikipedia, either from pages about the group (e.g. “A flight nurse is a registered”) or people representing it (e.g. “Over the years, Isaac Hayes was able”). Prompts are fed to different language models, whose outputs are automatically labelled for sentiment, regard, toxicity, emotion and gender polarity. These labels are also validated by human annotators hired on Amazon Mechanical Turk.
- **Affiliation of creators:** Amazon; University of California, Santa Barbara.
- **Domain:** linguistics.
- **Tasks in fairness literature:** bias evaluation in language models [136].
- **Data spec:** text.
- **Sample size:** ~ 20K prompts.
- **Year:** 2021.
- **Sensitive features:** gender, race, religion, profession, political leaning.
- **Link:** <https://github.com/amazon-research/bold>
- **Further info:** Dhamala et al. [136]

A.26 BookCorpus

- **Description:** this dataset was developed for the problem of learning general representations of text useful for different downstream tasks. It consist of text from 11,038 books from the web by unpublished authors available on <https://www.smashwords.com/> in 2015. The BookCorpus contains thousands of duplicate books (only 7,185 are unique) and many contain copyright restrictions. The GPT [415] and BERT [135] language models were trained on this dataset.
- **Affiliation of creators:** University of Toronto; Massachusetts Institute of Technology.
- **Domain:** linguistics.
- **Tasks in fairness literature:** data bias evaluation [487].
- **Data spec:** text.
- **Sample size:** ~ 1B words in ~74M sentences from ~11K books.
- **Year:** unknown.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** not available
- **Further info:** Bandy and Vincent [27], Zhu et al. [589]

A.27 BUPT Faces

- **Description:** this resource consists of two datasets, developed as a large scale collection, suitable for training face verification algorithms operating on diverse populations. The underlying data collection procedure mirrors the one from RFW (§ A.153), including sourcing from MS-Celeb-1M and automated annotation of so-called *race* into one of four categories: Caucasian, Indian, Asian and African. For categories where not enough images were readily available, the authors resort to the FreeBase celebrity list, downloading images of people from Google and cleaning them "both automatically and manually". The remaining images were obtained from MS-Celeb-1M (§ A.124), on which the BUPT Faces datasets are heavily based.
- **Affiliation of creators:** Beijing University of Posts and Telecommunications.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair reinforcement learning [529], fair classification [550], fair representation learning [196].

- **Data spec:** image.
- **Sample size:** ~ 2M images of ~ 40K celebrities (BUPT-Globalface); ~ 1M images of ~ 30K celebrities (BUPT-Balancedface).
- **Year:** 2019.
- **Sensitive features:** race.
- **Link:** <http://www.whdeng.cn/RFW/Trainingdataste.html>
- **Further info:** Wang and Deng [529]

A.28 Burst

- **Description:** Burst is a free provider of stock photography powered by Shopify. This dataset features a subset of Burst images used as a resource to test algorithms for fair image retrieval and ranking, aimed at providing, in response to a query, a collection of photos that is balanced across demographics. Images come with human-curated tags annotated internally by the Burst team.
- **Affiliation of creators:** Shopify.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair ranking [264].
- **Data spec:** image.
- **Sample size:** ~ 3K images.
- **Year:** present.
- **Sensitive features:** gender.
- **Link:** not available
- **Further info:** Karako and Manggala [264]; <https://burst.shopify.com/>

A.29 Business Entity Resolution

- **Description:** A proprietary Google dataset, where the task is to predict whether a pair of business descriptions describe the same real business.
- **Affiliation of creators:** Google.
- **Domain:** linguistics.
- **Tasks in fairness literature:** fair entity resolution [119].
- **Data spec:** text.
- **Sample size:** ~15K samples.
- **Year:** 2019.
- **Sensitive features:** geography, business size.
- **Link:** not available
- **Further info:** Cotter et al. [119]

A.30 Campus Recruitment

- **Description:** this dataset was published to Kaggle in 2020 by Ben Roshan, who was then enrolled in an MBA in Business Analytics at Jain University Bangalore. The provenance of this dataset is not clear. It was provided by a Jain University professor as a class resource to study and experiment with data analysis. It encodes information about students at an Indian institution, including their degree, their performance in school and placement information at the end of school, including salary.
- **Affiliation of creators:** Jain University Bangalore.
- **Domain:** education.
- **Tasks in fairness literature:** fair data generation [321].
- **Data spec:** tabular data.
- **Sample size:** ~ 200 students.
- **Year:** 2020.
- **Sensitive features:** gender.
- **Link:** <https://www.kaggle.com/datasets/benroshan/factors-affecting-campus-placement>
- **Further info:**

A.31 Cars3D

- **Description:** this dataset consists of CAD-generated models of 199 cars rendered from from 24 rotation angles. Originally devised for visual analogy making, it is also used for more general research on learning disentangled representation.
- **Affiliation of creators:** University of Michigan.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair representation learning [327].
- **Data spec:** image.

- **Sample size:** ~ 5K images.
- **Year:** 2020.
- **Sensitive features:** none.
- **Link:** https://github.com/google-research/disentanglement_lib/tree/master/disentanglement_lib/data/ground_truth
- **Further info:** Reed et al. [428]

A.32 CelebA

- **Description:** CelebFaces Attributes Dataset (CelebA) features images of celebrities from the CelebFaces dataset, augmented with annotations of landmark location and binary attributes. The attributes, ranging from highly subjective features (e.g. attractive, big nose) and potentially offensive (e.g. double chin) to more objective ones (e.g. black hair) were annotated by a “professional labeling company”.
- **Affiliation of creators:** Chinese University of Hong Kong.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair classification [108, 122, 255, 276, 328, 447], fair anomaly detection [567], bias discovery [11] fair anomaly detection [567], fairness evaluation of private classification [98], fairness evaluation of selective classification [253], fairness evaluation [451, 533], fair representation learning [414], fair data summarization [100], fair data generation [103, 424].
- **Data spec:** image.
- **Sample size:** ~ 200K face images of over ~ 10K unique individuals.
- **Year:** 2015.
- **Sensitive features:** gender, age, skin tone.
- **Link:** <http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>
- **Further info:** Liu et al. [326]

A.33 CheXpert

- **Description:** this dataset consists of chest X-ray images from patients that have been treated at the Stanford Hospital between October 2002 and July 2017. Each radiograph, either frontal or lateral, is annotated for the presence of 14 observations related to medical conditions. Most annotations were automatically extracted from free text radiology reports and validated against a set of 1,000 held-out reports, manually reviewed by a radiologist. For a subset of the X-ray images, high-quality labels are provided by a group of 3 radiologists. The task associated with this dataset is the automated diagnosis of medical conditions from radiographs.
- **Affiliation of creators:** Stanford University.
- **Domain:** radiology.
- **Tasks in fairness literature:** fairness evaluation of selective classification [253], fairness evaluation of private classification [98].
- **Data spec:** image.
- **Sample size:** ~ 200K chest radiographs from 60K patients.
- **Year:** 2019.
- **Sensitive features:** sex, age (of patient).
- **Link:** <https://stanfordmlgroup.github.io/competitions/chexpert/>
- **Further info:** Garbin et al. [178], Irvin et al. [244]

A.34 Chicago Ridesharing

- **Description:** this resource describes all trips reported by ridesharing companies to the City of Chicago, starting November 2018. It is the result of an ongoing transparency effort, following the introduction of a city-wide ordinance requiring the disclosure of trips and fares on part of transportation network providers. For each trip, this dataset reports geographical information (pickup and dropoff), duration and cost. To avoid individual re-identification, the granularity of times and locations is reduced to the nearest 15-minutes interval and census tract. Moreover, for rare combinations of census tract an interval, location data is provided at coarser granularity (community area).
- **Affiliation of creators:** City of Chicago.
- **Domain:** transportation.
- **Tasks in fairness literature:** fair pricing evaluation [392].
- **Data spec:** tabular data.
- **Sample size:** ~ 200M trips.
- **Year:** present.
- **Sensitive features:** geography.
- **Link:** <https://data.cityofchicago.org/Transportation/Transportation-Network-Providers-Trips/m6dm-c72p>
- **Further info:** <http://dev.cityofchicago.org/open%20data/data%20portal/2020/04/28/tnp-trips-2019-additional.html>; <http://dev.cityofchicago.org/open%20data/data%20portal/2019/04/12/tnp-taxi-privacy.html>

A.35 CIFAR

- **Description:** CIFAR-10 and CIFAR-100 are a labelled subset of the 80 million tiny images database. CIFAR consists of 32x32 colour images that students were paid to annotate. The project, aimed at advancing the effectiveness of supervised learning techniques in computer vision, was funded by the the Canadian Institute for Advanced Research, after which the dataset is named.
- **Affiliation of creators:** University of Toronto.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair classification [255, 533], fair incremental learning [579], robust fairness evaluation [378].
- **Data spec:** image.
- **Sample size:** ~ 6K images x 10 classes (CIFAR-10) or 600 images x 100 classes (CIFAR-100).
- **Year:** 2009.
- **Sensitive features:** none.
- **Link:** <https://www.cs.toronto.edu/~kriz/cifar.html>
- **Further info:** Krizhevsky [290]
- **Variants:** CIFAR-10S [533] is a modified version specifically aimed at studying biases in image classification across an artificial sensitive attribute (color/grayscale).

A.36 CiteSeer Papers

- **Description:** this dataset was created to study the problem of link-based classification of connected entities. The creators extracted a network of papers from CiteSeer, belonging to one of six categories: Agents, Artificial Intelligence, Database, Human Computer Interaction, Machine Learning and Information Retrieval. Each article is associated with a bag-of-word representation, and the associated task is classification into one of six topics.
- **Affiliation of creators:** University of Maryland.
- **Domain:** library and information sciences.
- **Tasks in fairness literature:** fair graph mining [312].
- **Data spec:** paper-paper pairs.
- **Sample size:** ~ 3K articles connected by ~ 5K citations.
- **Year:** 2016.
- **Sensitive features:** none.
- **Link:** <http://networkrepository.com/citeseer.php>
- **Further info:** Lu and Getoor [331]

A.37 Civil Comments

- **Description:** this dataset derives from an archive of the Civil Comments platform, a browser plugin for independent news sites, whose users peer-reviewed each other's comments with civility ratings. When the plugin shut down, they decided to make comments and metadata available, including the crowd-sourced toxicity ratings. A subset of this dataset was later annotated with a variety of sensitive attributes, capturing whether members of a certain group are mentioned in comments. This dataset powers the Jigsaw Unintended Bias in Toxicity Classification challenge.
- **Affiliation of creators:** Jigsaw; Civil Comments.
- **Domain:** social media.
- **Tasks in fairness literature:** fair toxicity classification [2, 108, 561], fairness evaluation of selective classification [253], fair robust toxicity classification [2], fairness evaluation of toxicity classification [241], fairness evaluation [19].
- **Data spec:** text.
- **Sample size:** ~ 2M comments.
- **Year:** 2019.
- **Sensitive features:** race/ethnicity, gender, sexual orientation, religion, disability.
- **Link:** <https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification>
- **Further info:** Borkan et al. [52]

A.38 Climate Assembly UK

- **Description:** this resource was curated to study the problem of subset selection for *sortition*, a political system where decisions are taken by a subset of the whole voting population selected at random. The data describes participants to Climate Assembly UK, a panel organized by the Sortition Foundation in 2020. With the goal of understanding public opinion on how the UK can meet greenhouse gas emission targets. The panel consisted of 110 UK residents selected from a pool of 1,715 who responded to an invitation from the Sortition Foundation reaching ~ 60K citizens. Features for each subject in the pool describe their demographics and climate concern level.
- **Affiliation of creators:** Carnegie Mellon University; Harvard University; Sortition Foundation.

- **Domain:** political science.
- **Tasks in fairness literature:** fair subset selection [171].
- **Data spec:** tabular data.
- **Sample size:** ~ 2K pool participants.
- **Year:** 2020.
- **Sensitive features:** gender, age, education, urban/rural, geography, ethnicity.
- **Link:** not available
- **Further info:** Flanigan et al. [171]; <https://www.climateassembly.uk/>

A.39 Columbia University Speed Dating

- **Description:** this dataset is a result of a speed dating experiment aimed at understanding preferences in mate selection in men and women. Subjects were recruited from students at Columbia University. Fourteen rounds were conducted with different proportions of male and female subjects, over the period 2002–2004, with participants meeting each potential mate for four minutes and rating them thereafter on six attributes. They also provide an overall evaluation of each potential mate and a binary decision indicating interest in meeting again. Before an event, each participant filled in a survey disclosing their preferences, expectations, and demographics. The inference task associated with this dataset is optimal recommendation in symmetrical two-sided markets.
- **Affiliation of creators:** Columbia University; Harvard University; Stanford University.
- **Domain:** sociology.
- **Tasks in fairness literature:** fair matching [586], preference-based fair ranking [394].
- **Data spec:** person-person pairs.
- **Sample size:** ~ 10K dating records involving ~ 400 people.
- **Year:** 2016.
- **Sensitive features:** gender, age, race, geography.
- **Link:** <https://data.world/annavmontoya/speed-dating-experiment>
- **Further info:** Fisman et al. [170]

A.40 Communities and Crime

- **Description:** this dataset was curated to develop a software tool supporting the work of US police departments. It was especially aimed at identifying similar precincts to exchange best practices and share experiences among departments. The creators were supported by the police departments of Camden (NJ) and Philadelphia (PA). The factors included in the dataset were the ones deemed most important to define similarity of communities from the perspective of law enforcement; they were chosen with the help of law enforcement officials from partner institutions and academics of criminal justice, geography and public policy. The dataset includes socio-economic factors (aggregate data on age, income, immigration, and racial composition) obtained from the 1990 US census, along with information about policing (e.g. number of police cars available) based on the 1990 Law Enforcement Management and Administrative Statistics survey, and crime data derived from the 1995 FBI Uniform Crime Reports. In its released version on UCI, the task associated with the dataset is predicting the total number of violent crimes per 100K population in each community. The most referenced version of this dataset was preprocessed with a normalization step; after receiving multiple requests, the creators also published an unnormalized version.
- **Affiliation of creators:** La Salle University; Rutgers University.
- **Domain:** law.
- **Tasks in fairness literature:** fair classification [118, 119, 122, 222, 328, 455, 551], fair regression evaluation [223], fair few-shot learning [466, 468], rich-subgroup fairness evaluation [270], rich-subgroup fair classification [269], fair regression [4, 41, 110, 111, 137, 286, 344, 386, 435], fair representation learning [441], robust fair classification [341], fair private classification [247], fairness evaluation of transfer learning [300], preference-based fair clustering [177].
- **Data spec:** tabular data.
- **Sample size:** ~ 2K communities.
- **Year:** 2009.
- **Sensitive features:** race, geography.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/communities+and+crime> and <http://archive.ics.uci.edu/ml/datasets/communities+and+crime+unnormalized>
- **Further info:** Redmond and Baveja [426]

A.41 COMPAS

- **Description:** this dataset was created for an external audit of racial biases in the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) risk assessment tool developed by Northpointe (now Equivant), which estimates the likelihood of a defendant becoming a recidivist. Instances represent defendants scored by COMPAS in Broward County, Florida, between 2013–2014,

reporting their demographics, criminal record, custody and COMPAS scores. Defendants' public criminal records were obtained from the Broward County Clerk's Office website matching them based on date of birth, first and last names. The dataset was augmented with jail records and COMPAS scores provided by the Broward County Sheriff's Office. Finally, public incarceration records were downloaded from the Florida Department of Corrections website. Instances are associated with two target variables (`is_recid` and `is_violent_recid`), indicating whether defendants were booked in jail for a criminal offense (potentially violent) that occurred after their COMPAS screening but within two years. See Appendix C for extensive documentation.

- **Affiliation of creators:** ProPublica.
- **Domain:** law.
- **Tasks in fairness literature:** fair classification [9, 41, 72, 73, 81, 83, 101, 119, 131, 137, 138, 143, 190, 221, 222, 328, 332, 340, 344, 386, 388, 398, 413, 433, 434, 442, 447, 458, 515, 521, 527, 549, 562], fairness evaluation [3, 75, 86, 105, 113, 176, 203, 245, 246, 267, 323, 356, 381, 405, 476, 490, 537, 569], fair risk assessment [116, 367, 376], fair *task assignment* [189], fair classification under unawareness [109, 272, 297, 299], data bias evaluation [40], fair representation learning [55, 441, 581], robust fair classification [44, 341, 429], dynamical fairness evaluation [572], fair reinforcement learning [362], fair ranking evaluation [259, 553], fair multi-stage classification [338], dynamical fair classification [512], preference-based fair classification [507, 563], fair regression [286], fair multi-stage classification [188], limited-label fair classification [104, 109, 528], robust fairness evaluation [467, 468], rich subgroup fairness evaluation [106, 577].
- **Data spec:** tabular data.
- **Sample size:** ~ 12K defendants.
- **Year:** 2016.
- **Sensitive features:** sex, age, race.
- **Link:** <https://github.com/propublica/compas-analysis>
- **Further info:** Angwin et al. [15], Larson et al. [301]

A.42 Cora Papers

- **Description:** this resource was produced within the wider development effort for *Cora*, an Internet portal for computer science research papers available in the early 2000s. The portal supported keyword search, topical categorization of articles, and citation mapping. This dataset consists of articles and citation links between them. It contains bag-of-word representations for the text of each article, and the associated task is classification into one of seven topics.
- **Affiliation of creators:** Just Research Carnegie Mellon University; Massachusetts Institute of Technology; Univeristy of Maryland; Lawrence Livermore National Laboratory.
- **Domain:** library and information sciences.
- **Tasks in fairness literature:** .
- **Data spec:** article-article pairs.
- **Sample size:** ~ 3K articles connected by ~ 5K citations.
- **Year:** 2019.
- **Sensitive features:** none.
- **Link:** <https://relational.fit.cvut.cz/dataset/CORA>
- **Further info:** McCallum et al. [350], Sen et al. [452]

A.43 Costarica Household Survey

- **Description:** this data comes from the national household survey of Costa Rica, performed by the national institute of statistics and census (Instituto Nacional de Estadística y Censos). The survey is aimed at measuring the socio-economical situation in the country and informing public policy. The data collection procedure is specially designed to allow for precise conclusions with respect to six different regions of the country and about differences in urban vs rural areas; stratification along these variables is deemed suitable. The 2018 survey contains a special section on the crimes suffered by respondents.
- **Affiliation of creators:** Instituto Nacional de Estadística y Censos.
- **Domain:** economics.
- **Tasks in fairness literature:** fair classification [384].
- **Data spec:** tabular data.
- **Sample size:** ~ 13K households.
- **Year:** 2018.
- **Sensitive features:** sex, age, birthplace, disability, geography, family size.
- **Link:** <https://www.inec.cr/encuestas/encuesta-nacional-de-hogares>
- **Further info:** <https://www.inec.cr/sites/default/files/documentos-biblioteca-virtual/enaho-2018.pdf>

A.44 Credit Card Default

- **Description:** this dataset was built to investigate automated mechanisms for credit card default prediction following a wave of defaults in Taiwan connected to patters of card over-issuing and over-usage. The dataset contains payment history of customers of an important Taiwanese bank, from April to October 2005. Demographics, marital status, and education of customers are also provided, along with the amount of credit and a binary variable encoding default on payment, which is the target variable of the associated task.
- **Affiliation of creators:** Chung-Hua University; Thompson Rivers University.
- **Domain:** finance.
- **Tasks in fairness literature:** fair classification [41, 101], fair clustering [39, 186, 214, 214], fair clustering under unawareness [157], fair classification under unawareness [531], fair data summarization [446, 489], fairness evaluation [320], fair anomaly detection [459].
- **Data spec:** tabular data.
- **Sample size:** ~ 30K credit card holders.
- **Year:** 2016.
- **Sensitive features:** gender, age.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/default+of+credit+card+clients>
- **Further info:** Yeh and hui Lien [557]

A.45 Credit Elasticities

- **Description:** this dataset stems from a randomized trial conducted by a consumer lender in South Africa to study loan price elasticity. Prior customers were contacted by mail with limited-time loan offers at variable and randomized interest rates. The aim of the study was understanding the relationship between interest rate and customer acceptance rates, along with the benefits for the lender. Customers who accepted and received formal approval, filled in a short survey with factors of interest for the study, including demographics, education, and prior borrowing history.
- **Affiliation of creators:** Yale University; Dartmouth College.
- **Domain:** finance.
- **Tasks in fairness literature:** fair pricing evaluation [260].
- **Data spec:** tabular data.
- **Sample size:** ~ 50K clients.
- **Year:** 2008.
- **Sensitive features:** gender, age, geography.
- **Link:** <http://doi.org/10.3886/E113240V1>
- **Further info:** Karlan and Zinman [266]

A.46 Crowd Judgement

- **Description:** this dataset was assembled to compare the performance of the COMPAS recidivism risk prediction system against that of non-expert human assessors [144]. A subset of 1,000 defendants were selected from the COMPAS dataset. Crowd-sourced assessors were recruited through Amazon Mechanical Turk. They were presented with a summary of each defendant, including demographics and previous criminal history, and asked to predict whether they would recidivate within 2 years of their most recent crime. These judgements, assembled via plain majority voting, ended up exhibiting accuracy and fairness levels comparable to that displayed by the COMPAS system. While this dataset was assembled for an experiment, it was later used to study the problem of fairness in crowdsourced judgements.
- **Affiliation of creators:** Dartmouth College.
- **Domain:** law.
- **Tasks in fairness literature:** fair *truth discovery* [316], fair *task assignment* [189, 316]
- **Data spec:** judge-defendant pair.
- **Sample size:** ~ 1K defendants from COMPAS and ~ 400 crowd-sourced labellers. Each defendant is judged by 20 different labellers.
- **Year:** 2018.
- **Sensitive features:** sex, age and race of defendants and crowd-sourced judges.
- **Link:** <https://farid.berkeley.edu/downloads/publications/scienceadvances17/>
- **Further info:** [144]
- **Variants:** a similar dataset was collected by Wang et al. [526].

A.47 Curatr British Library Digital Corpus

- **Description:** this dataset is a subset of English language digital texts from the British Library focused on volumes of 19th-century fiction, obtained through the Curatr platform. It was selected for the well-researched presence of stereotypical and binary concepts of gender in this literary production. The goal of the creators was studying gender biases in large text corpora and their relationship with biases in word embeddings trained on those corpora.

- **Affiliation of creators:** University College Dublin.
- **Domain:** literature.
- **Tasks in fairness literature:** data bias evaluation [304].
- **Data spec:** text.
- **Sample size:** ~ 20K books.
- **Year:** 2020.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <http://curatr.ucd.ie/>
- **Further info:** Leavy et al. [303]

A.48 CVs from Singapore

- **Description:** this dataset was developed to test demographic biases in resume filtering. In particular, the authors studied nationality bias in automated resume filtering in Singapore, across the three major ethnic groups of the city state: Chinese, Malaysian and Indian. The dataset consists of 135 resumes (45 per ethnic group) used for application to finance jobs in Singapore, collected by Jai Janyani. The dataset only includes resumes for which the origin of the candidates can be reliably inferred to be either Chinese, Malaysian, or Indian from education and initial employment. The dataset also comprises 9 finance job postings from China, Malaysia, and India (3 per country). All job-resume pairs are rated for relevance/suitability by three annotators.
- **Affiliation of creators:** University of Maryland.
- **Domain:** information systems, management information systems.
- **Tasks in fairness literature:** fair ranking [133].
- **Data spec:** text.
- **Sample size:** ~ 100 resumes.
- **Year:** 2020.
- **Sensitive features:** ethnic group.
- **Link:** not available
- **Further info:** Deshpande et al. [133]

A.49 Dallas Police Incidents

- **Description:** this dataset is due to the Dallas OpenData initiative⁹ and “reflects crimes as reported to the Dallas Police Department” beginning June 1, 2014. Each incident comes with rich spatio-temporal data, information about the victim, the officers involved and the type of crime. A subset of the dataset is available on Kaggle¹⁰.
- **Affiliation of creators:** Dallas Police Department.
- **Domain:** law.
- **Tasks in fairness literature:** fair spatio-temporal process learning [454].
- **Data spec:** tabular.
- **Sample size:** ~ 800K incidents.
- **Year:** present.
- **Sensitive features:** age, race, and gender (of victim), geography.
- **Link:** <https://www.dallasopendata.com/Public-Safety/Police-Incidents/qv6i-rr17>
- **Further info:**

A.50 Demographics on Twitter

- **Description:** this dataset was developed to test demographic classifiers on Twitter data. In particular, the tasks associated with this resource are the automatic inference of gender, age, location and political orientation of users. The true values for these attributes, which act as a ground truth for learning algorithms, were inferred from tweets and user bios, such as the ones containing the regexp "I'm a <gendered noun>", with gendered nouns including mother, woman, father, man.
- **Affiliation of creators:** Massachusetts Institute of Technology.
- **Domain:** social media.
- **Tasks in fairness literature:** fairness evaluation of sentiment analysis [460].
- **Data spec:** mixture.
- **Sample size:** ~ 80K profiles.
- **Year:** 2017.
- **Sensitive features:** gender, age, political orientation, geography.
- **Link:** not available

⁹<https://www.dallasopendata.com/>

¹⁰<https://www.kaggle.com/carrie1/dallaspoliceincidentreports>

- **Further info:** Vijayaraghavan et al. [517]

A.51 Diabetes 130-US Hospitals

- **Description:** this dataset contains 10 years of care data from 130 US hospitals extracted from Health Facts, a clinical database associated with a multi-institution data collection program. The dataset was extracted to study the association between the measurement of HbA1c (glycated hemoglobin) in human bloodstream and early hospital readmission, and was donated to UCI in 2014. The dataset includes patient demographics, in-hospital procedures, and diagnoses, along with information about subsequent readmissions.
- **Affiliation of creators:** Virginia Commonwealth University; University of Cordoba; Polish Academy of Sciences.
- **Domain:** endocrinology.
- **Tasks in fairness literature:** fair clustering [21, 39, 39, 99, 236, 339].
- **Data spec:** tabular data.
- **Sample size:** ~ 100K patients.
- **Year:** 2014.
- **Sensitive features:** age, race, gender.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/diabetes+130-us+hospitals+for+years+1999-2008>
- **Further info:** Strack et al. [480]

A.52 Diversity in Faces (DiF)

- **Description:** this large dataset was created to favour the development and evaluation of robust face analysis algorithms across diverse demographics and domain-specific features, such as craniofacial distances and facial contrast). One million images of people's faces from Flickr were labelled, mostly automatically, according to 10 different coding schemes, comprising, e.g., cranio-facial measurements, pose, and demographics. Age and gender were inferred both automatically and by human workers. Statistics about the diversity of this dataset along these coded measures are available in the accompanying report.
- **Affiliation of creators:** IBM.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair representation learning [414], fairness evaluation of private classification [22].
- **Data spec:** image.
- **Sample size:** ~ 1M images.
- **Year:** 2019.
- **Sensitive features:** skin color, age, and gender.
- **Link:** <https://www.ibm.com/blogs/research/2019/01/diversity-in-faces/>
- **Further info:** Merler et al. [361]

A.53 Drug Consumption

- **Description:** this dataset was collected by Elaine Fehrman between March 2011 and March 2012 after receiving approval from relevant ethics boards from the University of Leicester. The goal of this dataset is to seek patterns connecting an individual's risk of drug consumption with demographics and psychometric measurements of the Big Five personality traits (NEO-FFI-R), impulsivity (BIS-11), and sensation seeking (ImpSS). The study employed an online survey tool from Survey Gizmo to recruit participants world-wide; over 93% of the final usable sample reported living in an English-speaking country. Target variables summarize the consumption of 18 psychoactive substances on an ordinal scale ranging from never using the drug to using it over a decade ago, or in the last decade, year, month, week, or day. The 18 substances considered in the study are classified as central nervous system depressants, stimulants, or hallucinogens and comprise the following: alcohol, amphetamines, amyl nitrite, benzodiazepines, cannabis, chocolate, cocaine, caffeine, crack, ecstasy, heroin, ketamine, legal highs, LSD, methadone, magic mushrooms, nicotine, and Volatile Substance Abuse (VSA), along with one fictitious drug (Semeron) introduced to identify over-claimers. A version of the dataset donated to the UCI Machine Learning Repository is associated with 18 prediction tasks, i.e. one per substance.
- **Affiliation of creators:** Rampton Hospital; Nottinghamshire Healthcare NHS Foundation Trust; University of Leicester; University of Nottingham; University of Salahaddin.
- **Domain:** applied psychology.
- **Tasks in fairness literature:** fair classification [143, 344], evaluation of data bias [40], limited-label fair classification [109], robust fair classification [429].
- **Data spec:** tabular data.
- **Sample size:** ~ 2K respondents.
- **Year:** 2016.
- **Sensitive features:** age, gender, ethnicity, geography.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/Drug+consumption+%28quantified%29>
- **Further info:** Fehrman et al. [163, 164]

A.54 DrugNet

- **Description:** this dataset was collected to study drug consumption patterns in connection with social ties and behaviour of drug users. This work puts particular emphasis on situations at risk of disease transmission and to assess the opportunity for prevention via recruitment of peer educators to demonstrate, disseminate and support HIV prevention practices among their connections. Participants were recruited in Hartford neighbourhoods of high drug-use activity, mostly via street outreach and recruitment by early participants. Eligibility criteria included being at least 18 years old, using an illicit drug, and signing an informed consent form. Each participant provided data about their drug use, most common sites of usage, HIV risk practices associated with drug use and sexual behavior, and social ties deemed important by the respondent and their demographics.
- **Affiliation of creators:** Institute for Community Research of Hartford; Hispanic Health Council, Hartford; Boston College.
- **Domain:** social work, social networks.
- **Tasks in fairness literature:** fair graph clustering [281].
- **Data spec:** person-person pairs.
- **Sample size:** ~ 300 people.
- **Year:** 2016.
- **Sensitive features:** ethnicity, sex, age.
- **Link:** <https://sites.google.com/site/ucinetsoftware/datasets/covert-networks/drugnet>
- **Further info:** Weeks et al. [536]

A.55 dSprites

- **Description:** this dataset was assembled by researchers affiliated with Google DeepMind as an artificial benchmark for unsupervised methods aimed at learning disentangled data representations. Each image in the dataset consists of a black-and-white sprite with variable shape, scale, orientation and position. Together these are the *generative factors* underlying each image. Ideally, systems trained on this data should learn disentangled representations, such that latent image representations are clearly associated with changes in a single generative factor.
- **Affiliation of creators:** Google.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair representation learning [122, 327].
- **Data spec:** image.
- **Sample size:** ~ 700K images.
- **Year:** 2017.
- **Sensitive features:** none.
- **Link:** <https://github.com/deepmind/dsprites-dataset>
- **Further info:** Higgins et al. [226]

A.56 Dutch Census

- **Description:** this dataset was derived from the 2001 census carried out by the Dutch Central Bureau for Statistics to gather data about family composition, economic activities, levels of education, and occupation of Dutch citizens and foreigners from various countries of origin. A version of the dataset commonly employed in the fairness research literature has been preprocessed and made available online. The associated task is the classification of individuals into high-income and low-income professions.
- **Affiliation of creators:** Bournemouth University; TU Eindhoven.
- **Domain:** demography.
- **Tasks in fairness literature:** fair classification [3, 328, 549, 570], fairness evaluation [75].
- **Data spec:** tabular data.
- **Sample size:** ~ 60K respondents.
- **Year:** 2001.
- **Sensitive features:** sex, age, citizenship.
- **Link:** <https://sites.google.com/site/conditionaldiscrimination/>
- **Further info:** Žliobaite et al. [592]; https://microdata.worldbank.org/index.php/catalog/2102/data-dictionary/F2?file_name=NLD2001-P-H; <https://www.cbs.nl/nl-nl/publicatie/2004/31/the-dutch-virtual-census-of-2001>

A.57 EdGap

- **Description:** this dataset focuses on education performance in different US counties, with a focus on inequality of opportunity and its connection to socioeconomic factors. Along with average SAT and ACT test scores by county, this dataset reports socioeconomic data from the American Community Survey by the Bureau of Census, including household income, unemployment, adult educational attainment, and family structure. Importantly, some states require all students to take ACT or SAT tests, while others do not. As a

result, average test scores are inherently higher in states that do not require all students to test, and they are not directly comparable to average scores in states where testing is mandatory.

- **Affiliation of creators:** Memphis Teacher Residency.
- **Domain:** education.
- **Tasks in fairness literature:** fair risk assessment [220].
- **Data spec:** tabular data.
- **Sample size:** ~ 2K counties.
- **Year:** 2019.
- **Sensitive features:** geography.
- **Link:** <https://www.edgap.org/>
- **Further info:**

A.58 Epileptic Seizures

- **Description:** this dataset was curated to study electroencephalographic (EEG) time series in relation to epilepsy. The dataset consists of EEG recordings from healthy volunteers with eyes closed and eyes open, and from epilepsy patients during seizure-free intervals and during epileptic seizures. Volunteers and patients are recorded for 23.6-sec. A version of this dataset, used in fairness research, was donated to UCI Machine Learning Repository by researchers affiliated with Rochester Institute of Technology in 2017, with a classification task based on the patients' condition and state at the time of recording. The data was later removed from UCI at the original curators' request.
- **Affiliation of creators:** University of Bonn.
- **Domain:** neurology.
- **Tasks in fairness literature:** robust fairness evaluation [45].
- **Data spec:** time series.
- **Sample size:** ~ 500 individuals, each summarized by ~ 4K-points time series.
- **Year:** 2017.
- **Sensitive features:** none.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/Epileptic+Seizure+Recognition>; <http://epileptologie-bonn.de/cms/upload/workgroup/lehnertz/eegdata.html>
- **Further info:** Andrzejak et al. [14]

A.59 Equitable School Access in Chicago

- **Description:** this resource was assembled from disparate sources to evaluate school access in Chicago for different race groups. A transportation network was inferred from data on public bus lines available on the Chicago Transit Authority website. Data on school location and quality evaluation was obtained from the Chicago Public School data portal. Finally, demographic information on race representation in different tracts was retrieved from the 2010 US census.
- **Affiliation of creators:** Salesforce.
- **Domain:** transportation.
- **Tasks in fairness literature:** fair graph augmentation [423].
- **Data spec:** location-location pairs.
- **Sample size:** ~ 2K nodes (locations), connected by ~ 8K edges (bus lines).
- **Year:** 2020.
- **Sensitive features:** race.
- **Link:** <https://github.com/salesforce/GAEA>
- **Further info:** Ramachandran et al. [423]

A.60 Equity Evaluation Corpus (EEC)

- **Description:** this dataset was compiled to audit sentiment analysis systems for gender and race bias. It is based on 11 short sentence templates; 7 templates include emotion words, while the remaining 4 do not. Moreover, each sentence includes one gender- or race-associated word, such as names predominantly associated with African American or European American people. Gender-related words consist of names, nouns, and pronouns.
- **Affiliation of creators:** National Research Council Canada.
- **Domain:** linguistics.
- **Tasks in fairness literature:** fair sentiment analysis evaluation [318].
- **Data spec:** text.
- **Sample size:** ~ 9K sentences.
- **Year:** 2018.

- **Sensitive features:** race, gender.
- **Link:** <https://saifmohammad.com/WebPages/Biases-SA.html>
- **Further info:** Kiritchenko and Mohammad [277]

A.61 Facebook Ego-networks

- **Description:** this dataset was collected to study the problem of identifying users' social circles, i.e. categorizing links between nodes in a social network. The data represents ten ego-networks whose central user was asked to fill in a survey and manually identify the circles to which their friends belonged. Features from each profile, including education, work and location are anonymized.
- **Affiliation of creators:** Stanford University.
- **Domain:** social networks.
- **Tasks in fairness literature:** fair graph mining [312].
- **Data spec:** user-user pairs.
- **Sample size:** ~ 4K people connected by ~ 90K friend relations.
- **Year:** 2012.
- **Sensitive features:** geography, gender.
- **Link:** <https://snap.stanford.edu/data/egonets-Facebook.html>
- **Further info:** Leskovec and Mcauley [310]

A.62 Facebook Large Network

- **Description:** this dataset was developed to study the effectiveness of node embeddings for learning tasks defined on graphs. The dataset concentrates on verified Facebook pages of politicians, governmental organizations, television shows, and companies, represented as nodes, while edges represent mutual likes. In addition, each page comes with node embeddings which are extracted from the textual description of each page. The original task on this dataset is page category classification.
- **Affiliation of creators:** University of Edinburgh.
- **Domain:** social networks.
- **Tasks in fairness literature:** fair graph mining evaluation [262].
- **Data spec:** page-page pairs.
- **Sample size:** ~20K nodes (pages) connected by ~ 200K edges (mutual likes).
- **Year:** 2019.
- **Sensitive features:** none.
- **Link:** <http://snap.stanford.edu/data/facebook-large-page-page-network.html>
- **Further info:** Rozemberczki et al. [438]

A.63 FACES

- **Description:** this resource contains images of Caucasian individuals of variable age and gender under six predefined facial expressions (neutrality, sadness, disgust, fear, anger, and happiness). This dataset is described as a database of emotion-related stimuli for scientific research. Subjects were hired through a model agency in Berlin, and suitably informed about the purpose of the photo-shooting session, thereafter signing an informed consent document. Each model reported their own age and gender. The necessary facial expressions were carefully explained with the help of a manual, with attention to the position of muscles. Photographs were obtained and post-processed in a standardized fashion, and later validated by raters of different ages with respect to the perceived expression and age of subjects. At a later stage, images were also annotated for attractiveness and distinctiveness. Currently, a small subset of the images is publicly available, while the full dataset is available after registration.
- **Affiliation of creators:** Max Planck Institute for Human Development.
- **Domain:** computer vision, experimental psychology.
- **Tasks in fairness literature:** fairness evaluation [273].
- **Data spec:** image.
- **Sample size:** ~ 2K images of ~ 200 people.
- **Year:** 2010.
- **Sensitive features:** age, gender.
- **Link:** <https://faces.mpib-berlin.mpg.de/imeji/>
- **Further info:** Ebner et al. [149]

A.64 FairFace

- **Description:** this dataset was developed as a balanced resource for face analysis with diverse race, gender and age composition. The associated task is race, gender and age classification. Starting from a large public image dataset (Yahoo YFCC100M), the authors sampled images incrementally to ensure diversity with respect to race, for which they considered seven categories: White, Black,

Indian, East Asian, Southeast Asian, Middle East, and Latino. Sensitive attributes were annotated by workers on Amazon Mechanical Turk, and also through a model based on these annotations. Faces with low agreement between model and annotators were manually re-verified by the dataset curators. This dataset was annotated automatically with a binary Fitzpatrick skin tone label [98].

- **Affiliation of creators:** University of California, Los Angeles.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fairness evaluation of private classification [98].
- **Data spec:** image.
- **Sample size:** ~ 100K images.
- **Year:** 2019.
- **Sensitive features:** race, age, gender, skin tone.
- **Link:** <https://github.com/joojs/fairface>
- **Further info:** Karkkainen and Joo [265]

A.65 Fantasy Football

- **Description:** this resource was curated to study the problem of fair ranking aggregation. The creators collected rankings of National Football League players from the top 25 experts on the popular fantasy sports website FantasyPros. The data covers 16 weeks during the 2019 football season. Players are assigned to different sensitive groups based on the conference of their team (American Football Conference or National Football Conference). The data available online concentrates on wide receivers.
- **Affiliation of creators:** Worcester Polytechnic Institute.
- **Domain:** sports.
- **Tasks in fairness literature:** fair ranking evaluation [291].
- **Data spec:** player-expert pairs.
- **Sample size:** ~ 50 players, ranked by 25 experts (on a weekly basis), over 16 weeks.
- **Year:** 2020.
- **Sensitive features:** football conference.
- **Link:** <https://arcgit.wpi.edu/cakuhlman/VLDB2020/tree/master/charts/data>
- **Further info:** Kuhlman and Rundensteiner [292]

A.66 Fashion MNIST

- **Description:** this dataset is based on product assortment from the Zalando website. It contains gray-scale resized versions of thumbnail images of unique clothing products, labeled by in-house fashion experts according to their category, including e.g. trousers, coat and shirt. The envisioned task is object classification. The dataset, sharing the same size and structure as MNIST, was developed to provide a harder and more representative task, and to replace MNIST as a popular computer vision benchmark.
- **Affiliation of creators:** Zalando.
- **Domain:** computer vision.
- **Tasks in fairness literature:** robust fairness evaluation [45].
- **Data spec:** image.
- **Sample size:** ~ 70K images across 10 product categories.
- **Year:** 2017.
- **Sensitive features:** none.
- **Link:** <https://github.com/zalando-research/fashion-mnist>
- **Further info:** Xiao et al. [545]

A.67 FICO

- **Description:** based on a sample of 301,536 TransUnion TransRisk scores from 2003, this dataset was created to study the problem of adjusting predictors for compliance with the equality of opportunity fairness metric. The TransUnion data was preprocessed and aggregated to summarize the CDF of risk scores by race (Non-Hispanic white, Black, Hispanic, Asian). The original data comes from a 2007 report to the US Congress on credit scoring and its effects on the availability and affordability of credit carried out by a dedicated Federal Reserve working group. The collection, creation, processing, and aggregation was carried out by the working group; the data was later scraped by the creators, who made it available without any modification.
- **Affiliation of creators:** Google; University of Texas at Austin; Toyota Technological Institute at Chicago.
- **Domain:** finance.
- **Tasks in fairness literature:** fairness evaluation [215], dynamical fair classification [324], dynamical fairness evaluation [123, 322, 572], fair resource allocation [193].
- **Data spec:** tabular data.
- **Sample size:** N/As. CDFs are provided over risk scores which are normalized (0-100%) and quantized with step 0.5%.

- **Year:** 2016.
- **Sensitive features:** race.
- **Link:** <https://github.com/fairmlbook/fairmlbook.github.io/tree/master/code/creditscore/data>
- **Further info:** Barocas et al. [33], Hardt et al. [215], US Federal Reserve [506]

A.68 FIFA 20 Players

- **Description:** this dataset was scraped by Stefano Leone and made available on Kaggle. It includes the players' data for the Career Mode from FIFA 15 to FIFA 20, a popular football game. Several tasks are envisioned for this dataset, including a historical comparison of players.
- **Affiliation of creators:** unknown.
- **Domain:** sports.
- **Tasks in fairness literature:** fairness evaluation under unawareness [18].
- **Data spec:** tabular data.
- **Sample size:** ~ 20K players.
- **Year:** 2019.
- **Sensitive features:** geography.
- **Link:** <https://www.kaggle.com/stefanoleone992/fifa-20-complete-player-dataset>
- **Further info:**

A.69 FilmTrust

- **Description:** this dataset was crawled from the entire FilmTrust website, a movie recommendation service with a social network component. The dataset comprises user-movie ratings on a 5-star scale and user-user indications of trust about movie taste. This resource can be used to train and evaluate recommender systems.
- **Affiliation of creators:** Northeastern University; Nanyang Technological University; American University of Beirut; University of Cambridge.
- **Domain:** information systems, movies.
- **Tasks in fairness literature:** fair ranking [325].
- **Data spec:** user-movie pairs and user-user pairs.
- **Sample size:** ~ 40K ratings by ~ 2K users over ~ 2K movies.
- **Year:** 2011.
- **Sensitive features:** none.
- **Link:** <https://guoguibing.github.io/librec/datasets.html>
- **Further info:** Guo et al. [207]

A.70 Framingham

- **Description:** the Framingham Heart Study began in 1948 under the direction of the National Heart, Lung, and Blood Institute (NHLBI), with the goal of identifying key factors that contribute to cardiovascular disease, given a mounting epidemic of cardiovascular disease whose etiology was mostly unknown at the time. Six different cohorts have been recruited over the years among citizens of Framingham, Massachusetts, without symptoms of cardiovascular disease. After the original cohort, two more were enrolled from the children and grandchildren of the first one. Additional cohorts were also started to reflect the increased racial and ethnic diversity in the town of Framingham. Participants in the study report on their habits (e.g. physical activity, smoking) and undergo regular physical examination and laboratory tests.
- **Affiliation of creators:** National Heart, Lung, and Blood Institute (NHLBI); Boston University.
- **Domain:** cardiology.
- **Tasks in fairness literature:** fair ranking evaluation [259].
- **Data spec:** mixture.
- **Sample size:** ~ 15K respondents.
- **Year:** present.
- **Sensitive features:** age, sex, race.
- **Link:** <https://framinghamheartstudy.org/>
- **Further info:** Kannel and McGee [263], Tsao and Vasan [498]

A.71 Freebase15k-237

- **Description:** Freebase was a collaborative knowledge base which allowed its community members to fill in structured data about diverse entities and relations between them. This database was developed from a prior Freebase dataset [50], pruning it from redundant relations and augmenting it with textual relationships from the ClueWeb12 corpus. The creators of this dataset worked on the joint

optimization of entity knowledge base and representations of the entities' textual relations, with the goal of providing representations of entities suited for knowledge base completion.

- **Affiliation of creators:** Microsoft; Stanford University.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair graph mining [53], fairness evaluation in graph mining [168].
- **Data spec:** entity-relation-entity triples.
- **Sample size:** ~ 15K entities connected by 170K edges (relations).
- **Year:** 2016.
- **Sensitive features:** demographics of people featured in entities and their relations.
- **Link:** <https://www.microsoft.com/en-us/download/details.aspx?id=52312>
- **Further info:** Toutanova et al. [496]

A.72 GAP Coreference

- **Description:** this resource was developed as a gender-balanced coreference resolution dataset, useful for auditing gender-dependent differences in the accuracy of existing pronoun resolution algorithms and for training new algorithms that are less gender-biased. The dataset consists of thousands of ambiguous pronoun-name pairs in sentences extracted from Wikipedia. Several measures are taken to avoid the success of naïve heuristics and to favour diversity. Most notably, while the initial (automated) stage of the data collection pipeline extracts contexts with a female:male ratio of 1:9, feminine pronouns are oversampled to achieve a 1:1 ratio. Each example is presented to and annotated for coreference by three in-house workers.
- **Affiliation of creators:** Google.
- **Domain:** linguistics.
- **Tasks in fairness literature:** data bias evaluation [284].
- **Data spec:** text.
- **Sample size:** ~ 9K sentences.
- **Year:** 2018.
- **Sensitive features:** gender.
- **Link:** <https://github.com/google-research-datasets/gap-coreference>
- **Further info:** Webster et al. [535]

A.73 German Credit

- **Description:** the German Credit dataset was created to study the problem of automated credit decisions at a regional Bank in southern Germany. Instances represent loan applicants from 1973 to 1975, who were deemed creditworthy and were granted a loan, bringing about a natural selection bias. The data summarizes their financial situation, credit history and personal situation, including housing and number of liable people. A binary variable encoding whether each loan recipient punctually payed every installment is the target of a classification task. Among covariates, marital status and sex are jointly encoded in a single variable. Many documentation mistakes are present in the UCI entry associated with this resource [501]. Due to one of these mistakes, users of this dataset are led to believe that the variable sex can be retrieved from the joint marital_status-sex variable, however this is false. A revised version with correct variable encodings, called South German Credit, was donated to UCI Machine Learning Repository [503] with an accompanying report [204]. See Appendix D for extensive documentation.
- **Affiliation of creators:** Hypo Bank (OP/EDV-VP); Universität Hamburg; Strathclyde University (German Credit); Beuth University of Applied Sciences Berlin (South German Credit).
- **Domain:** finance.
- **Tasks in fairness literature:** fair classification [23, 81, 131, 143, 221, 328, 343, 344, 398, 418, 419, 456, 458, 515, 551], fairness evaluation [165, 176], fair active resource allocation [67], preference-based fair classification [574], fair active classification [383], fair classification under unawareness [272], robust fairness evaluation [45], fair representation learning [329, 441], fair reinforcement learning [362], fair ranking evaluation [259, 543, 553], fair ranking [54, 463], fair multi-stage classification [188], limited-label fair classification [104, 109, 528], limited-label fairness evaluation [248].
- **Data spec:** tabular data.
- **Sample size:** ~ 1K.
- **Year:** 1994 (German Credit); 2020 (South German Credit).
- **Sensitive features:** age, geography.
- **Link:** [https://archive.ics.uci.edu/ml/datasets/statlog+\(german+credit+data\)](https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data)) (German Credit); <https://archive.ics.uci.edu/ml/datasets/South+German+Credit+%28UPDATE%29> (South German Credit)
- **Further info:** Grömping [204]

A.74 German Political Posts

- **Description:** this dataset was used as a training set for German word embeddings, with the goal of investigating biases in word representations. The authors used the Facebook and Twitter APIs to collect posts and comments from the social media channels of six main political parties in Germany (CDU/CSU, SPD, Bündnis90/Die Grünen, FDP, Die Linke, AfD). Facebook posts are from the period 2015–2018, while tweets were collected between January and October 2018. Overall, the dataset consists of millions of posts, for a total of half a billion tokens. A subset of the Facebook comments (100,000) were labeled by human annotators based on whether they contain sexist content, with four sub-labels indicating sexist comments, sexist buzzwords, gender-related compliments, statements against gender equality and assignment of gender stereotypical roles to people.
- **Affiliation of creators:** Technical University of Munich.
- **Domain:** social media.
- **Tasks in fairness literature:** bias evaluation in WEs [393].
- **Data spec:** text.
- **Sample size:** ~ 20M posts comments and tweets.
- **Year:** 2020.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** not available
- **Further info:** Papakyriakopoulos et al. [393]

A.75 GLUE

- **Description:** this benchmark was assembled to reliably evaluate the progress of natural language processing models. It consists of multiple datasets and associated tasks from the natural language processing domain, including paraphrase detection, textual entailment, sentiment analysis and question answering. Given the quick progress registered by language models on GLUE, a similar benchmark called SuperGLUE was subsequently released comprising more challenging and diverse tasks [524].
- **Affiliation of creators:** New York University; University of Washington; DeepMind.
- **Domain:** linguistics.
- **Tasks in fairness literature:** fairness evaluation [19, 439], bias evaluation in language models [97], fairness evaluation of selective classification [253].
- **Data spec:** text.
- **Sample size:** ~ 100 – 400K samples. Datasets have variable sizes spanning three orders of magnitude.
- **Year:** 2018.
- **Sensitive features:** none.
- **Link:** <https://gluebenchmark.com/>
- **Further info:** Wang et al. [525]

A.76 Goodreads Reviews

- **Description:** there are several versions of this dataset, corresponding to different crawls. Here we refer to the most well documented one by Wan and McAuley [523]. This resource consists of anonymized reviews collected from public user *book shelves*. Rich metadata is available for books and reviews, including authors, country code, publisher, userid, rating, timestamp, and text. A few medium-size subsamples focused on specific book genres are available. The task typically associated with this resource is book recommendation.
- **Affiliation of creators:** University of California, San Diego.
- **Domain:** literature, information systems.
- **Tasks in fairness literature:** fair ranking evaluation [421], fairness evaluation [93].
- **Data spec:** user-book pairs.
- **Sample size:** ~ 200M records from ~ 900K users over ~ 2M books.
- **Year:** 2019.
- **Sensitive features:** author.
- **Link:** <https://sites.google.com/eng.ucsd.edu/ucsdbookgraph/>
- **Further info:** Wan and McAuley [523]

A.77 Google Local

- **Description:** this dataset contains reviews and ratings from millions of users on local businesses from five different continents. Businesses are labelled with nearly 50 thousand categories. This resource was collected as a real world example of interactions between users and ratable items, with the goal of testing novel recommendation approaches. The dataset comprises data that is specific to users (e.g. places lived), businesses (e.g. GPS coordinates), and reviews (e.g. timestamps).
- **Affiliation of creators:** University of California, San Diego.
- **Domain:** information systems.

- **Tasks in fairness literature:** fair ranking [395].
- **Data spec:** user-business pairs.
- **Sample size:** ~ 10M reviews and ratings from ~ 5M users on ~ 3M local businesses.
- **Year:** 2018.
- **Sensitive features:** geography.
- **Link:** https://cseweb.ucsd.edu/~jmcauley/datasets.html#google_local
- **Further info:** He et al. [218]

A.78 Greek Websites

- **Description:** this dataset was created to demonstrate the *bias goggles* tools, which enables users to explore diverse bias aspects connected with popular Greek web domains. The dataset is a subset of the Greek web, crawled from Greek websites that cover politics and sports, represent big industries, or are generally popular. Starting from a seed of hundreds of websites, crawlers followed the links up to depth 7, avoiding popular sites such as Facebook and Twitter. The final dataset has a graph structure, comprising pages and links between them.
- **Affiliation of creators:** FORTH-ICS, University of Crete.
- **Domain:** .
- **Tasks in fairness literature:** bias discovery[287].
- **Data spec:** page-page pairs.
- **Sample size:** ~ 900k pages from ~ 90k domains.
- **Year:** 2020.
- **Sensitive features:** none.
- **Link:** <https://pangaia.ics.forth.gr/bias-goggles/about.html#Dataset>
- **Further info:** Konstantakis et al. [287]

A.79 Guardian Articles

- **Description:** this dataset consists of articles from *The Guardian*, retrieved from The Guardian Open Platform API. In particular, the authors crawled every article that appeared on the website between 2009 and 2018. They created this dataset to demonstrate a framework for the identification of gender biases in training data for machine learning.
- **Affiliation of creators:** University College Dublin.
- **Domain:** news.
- **Tasks in fairness literature:** data bias evaluation [304].
- **Data spec:** text.
- **Sample size:** unknown.
- **Year:** 2020.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** not available
- **Further info:** Leavy et al. [304]

A.80 HAM10000

- **Description:** the dataset comprises 10,015 dermatoscopic images collected over a period of 20 years the Department of Dermatology at the Medical University of Vienna, Austria and the skin cancer practice of Cliff Rosendahl in Queensland, Australia. Images were acquired and stored through different modalities; each image depicts a lesion and comes with metadata detailing the region of skin lesion, patient demographics, and diagnosis, which is the target variable. The dataset was employed for the lesion disease classification of the ISIC 2018 challenge.
- **Affiliation of creators:** Medical University of Vienna; University of Queensland.
- **Domain:** dermatology.
- **Tasks in fairness literature:** fair classification [343].
- **Data spec:** image.
- **Sample size:** ~10K images.
- **Year:** 2018.
- **Sensitive features:** age, sex.
- **Link:** <https://doi.org/10.7910/DVN/DBW86T>
- **Further info:** Tschandl et al. [499]

A.81 Harvey Rescue

- **Description:** this dataset is the result of crowdsourced efforts to connect rescue parties with people requesting help in the Houston area, mostly due to the flooding caused by Hurricane Harvey. Most requests are from August 28, 2017, and were sent via social media; they are timestamped and associated with the location of the people seeking help.
- **Affiliation of creators:** Harvey Relief Handiworks; Harvey Relief Coalition.
- **Domain:** social work.
- **Tasks in fairness literature:** fair spatio-temporal process learning [454].
- **Data spec:** tabular data.
- **Sample size:** ~1K help requests.
- **Year:** 2017.
- **Sensitive features:** geography.
- **Link:** not available
- **Further info:** <http://harveyrelief.handiworks.co/>

A.82 Heart Disease

- **Description:** this dataset is a collection of medical data from separate groups of patients referred for cardiac catheterisation and coronary angiography at 5 different medical centers, namely the Cleveland Clinic (data from 1981–1984), the Hungarian Institute of Cardiology in Budapest (1983–1987), the Long Beach Veterans Administration Medical Center (1984–1987) and the University Hospitals of Basel and Zurich (1985). The binary target variable in this dataset encodes a diagnosis of Coronary artery disease. Covariates relate to patient demographics, exercise data (e.g. maximum heart rate) and routine test data (e.g. resting blood pressure). Overall, 76 covariates are available but 14 are recommended. Names and social security numbers of the patients were initially available, but have been removed from the publicly available dataset.
- **Affiliation of creators:** Veterans Administration Medical Center, Long Beach; Hungarian Institute of Cardiology, Budapest; University Hospital, Zurich; University Hospital, Basel; Studer Corporation; Stanford University.
- **Domain:** cardiology.
- **Tasks in fairness literature:** fairness evaluation [405], fair active classification [383].
- **Data spec:** tabular data.
- **Sample size:** ~ 1K patients.
- **Year:** 1988.
- **Sensitive features:** age, sex.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/heart+disease>
- **Further info:** Detrano et al. [134]

A.83 Heritage Health

- **Description:** this dataset was developed as part of the Heritage Health Prize competition with the goal of reducing the cost of health care by decreasing the number of avoidable hospitalizations. The competition requires predicting the number of days a patient will spend in hospital during the 12 months following a cutoff date. The dataset features basic demographic information about patients, along with data about prior hospitalizations (e.g. length of stay and diagnosis), laboratory tests and prescriptions.
- **Affiliation of creators:** CHEO Research Institute, Inc; University of Ottawa; University of Maryland; Privacy Analytics, Inc; Kaggle; Heritage Provider Network.
- **Domain:** health policy.
- **Tasks in fairness literature:** fair multi-stage classification [338], fair representation learning [329], fair classification [418, 419], fair transfer learning [336], fairness evaluation [245].
- **Data spec:** tabular data.
- **Sample size:** ~ 150K patients.
- **Year:** 2011.
- **Sensitive features:** age, sex.
- **Link:** <https://www.kaggle.com/c/hhp/data>
- **Further info:** El Emam et al. [152]

A.84 High School Contact and Friendship Network

- **Description:** this dataset was developed to compare and contrast different methods commonly employed to measure human interaction and build the underlying social network. Data corresponds to interactions and friendship relations between students of a French high school in Marseilles. The authors consider four different methods of network data collection, namely face-to-face contacts measured by two concurrent methods (sensors and diaries), self-reported friendship surveys, and Facebook links.
- **Affiliation of creators:** Aix Marseille Université; Université de Toulon; Centre national de la recherche scientifique; ISI Foundation.

- **Domain:** social networks.
- **Tasks in fairness literature:** fair graph clustering [281].
- **Data spec:** student-student pairs.
- **Sample size:** ~ 300 students.
- **Year:** 2015.
- **Sensitive features:** gender.
- **Link:** <http://www.sociopatterns.org/datasets/high-school-contact-and-friendship-networks/>
- **Further info:** Mastrandrea et al. [345]

A.85 HMDA

- **Description:** The Home Mortgage Disclosure Act (HMDA) is a US federal law from 1975 mandating that financial institutions maintain and disclose information about mortgages to the public. Companies submit a Loan Application Register (LAR) to the Federal Financial Institutions Examination Council FFIEC who maintain and disclose the data. The LAR format is subject to changes, such as the one which happened in 2017. From 2018 onward, entries to the LAR comprise information about the financial institution (e.g. geography, id), the applicants (e.g. demographics, income), the house (e.g. value, construction method), the mortgage conditions (type, interest rate, amount) and the outcome. Ethnicity, race, and sex of applicants are self-reported.
- **Affiliation of creators:** Federal Financial Institutions Examination Council.
- **Domain:** finance.
- **Tasks in fairness literature:** fairness evaluation under unawareness [94, 256].
- **Data spec:** tabular data.
- **Sample size:** ~ 200M records.
- **Year:** present.
- **Sensitive features:** sex, geography, race, ethnicity.
- **Link:** <https://ffiec.cfbp.gov/data-browser/>
- **Further info:** <https://ffiec.cfbp.gov/>; <https://www.consumerfinance.gov/data-research/hmda/>

A.86 Homeless Youths' Social Networks

- **Description:** this dataset was collected to study methamphetamine use norms among homeless youth in association with their social networks. A sample of homeless youth aged 13–25 years was recruited between 2011–2012 from two drop-in centers in California. After obtaining informed consent/assent, participants filled in a survey and answered questions from an interview. The survey included questions on demographics, migratory status, educational status and housing. To reconstruct the social network between them, each participant provided information for up to 50 people with whom they had interacted during the previous 30 days.
- **Affiliation of creators:** University of Denver; University of Southern California.
- **Domain:** social work.
- **Tasks in fairness literature:** fair graph diffusion [420].
- **Data spec:** person-person pairs.
- **Sample size:** ~ 300 youth.
- **Year:** 2015.
- **Sensitive features:** age, gender, sexual orientation, race and ethnicity.
- **Link:** not available
- **Further info:** Barman-Adhikari et al. [32]

A.87 IBM HR Analytics

- **Description:** based on the information available on Kaggle, this is a fictional dataset created by IBM data scientists. It describes employees along dimensions that may be relevant for attrition, the target variable encoding employee departure. Available covariates include information on employee background (education, number of prior companies), work satisfaction (recent promotions, environment and job satisfaction) and seniority (years at the company, years in current role, job level).
- **Affiliation of creators:** IBM.
- **Domain:** information systems, management information systems.
- **Tasks in fairness literature:** fair data generation [321].
- **Data spec:** tabular data.
- **Sample size:** ~ 1K employees.
- **Year:** 2019.
- **Sensitive features:** gender.
- **Link:** <https://www.kaggle.com/datasets/pavansubhasht/ibm-hr-analytics-attrition-dataset>
- **Further info:** <https://github.com/IBM/employee-attrition-aif360>

A.88 IIT-JEE

- **Description:** this dataset was released in response to a Right to Information application filed in June 2009, and contains country-wide results for the Joint Entrance Exam (EET) to Indian Institutes of Technology (IITs), a group of prestigious engineering schools in India. The dataset contains the marks obtained by every candidate who took the test in 2009, divided according to the specific Math, Physics, and Chemistry sections of the test. Demographics such as ZIP code, gender, and birth categories (ethnic categories relating to the caste system) are also included.
- **Affiliation of creators:** Indian Institute of Technology, Kharagpur.
- **Domain:** education.
- **Tasks in fairness literature:** fair ranking [84].
- **Data spec:** tabular data.
- **Sample size:** ~ 400K students.
- **Year:** 2009.
- **Sensitive features:** gender, birth category.
- **Link:** not available
- **Further info:** Celis et al. [84]

A.89 IJB-A

- **Description:** the IARPA Janus Benchmark A (IJB-A) dataset was proposed as a face recognition benchmark with wide geographic representation and pose variation for subjects. It consists of *in-the-wild* images and videos of 500 subjects, obtained through internet searches over Creative Commons licensed content. The subjects were manually specified by the creators of the dataset to ensure broad geographic representation. The tasks associated with the dataset are face identification and verification. The dataset curators also collected the subjects' skin color and gender, through an unspecified annotation procedure. Similar protected attributes (gender and Fitzpatrick skin type) were labelled by one author of Buolamwini and Gebru [63].
- **Affiliation of creators:** Noblis; National Institute of Standards and Technology (NIST); Intelligence Advanced Research Projects Activity (IARPA); Michigan State University.
- **Domain:** computer vision.
- **Tasks in fairness literature:** data bias evaluation [63].
- **Data spec:** image.
- **Sample size:** ~ 6K images of ~ 500 subjects.
- **Year:** 2015.
- **Sensitive features:** gender, skin color.
- **Link:** <https://www.nist.gov/itl/iad/image-group/ijb-dataset-request-form>
- **Further info:** Klare et al. [279]

A.90 ILEA

- **Description:** this dataset was created by the Inner London Education Authority (ILEA) considering data from 140 British schools. It comprises the results of public examinations taken by students of age 16 over the period 1985–1987. These values are used as a measurement of school effectiveness, with emphasis on quality of education and equality of opportunity for students of different backgrounds and ethnicities. Student-level records report their sex and ethnicity, while school-level factors include the percentage of students eligible for free meals and the percentage of girls in each institute.
- **Affiliation of creators:** Inner London Education Authority (ILEA).
- **Domain:** education.
- **Tasks in fairness literature:** fair representation learning [389, 390].
- **Data spec:** unknown.
- **Sample size:** ~ 30K students from 140 secondary schools.
- **Year:** unknown.
- **Sensitive features:** age, sex, ethnicity.
- **Link:** not available
- **Further info:** [195, 385]

A.91 Image Embedding Association Test (iEAT)

- **Description:** the Image Embedding Association Test (iEAT) is a resource for quantifying biased associations between representations of social concepts and attributes in images. It mimics seminal work on biases in WEs [71], following the Implicit Association Test (IAT) from social psychology [202]. The curators identified several combinations of target concepts (e.g. young) and attributes (e.g. pleasant), testing similarities between representations of these concepts learnt by unsupervised computer vision models. For each

attribute/concept they obtained a set of images from the IAT, the CIFAR-100 dataset or Google Image Search, which act as the source of images and the associated sensitive attribute labels.

- **Affiliation of creators:** Carnegie Mellon University; George Washington University.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fairness evaluation of learnt representations [479].
- **Data spec:** image.
- **Sample size:** ~ 200 image for 15 iEATs.
- **Year:** 2021.
- **Sensitive features:** religion, gender, age, race, sexual orientation, disability, skin tone, weight.
- **Link:** <https://github.com/ryansteed/ieat/tree/master/data>
- **Further info:** Steed and Caliskan [479]

A.92 ImageNet

- **Description:** Imagenet is one of the most influential machine learning dataset of the 2010s. Much important work on computer vision, including early breakthroughs in deep learning has been sparked by ImageNet Large Scale Visual Recognition Challenge (ILSVRC), a competition held yearly from 2010 to 2017. The most used portion of ImageNet is indeed the data powering the classification task in ILSVRC 2012, featuring 1,000 classes, over 100 of which represent different dog breeds. Recently, several problematic biases were found in the person subtree of ImageNet, tracing their causes and proposing approaches to remove them [120, 406, 552].
- **Affiliation of creators:** Princeton University.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair classification [148], bias discovery [11], data bias evaluation [552], fair incremental learning [579], fairness evaluation [147].
- **Data spec:** image.
- **Sample size:** ~ 14M images depicting ~ 20K categories (synsets).
- **Year:** 2021.
- **Sensitive features:** people's gender and other sensitive annotations may be present in synsets from the person subtree.
- **Link:** <https://image-net.org/>
- **Further info:** Barocas et al. [33], Crawford and Paglen [120], Deng et al. [132], Prabhu and Birhane [406], Yang et al. [552]

A.93 In-Situ

- **Description:** this dataset was curated to measure biases in named entity recognition algorithms, based on gender, race and religion of people represented by entities. The authors exploit census data to build a list of 123 names typical of men and women of different race and religion. Next, they extract 289 sentences mentioning people from the CoNLL 2003 NER test data [494], itself derived from Reuters 1990s news stories. Finally, they substitute the unigram person entity from the CoNLL 2003 shared task with each of names obtained previously as specific to a demographic group.
- **Affiliation of creators:** Twitter.
- **Domain:** linguistics.
- **Tasks in fairness literature:** fairness evaluation in entity recognition [368].
- **Data spec:** text.
- **Sample size:** ~ 50K sentences.
- **Year:** 2020.
- **Sensitive features:** gender, race and religion.
- **Link:** https://github.com/napsternxg/NER_bias
- **Further info:** Mishra et al. [368]

A.94 iNaturalist Datasets

- **Description:** these datasets were curated as challenging real-world benchmarks for large-scale fine-grained visual classification and feature visually similar classes with large class imbalance. They consist of images of plants and animals from iNaturalist, a social network where nature enthusiasts share information and observations about biodiversity. There are four different releases of the dataset: 2017, 2018, 2019, and 2021. A subset of the images are also annotated with bounding boxes and have additional metadata such as where and when the images were captured.
- **Affiliation of creators:** California Institute of Technology; University of Edinburgh; Google; Cornell University; iNaturalist.
- **Domain:** biology.
- **Tasks in fairness literature:** fairness evaluation of private classification [22].
- **Data spec:** image.
- **Sample size:** ~ 3M images from ~ 10K different species of plants and animals.

- **Year:** 2021.
- **Sensitive features:** none.
- **Link:** https://github.com/visipedia/inat_comp
- **Further info:** [513, 514]

A.95 Indian Census

- **Description:** very little information seems to be available on this dataset. It represents a count of residents of 35 Indian states, repeated every ten years between 1951 and 2001.
- **Affiliation of creators:** Office of the Registrar General of India.
- **Domain:** demography.
- **Tasks in fairness literature:** fairness evaluation of private resource allocation [410].
- **Data spec:** tabular data.
- **Sample size:** ~ 30 state.
- **Year:** unknown.
- **Sensitive features:** geography.
- **Link:** https://www.indiabudget.gov.in/budget_archive/es2006-07/chapt2007/tab97.pdf
- **Further info:**

A.96 Indian Student Performance

- **Description:** this dataset was curated to support educational data mining algorithms. The creators collected data from three colleges of Assam, India (Duliajan College, Doomdooma College, and Digboi College). Each data point represents a student, summarizing information on their demographics (gender, caste), family (occupation and qualification of parents), and school fruition (study hours, attendance, home-to-school travel). Among the latter there are four variables summarizing student performance in different classes and examinations, which represent the response variable of a prediction task.
- **Affiliation of creators:** Dibrugarh University; Sana'a University; Abdelmalek Essaâdi University.
- **Domain:** education.
- **Tasks in fairness literature:** fair data summarization [36].
- **Data spec:** tabular data.
- **Sample size:** ~ 300 students.
- **Year:** 2018.
- **Sensitive features:** gender, caste, geography.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/Student+Academics+Performance>
- **Further info:** Hussain et al. [240]

A.97 Infant Health and Development Program (IHDP)

- **Description:** this dataset is the result of the IHDP program carried out between 1985 and 1988 in the US. A longitudinal randomized trial was conducted to evaluate the effectiveness of comprehensive early intervention in reducing developmental and health problems in low birth weight premature infants. Families in the experimental group received an intervention based on an educational program delivered through home visits, a daily center-based program and a parent supporting group. Children in the study were assessed across multiple cognitive, behavioral, and health dimensions longitudinally in four phases at ages 3, 5, 8, and 18. The dataset also contains information on household composition, source of health care, parents' demographics and employment.
- **Affiliation of creators:** unknown.
- **Domain:** pediatrics.
- **Tasks in fairness literature:** fair risk assessment [337, 558].
- **Data spec:** mixture.
- **Sample size:** ~ 1K infants.
- **Year:** 1993.
- **Sensitive features:** race and ethnicity (of parents), age (maternal), gender (of infant).
- **Link:** <https://www.icpsr.umich.edu/web/HMCA/studies/9795>
- **Further info:** Brooks-Gunn et al. [58]

A.98 Instagram Photos

- **Description:** this dataset was crawled from Instagram to explore trade-offs between fairness and revenue in platforms that serve ads to their users. The authors crawled metadata from photos (location and tags) and users (names), using Kevin Systrom as a seed user and cascading into profiles that like or comment photos. The curators concentrated on cities with enough geotagged data, namely New York and Los Angeles. Moreover, they labeled the users with gender and race. Gender was labeled via US social security data,

using the proportion of babies with a given name registered with either gender. Gender was only assigned to users with a first name for which there were both at least 50 births and 95% of recorded births were one gender. Race were labeled using the Face++ API on a subset of photos. Photos were not downloaded, rather they were fed to Face++ via their publicly available URL. Finally, the ground truth labels were validated by two research assistants. To emulate a location-based advertisement model, the creators devised a task aimed at predicting what topics a user will be interested in, given their locations from previous check-ins.

- **Affiliation of creators:** Columbia University.
- **Domain:** social media.
- **Tasks in fairness literature:** fair advertising [430].
- **Data spec:** unknown.
- **Sample size:** ~ 1M photos from ~ 40K users.
- **Year:** 2017.
- **Sensitive features:** race, gender, geography.
- **Link:** not available
- **Further info:** Riederer and Chaintreau [430]

A.99 Internet Ads

- **Description:** this dataset was assembled to study the problem of automated advertisement removal in browsers. It consists of images crawled from randomly generated urls, manually classified as ad/no-ad. Image encodings are derived from raw html, thus containing no information about pixel values, but rather encoding width, height, anchor text and image source. The associated task is classifying each image encoding as an ad or a no-ad image.
- **Affiliation of creators:** University College Dublin.
- **Domain:** pattern recognition.
- **Tasks in fairness literature:** fair anomaly detection [459].
- **Data spec:** tabular data.
- **Sample size:** ~ 3K image encodings.
- **Year:** 1998.
- **Sensitive features:** none.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/internet+advertisements>
- **Further info:** Kushmerick [294]

A.100 Iris

- **Description:** the most popular dataset on the UCI Machine Learning Repository was created by E. Anderson and popularized by R.A. Fisher in the pattern recognition community in the 1930s. The measurements in this collection represent the length and width of sepal and petals of different Iris flowers, collected to evaluate the morphological variation of different Iris species. The typical learning task associated with this dataset is labelling the species based on the available measurements.
- **Affiliation of creators:** Missouri Botanical Garden; Washington University.
- **Domain:** plant science.
- **Tasks in fairness literature:** fair clustering [1, 95].
- **Data spec:** tabular data.
- **Sample size:** ~ 100 samples from three species of Iris.
- **Year:** 1988.
- **Sensitive features:** none.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/iris>
- **Further info:** [12, 169]

A.101 Italian Car Insurance

- **Description:** this resource was curated to study discriminatory practices in the Italian car insurance market. More specifically, the data was collected to estimate the direct effect of gender and birthplace on yearly quoted premiums. It was collected in 2020 from a popular Italian car insurance comparison website, where the curators tried different hypothetical driver profiles and collected the quotes provided by nine companies. Along with gender and birthplace, additional driver features include age, city of residence, insured vehicle, mileage, and a summary of claim history.
- **Affiliation of creators:** University of Padua; Carnegie Mellon University; University of Udine.
- **Domain:** economics.
- **Tasks in fairness literature:** fair pricing evaluation [159].
- **Data spec:** tabular data.
- **Sample size:** ~ 2K driver profiles.

- **Year:** 2021.
- **Sensitive features:** gender, birthplace.
- **Link:** not available
- **Further info:** Fabris et al. [159]

A.102 KDD Cup 99

- **Description:** this dataset was developed for a data mining competition on cybersecurity, focused on building an automated network intrusion detector based on TCP dump data. The task is predicting whether a connection is legitimate and inoffensive or symptomatic of an attack, such as denial-of-service or user-to-root; tens of attack classes have been simulated and annotated within this dataset. The available features include basic TCP/IP information, network traffic and contextual features, such as number of failed login attempts.
- **Affiliation of creators:** Massachusetts Institute of Technology.
- **Domain:** computer networks.
- **Tasks in fairness literature:** fair clustering [95].
- **Data spec:** tabular data.
- **Sample size:** ~ 7M connections.
- **Year:** 1999.
- **Sensitive features:** none.
- **Link:** <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>
- **Further info:** Tavallae et al. [492]

A.103 Kidney Exchange Program

- **Description:** this dataset is based on data of the Canadian Kidney Paired Donation Program (KPD) to study strategic behavior among entities controlling part of the incompatible patient-donor pairs. Based on data from the Canadian Blood Services on the KPD and census, these instances were generated. The random instance generator is available upon request. The instances are weighted graphs. The incompatible patient-donor pairs represent the vertices of the graph, an arc means that the donor of a vertex is compatible with the patient of another vertex, and weights represent the benefit of the donation. Compatibility is encoded based on true blood type distribution and risk of transplant rejection.
- **Affiliation of creators:** Université de Montréal; Polytechnique de Montréal.
- **Domain:** public health.
- **Tasks in fairness literature:** fair matching evaluation [161].
- **Data spec:** patient-donor pairs.
- **Sample size:** 180.
- **Year:** 2020.
- **Sensitive features:** blood type, geography.
- **Link:** <https://github.com/mxmmargarida/KEG>
- **Further info:** Carvalho and Lodi [76]

A.104 Kidney Matching

- **Description:** this dataset was created via a simulator based on real data provided by the Organ and Tissue Authority of Australia. The data was validated against additional information from the Australian Bureau of Statistics, the Public and Research sets, and Wikipedia. The simulator models the probability distribution over the Blood Type and State of donors and patients, along with the quality of a donated organ (summarized by Kidney Donor Patient Index) and of a patient (quantified by the Expected Post-Transplant Survival). The envisioned task for this data is optimal matching of organs and patients.
- **Affiliation of creators:** unknown.
- **Domain:** public health.
- **Tasks in fairness literature:** fairness matching evaluation [347].
- **Data spec:** tabular data.
- **Sample size:** unknown.
- **Year:** 2018.
- **Sensitive features:** age, geography, blood type.
- **Link:** not available
- **Further info:** Mattei et al. [346]

A.105 Kiva

- **Description:** this dataset was obtained from kiva.org, a non-profit organization allowing low-income entrepreneurs and students to borrow money through loan crowdfunding. The data summarizes all transactions occurred in 2017. Transactions are typically between 25\$ to 50\$ and range from 5\$ to 10,000\$. Features include information about the loan, such as its purpose, sector and amount, and data specific to the borrower and their demographics. Women are prevalent in this dataset, probably due to the priorities of partner organizations and the easier access to capital enjoyed by men in many countries.
- **Affiliation of creators:** Kiva; DePaul University.
- **Domain:** finance.
- **Tasks in fairness literature:** fair ranking [65, 325, 473], bias discovery [472].
- **Data spec:** tabular data.
- **Sample size:** ~ 1M transactions involving ~ 100K loans and ~ 200K users.
- **Year:** 2018.
- **Sensitive features:** gender, geography, activity.
- **Link:** not available
- **Further info:** Sonboli and Burke [472]

A.106 Labeled Faces in the Wild (LFW)

- **Description:** LFW is a public benchmark for face verification, maintained by researchers affiliated with the University of Massachusetts. It was built to measure the progress of face verification systems in unconstrained settings (e.g. variable pose, illumination, resolution). The dataset consists of images of people who appeared in the news, labelled with the name of the respective individual. According to perception of human coders who were later asked to annotate this dataset, images mostly skew white, male and below 60.
- **Affiliation of creators:** University of Massachusetts, Amherst; Stony Brook University.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair data summarization [446], fair clustering [186], robust fairness evaluation [45], fairness evaluation [451].
- **Data spec:** image.
- **Sample size:** ~ 13K face images of ~ 6K individuals.
- **Year:** 2007.
- **Sensitive features:** gender, age, race.
- **Link:** <http://vis-www.cs.umass.edu/lfw/>
- **Further info:** Gebru et al. [182], Han and Jain [211], Huang et al. [235]

A.107 Large Movie Review

- **Description:** a set of reviews from IMDB, collected, filtered and preprocessed by researchers affiliated with Stanford University. Polarity judgements are balanced in terms of positive and negative reviews and automatically inferred from star-based ratings, so that 7 or more is positive, while 4 or less is considered negative. The dataset was collected to provide a large benchmark for sentiment analysis algorithms.
- **Affiliation of creators:** Stanford University.
- **Domain:** linguistics.
- **Tasks in fairness literature:** fair sentiment analysis evaluation [318].
- **Data spec:** text.
- **Sample size:** ~ 50K reviews.
- **Year:** 2011.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <https://ai.stanford.edu/~amaas/data/sentiment/>
- **Further info:** Maas et al. [334]

A.108 Last.fm

- **Description:** the Last.fm datasets were collected via the Last.fm API with the purpose of studying music consumption, discovery and recommendation on the web. Two datasets are provided: LFM1K, comprising timestamped listening habits of a limited user sample (~1K) at song granularity, and LFM360K, containing the top 50 most played artists of a wider user population (~360K).
- **Affiliation of creators:** Barcelona Music and Audio Technologies; Universitat Pompeu Fabra.
- **Domain:** music, information systems.
- **Tasks in fairness literature:** fair ranking evaluation [151].
- **Data spec:** user-song pairs (LFM1K); user-artist pairs (LFM360K).

- **Sample size:** ~19M timestamped records of ~1K users playing songs from ~170K artists (LFM1K); ~ 20M play counts (user-artist pairs) for ~400K users over ~300K artists (LFM360K).
- **Year:** 2010.
- **Sensitive features:** user age, gender, geography; artist.
- **Link:** <http://ocelma.net/MusicRecommendationDataset/>
- **Further info:** Celma [85]

A.109 Latin Newspapers

- **Description:** this dataset was built to study gender bias in language models and their connection with the corpora they have been trained on. It was built crawling articles from the websites of three newspapers from Chile, Peru, and Mexico. More detailed information about this resource seems to be missing.
- **Affiliation of creators:** Capital One.
- **Domain:** news.
- **Tasks in fairness literature:** data bias evaluation [172].
- **Data spec:** text.
- **Sample size:** ~ 60K articles.
- **Year:** 2019.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** not available
- **Further info:** Florez [172]

A.110 Law School

- **Description:** This dataset was collected to study performance in law school and bar examination of minority examinees in connection with affirmative action programs established after 1967 and subsequent anecdotal reports suggesting low bar passage rates for black examinees. Students, law schools, and state boards of bar examiners contributed to this dataset. The study tracks students who entered law school in fall 1991 through three or more years of law school and up to five administrations of the bar examination. Variables include demographics of candidates (e.g. age, race, sex), their academic performance (undergraduate GPA, law school admission test, and GPA), personal condition (e.g. financial responsibility for others during law school) along with information about law schools and bar exams (e.g. geographical area where it was taken). The associated task in machine learning is prediction of passage of the bar exam.
- **Affiliation of creators:** Law School Admission Council (LSAC).
- **Domain:** education.
- **Tasks in fairness literature:** fair classification [3, 41, 101, 442, 551], rich-subgroup fairness evaluation [270], fair classification under unawareness [297, 299], fairness evaluation [46, 295], fair regression [4, 110, 111, 286], fair representation learning [441], robust fair classification [341], limited-label fair classification [528].
- **Data spec:** tabular data.
- **Sample size:** ~ 20K examinees.
- **Year:** 1998.
- **Sensitive features:** sex, race, age.
- **Link:** not available
- **Further info:** Wightman et al. [538]

A.111 Libimseti

- **Description:** this dataset was collected to explore the effectiveness of recommendations in online dating services based on collaborative filtering. It was collected in collaboration with employees of the dating platform libimseti.cz, one of the largest Czech dating websites at the time. The data consists of anonymous ratings provided by (and to) users of the web service on a 10-point scale.
- **Affiliation of creators:** Charles University in Prague; Libimseti.
- **Domain:** sociology, information systems.
- **Tasks in fairness literature:** fair matching [500].
- **Data spec:** user-user pairs.
- **Sample size:** ~10M ratings over ~200K users.
- **Year:** 2007.
- **Sensitive features:** gender.
- **Link:** <http://colfi.wz.cz/>
- **Further info:** Brozovsky and Petricek [59], Brožovský [60]

A.112 Los Angeles City Attorney's Office Records

- **Description:** this dataset was extracted from the Los Angeles City Attorney's case management system. It consists of a collection of records aimed at powering data-driven approaches to decision making and resource allocation for misdemeanor recidivism reduction via individually tailored social service interventions. Focusing on cases handled by the office between 1995–2017, the data includes information about jail bookings, charges, court appearances, outcomes, and demographics.
- **Affiliation of creators:** Los Angeles City Attorney's Office; University of Chicago.
- **Domain:** law.
- **Tasks in fairness literature:** fair classification [432].
- **Data spec:** tabular data.
- **Sample size:** ~ 1M unique individuals associated with ~ 2M cases.
- **Year:** 2020.
- **Sensitive features:** race, ethnicity.
- **Link:** not available
- **Further info:** [432]

A.113 MEPS-HC

- **Description:** the Medical Expenditure Panel Survey (MEPS) data is collected by the US Department of Health and Human Services, to survey healthcare spending and utilization by US citizens. Overall, this is a set of large-scale surveys of families and individuals, their employers, and medical providers (e.g. doctors, hospitals, pharmacies). The Household Component (HC) focuses on households and individuals, who provide information about their demographics, medical conditions and expenses, health insurance coverage, and access to care. Individuals included in a panel undergo five rounds of interviews over two years. Healthcare expenditure is often regarded as a target variable in machine learning applications, where it has been used as a proxy for healthcare utilization, with the goal of identifying patients in need.
- **Affiliation of creators:** Agency for Healthcare Research and Quality.
- **Domain:** health policy.
- **Tasks in fairness literature:** fair transfer learning [117], fair regression [435], fairness evaluation [465], robust fair classification [44], fair classification [456].
- **Data spec:** tabular data.
- **Sample size:** ~ 30K, variable on a yearly basis.
- **Year:** present.
- **Sensitive features:** gender, ethnicity, age.
- **Link:** https://meps.ahrq.gov/mepsweb/data_stats/download_data_files.jsp
- **Further info:** <https://www.ahrq.gov/data/meps.html>

A.114 MGGG States

- **Description:** developed by the Metric Geometry and Gerrymandering Group¹¹, this dataset contains precinct-level aggregated information about demographics and political leaning of voters in each district. The data hinges on several distinct sources of data, including GIS mapping files from the US Census Bureau¹², demographic data from IPUMS¹³ and election data from MIT Election and Data Science¹⁴. Source and precise data format vary by state.
- **Affiliation of creators:** Tufts University.
- **Domain:** political science.
- **Tasks in fairness literature:** fair districting for electoral precincts [450].
- **Data spec:** mixture.
- **Sample size:** variable number of precincts (thousands) per state.
- **Year:** 2021.
- **Sensitive features:** race, political affiliation (representation in different precincts).
- **Link:** <https://github.com/mggg-states>
- **Further info:** <https://mggg.org/>

¹¹<https://mggg.org/>

¹²<https://www.census.gov/geographies/mapping-files.html>

¹³<https://www.nhgis.org/>

¹⁴<https://electionlab.mit.edu/>

A.115 Microsoft Learning to Rank

- **Description:** this dataset was released to spur advances in learning to rank algorithms, capable of producing a list of documents in response to a text query, ranked according to their relevance for the query. The dataset contains relevance judgements for query-document pairs, obtained “from a retired labeling set” of the Bing search engine. Over 100 numerical features are provided for each query-document pair, summarizing the salient lexical properties of the pair and the quality of the webpage, including its page rank.
- **Affiliation of creators:** Microsoft.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair ranking [54].
- **Data spec:** query document pairs.
- **Sample size:** ~ 30K queries.
- **Year:** 2013.
- **Sensitive features:** none.
- **Link:** <https://www.microsoft.com/en-us/research/project/mslr/>
- **Further info:** [412]

A.116 Million Playlist Dataset (MPD)

- **Description:** this dataset powered the 2018 RecSys Challenge on automatic playlist continuation. It consists of a sample of public Spotify playlists created by US Spotify users between 2010–2017. Each playlist consists of a title, track list and additional metadata. For each track, MPD provides the title, artist, album, duration and Spotify pointers. User data is anonymized. The dataset was augmented with record label information crawled from the web [282].
- **Affiliation of creators:** Spotify; Johannes Kepler University; University of Massachusetts.
- **Domain:** music, information systems.
- **Tasks in fairness literature:** data bias evaluation [282].
- **Data spec:** tabular data.
- **Sample size:** ~ 1M playlists containing ~ 2M unique tracks by ~ 300K artists.
- **Year:** 2018.
- **Sensitive features:** artist, record label.
- **Link:** <https://www.aicrowd.com/challenges/spotify-million-playlist-dataset-challenge>
- **Further info:** Chen et al. [92]

A.117 Million Song Dataset (MSD)

- **Description:** this dataset was created as a large-scale benchmark for algorithms in the musical domain. Song data was acquired through The Echo Nest API, capturing a wide array of information about the song (duration, loudness, key, tempo, etc.) and the artist (name, id, location, etc.). In total the dataset creators retrieved one million songs, and for each song 55 fields are provided as metadata. This dataset also powers the Million Song Dataset Challenge, integrating the MSD with implicit feedback from taste profiles gathered from an undisclosed set of applications.
- **Affiliation of creators:** Columbia University; The Echo Nest.
- **Domain:** music, information systems.
- **Tasks in fairness literature:** dynamical evaluation of fair ranking [166].
- **Data spec:** user-song pairs.
- **Sample size:** ~ 50M play counts over ~ 1M users and ~ 400K songs.
- **Year:** 2012.
- **Sensitive features:** artist; geography.
- **Link:** <http://millionsongdataset.com/>; <https://www.kaggle.com/c/msdchallenge>
- **Further info:** Bertin-Mahieux et al. [42], McFee et al. [352]

A.118 MIMIC-CXR-JPG

- **Description:** this dataset was curated to encourage research in medical computer vision. It consists of chest x-rays sourced from the Beth Israel Deaconess Medical Center between 2011–2016. Each image is tagged with one or more of fourteen labels, derived from the corresponding free-text radiology reports via natural language processing tools. A subset of 687 report-label pairs have been validated by a board of certified radiologists with 8 years of experience.
- **Affiliation of creators:** Massachusetts Institute of Technology; Beth Israel Deaconess Medical Center; Stanford University; Harvard Medical School; National Library of Medicine.
- **Domain:** radiology.
- **Tasks in fairness literature:** fairness evaluation of private classification [98].
- **Data spec:** images.

- **Sample size:** ~ 400K images of ~ 70K patients.
- **Year:** 2019.
- **Sensitive features:** sex.
- **Link:** <https://physionet.org/content/mimic-cxr-jpg/2.0.0/>
- **Further info:** [251]

A.119 MIMIC-III

- **Description:** this dataset was extracted from a database of patients admitted to critical care units at the Beth Israel Deaconess Medical Center in Boston (MA), following the widespread adoption of digital health records in US hospitals. Data comprises vital signs, medications, laboratory measurements, notes and observations by care providers, fluid balance, procedure codes, diagnostic codes, imaging reports, length of stay, survival data, and demographics. The dataset spans over a decade of intensive care unit stays for adult and neonatal patients.
- **Affiliation of creators:** Massachusetts Institute of Technology; Beth Israel Deaconess Medical Center; A*STAR.
- **Domain:** critical care medicine.
- **Tasks in fairness literature:** fair classification [343], fairness evaluation [93, 568], robust fair classification [464].
- **Data spec:** mixture.
- **Sample size:** ~ 60K patients.
- **Year:** 2016.
- **Sensitive features:** age, ethnicity, gender.
- **Link:** <https://mimic.mit.edu/>
- **Further info:** Johnson et al. [252]

A.120 ML Fairness Gym

- **Description:** this resource was developed to study the long-term behaviour and emergent properties of fair ML systems. It is an extension of OpenAI Gym [57], simulating the actions of agents within environments as Markov Decision Processes. As of 2021, four environments have been released. (1) *Lending* emulates the decisions of a bank, based on perceived credit-worthiness of individuals, which is distributed according to an artificial sensitive feature. (2) *Attention allocation* concentrates on agents tasked with monitoring sites for incidents. (3) *College admission* relates to sequential game theory, where agents represent colleges and environments contain students capable of strategically manipulating their features at different costs, for instance through preparation courses. (4) *Infectious disease* models the problem of vaccine allocation and its long-term consequences on people in different demographic groups.
- **Affiliation of creators:** Google.
- **Domain:** N/A.
- **Tasks in fairness literature:** dynamical fair resource allocation [17, 124], dynamical fair classification [124].
- **Data spec:** time series.
- **Sample size:** variable.
- **Year:** 2020.
- **Sensitive features:** synthetic.
- **Link:** <https://github.com/google/ml-fairness-gym>
- **Further info:** D'Amour et al. [124]

A.121 MNIST

- **Description:** one of the most famous resources in computer vision, this dataset was created from an earlier database released by the National Institute of Standards and Technology (NIST). It consists of hand-written digits collected among high-school students and Census Bureau employees, which have to be correctly labelled by image processing systems. Several augmentations have also been used in the fairness literature, discussed at the end of this section.
- **Affiliation of creators:** AT&T Labs.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair clustering [213, 313], fair anomaly detection [567], fair classification [121], fairness evaluation [451].
- **Data spec:** image.
- **Sample size:** ~ 70K images across 10 digits.
- **Year:** 1998.
- **Sensitive features:** none.
- **Link:** <http://yann.lecun.com/exdb/mnist/>
- **Further info:** Barocas et al. [33], Lecun et al. [305]
- **Variants:**

- MNIST-USPS [313]: merge with USPS dataset of handwritten digits [239].
- Color-reverse MNIST [313] or MNIST-Invert [567]: images from MNIST, reversed via $p = 255 - p$ for each pixel p .
- Color MNIST [16]: images from MNIST colored red or green based on class label.
- C-MNIST: images from MNIST, such that both digits and background are colored.

A.122 Mobile Money Loans

- **Description:** this dataset captures the ongoing collaboration between some banks and mobile network operators in East Africa. Phone data, including mobile money transactions, is used as “soft” financial data to create a credit score. Mobile money (bank-less) transactions represent a low-barrier tool for the financial inclusion of the poor and are fairly popular in some African countries.
- **Affiliation of creators:** unknown.
- **Domain:** finance.
- **Tasks in fairness literature:** fair transfer learning [117].
- **Data spec:** tabular data.
- **Sample size:** ~ 200K people.
- **Year:** unknown.
- **Sensitive features:** age, gender.
- **Link:** not available
- **Further info:** Speakman et al. [474]

A.123 MovieLens

- **Description:** first released in 1998, MovieLens datasets represent user ratings from the movie recommender platform run by the GroupLens research group from the University of Minnesota. While different datasets have been released by GroupLens, in this section we concentrate on MovieLens 1M, the one predominantly used in fairness research. User-system interactions take the form of a quadruple (UserID, MovieID, Rating, Timestamp), with ratings expressed on a 1-5 star scale. The dataset also reports user demographics such as age and gender, which is voluntarily provided by the users.
- **Affiliation of creators:** University of Minnesota.
- **Domain:** information systems, movies.
- **Tasks in fairness literature:** fair ranking [65, 139, 162, 325, 473], fair ranking evaluation [151, 555, 556], fair data summarization [153], fair representation learning [389, 390], fair graph mining [53, 66], fair data generation [64].
- **Data spec:** user-movie pairs.
- **Sample size:** ~ 1M reviews by ~ 6K users over ~ 4K movies.
- **Year:** 2003.
- **Sensitive features:** gender, age.
- **Link:** <https://grouplens.org/datasets/movielens/1m/>
- **Further info:** Harper and Konstan [216]

A.124 MS-Celeb-1M

- **Description:** this dataset was created as a large scale public benchmark for face recognition. The creators cover a wide range of countries and emphasizes diversity echoing outdated notions of race: “We cover all the major races in the world (Caucasian, Mongoloid, and Negroid)” [209]. While (in theory) containing only images of celebrities, the dataset was found to feature people who simply must maintain an online presence, and was retracted for this reason. Despite termination of the hosting website, the dataset is still searched for, available and used to build new fairness datasets, such as RFW (§ A.153) and BUPT Faces (§ A.27). The dataset was recently augmented with gender and nationality data automatically inferred from biographies of people [351]. From nationality, a race-related attribute was also annotated on a subset of 20,000 images.
- **Affiliation of creators:** Microsoft.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fairness evaluation through artificial data generation [351].
- **Data spec:** image.
- **Sample size:** ~ 10M images representing ~ 100K people.
- **Year:** 2016.
- **Sensitive features:** gender, race, geography.
- **Link:** not available
- **Further info:** Guo et al. [209], McDuff et al. [351], Murgia [375]

A.125 MS-COCO

- **Description:** this dataset was created with the goal of improving the state of the art in object recognition. The dataset consists of over 300,000 labeled images collected from Flickr. Each image was annotated based on whether it contains one or more of the 91 object types proposed by the authors. Segmentations are also provided to indicate the region where objects are located in each image. Finally, five human-generated captions are provided for each image. Annotation, segmentation and captioning were performed by human annotators hired on Amazon Mechanical Turk. A subset of the images depicting people have been augmented with gender labels “man” and “woman” based on whether captions mention one word but not the other [225, 583].
- **Affiliation of creators:** Cornell University; Toyota Technological Institute; Facebook; Microsoft; Brown University; California Institute of Technology; University of California at Irvine.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair representation learning [127], fair classification [225].
- **Data spec:** image.
- **Sample size:** ~ 300K images.
- **Year:** 2014.
- **Sensitive features:** gender.
- **Link:** <https://cocodataset.org/>
- **Further info:** Lin et al. [319]

A.126 Multi-task Facial Landmark (MTFL)

- **Description:** this dataset was developed to evaluate the effectiveness of multi-task learning in problems of facial landmark detection. The dataset builds upon an existing collection of outdoor face images sourced from the web already labelled with bounding boxes and landmarks [559], by annotating whether subjects are smiling or wearing glasses, along with their gender and pose. These annotations, whose provenance is not documented, allow researchers to define additional classification tasks for their multi-task learning pipeline.
- **Affiliation of creators:** The Chinese University of Hong Kong.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair clustering [313].
- **Data spec:** image.
- **Sample size:** ~ 10K images.
- **Year:** 2014.
- **Sensitive features:** gender.
- **Link:** <http://mmlab.ie.cuhk.edu.hk/projects/TCDCN.html>
- **Further info:** Zhang et al. [575, 576]

A.127 National Longitudinal Survey of Youth

- **Description:** the National Longitudinal Surveys from the US Bureau of Labor Statistics follow the lives of representative samples of US citizens, focusing on their labor market activities and other significant life events. Subjects periodically provide responses to questions about their education, employment, housing, income, health, and more. Two different cohorts were started in 1979 (NLSY79) and (NLSY97), which have been associated with machine learning tasks of income prediction and GPA prediction respectively.
- **Affiliation of creators:** US Bureau of Labor Statistics.
- **Domain:** demography.
- **Tasks in fairness literature:** fair regression [110, 111, 286].
- **Data spec:** tabular data.
- **Sample size:** ~ 10K respondents (NLSY79); ~ 9K respondents (NLSY97).
- **Year:** present.
- **Sensitive features:** age, race, sex.
- **Link:** <https://www.bls.gov/nls/nlsy79.htm> (NLSY79); <https://www.bls.gov/nls/nlsy97.htm> (NLSY97)
- **Further info:**

A.128 National Lung Screening Trial (NLST)

- **Description:** the NLST was a randomized controlled trial aimed at understanding whether imaging through low-dose helical computed tomography reduces lung cancer mortality relative to chest radiography. Participants were recruited at 33 screening centers across the US, among subjects deemed at risk of lung cancer based on age and smoking history, and were made aware of the trial. A breadth of features about participants is available, including demographics, disease history, smoking history, family history of lung cancer, type, and results of screening exams.
- **Affiliation of creators:** National Cancer Institute’s Division of Cancer Prevention, Division of Cancer Treatment and Diagnosis.
- **Domain:** radiology.

- **Tasks in fairness literature:** fair preference-based classification [507].
- **Data spec:** image.
- **Sample size:** ~ 50K participants.
- **Year:** 2020.
- **Sensitive features:** age, ethnicity, race, sex.
- **Link:** <https://cdas.cancer.gov/nlst/>
- **Further info:** NLST Trial Research Team [382]; <https://www.cancer.gov/types/lung/research/nlst>

A.129 New York Times Annotated Corpus

- **Description:** this corpus contains nearly two million articles published in The New York Times over the period 1987–2007. For some articles, annotations by library scientists are available, including topics, mentioned entities, and summaries. The data is provided in News Industry Text Format (NITF).
- **Affiliation of creators:** The New York Times.
- **Domain:** news.
- **Tasks in fairness literature:** bias evaluation in WEs [62].
- **Data spec:** text.
- **Sample size:** ~ 2M articles.
- **Year:** 2008.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <https://catalog ldc.upenn.edu/LDC2008T19>
- **Further info:**

A.130 Nominees Corpus

- **Description:** this corpus was curated to study gender-related differences in literary production, with attention to perception of quality. It consists of fifty Dutch-language fiction novels nominated for either the AKO Literatuurprijs(shortlist) or the Libris Literatuur Prijs (longlist) in the period 2007–2012. The corpus was curated to control for nominee gender and country of origin. Word counts, LIWC counts, and metadata for this dataset are available at <http://dx.doi.org/10.17632/tmp32v54ss.2>.
- **Affiliation of creators:** University of Amsterdam.
- **Domain:** literature.
- **Tasks in fairness literature:** fairness evaluation [289].
- **Data spec:** text.
- **Sample size:** ~ 50 novels.
- **Year:** 2017.
- **Sensitive features:** gender, geography (of author).
- **Link:** not available
- **Further info:** Koolen [288], Koolen and van Cranenburgh [289]

A.131 North Carolina Voters

- **Description:** US voter data is collected, curated, and maintained for multiple reasons. Data about voters in North Carolina is collected publicly as part of voter registration requirements and also privately. Private companies curating these datasets sell voter data as part of products, which include outreach lists and analytics. These datasets include voters' full names, address, demographics, and party affiliation.
- **Affiliation of creators:** North Carolina State Board of Elections.
- **Domain:** political science.
- **Tasks in fairness literature:** data bias evaluation [115], fair clustering [1], fairness evaluation of advertisement [475].
- **Data spec:** tabular data.
- **Sample size:** ~ 8M voters.
- **Year:** present.
- **Sensitive features:** race, ethnicity, age, geography.
- **Link:** <https://www.ncsbe.gov/results-data/voter-registration-data>
- **Further info:**
- **Variants:** a privately curated version of this dataset is maintained by L2.¹⁵

¹⁵<https://l2-data.com/states/north-carolina/>

A.132 Nursery

- **Description:** this dataset encodes applications for a nursery school in Ljubljana, Slovenia. To favour transparent and objective decision-making, a computer-based decision support system was developed for the selection and ranking of applications. The target variable reported is thus the output of an expert systems based on a set of rules, taking as an input information about the family, including housing, occupation and financial status, included in the dataset. The variables were reportedly constructed in a careful manner, taking into account laws that were in force at that time and following advice given by leading experts in that field. However, the variables also appear to be coded rather subjectively. For example, the variable *social condition* admits as a value *Slightly problematic*, allegedly reserved for “When education ability of parents is low (unequal, inconsistent education, exaggerated pretentiousness or indulgence, neurotic reactions of parents), or there are improper relations in family (easier forms of parental personality disturbances, privileged or ignored children, conflicts in the family)”. Given that the true map between inputs and outputs is known, this resource is mostly useful to evaluate methods of structure discovery.
- **Affiliation of creators:** University of Maribor; Jožef Stefan Institute; University of Ljubljana; Center for Public Enterprises in Developing Countries.
- **Domain:** education.
- **Tasks in fairness literature:** fair classification [435].
- **Data spec:** tabular data.
- **Sample size:** ~ 10K combinations of input data (hypothetical applicants).
- **Year:** 1997.
- **Sensitive features:** family wealth.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/nursery>
- **Further info:** Olave et al. [387]

A.133 NYC Taxi Trips

- **Description:** this dataset was collected through a Freedom of Information Law request from the NYC Taxi and Limousine Commission. Data points represent New York taxi trips over 4 years (2010–2013), complete with spatio-temporal data, trip duration, number of passengers, and cost. Reportedly, the dataset contains a large number of errors, including misreported trip distance, duration, and GPS coordinates. Overall, these errors account for 7% of all trips in the dataset.
- **Affiliation of creators:** University of Illinois.
- **Domain:** transportation.
- **Tasks in fairness literature:** fair matching [311, 379].
- **Data spec:** tabular data.
- **Sample size:** ~ 700M taxi trips.
- **Year:** 2016.
- **Sensitive features:** none.
- **Link:** <https://experts.illinois.edu/en/datasets/new-york-city-taxi-trip-data-2010-2013-2>
- **Further info:** <https://bit.ly/3yrT8jt>
- **Variants:** a similar, smaller dataset was obtained by Chris Whong from the NYC Taxi and Limousine Commission under the Freedom of Information Law.¹⁶

A.134 Occupations in Google Images

- **Description:** this dataset was collected to study gender and skin tone diversity in image search results for jobs, and its relation with gender and race concentration in different professions. The dataset consists of the top 100 results for 96 occupations from Google Image Search, collected in December 2019. The creators hired workers on Amazon Mechanical Turk to label the gender (male, female) and Fitzpatrick skin tone (Type 1–6) of the primary person in each image, adding “Not applicable” and “Cannot determine” as possible options. Three labels were collected for each image, to which the majority label was assigned where possible.
- **Affiliation of creators:** Yale University.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair subset selection under unawareness [359].
- **Data spec:** image.
- **Sample size:** ~ 10K images of ~100 occupations.
- **Year:** 2019.
- **Sensitive features:** gender, skin tone (inferred).
- **Link:** <https://drive.google.com/drive/u/0/folders/1j9I5ESc-7NRCZ-zSD0C6LHjeNp42Rjkj>
- **Further info:** Celis and Keswani [82]

¹⁶<http://www.andresmh.com/nyctaxitrips/>

A.135 Office31

- **Description:** this dataset was curated to support domain adaptation algorithms for computer vision systems. It features images of 31 different office tools (e.g. chair, keyboard, printer) from 3 different domains: listings on Amazon, high quality camera images, low quality webcam shots.
- **Affiliation of creators:** University of California, Berkeley.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair clustering [313].
- **Data spec:** image.
- **Sample size:** ~ 4K images.
- **Year:** 2011.
- **Sensitive features:** none.
- **Link:** <https://paperswithcode.com/dataset/office-31>
- **Further info:** Saenko et al. [444]

A.136 Olympic Athletes

- **Description:** this is a historical sports-related dataset on the modern Olympic Games from their first edition in 1896 to the 2016 Rio Games. The dataset was consolidated by Randi H Griffin utilizing SportsReference as the primary source of information. For each athlete, the dataset comprises demographics, height, weight, competition, and medal.
- **Affiliation of creators:** unknown.
- **Domain:** sports.
- **Tasks in fairness literature:** fair clustering [236].
- **Data spec:** tabular data.
- **Sample size:** ~ 300K athletes.
- **Year:** 2018.
- **Sensitive features:** sex, age.
- **Link:** <https://www.kaggle.com/heesoo37/120-years-of-olympic-history-athletes-and-results>
- **Further info:** <https://www.sports-reference.com/>

A.137 Omniglot

- **Description:** this dataset was designed to study the problem of automatically learning basic visual concepts. It consists of handwritten characters from different alphabets drawn online via Amazon Mechanical Turk by 20 different people.
- **Affiliation of creators:** New York University; University of Toronto; Massachusetts Institute of Technology.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair few-shot learning [314].
- **Data spec:** image.
- **Sample size:** ~ 2K images from 50 different alphabets.
- **Year:** 2019.
- **Sensitive features:** none.
- **Link:** <https://github.com/brendenlake/omniglot>
- **Further info:** Lake et al. [298]

A.138 One billion word benchmark

- **Description:** this dataset was proposed in 2014 as a benchmark for language models. The authors sourced English textual data from the EMNLP 6th workshop on Statistical Machine Translation¹⁷, more specifically the Monolingual language model training data, comprising a news crawl from 2007–2011 and data from the European Parliament website. Preprocessing includes removal of duplicate sentences, rare words (appearing less than 3 times) and mapping out-of-vocabulary words to the <UNK> token. The ELMo contextualized WEs [400] were trained on this benchmark.
- **Affiliation of creators:** Google; University of Edinburgh; Cantab Research Ltd.
- **Domain:** linguistics.
- **Tasks in fairness literature:** data bias evaluation [487].
- **Data spec:** text.
- **Sample size:** ~ 800M words.
- **Year:** 2014.
- **Sensitive features:** textual references to people and their demographics.

¹⁷<http://statmt.org/wmt11/training-monolingual.tgz>

- **Link:** <https://opensource.google/projects/lm-benchmark>
- **Further info:** Chelba et al. [90]

A.139 Online Freelance Marketplaces

- **Description:** this dataset was created to audit racial and gender biases on TaskRabbit and Fiverr, two popular online freelancing marketplaces. The dataset was built by crawling workers' profiles from both websites, including metadata, activities, and past job reviews. Profiles were later annotated with perceived demographics (gender and race) by Amazon Mechanical Turk based on profile images. On TaskRabbit, the authors executed search queries for all task categories in the 10 largest cities where the service is available, logging workers' ranking in search results. On Fiverr, they concentrated on 9 tasks of diverse nature. The total number of queries that were issued on each platform, resulting in as many search result pages, is not explicitly stated.
- **Affiliation of creators:** Northeastern University, GESIS Leibniz Institute for the Social Sciences, University of Koblenz-Landau, ETH Zürich.
- **Domain:** information systems.
- **Tasks in fairness literature:** fairness evaluation [212].
- **Data spec:** query-result pairs.
- **Sample size:** ~ 10K workers (Fiverr); ~ 4K (TaskRabbit).
- **Year:** 2017.
- **Sensitive features:** gender, race.
- **Link:** not available
- **Further info:** Hannák et al. [212]

A.140 Open Images Dataset

- **Description:** this dataset was curated to improve and measure the performance of computer vision algorithms. Images with CC-BY license were downloaded from Flickr, and further filtered to remove near-duplicates, inappropriate content, and images appearing elsewhere in the internet. Different versions of this dataset were released, progressively adding a wealth of information on these images, including labels, bounding boxes, segmentation masks, visual relationships, and localized narratives. Bounding boxes relate to 600 classes, including "person", which admits "girl", "boy", "woman", and "man" as a subclass. Labels were generated automatically and later verified by humans with majority voting, with a variable annotator pool size depending on the estimated complexity of verifying each label.
- **Affiliation of creators:** Google.
- **Domain:** computer vision.
- **Tasks in fairness literature:** data bias evaluation [449], fairness evaluation [7].
- **Data spec:** image.
- **Sample size:** ~ 9M images.
- **Year:** 2020.
- **Sensitive features:** gender, age.
- **Link:** <https://storage.googleapis.com/openimages/web/index.html>
- **Further info:** [296]

A.141 Paper-Reviewer Matching

- **Description:** this dataset summarizes the peer review assignment process of 3 different conferences, namely one edition of Medical Imaging and Deep Learning (MIDL) and two editions of the Conference on Computer Vision and Pattern Recognition (called CVPR and CVPR2018). The data, provided by OpenReview and the Computer Vision Foundation, consist of a matrix of paper-reviewer affinities, a set of coverage constraints to ensure each paper is properly reviewed, and a set of upper bound constraints to avoid imposing an excessive burden on reviewers.
- **Affiliation of creators:** unknown.
- **Domain:** library and information sciences.
- **Tasks in fairness literature:** fair matching [283].
- **Data spec:** paper-reviewer pairs.
- **Sample size:** ~ 200 reviewers for ~ 100 papers (MIDL); ~ 1K reviewers for ~ 3K papers (CVPR). ~ 3K reviewers for ~ 5K papers (CVPR2018).
- **Year:** 2019.
- **Sensitive features:** none.
- **Link:** not available
- **Further info:** Kobren et al. [283]

A.142 Philadelphia Crime Incidents

- **Description:** this dataset is provided as part of OpenDataPhilly initiative. It summarizes hundreds of thousands of crime incidents handled by the Philadelphia Police Department over a period of ten years (2006–2016). The dataset comes with fine spatial and temporal granularity and has been used to monitor seasonal and historical trends and measure the effect of police strategies.
- **Affiliation of creators:** Philadelphia Police Department.
- **Domain:** law.
- **Tasks in fairness literature:** fair resource allocation [154].
- **Data spec:** tabular data.
- **Sample size:** ~ 1M crime incidents.
- **Year:** present.
- **Sensitive features:** geography.
- **Link:** <https://www.opendataphilly.org/dataset/crime-incidents>
- **Further info:**

A.143 Pilot Parliaments Benchmark (PPB)

- **Description:** this dataset was developed as a benchmark with a balanced representation of gender and skin type to evaluate the performance of face analysis technology. The dataset features images of parliamentary representatives from three African countries (Rwanda, Senegal, South Africa) and three European countries (Iceland, Finland, Sweden) to achieve a good balance between skin type and gender while reducing potential harms connected with lack of consent from the people involved. Three annotators provided gender and Fitzpatrick labels. A certified surgical dermatologist provided the definitive Fitzpatrick skin type labels. Gender was annotated based on name, gendered title, and photo appearance.
- **Affiliation of creators:** Massachusetts Institute of Technology; Microsoft.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair classification [11, 276], fairness evaluation [63, 422], bias discovery [11, 276].
- **Data spec:** image.
- **Sample size:** ~ 1K images of ~ 1K individuals.
- **Year:** 2018.
- **Sensitive features:** gender, skin type.
- **Link:** <http://gendershades.org/>
- **Further info:** Buolamwini and Gebru [63]

A.144 Pima Indians Diabetes Dataset (PIDD)

- **Description:** this resource owes its name to the respective entry on the UCI repository (now unavailable), and was derived from a medical study of Native Americans from the Gila River Community, often called Pima. The study was initiated in the 1960s by the National Institute of Diabetes and Digestive and Kidney Diseases and found a large prevalence of *diabetes mellitus* in this population. The dataset commonly available nowadays represents a subset of the original study, focusing on women of age 21 or older. It reports whether they tested positive for diabetes, along with eight covariates that were found to be significant risk factors for this population. These include the number of pregnancies, skin thickness, and body mass index, based on which algorithms should predict the test results.
- **Affiliation of creators:** Logistics Management Institute; National Institute of Diabetes Digestive and Kidney Diseases; John Hopkins University.
- **Domain:** endocrinology.
- **Tasks in fairness literature:** fairness evaluation [457], fair clustering [95].
- **Data spec:** tabular data.
- **Sample size:** ~ 800 subjects.
- **Year:** 2016.
- **Sensitive features:** age.
- **Link:** <https://www.kaggle.com/uciml/pima-indians-diabetes-database>
- **Further info:** Radin [417], Smith et al. [470]

A.145 Pokec Social Network

- **Description:** this graph dataset summarizes the networks of Pokec users, a social network service popular in Slovakia and Czech Republic. Due to default privacy settings being predefined as public, a wealth of information for each profile was collected by curators including information on demographics, politics, education, marital status, and children wherever available. This resource was collected to perform data analysis in social networks.
- **Affiliation of creators:** University of Zilina.

- **Domain:** social networks.
- **Tasks in fairness literature:** fair data summarization [153].
- **Data spec:** user-user pairs.
- **Sample size:** ~ 2M nodes (profiles) connected by ~ 30M edges (friendship relations).
- **Year:** 2013.
- **Sensitive features:** gender, geography, age.
- **Link:** <https://snap.stanford.edu/data/soc-pokec.html>
- **Further info:** Takac and Zabovsky [486]

A.146 Popular Baby Names

- **Description:** this dataset summarizes birth registration in New York City, focusing on names sex and race of newborns, providing a reliable source of data to assess naming trends in New York. A similar nation-wide database is maintained by the US Social Security Administration.
- **Affiliation of creators:** City of New York, Department of Health and Mental Hygiene (NYC names); United States Social Security Administration (US names).
- **Domain:** linguistics.
- **Tasks in fairness literature:** fair sentiment analysis [374, 560], bias discovery in WEs [485].
- **Data spec:** tabular data.
- **Sample size:** ~ 3K unique names (NYC names); ~ 30K unique names (US names).
- **Year:** 2021.
- **Sensitive features:** sex, race.
- **Link:** <https://catalog.data.gov/dataset/popular-baby-names> (NYC names); <https://www.ssa.gov/oact/babynames/limits.html> (US names)
- **Further info:**

A.147 Poverty in Colombia

- **Description:** this dataset stems from an official survey of households performed yearly by the Colombian national statistics department (Departamento Administrativo Nacional de Estadística). The survey is aimed at soliciting information about employment, income, and demographics. The data serves as an input for studies on poverty in Colombia.
- **Affiliation of creators:** Departamento Administrativo Nacional de Estadística.
- **Domain:** economics.
- **Tasks in fairness literature:** fair classification [384].
- **Data spec:** tabular data.
- **Sample size:** unknown.
- **Year:** 2018.
- **Sensitive features:** age, sex, geography.
- **Link:** <https://www.dane.gov.co/index.php/estadisticas-por-tema/pobreza-y-condiciones-de-vida/pobreza-y-desigualdad/pobreza-monetaria-y-multidimensional-en-colombia-2018>
- **Further info:** https://www.dane.gov.co/files/investigaciones/condiciones_vida/pobreza/2018/bt_pobreza_monetaria_18.pdf

A.148 PP-Pathways

- **Description:** this dataset represents a network of physical interactions between proteins that are experimentally documented in humans. The dataset was assembled to study the problem of automated discovery of the proteins (nodes) associated with a given disease. Starting from a few known disease-associated proteins and a map of protein-protein interactions (edges), the task is to find the full list of proteins associated with said disease.
- **Affiliation of creators:** Stanford University; Chan Zuckerberg Biohub.
- **Domain:** biology.
- **Tasks in fairness literature:** fair graph mining [262].
- **Data spec:** protein-protein pairs.
- **Sample size:** ~ 20K proteins (nodes) linked by ~ 300K physical interactions.
- **Year:** 2018.
- **Sensitive features:** none.
- **Link:** <http://snap.stanford.edu/biodata/datasets/10000/10000-PP-Pathways.html>
- **Further info:** Agrawal et al. [5]

A.149 Prosper Loans Network

- **Description:** this dataset represents transactions on the Prosper marketplace, a famous peer-to-peer lending service where US-based users can register as lenders or borrowers. This resource has a graph structure and covers the period 2005–2011. Loan records include user ids, timestamps, loan amount, and rate. The dataset was first associated with a study of arbitrage and its profitability in a peer-to-peer lending system.
- **Affiliation of creators:** Prosper; University College Dublin.
- **Domain:** finance.
- **Tasks in fairness literature:** fair classification [315].
- **Data spec:** lender-borrower pairs.
- **Sample size:** ~ 3M loan records involving ~ 100K people.
- **Year:** 2015.
- **Sensitive features:** none.
- **Link:** <http://mlg.ucd.ie/datasets/prosper.html>
- **Further info:** Redmond and Cunningham [427]

A.150 PubMed Diabetes Papers

- **Description:** this dataset was created to study the problem of classification of connected entities via active learning. The creators extracted a set of articles related to diabetes from PubMed, along with their citation network. The task associated with the dataset is inferring a label specifying the type of diabetes addressed in each publication. For this task, TF/IDF-weighted term frequencies of every article are available.
- **Affiliation of creators:** University of Maryland.
- **Domain:** library and information sciences.
- **Tasks in fairness literature:** fair graph mining [312].
- **Data spec:** article-article pairs.
- **Sample size:** ~ 20K articles connected by ~ 40K citations.
- **Year:** 2020.
- **Sensitive features:** none.
- **Link:** <https://linqs.soe.ucsc.edu/data>
- **Further info:** Namata et al. [377]

A.151 Pymetrics Bias Group

- **Description:** Pymetrics is a company that offers a candidate screening tool to employers. Candidates play a core set of twelve games, derived from psychological studies. The resulting gamified psychological measurements are exploited to build predictive models for hiring, where positive examples are provided by high-performing employees from the employer. Pymetrics staff maintain a *Pymetrics Bias Group* dataset for internal fairness audits by asking players to fill in an optional demographic survey after they complete the games.
- **Affiliation of creators:** Pymetrics.
- **Domain:** information systems, management information systems.
- **Tasks in fairness literature:** fairness evaluation [542].
- **Data spec:** tabular data.
- **Sample size:** ~ 10K users.
- **Year:** 2021.
- **Sensitive features:** gender, race.
- **Link:** not available
- **Further info:** Wilson et al. [542]

A.152 Race on Twitter

- **Description:** this dataset was collected to power applications of user-level race prediction on Twitter. Twitter users were hired through Qualtrics, were they filled in a survey providing their Twitter handle and demographics, including race, gender, age, education, and income. The dataset creators downloaded the most recent 3,200 tweets by the users who provided their handle. The data, allegedly released in an anonymized and aggregated format, appears to be unavailable.
- **Affiliation of creators:** University of Pennsylvania.
- **Domain:** social media.
- **Tasks in fairness literature:** fairness evaluation [26].
- **Data spec:** text.
- **Sample size:** ~ 5M tweets from ~ 4K users.

- **Year:** 2018.
- **Sensitive features:** race, gender, age.
- **Link:** <http://www.preotiuc.ro/>
- **Further info:** Preoțiuc-Pietro and Ungar [407]

A.153 Racial Faces in the Wild (RFW)

- **Description:** this dataset was developed as a benchmark for face verification algorithms operating on diverse populations. The dataset comprises 4 clusters of images extracted from MS-Celeb-1M (§ A.124), a dataset that was discontinued by Microsoft due to privacy violations. Clusters are of similar size and contain individuals labelled Caucasian, Asian, Indian and African. Half of the labels (Asian, Indian) are derived from the “Nationality attribute of FreeBase celebrities”; the remaining half (Caucasian, African) is automatically estimated via the Face++ API. This attribute is referred to as “race” by the authors, who also assert “carefully and manually” cleaning every image. Clusters feature multiple images of each individual to allow for face verification applications.
- **Affiliation of creators:** Beijing University of Posts; Telecommunications and Canon Information Technology (Beijing).
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair reinforcement learning [529], fair representation learning [196].
- **Data spec:** image.
- **Sample size:** ~ 50K images of ~ 10K individuals.
- **Year:** 2019.
- **Sensitive features:** race (inferred).
- **Link:** <http://www.whdeng.cn/RFW/testing.html>
- **Further info:** Wang et al. [530]

A.154 Real-Time Crime Forecasting Challenge

- **Description:** this dataset was assembled and released by the US National Institute of Justice in 2017 with the goal of advancing the state of automated crime forecasting. It consists of calls-for-service (CFS) records provided by the Portland Police Bureau for the period 2012–2017. Each CFS record contains spatio-temporal data and crime-related categories. The dataset was released as part of a challenge with a total prize of 1,200,000\$.
- **Affiliation of creators:** National Institute of Justice.
- **Domain:** law.
- **Tasks in fairness literature:** fair spatio-temporal process learning [454].
- **Data spec:** tabular data.
- **Sample size:** ~ 700K CFS records.
- **Year:** 2017.
- **Sensitive features:** geography.
- **Link:** <https://nij.ojp.gov/funding/real-time-crime-forecasting-challenge-posting#data>
- **Further info:** Team Conduent Public Safety Solutions [493]

A.155 Recidivism of Felons on Probation

- **Description:** this dataset covers probation cases of persons who were sentenced in 1986 in 32 urban and suburban US jurisdictions. It was assembled to study the behaviour of individuals on probation and their compliance with court orders across states. Possible outcomes include successful discharge, new felony rearrest, and absconding. The information on probation cases was frequently obtained through manual reviews and transcription of probation files, mostly by college students. Variables include probationer’s demographics, educational level, wage, history of convictions, disciplinary hearings and probation sentences. The final dataset consists of ~ 10K probation cases “representative of 79,043 probationers”.
- **Affiliation of creators:** US Department of Justice; National Association of Criminal Justice Planners.
- **Domain:** law.
- **Tasks in fairness literature:** limited-label fair classification [532].
- **Data spec:** tabular data.
- **Sample size:** ~ 10K probation cases.
- **Year:** 2005.
- **Sensitive features:** sex, race, ethnicity, age.
- **Link:** <https://www.icpsr.umich.edu/web/NACJD/studies/9574>
- **Further info:** <https://bjs.ojp.gov/data-collection/recidivism-survey-felons-probation>

A.156 Reddit Comments

- **Description:** this resource consists of Reddit comments and relative metadata, crawled and made available online for research purposes. While the available dumps cover the period 2006-2021, below the “sample size” field refers to comments from 2014 used in one surveyed work.
- **Affiliation of creators:** Pushshift data.
- **Domain:** social media, linguistics.
- **Tasks in fairness literature:** bias evaluation in language models [208].
- **Data spec:** text.
- **Sample size:** ~ 500M comments.
- **Year:** 2021.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <https://files.pushshift.io/reddit/comments/>
- **Further info:** Guo and Caliskan [208]

A.157 Renal Failure

- **Description:** the dataset was created to compare the performance of two different algorithms for automated renal failure risk assessment. Considering patients who received care at NYU Langone Medical Center, each entry encodes their health records, demographics, disease history, and lab results. The final version of the dataset has a cutoff date, considering only patients who did not have kidney failure by that time, and reporting, as a target ground truth, whether they proceeded to have kidney failure within the next year.
- **Affiliation of creators:** New York University; New York University Langone Medical Center.
- **Domain:** nephrology.
- **Tasks in fairness literature:** fairness evaluation [540].
- **Data spec:** tabular data.
- **Sample size:** ~ 2M patients.
- **Year:** 2019.
- **Sensitive features:** age, gender, race.
- **Link:** not available
- **Further info:** Williams and Razavian [540]

A.158 Reuters 50 50

- **Description:** this dataset was extracted from the Reuters Corpus Volume 1 (RCV1), a large corpus of newswire stories, to study the problem of authorship attribution. The 50 most prolific authors were selected from RCV1, considering only texts labeled corporate/industrial. The dataset consists of short news stories from these authors, labelled with the name of the author.
- **Affiliation of creators:** University of the Aegean.
- **Domain:** news.
- **Tasks in fairness literature:** fair clustering [214].
- **Data spec:** text.
- **Sample size:** ~ 5K articles.
- **Year:** 2011.
- **Sensitive features:** author, textual references to people and their demographics.
- **Link:** http://archive.ics.uci.edu/ml/datasets/Reuter_50_50
- **Further info:** Houvardas and Stamatatos [231]

A.159 Ricci

- **Description:** this dataset relates to the US supreme court labor case on discrimination *Ricci vs DeStefano* (2009), connected with the disparate impact doctrine. It represents 118 firefighter promotion tests, providing the scores and race of each test taker. Eighteen firefighters from the New Haven Fire Department claimed “reverse discrimination” after the city refused to certify a promotion examination where they had obtained high scores. The reasons why city officials avoided certifying the examination included concerns of potential violation of the ‘four-fifths’ rule, as, given the vacancies at the time, no black firefighter would be promoted. The dataset was published and popularized by Weiwen Miao for pedagogical use.
- **Affiliation of creators:** Haverford College.
- **Domain:** law.
- **Tasks in fairness literature:** fairness evaluation [165, 176], limited-label fairness evaluation [248].
- **Data spec:** tabular data.
- **Sample size:** ~ 100 test takers.

- **Year:** 2018.
- **Sensitive features:** race.
- **Link:** http://jse.amstat.org/jse_data_archive.htm; <https://github.com/algofairness/fairness-comparison/tree/master/fairness/data/raw>
- **Further info:** Gastwirth and Miao [180], Miao [364]

A.160 Rice Facebook Network

- **Description:** this dataset represents the Facebook sub-network of students and alumni of Rice University. It consists of a crawl of reachable profiles in the Rice Facebook network, augmented with academic information obtained from Rice University directories. This collection was created to study the problem of inferring unknown attributes in a social network based on the network graph and attributes that are available for a fraction of users.
- **Affiliation of creators:** MPI-SWS; Rice University; Northeastern University.
- **Domain:** social networks.
- **Tasks in fairness literature:** fair graph diffusion [8].
- **Data spec:** user-user pairs.
- **Sample size:** ~ 1K profiles connected by 40K edges.
- **Year:** 2010.
- **Sensitive features:** none.
- **Link:** not available
- **Further info:** Mislove et al. [369]

A.161 Riddle of Literary Quality

- **Description:** this text corpus was assembled to study the factors that correlate with the acceptance of a text as literary (or non-literary) and good (or bad). It consists of 401 Dutch-language novels published between 2007–2012. These works were selected for being bestsellers or often lent from libraries in the period 2009–2012. Due to copyright reasons, the data is not publicly available.
- **Affiliation of creators:** Huygens ING – KNAW; University of Amsterdam; Fryske Akademy.
- **Domain:** literature.
- **Tasks in fairness literature:** fairness evaluation [289].
- **Data spec:** text.
- **Sample size:** ~ 400 novels.
- **Year:** 2017.
- **Sensitive features:** gender (of author).
- **Link:** not available
- **Further info:** Koolen and van Cranenburgh [289]; <https://literaryquality.huygens.knaw.nl/>

A.162 Ride-hailing App

- **Description:** this dataset was gathered from a ride-hailing app operating in an undisclosed major Asian city. It summarizes spatio-temporal data about ride requests (jobs) and assignments to drivers during 29 consecutive days. The data tracks the position and status of taxis logging data every 30-90 seconds.
- **Affiliation of creators:** Max Planck Institute for Software Systems; Max Planck Institute for Informatics.
- **Domain:** transportation.
- **Tasks in fairness literature:** fair matching [481].
- **Data spec:** driver-job pairs.
- **Sample size:** ~ 1K drivers handling ~ 200K job requests.
- **Year:** 2019.
- **Sensitive features:** geography.
- **Link:** not available
- **Further info:** Sühr et al. [481]

A.163 RtGender

- **Description:** this dataset captures differences in online commenting behaviour to posts and videos of female and male users. It was created by collecting posts and top-level comments from four platforms: Facebook, Reddit, Fitocracy, TED talks. For each of the four sources, the possibility to reliably report the gender of the poster or presenter shaped the data collection procedure. Authors of posts and videos were selected among users self-reporting their gender or public figures for which gender annotations were available. For instance, the authors created two Facebook-based datasets: one containing all posts and associated top-level comments for all 412 members of US parliament who have public Facebook pages, and a similar one for 105 American public figures (journalists, novelists, actors, actresses, etc.). The gender of these figures was derived based on their presence on Wikipedia category pages relevant for

gender.¹⁸ The gender of commenters and a reliable ID to identify them across comments may be useful for some analyses. The authors report commenters' first names and a randomized ID, which should support these goals, while reducing chances of re-identification based on last name and Facebook ID.

- **Affiliation of creators:** Stanford University; University of Michigan; Carnegie Mellon University.
- **Domain:** social media, linguistics.
- **Tasks in fairness literature:** fairness evaluation [19].
- **Data spec:** text.
- **Sample size:** ~ 2M posts with ~ 25M comments.
- **Year:** 2018.
- **Sensitive features:** gender.¹⁹
- **Link:** <https://nlp.stanford.edu/robvoigt/rtgender/>
- **Further info:** [518]

A.164 SafeGraph Research Release

- **Description:** this dataset captures mobility patterns in the US and Canada. It is maintained by SafeGraph, a data company powering analytics about access to Points-of-Interest (POI) and mobility, including pandemic research. SafeGraph data is sourced from millions of mobile devices, whose users allow location tracking by some apps. The *Research Release* dataset consists of aggregated estimates of hourly visit counts to over 6 million POI. Given the increasing importance of SafeGraph data, directly influencing not only private initiative but also public policy, audits of data representativeness are being carried out both internally [477] and externally [115].
- **Affiliation of creators:** Safegraph.
- **Domain:** urban studies.
- **Tasks in fairness literature:** data bias evaluation [115].
- **Data spec:** mixture.
- **Sample size:** ~ 7M POI.
- **Year:** present.
- **Sensitive features:** geography.
- **Link:** <https://www.safegraph.com/academics>
- **Further info:** <https://docs.safegraph.com/v4.0/docs>

A.165 Scientist+Painter

- **Description:** this resource was crawled to study the problem of fair and diverse representation in subsets of instances selected from a large dataset, with a focus on gender concentration in professions. The dataset consists of approximately 800 images that equally represent male scientists, female scientists, male painters, and female painters. These images were gathered from Google image search, selecting the top 200 medium sized JPEG files that passed the strictest level of Safe Search filtering. Then, each image was processed to obtain sets of 128-dimensional SIFT descriptors. The descriptors are combined, subsampled and then clustered using k-means into 256 clusters.
- **Affiliation of creators:** École Polytechnique Fédérale de Lausanne (EPFL); Microsoft; University of California, Berkeley.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair data summarization [78, 80].
- **Data spec:** image.
- **Sample size:** ~ 800 images.
- **Year:** 2016.
- **Sensitive features:** male/female.
- **Link:** goo.gl/hNukfP
- **Further info:** Celis et al. [80]

A.166 Section 203 determinations

- **Description:** this dataset is created in support of the language minority provisions of the Voting Rights Act, Section 203. The data contains information about limited-English proficient voting population by jurisdiction, which is used to determine whether election materials must be printed in minority languages. For each combination of language protected by Section 203 and US jurisdiction, the dataset provides information about total population, population of voting age, US citizen population of voting age, combining this information with language spoken at home and overall English proficiency.
- **Affiliation of creators:** US Census Bureau.
- **Domain:** demography.

¹⁸ e.g. https://en.wikipedia.org/wiki/Category:American_female_tennis_players

¹⁹ Annotations for Facebook and TED come from Wikipedia and Mirkin et al. [366] respectively. Reddit and Fitocracy rely on self-reported labels.

- **Tasks in fairness literature:** fairness evaluation of private resource allocation [410].
- **Data spec:** tabular data.
- **Sample size:** ~ 600K combinations of jurisdictions and languages potentially spoken therein.
- **Year:** 2017.
- **Sensitive features:** geography, language.
- **Link:** <https://www.census.gov/data/datasets/2016/dec/rdo/section-203-determinations.html>
- **Further info:** <https://www.census.gov/programs-surveys/decennial-census/about/voting-rights/voting-rights-determination-file-2016.html>

A.167 Sentiment140

- **Description:** this dataset was created to study the problem of sentiment analysis in social media, envisioning applications of product quality and brand reputation analysis via Twitter monitoring. The sentiment of tweets, retrieved via Twitter API, is automatically inferred based on the presence of emoticons conveying joy or sadness. This dataset is part of the LEAF benchmark for federated learning. In federated learning settings, devices correspond to accounts.
- **Affiliation of creators:** Stanford University.
- **Domain:** social media.
- **Tasks in fairness literature:** fair federated learning [314].
- **Data spec:** text.
- **Sample size:** ~ 2M tweets by ~ 600K accounts.
- **Year:** 2012.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <http://help.sentiment140.com/home>
- **Further info:** Go et al. [187]

A.168 Seoul Bike Sharing

- **Description:** this resource, summarizing hourly public rental history of *Seoul Bikes*, was curated to study the problem of bike sharing demand prediction. The data was downloaded from the Seoul Public Data Park website of South Korea and spans one year of utilization (December 2017 to November 2018) of Seoul Bikes, a bike sharing system that started in 2015. This dataset consists of hourly information about weather (e.g. temperature, solar radiation, rainfall) and time (date, time, season, holiday), along with the number of bikes rented at each hour, which is the target of a prediction task.
- **Affiliation of creators:** Sunchon National University.
- **Domain:** transportation.
- **Tasks in fairness literature:** fair regression [137].
- **Data spec:** time series.
- **Sample size:** ~ 9K hourly points.
- **Year:** 2020.
- **Sensitive features:** none.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/Seoul+Bike+Sharing+Demand>
- **Further info:** V E and Cho [509], V E et al. [510], <https://data.seoul.go.kr/index.do>

A.169 Shakespeare

- **Description:** this dataset is available as part of the LEAF benchmark for federated learning [68]. It is built from “The Complete Works of William Shakespeare”, where each speaking role represents a different device. The task envisioned for this dataset is next character prediction.
- **Affiliation of creators:** Google; Carnegie Mellon University; Determined AI.
- **Domain:** literature.
- **Tasks in fairness literature:** fair federated learning [314].
- **Data spec:** text.
- **Sample size:** ~ 4M tokens over ~ 1K speaking roles.
- **Year:** 2020.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** https://www.tensorflow.org/federated/api_docs/python/tff/simulation/datasets/shakespeare
- **Further info:** Caldas et al. [68], McMahan et al. [355]

A.170 Shanghai Taxi Trajectories

- **Description:** this semi-synthetic dataset represents the road network and traffic patterns of Shanghai. Trajectories were collected from thousands of taxis operating in Shanghai. Spatio-temporal traffic patterns were extracted from these trajectories and used to build the dataset.
- **Affiliation of creators:** Shanghai Jiao Tong University; CITI-INRIA Lab.
- **Domain:** transportation.
- **Tasks in fairness literature:** fair routing [411].
- **Data spec:** unknown.
- **Sample size:** unknown.
- **Year:** 2015.
- **Sensitive features:** geography.
- **Link:** not available
- **Further info:** Qian et al. [411]

A.171 shapes3D

- **Description:** this dataset is an artificial benchmark for unsupervised methods aimed at learning disentangled data representations. It consists of images of 3D shapes in a walled environment, with variable floor colour, wall colour, object colour, scale, shape and orientation.
- **Affiliation of creators:** DeepMind; Wayve.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair representation learning [327], fair data generation [103].
- **Data spec:** image.
- **Sample size:** ~ 500K images.
- **Year:** 2018.
- **Sensitive features:** none.
- **Link:** <https://github.com/deepmind/3d-shapes>
- **Further info:** Kim and Mnih [274]

A.172 SIIM-ISIC Melanoma Classification

- **Description:** this dataset was developed to advance the study of automated melanoma classification. The resource consists of dermoscopy images from six medical centers. Images in the dataset are tagged with a patient identifier, allowing lesions from the same patient to be mapped to one another. Images were queried from medical databases among patients with dermoscopy imaging from 1998 to 2019, ranging in quality from 307,200 to 24,000,000 pixels. A curated subset is employed for the 2020 ISIC Grand Challenge.²⁰ This dataset was annotated automatically with a binary Fitzpatrick skin tone label [98].
- **Affiliation of creators:** Memorial Sloan Kettering Cancer Center; University of Queensland; University of Athens; IBM; Universitat de Barcelona; Melanoma Institute Australia; Sydney Melanoma Diagnostic Center; Emory University; Medical University of Vienna; Mayo Clinic; SUNY Downstate Medical School; Stony brook Medical School; Rabin Medical Center; Weill Cornell Medical College.
- **Domain:** dermatology.
- **Tasks in fairness literature:** fairness evaluation of private classification [98].
- **Data spec:** image.
- **Sample size:** ~ 30K images of ~ 2K patients.
- **Year:** 2020.
- **Sensitive features:** skin type.
- **Link:** [urlhttps://doi.org/10.34970/2020-ds01](https://doi.org/10.34970/2020-ds01)
- **Further info:** Rotemberg et al. [437]

A.173 SmallORB

- **Description:** this dataset was assembled by researchers affiliated with New York University as a benchmark for robust object recognition under variable pose and lighting conditions. It consists of images of 50 different toys belonging to 5 categories (four-legged animals, human figures, airplanes, trucks, and cars) obtained by 2 different cameras.
- **Affiliation of creators:** New York University; NEC Labs America.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fair representation learning [327].
- **Data spec:** image.

²⁰<https://www.kaggle.com/c/siim-isic-melanoma-classification>

- **Sample size:** ~ 100K images.
- **Year:** 2005.
- **Sensitive features:** none.
- **Link:** <https://cs.nyu.edu/~yclab/data/norb-v1.0-small/>
- **Further info:** LeCun et al. [306]

A.174 Spliddit Divide Goods

- **Description:** this dataset summarizes instances of usage of the *divide goods* feature of Spliddit, a not-for-profit academic endeavor providing easy access to fair division methods. A typical use case for the service is inheritance division. Participants express their preferences by dividing 1,000 points between the available goods. In response, the service provides suggestions that are meant to maximize the overall satisfaction of all stakeholders.
- **Affiliation of creators:** Spliddit.
- **Domain:** economics.
- **Tasks in fairness literature:** fair preference-based resource allocation [20].
- **Data spec:** tabular data.
- **Sample size:** ~ 1K division instances.
- **Year:** 2016.
- **Sensitive features:** none.
- **Link:** not available
- **Further info:** Caragiannis et al. [74]; <http://www.spliddit.org/apps/goods>

A.175 Stanford Medicine Research Data Repository

- **Description:** this is a data lake/repository developed at Stanford University, supporting a number of data sources and access pipelines. The aim of the underlying project is favouring access to clinical data for research purposes through flexible and robust management of medical data. The data comes from Stanford Health Care, the Stanford Children's Hospital, the University Healthcare Alliance and Packard Children's Health Alliance clinics.
- **Affiliation of creators:** Stanford University.
- **Domain:** medicine.
- **Tasks in fairness literature:** fair risk assessment [402].
- **Data spec:** mixture.
- **Sample size:** ~3M individuals.
- **Year:** present.
- **Sensitive features:** race, ethnicity, gender, age.
- **Link:** <https://starr.stanford.edu/>
- **Further info:** Datta et al. [126], Lowe et al. [330]

A.176 State Court Processing Statistics (SCPS)

- **Description:** this resource was curated as part of the SCPS program. The program tracked felony defendants from charging by the prosecutor until disposition of their cases for a maximum of 12 months (24 months for murder cases). The data represents felony cases filed in approximately 40 populous US counties in the period 1990-2009. Defendants are summarized by 106 variables summarizing demographics, arrest charges, criminal history, pretrial release and detention, adjudication, and sentencing.
- **Affiliation of creators:** US Department of Justice.
- **Domain:** law.
- **Tasks in fairness literature:** fairness evaluation of multi-stage classification [201].
- **Data spec:** tabular data.
- **Sample size:** ~ 200K defendants.
- **Year:** 2014.
- **Sensitive features:** gender, race, age, geography.
- **Link:** <https://www.icpsr.umich.edu/web/NACJD/studies/2038/datadocumentation>
- **Further info:** <https://bjs.ojp.gov/data-collection/state-court-processing-statistics-scps>

A.177 Steemit

- **Description:** this resource was collected to test novel approaches for personalized content recommendation in social networks. It consists of two separate datasets summarizing interactions in the Spanish subnetwork and the English subnetwork of Steemit, a blockchain-based social media website. The datasets summarize user-post interactions in a binary fashion, using comments as a proxy

for positive engagement. The datasets cover a whole year of commenting activities over the period 2017–2018 and comprise the text of posts.

- **Affiliation of creators:** Hong Kong University of Science and Technology; WeBank.
- **Domain:** social media.
- **Tasks in fairness literature:** fairness evaluation [546].
- **Data spec:** user-post pairs.
- **Sample size:** ~ 50K users interacting over ~ 200K posts.
- **Year:** 2019.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <https://github.com/HKUST-KnowComp/Social-Explorative-Attention-Networks>
- **Further info:** Xiao et al. [546]

A.178 Stop, Question and Frisk

- **Description:** Stop, Question and Frisk (SQF) is an expression that commonly refers to a New York City policing program under which officers can briefly detain, question, and search a citizen if the officer has a reasonable suspicion of criminal activity. Concerns about race-based disparities in this practice have been expressed multiple times, especially in connection with the subjective nature of “reasonable suspicion” and the fact that being in a “high-crime area” lawfully lowers the bar of what may constitute reasonable suspicion. The NYPD has a policy of keeping track of most stops, recording them in UF-250 forms which are maintained centrally and distributed by the NYPD. The form includes several information such as place and time of a stop, the duration of the stop and its outcome along with data on demographics and physical appearance of the suspect. Currently available data pertains to years 2003–2020.
- **Affiliation of creators:** New York Police Department.
- **Domain:** law.
- **Tasks in fairness literature:** preference-based fair classification [563], robust fair classification [257], fair classification under unawareness [272], fairness evaluation [191], fair classification [9].
- **Data spec:** tabular data.
- **Sample size:** ~ 1M records.
- **Year:** 2021.
- **Sensitive features:** race, age, sex, geography.
- **Link:** <https://www1.nyc.gov/site/nypd/stats/reports-analysis/stopfrisk.page>
- **Further info:** Gelman et al. [184], Goel et al. [192]

A.179 Strategic Subject List

- **Description:** this dataset was funded through a Bureau of Justice Assistance grant and leveraged by the Illinois Institute of Technology to develop the Chicago Police Department’s Strategic Subject Algorithm. The algorithm provides a risk score which reflects an individual’s probability of being involved in a shooting incident either as a victim or an offender. For each individual, the dataset provides information about the circumstances of their arrest, their demographics and criminal history. The dataset covers arrest data from the period 2012–2016; the associated program was discontinued in 2019.
- **Affiliation of creators:** Chicago Police Department; Illinois Institute of Technology.
- **Domain:** law.
- **Tasks in fairness literature:** fairness evaluation [46].
- **Data spec:** tabular data.
- **Sample size:** ~ 400K individuals.
- **Year:** 2020.
- **Sensitive features:** ace, sex, age.
- **Link:** <https://data.cityofchicago.org/Public-Safety/Strategic-Subject-List-Historical/4aki-r3np>
- **Further info:** Hollywood et al. [228]

A.180 Student

- **Description:** the data was collected from two Portuguese public secondary schools in the Alentejo region, to investigate student achievement prediction and identify decisive factors in student success. The data tracks student performance in Mathematics and Portuguese through school year 2005-2006 and is complemented by demographic, socio-economical, and personal data obtained through a questionnaire. Numerical grades (20-point scale) collected by students over three terms are typically the target of the associated prediction task.
- **Affiliation of creators:** University of Minho.
- **Domain:** education.

- **Tasks in fairness literature:** fair regression [110, 111, 224], rich-subgroup fairness evaluation [270], fair data summarization [36, 254].
- **Data spec:** tabular data.
- **Sample size:** ~ 600 students.
- **Year:** 2014.
- **Sensitive features:** sex, age.
- **Link:** <https://archive.ics.uci.edu/ml/datasets/student+performance>
- **Further info:** Cortez and Silva [114]

A.181 Sushi

- **Description:** this dataset was sourced online via a commercial survey service to evaluate rank-based approaches to solicit preferences and provide recommendations. The dataset captures the preferences for different types of sushi held by people in different areas of Japan. These are encoded both as ratings in a 5-point scale and ordered lists of preferences, which recommenders should learn via collaborative filtering. Demographic data was also collected to study geographical preference patterns.
- **Affiliation of creators:** Japanese National Institute of Advanced Industrial Science and Technology (AIST).
- **Domain:** .
- **Tasks in fairness literature:** fair data summarization [100].
- **Data spec:** user-sushi pairs.
- **Sample size:** ~ 5K respondents.
- **Year:** 2016.
- **Sensitive features:** gender, age, geography.
- **Link:** <https://www.kamishima.net/%20sushi/>
- **Further info:** Kamishima [261]

A.182 Symptoms in Queries

- **Description:** the purpose of this dataset is to study, using only aggregate statistics, the fairness and accuracy of a classifier that predicts whether an individual has a certain type of cancer based on their Bing search queries. The dataset does not include individual data points. It provides, for each US state, and for 18 types of cancer, the proportion of individuals who have this cancer in the state according to CDC 2019 data,²¹ and the proportion of individuals who are predicted to have this cancer according to the classifier that was calculated using Bing queries.
- **Affiliation of creators:** Microsoft; Ben-Gurion University of the Negev.
- **Domain:** information systems, public health.
- **Tasks in fairness literature:** limited-label fairness evaluation [443].
- **Data spec:** tabular data.
- **Sample size:** statistics for ~ 20 cancer types across ~ 50 US states.
- **Year:** 2020.
- **Sensitive features:** geography.
- **Link:** https://github.com/sivansabato/bfa/blob/master/cancer_data.m
- **Further info:** Sabato and Yom-Tov [443]

A.183 TAPER Twitter Lists

- **Description:** this resource was collected to study the problem of personalized expert recommendation, leveraging Twitter lists where users labelled other users as relevant for (or expert in) a given topic. The creators started from a seed dataset of over 12 million geo-tagged Twitter lists, which they filtered to only keep US-based users in topics: news, music, technology, celebrities, sports, business, politics, food, fashion, art, science, education, marketing, movie, photography, and health. A subset of this dataset was annotated with user race (whites and non-whites) via Face++ [590].
- **Affiliation of creators:** Texas A&M University.
- **Domain:** social media.
- **Tasks in fairness literature:** fair ranking [590].
- **Data spec:** user-topic pairs.
- **Sample size:** ~ 10K Twitter lists featuring ~ 8K list members.
- **Year:** 2016.
- **Sensitive features:** race.
- **Link:** not available
- **Further info:** Ge et al. [181]

²¹<https://gis.cdc.gov/Cancer/USCS/DataViz.html>

A.184 TaskRabbit

- **Description:** this resource was assembled to study the effectiveness of fair ranking approaches in improving outcomes for protected groups in online hiring. It consists of the top 10 results returned by the online freelance marketplace TaskRabbit for three queries: “Shopping”, “Event staffing”, and “Moving Assistance”. The geographic location for a query was especially selected to yield a ranking with 3 female candidates among the top 10, with most of them appearing in the bottom 5, which may be a motivating condition for a fairness intervention. Candidates’ gender was manually labelled by creators based on pronoun usage and profile pictures. For each profile, the authors extracted information on job suitability, including TaskRabbit relevance scores, number of completed tasks and positive reviews.
- **Affiliation of creators:** Technische Universität Berlin; Harvard University.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair ranking evaluation [482], multi-stage fairness evaluation [482].
- **Data spec:** query-worker pairs.
- **Sample size:** 3 rankings (one per query) of ~ 10 workers.
- **Year:** 2021.
- **Sensitive features:** gender.
- **Link:** not available
- **Further info:** Sühr et al. [482]

A.185 TIMIT

- **Description:** this resource was curated to power studies of phonetics and to evaluate systems of automated speech recognition. The dataset features speakers of different American English dialects, and includes time-aligned orthographic, phonetic and word transcriptions. Utterances are sampled at a 16kHz frequency.
- **Affiliation of creators:** University of Pennsylvania; National Institute of Standards and Technology; Massachusetts Institute of Technology; SRI International; Texas Instruments.
- **Domain:** linguistics.
- **Tasks in fairness literature:** fairness evaluation of speech recognition [451].
- **Data spec:** time series.
- **Sample size:** ~ 600 speakers, each uttering ~ 10 sentences.
- **Year:** 1993.
- **Sensitive features:** dialect, gender.
- **Link:** <https://catalog.ldc.upenn.edu/LDC93S1>
- **Further info:** <https://en.wikipedia.org/wiki/TIMIT>

A.186 Toy Dataset 1

- **Description:** this dataset consists of $\sim 4K$ points generated as follows. Binary class labels y are generated at random for each point. Next, two-dimensional features x are assigned to each point, sampling from gaussian distributions whose mean and variance depend on y , so that $p(x|y = 1) = \mathcal{N}([2; 2], [5, 1; 1, 5])$; $p(x|y = -1) = \mathcal{N}([-2; -2], [10, 1; 1, 3])$. Finally, each point’s sensitive attribute z is sampled from a Bernoulli distribution so that $p(z = 1) = p(x'|y = 1)/(p(x'|y = 1) + p(x'|y = -1))$, where x' is a rotated version of x : $x' = [\cos(\phi), -\sin(\phi); \sin(\phi), \cos(\phi)]x$. Parameter ϕ controls the correlation between class label y and sensitive attribute z .
- **Affiliation of creators:** Max Planck Institute for Software Systems.
- **Domain:** N/A.
- **Tasks in fairness literature:** fair classification [433, 564], fair preference-based classification [10, 563], fair few-shot learning [466, 468], fair classification under unawareness [272].
- **Data spec:** tabular data.
- **Sample size:** $\sim 4K$ points.
- **Year:** 2017.
- **Sensitive features:** N/A.
- **Link:** https://github.com/mbilalzafar/fair-classification/tree/master/disparate_impact/synthetic_data_demo
- **Further info:** Zafar et al. [564]

A.187 Toy Dataset 2

- **Description:** this dataset contains synthetic relevance judgements over pairs of queries and documents that are biased against a minority group. For each query, there are 10 candidate documents, 8 from group G_0 and 2 from minority group G_1 . Each document is associated with a feature vector (x_1, x_2) , with both components sampled uniformly at random from the interval $(0, 3)$. The relevance of documents is set to $y = x_1 + x_2$ and clipped between 0 and 5. Feature x_2 is then corrupted and replaced by zero for group G_1 , leading to a biased representation between groups, such that any use of x_2 should lead to unfair rankings.

- **Affiliation of creators:** Cornell University.
- **Domain:** N/A.
- **Tasks in fairness literature:** fair ranking [54, 463].
- **Data spec:** query-document pairs.
- **Sample size:** ~ 1 K relevance judgements overs ~ 100 queries with ~ 10 candidate documents.
- **Year:** 2019.
- **Sensitive features:** N/A.
- **Link:** <https://github.com/ashudeep/Fair-PGRank>
- **Further info:** Singh and Joachims [463]

A.188 Toy Dataset 3

- **Description:** this dataset was created to demonstrate undesirable properties of a family of fair classification approaches. Each instance in the dataset is associated with a sensitive attribute z , a target variable y encoding employability, one feature that is important for the problem at hand and correlated with z (work_experience) and a second feature which is unimportant yet also correlated with z (hair_length). The data generating process is the following:

$$\begin{aligned}
 z_i &\sim \text{Bernoulli}(0.5) \\
 \text{hair_length}_i | z_i = 1 &\sim 35 \cdot \text{Beta}(2, 2) \\
 \text{hair_length}_i | z_i = 0 &\sim 35 \cdot \text{Beta}(2, 7) \\
 \text{work_exp}_i | z_i &\sim \text{Poisson}(25 + 6z_i) - \text{Normal}(20, 0.2) \\
 y_i | \text{work_exp}_i &\sim 2 \cdot \text{Bernoulli}(p_i) - 1, \\
 \text{where } p_i &= 1 / (1 + \exp[-(-25.5 + 2.5\text{work_exp})])
 \end{aligned}$$

- **Affiliation of creators:** Carnegie Mellon University; University of California, San Diego.
- **Domain:** N/A.
- **Tasks in fairness literature:** fairness evaluation [46, 320].
- **Data spec:** tabular data.
- **Sample size:** ~ 2 K points.
- **Year:** 2018.
- **Sensitive features:** N/A.
- **Link:** not available
- **Further info:** Lipton et al. [320]

A.189 Toy Dataset 4

- **Description:** in this toy example, features are generated according to four 2-dimensional isotropic Gaussian distributions with different mean μ and variance σ^2 . Each of the four distributions corresponds to a different combination of binary label y and protected attribute s as follows: (1) $s = a, y = +1 : \mu = (-1, -1), \sigma^2 = 0.8$; (2) $s = a, y = -1 : \mu = (1, 1), \sigma^2 = 0.8$; (3) $s = b, y = +1 : \mu = (0.5, -0.5), \sigma^2 = 0.5$; (4) $s = b, y = -1 : \mu = (0.5, 0.5), \sigma^2 = 0.5$.
- **Affiliation of creators:** Istituto Italiano di Tecnologia; University of Genoa; University of Waterloo; University College London.
- **Domain:** N/A.
- **Tasks in fairness literature:** fair classification [143], fairness evaluation [541].
- **Data spec:** tabular data.
- **Sample size:** ~ 6 K points.
- **Year:** 2018.
- **Sensitive features:** N/A.
- **Link:** https://github.com/jmikko/fair_ERM
- **Further info:** Donini et al. [143]

A.190 TREC Robust04

- **Description:** this classic information retrieval collection is a set of topics, documents and relevance judgements collected as part of the Text REtrieval Conference (TREC) 2004 Robust Retrieval Track to catalyze research improving the consistency of information retrieval technology. Documents are taken from articles published during the 1990s in the Financial Times Limited, the Federal Register, the Foreign Broadcast Information Service, and the Los Angeles Times. Graded relevance (not relevant, relevant, highly

relevant) was judged by human assessors for a subset of all possible topic-document combinations, which were selected as “promising” by the automated systems that entered the TREC initiative. The associated task is predicting the relevance of documents for various textual queries.

- **Affiliation of creators:** National Institute of Standards and Technology.
- **Domain:** news, information systems.
- **Tasks in fairness literature:** fair ranking evaluation [185].
- **Data spec:** query-document pairs.
- **Sample size:** ~ 300K relevance judgements over ~ 200 queries and ~ 500K documents.
- **Year:** 2005.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** https://trec.nist.gov/data/t13_robust.html
- **Further info:** Voorhees [520]

A.191 Twitch Social Networks

- **Description:** this dataset was developed to study the effectiveness of node embeddings for learning tasks defined on graphs. This resource concentrates on Twitch content creators streaming in 6 different languages. The dataset has users as nodes, mutual friendships as edges, and node embeddings summarizing games liked, location and streaming habits. The original task on this dataset is predicting whether a streamer uses explicit language.
- **Affiliation of creators:** University of Edinburgh.
- **Domain:** social networks.
- **Tasks in fairness literature:** fair graph mining [262].
- **Data spec:** user-user pairs.
- **Sample size:** ~30K nodes (users) connected by ~ 400K edges (mutual friendship).
- **Year:** 2019.
- **Sensitive features:** none.
- **Link:** <http://snap.stanford.edu/data/twitch-social-networks.html>
- **Further info:** Rozemberczki et al. [438]

A.192 Twitter Abusive Behavior

- **Description:** this dataset is the result of an eight-month crowdsourced study of various forms of abusive behavior on Twitter. The authors began by considering a wide variety of inappropriate speech categories, analyzing how they are used by amateur annotators hired on CrowdFlower. After two exploratory rounds, they merged some labels and eliminated others, converging to a final four-class categorization into (normal, spam, abusive, hateful), requiring five crowdsourced judgements per tweet. Tweets were sampled according to a boosted random sampling technique. A large part of the dataset is randomly sampled, with the addition of tweets that are likely to belong to one or more of the minority (non-normal) classes. The dataset is available as a table mapping tweet IDs to behavior category, making it possible to identify Twitter users in this dataset.
- **Affiliation of creators:** Aristotle University of Thessaloniki; Cyprus University of Technology; Telefonica; University of Alabama at Birmingham; University College London.
- **Domain:** social media.
- **Tasks in fairness literature:** fairness evaluation of harmful content detection [26].
- **Data spec:** text.
- **Sample size:** ~ 100K tweets.
- **Year:** 2018.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <https://github.com/ENCASEH2020/hatespeech-twitter>
- **Further info:** Founta et al. [174]

A.193 Twitter Hate Speech Detection

- **Description:** this dataset was developed to study the problem of automated hate speech detection. The creators used the Twitter API to search for tweets containing racist and sexist terms and hashtags. The annotation was carried out by the authors, with an external review by a 25-year-old woman studying gender studies. After identifying a list of eleven criteria to identify hate speech against a minority, each tweet was labelled as sexism, racism or none. The task associated with this resource is hate speech detection. The dataset is available as a table mapping tweet IDs to hate speech category, making it possible to identify Twitter users in this dataset.
- **Affiliation of creators:** University of Copenhagen.
- **Domain:** social media.
- **Tasks in fairness literature:** fairness evaluation [26].

- **Data spec:** text.
- **Sample size:** ~ 20K tweets.
- **Year:** 2016.
- **Sensitive features:** textual references to people and their demographics.
- **Link:**
- **Further info:** Waseem and Hovy [534]

A.194 Twitter Offensive Language

- **Description:** this dataset was developed to study the problem of automated hate speech detection, and to distinguish between hate speech and other kinds of offensive language. The creators used the Twitter API to search for tweets containing terms from a hate speech lexicon compiled by *Hatebase.org*. Workers on CrowdFlower annotated a random subset of these tweets as hate speech, offensive but not hate speech, or neither offensive nor hate speech. Workers were explicitly told that the mere presence of a slur word does not amount to hate speech. Three of more workers annotated each tweet.
- **Affiliation of creators:** Cornell University; Qatar Computing Research Institute.
- **Domain:** social media.
- **Tasks in fairness literature:** fairness evaluation [26], fair multi-stage classification [271].
- **Data spec:** text.
- **Sample size:** ~ 20K tweets.
- **Year:** 2017.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <https://github.com/t-davidson/hate-speech-and-offensive-language/tree/master/data>
- **Further info:** Davidson et al. [129]

A.195 Twitter Online Harrassment

- **Description:** this dataset was developed as multidisciplinary resource to study online harrassment. The authors searched a stream of tweets for keywords likely to denote violent, offensive, threatening or hateful content based on race, gender, religion and sexual orientation. They developed coding guidelines to label a tweet as harrassing or non/harrassing and spent three weeks reviewing and refining it, annotating sample tweets as a group, and discussing the results. The curators are not publicly sharing the dataset due to Twitter terms of service restrictions and privacy concerns about individuals whose tweets are included; researchers can request access.
- **Affiliation of creators:** University of Maryland.
- **Domain:** social media.
- **Tasks in fairness literature:** fairness evaluation [26].
- **Data spec:** text.
- **Sample size:** ~ 40K tweets.
- **Year:** 2017.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** not available
- **Further info:** Golbeck et al. [194]

A.196 Twitter Political Searches

- **Description:** this dataset was collected to study political biases in Twitter search results, due to political leaning of tweets and biases in the Twitter ranking algorithm. The authors identified 25 popular political queries in December 2015, and collected relevant tweets during a week in which two presidential debates occurred, via the Twitter streaming API. Tweets were annotated based on users' political leaning. Users' leaning was automatically inferred from their topics of interest, via a classifier trained on representative sets of democratic and republican users. Both the accuracy of classifiers and the validity of user leaning as a proxy for tweet leaning was validated by workers recruited on Amazon Mechanical Turk.
- **Affiliation of creators:** Max Planck Institute for Software Systems; University of Illinois at Urbana-Champaign; Indian Institute of Engineering Science and Technology, Shibpur; Adobe Research.
- **Domain:** social media.
- **Tasks in fairness literature:** social media.
- **Data spec:** query-result pairs.
- **Sample size:** ~ 30K search results containing ~ 30K distinct tweets from ~ 20K users.
- **Year:** 2016.
- **Sensitive features:** political leaning.
- **Link:** not available
- **Further info:** Kulshrestha et al. [293]

A.197 Twitter Presidential Politics

- **Description:** this dataset was created by collecting tweets, through the Twitter API, from 576 accounts linked to presidential candidates and members of congress, from the entire account history until December 2019. Out of all the accounts considered, 258 accounts were classified as Republican and 318 as Democratic. The dataset was collected to build a political bias subspace from word embeddings, which could be a flexible tool to quantitatively investigate political leaning in text-based media.
- **Affiliation of creators:** Clarkson University.
- **Domain:** social media.
- **Tasks in fairness literature:** bias audit [198].
- **Data spec:** text.
- **Sample size:** ~ 1M tweets from ~ 500 accounts.
- **Year:** 2020.
- **Sensitive features:** political leaning.
- **Link:** not available
- **Further info:** Gordon et al. [198]

A.198 Twitter Trending Topics

- **Description:** this dataset was used to study the problem of fair recommendation. It comprises a random sample (1%) of all tweets posted in the US between February and July 2017, obtained through the Twitter Streaming API. This sample is paired with a collection of trending Twitter topics queried every 15-minutes through the Twitter REST API in July 2017. User interest in each topic was inferred using Twitter lists and follower-followee graphs. Finally, user demographics were also annotated to evaluate how user interest in different topics skews with respect to race, age, and gender. These attributes were obtained feeding user profile images to Face++.
- **Affiliation of creators:** Indian Institute of Technology Kharagpur; Max Planck Institute for Software Systems; Grenoble INP.
- **Domain:** social media.
- **Tasks in fairness literature:** fair ranking [87].
- **Data spec:** text.
- **Sample size:** ~ 200M tweets by ~ 10M users and ~ 10K trending topics.
- **Year:** 2018.
- **Sensitive features:** race, age, and gender.
- **Link:** not available
- **Further info:** Chakraborty et al. [87]

A.199 TwitterAAE

- **Description:** this resource was developed to study the use of dialect language on social media. The authors used Twitter APIs to collect public tweets sent on mobile phones from US users in 2013. They devise a distant supervision approach based on geolocation to annotate the probable language/dialect of the tweet, distinguishing between African American English (AAE) and Standard American English (SAE). To validate their approach, the creators studied the phonological and syntactic divergence of AAE tweets vs. SAE tweets, ensuring they align with linguistic phenomena that typically distinguish these variants of English.
- **Affiliation of creators:** University of Massachusetts Amherst.
- **Domain:** social media, linguistics.
- **Tasks in fairness literature:** fairness evaluation of sentiment analysis [460], fairness evaluation of private classification [22], fairness evaluation [26], robust fair language model [217], fairness evaluation of language identification [47].
- **Data spec:** text.
- **Sample size:** ~ 8M tweets.
- **Year:** 2016.
- **Sensitive features:** dialect (related to race).
- **Link:** <http://slanglab.cs.umass.edu/TwitterAAE/>
- **Further info:** Blodgett et al. [48]

A.200 U.S. Harmonized Tariff Schedules (HTS)

- **Description:** this resource represents a comprehensive classification system for goods imported in the US, which defines the applicable tariffs. It defines a fine-grained categorization for goods, based e.g. on their material and shape. The chapter on apparel was explicitly criticized for its differential treatment of men's and women's clothing, effectively resulting in discriminatory tariffs for consumers.
- **Affiliation of creators:** US International Trade Commission.
- **Domain:** economics.
- **Tasks in fairness literature:** fairness evaluation [333].
- **Data spec:** tabular data.

- **Sample size:** unknown.
- **Year:** present.
- **Sensitive features:** gender.
- **Link:** <https://hts.usitc.gov/current>
- **Further info:** Barbaro [30]

A.201 UniGe

- **Description:** this dataset is connected with the *DROP@UNIGE* project, aimed at studying the dynamics of university dropout, focusing on the University of Genoa as a case study. In ML fairness literature, the most common version of the dataset focuses on students who enrolled in 2017. Students are associated with attributes describing their ethnicity, gender, financial status, and prior school experience. The target variable encodes early academic success, as summarized by students' grades at the end of the first semester.
- **Affiliation of creators:** University of Genoa.
- **Domain:** education.
- **Tasks in fairness literature:** fair regression [110, 111], fair representation learning [389, 390].
- **Data spec:** tabular data.
- **Sample size:** ~ 5K students.
- **Year:** unknown.
- **Sensitive features:** ethnicity, gender, financial status.
- **Link:** not available
- **Further info:** Oneto et al. [391]

A.202 University Facebook Networks

- **Description:** a collection of 100 datasets shared with researchers in anonymized format by Adam D'Angelo of Facebook. The datasets used in the fairness literature consist of a 2005 snapshot from the Facebook network of the Universities of Oklahoma (Oklahoma97), North Carolina (UNC28), Caltech (Caltech36), Reed College (Reed98), and Michigan State (Michigan23), and links between them. User data comprises gender, class year, and anonymized data fields representing high school, major, and dormitory residences.
- **Affiliation of creators:** Facebook; University of North Carolina; Harvard University; University of Oxford.
- **Domain:** social networks.
- **Tasks in fairness literature:** fair graph mining [312], fair graph augmentation [423].
- **Data spec:** user-user pairs.
- **Sample size:** ~ 20K people connected by ~ 1M friend relations (Oklahoma97); ~ 20K people connected by ~ 1M friend relations (UNC28); ~ 30K people connected by ~ 1M friend relations (Michigan23); ~ 1K people connected by ~ 20K friend relations (Reed98); ~ 1K people connected by ~ 20K friend relations (Caltech36).
- **Year:** 2017.
- **Sensitive features:** gender.
- **Link:** <http://networkrepository.com/socfb-Oklahoma97.php> (Oklahoma97); <http://networkrepository.com/socfb-UNC28.php> (UNC28); <https://networkrepository.com/socfb-Michigan23.php> (Michigan23); <https://networkrepository.com/socfb-Reed98.php> (Reed98); <https://networkrepository.com/socfb-Caltech36.php> (Caltech36)
- **Further info:** Red et al. [425]

A.203 US Census Data (1990)

- **Description:** this resource is a one percent sample extracted from the 1990 US census data as a benchmark for clustering algorithms on large datasets. It contains a variety of features about different aspects of participants' lives, including demographics, wealth, and military service.
- **Affiliation of creators:** Microsoft.
- **Domain:** demography.
- **Tasks in fairness literature:** fair clustering [21, 39, 236], fair clustering under unawareness [157], limited-label fairness evaluation [443].
- **Data spec:** tabular data.
- **Sample size:** ~ 2M respondents.
- **Year:** 1999.
- **Sensitive features:** age, sex.
- **Link:** [https://archive.ics.uci.edu/ml/datasets/US+Census+Data+\(1990\)](https://archive.ics.uci.edu/ml/datasets/US+Census+Data+(1990))
- **Further info:** Meek et al. [357]

A.204 US Family Income

- **Description:** this resource was compiled from the Current Population Survey (CPS) Annual Social and Economic (ASEC) Supplement. It contains income data for over 80,000 thousand US families, broken down by age and race (White, Black, Asian, and Hispanic).
- **Affiliation of creators:** US Bureau of Labor Statistics; US Census Bureau.
- **Domain:** economics.
- **Tasks in fairness literature:** fair subset selection under unawareness [359].
- **Data spec:** tabular data.
- **Sample size:** 4 races x 12 age categories x 41 income categories.
- **Year:** 2020.
- **Sensitive features:** age, race.
- **Link:** <https://www.census.gov/data/tables/time-series/demo/income-poverty/cps-finc/finc-02.html>
- **Further info:** <https://www2.census.gov/programs-surveys/cps/techdocs/cpsmar20.pdf>

A.205 US Federal Judges

- **Description:** this dataset was extracted from Epstein et al. [155] to study the problem of judicial subset selection from the point of view of justice, fairness and interpretability. Given the fact that in several judicial systems a subset of judges is selected from the whole judicial body to decide the outcome of appeals, the creators extract cases where three judges are required from Epstein et al. [155], covering the period 2000–2004. They emulate prior probabilities of affirmance/reversal for specific judges based on their past decisions. The task associated with this dataset is the optimal selection of a subset of judges, so that the procedure is interpretable, the subset contains at least one female (junior) judge and the decision of the subset coincides with the decision of the whole judicial body.
- **Affiliation of creators:** Yale University.
- **Domain:** law.
- **Tasks in fairness literature:** fair subset selection [238].
- **Data spec:** judge-case pairs.
- **Sample size:** ~300 judges selected for ~ 2K cases.
- **Year:** 2020.
- **Sensitive features:** gender.
- **Link:** not available
- **Further info:** Huang et al. [238]

A.206 US Student Performance

- **Description:** this resource represents students at an undisclosed US research university, spanning the Fall 2014 to Spring 2019 terms. The associated task is predicting student success based on university administrative records. Student features include demographics and academic information on prior achievement and standardized test scores.
- **Affiliation of creators:** Cornell University.
- **Domain:** education.
- **Tasks in fairness literature:** fairness evaluation [307].
- **Data spec:** tabular data.
- **Sample size:** unknown.
- **Year:** 2020.
- **Sensitive features:** gender, racial-ethnic group.
- **Link:** not available
- **Further info:** Lee and Kizilcec [307]

A.207 UTK Face

- **Description:** the dataset was developed as a diverse resource for face regression and progression (models of aging), where diversity is intended with respect to age, gender and race. The creators sourced part of the images from two existing datasets (Morph and CACD datasets). To increase the representation of some age groups, additional images were crawled from major search engines based on specific keywords (e.g., baby). Age, gender, and race were estimated through an algorithm and validated by a human annotator.
- **Affiliation of creators:** University of Tennessee.
- **Domain:** computer vision.
- **Tasks in fairness literature:** robust fairness evaluation [378], fairness evaluation of private classification [22], fairness evaluation [451], fair classification [255].
- **Data spec:** image.
- **Sample size:** ~ 20K face images.
- **Year:** 2017.

- **Sensitive features:** age, gender, race (inferred).
- **Link:** <https://susanqq.github.io/UTKFace/>
- **Further info:** Zhang et al. [578]

A.208 Vehicle

- **Description:** this dataset comprises measurements from a distributed network of acoustic, seismic, and infrared sensors, as different types of military vehicles are driven in their proximity. This dataset was developed as part of a project supported by DARPA for the task of vehicle detection and type classification.
- **Affiliation of creators:** University of Wisconsin-Madison.
- **Domain:** signal processing.
- **Tasks in fairness literature:** fair federated learning [314].
- **Data spec:** time series.
- **Sample size:** unknown.
- **Year:** 2013.
- **Sensitive features:** none.
- **Link:** <http://www.ecs.umass.edu/mduarte/Software.html>
- **Further info:** Duarte and Hu [145]

A.209 Victorian Era Authorship Attribution

- **Description:** this resource was developed to benchmark different authorship attribution techniques. Querying the Gdelt database, the creators focus on English language authors from the 19th century with at least five books available. The corpus was split into text fragments of 1,000 words each. Only the most frequent 10,000 words were kept, while the remaining ones were removed.
- **Affiliation of creators:** Purdue University.
- **Domain:** literature.
- **Tasks in fairness literature:** fair clustering [214].
- **Data spec:** text.
- **Sample size:** ~ 100K text fragments.
- **Year:** 2018.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <http://archive.ics.uci.edu/ml/datasets/Victorian+Era+Authorship+Attribution>
- **Further info:** Gungor [206]

A.210 Visual Question Answering (VQA)

- **Description:** this dataset is curated as a benchmark for open-ended visual question answering. The collection features both real images from MS-COCO [319] and abstract scenes with human figures. Questions and answers were compiled by workers on Mechanical Turk who were instructed to formulate questions that require seeing the associated image for a correct answer.
- **Affiliation of creators:** Georgia Institute of Technology; Carnegie Mellon University; Army Research Lab; Facebook AI Research.
- **Domain:** computer vision.
- **Tasks in fairness literature:** bias discovery [342].
- **Data spec:** mixture (image, text).
- **Sample size:** ~ 1M questions over ~ 300K images.
- **Year:** 2017.
- **Sensitive features:** visual and textual references to gender.
- **Link:** <https://visualqa.org/>
- **Further info:** Goyal et al. [199]

A.211 Warfarin

- **Description:** this dataset was collected as part of a study about algorithmic estimation of optimal warfarin dosage as an oral anticoagulation treatment. The study was carried out by the International Warfarin Pharmacogenetics Consortium, comprising 21 research groups from 9 countries and 4 continents. The dataset was co-curated by staff at the Pharmacogenomics Knowledge Base (PharmGKB) including, for thousands of patients at centers around the world, their demographics, comorbidities, other medications and genetic factors, along with the steady-state dose of warfarin that led to stable levels of anticoagulation without adverse events.
- **Affiliation of creators:** PharmGKB; International Warfarin Pharmacogenetics Consortium.
- **Domain:** pharmacology.
- **Tasks in fairness literature:** fairness evaluation under unawareness [256].
- **Data spec:** tabular data.

- **Sample size:** ~ 6K patients.
- **Year:** 2009.
- **Sensitive features:** sex, ethnicity, age.
- **Link:** <https://www.pharmgkb.org/downloads>
- **Further info:** International Warfarin Pharmacogenetics Consortium [243]

A.212 Waterbirds

- **Description:** this computer vision dataset consists of photos where subjects and backgrounds are carefully paired to induce spurious correlations. Subjects are birds, taken from the CUB dataset [522], divided into waterbirds and landbirds. Pixel-level segmentation masks are exploited to cut out subjects and paste them onto land or water backgrounds from the Places dataset [588]. While in the provided validation and test splits both landbirds and waterbirds appear with the same frequency on either background, the training split is imbalanced so that 95% of all waterbirds are placed against a water background and 95% of all landbirds are depicted against a land background.
- **Affiliation of creators:** Stanford University; Microsoft.
- **Domain:** computer vision.
- **Tasks in fairness literature:** fairness evaluation of selective classification [253].
- **Data spec:** image.
- **Sample size:** ~ 10K images.
- **Year:** 2021.
- **Sensitive features:** none.
- **Link:** <https://github.com/ejones313/worst-group-sc/tree/main/src/data>
- **Further info:** Sagawa et al. [445]

A.213 WebText

- **Description:** this resource is a web scrape collected to train the GPT-2 language model. The authors considered all outbound links from Reddit which collected at least 3 *karma*. This inclusion criterion signals that the link received some upvotes by redditors and is treated as a quality heuristic for the webpage. To extract text data from each link, a combination of Dagnet [401] and Newspaper²² extractors was exploited. The curators performed deduplication and removed all Wikipedia pages to reduce text overlap with Wikipedia-based datasets.
- **Affiliation of creators:** OpenAI.
- **Domain:** linguistics.
- **Tasks in fairness literature:** data bias evaluation [487].
- **Data spec:** text.
- **Sample size:** ~ 8M documents.
- **Year:** 2019.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <https://github.com/openai/gpt-2-output-dataset> (partial)
- **Further info:** Radford et al. [416]

A.214 Wholesale

- **Description:** this dataset represents Portuguese businesses from the catering industry purchasing goods from the same wholesaler. The businesses are located in Lisbon, Oporto, and a third undisclosed area; 298 are from the Horeca (Hotel/Restaurant/Café) channel and 142 from the Retail channel. Each data point comprises this information along with yearly expenditures on different categories of products (e.g. milk, frozen goods, delicatessen). Collection of this data was presumably carried out by the wholesaler in a business intelligence initiative primarily aimed at customer segmentation and targeted marketing.
- **Affiliation of creators:** Université Pierre et Marie Curie; University Institute of Lisbon; INRIA.
- **Domain:** marketing.
- **Tasks in fairness literature:** fair data summarization [254].
- **Data spec:** tabular data.
- **Sample size:** ~ 400 businesses.
- **Year:** 2014.
- **Sensitive features:** geography.
- **Link:** <http://archive.ics.uci.edu/ml/datasets/wholesale+customers>
- **Further info:** Baudry et al. [34]

²²<https://github.com/codelucas/newspaper>

A.215 Wikidata

- **Description:** founded in 2012, Wikidata is a free, collaborative, multilingual knowledge base, maintained by editors and partly automated. It consists of items linked by properties. The most common items include humans, administrative territorial entities, architectural structures, chemical compounds, films, and scholarly articles.
- **Affiliation of creators:** Wikimedia Foundation.
- **Domain:** information systems.
- **Tasks in fairness literature:** fairness evaluation in graph mining [168].
- **Data spec:** item-property-value triples.
- **Sample size:** ~ 90M items.
- **Year:** present.
- **Sensitive features:** demographics of people featured in entities (age, sex, geography) and their relations.
- **Link:** https://www.wikidata.org/wiki/Wikidata:Data_access
- **Further info:** https://www.wikidata.org/wiki/Wikidata:Main_Page

A.216 Wikipedia dumps

- **Description:** Wikipedia dumps are maintained and updated regularly by the Wikimedia Foundation. Typically, they contain every article available in a language at a given time. As a large source of curated text, they have often been used by the natural language processing and computational linguistics communities to extract models of human language. We find usage of German, English, Mandarin Chinese, Spanish, Arabic, French, Farsi, Urdu, and Wolof dumps in the surveyed articles.
- **Affiliation of creators:** Wikimedia Foundation.
- **Domain:** linguistics.
- **Tasks in fairness literature:** bias evaluation in WEs [62, 96, 318, 393].
- **Data spec:** text.
- **Sample size:** ~ 6M articles (EN), ~ 3M articles (DE) as of May 2021.
- **Year:** present.
- **Sensitive features:** textual references to people and their demographics.
- **Link:** <https://dumps.wikimedia.org/enwiki/>; <https://dumps.wikimedia.org/dewiki/>
- **Further info:** https://meta.wikimedia.org/wiki/Data_dumps

A.217 Wikipedia Toxic Comments

- **Description:** this dataset was developed as a resource to analyze discourse and personal attacks on Wikipedia talk pages, which are used by editors to discuss improvements. It is aimed at using ML for better online conversations and flag posts that are likely to make other participants leave. The data consists of Wikipedia comments labelled by 5,000 crowd-workers according to their toxicity level (toxic, severe_toxic) and type (obscene, threat, insult, identity_hate). This resource powers a public Kaggle competition.
- **Affiliation of creators:** Wikimedia foundation; Google.
- **Domain:** social media.
- **Tasks in fairness literature:** fair classification, [179, 453], fairness evaluation [142].
- **Data spec:** text.
- **Sample size:** ~ 160K comments.
- **Year:** 2017.
- **Sensitive features:** textual reference to people and their demographics.
- **Link:** <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge>
- **Further info:** <https://www.perspectiveapi.com/research/>

A.218 Willingness-to-Pay for Vaccine

- **Description:** this dataset resulted from a study of willingness to pay for a vaccine against tick-borne encephalitis in Sweden. Thousands of citizens from different areas of the country filled in a survey about exposure, risk perception, knowledge, and protective behavior related to ticks and tick-borne diseases, along with socioeconomic information. The central question of the survey asks how much respondents would be willing to pay for a vaccine that provides a three-year protection against tick-borne encephalitis.
- **Affiliation of creators:** University of Gothenburg.
- **Domain:** public health.
- **Tasks in fairness literature:** fair pricing evaluation [260].
- **Data spec:** tabular data.
- **Sample size:** ~ 2K respondents.
- **Year:** 2015.
- **Sensitive features:** age, gender, geography.

- **Link:** <https://snd.gu.se/sv/catalogue/study/snd0987/1#dataset>
- **Further info:** Slunge [469]

A.219 Winobias

- **Description:** similarly to Winogender, this benchmark was built to study coreference resolution and gender bias, focusing on words that relate to professions with diverse gender representation. Example: “The physician hired the secretary because he (she) was overwhelmed with clients”. The correct pronoun resolution is clear from the syntax or semantics of the sentence and can be either stereotypical or counter-stereotypical. The accuracy of biased coreference resolution systems will vary accordingly.
- **Affiliation of creators:** University of California Los Angeles; University of Virginia; Allen Institute for Artificial Intelligence.
- **Domain:** linguistics.
- **Tasks in fairness literature:** fair entity resolution evaluation. [516].
- **Data spec:** text.
- **Sample size:** ~ 3K sentences.
- **Year:** 2020.
- **Sensitive features:** gender.
- **Link:** <https://github.com/uclanlp/corefBias/tree/master/WinoBias/wino>
- **Further info:** Zhao et al. [584]

A.220 Winogender

- **Description:** this dataset was crafted to systematically study gender bias in systems for coreference resolution, the task of resolving whom pronouns refer to in a sentence. This resource follows the Winograd schemas, with sentence templates mentioning a profession (nurse), a participant (patient), and a pronoun referring to either one of them: “The nurse notified the patient that her/his/their shift would be ending in an hour.” Sentence templates have been crafted so that the pronoun resolution can be done unambiguously based on contextual information, hence unbiased systems should display similar error rates, regardless of gender concentrations in different professions. The ground truth for each sentence has been validated by workers on Mechanical Turk with accuracy over 99%.
- **Affiliation of creators:** Johns Hopkins University.
- **Domain:** linguistics.
- **Tasks in fairness literature:** fair entity resolution evaluation [516], fairness evaluation in entity recognition [368].
- **Data spec:** text.
- **Sample size:** ~ 700 sentences.
- **Year:** 2018.
- **Sensitive features:** gender.
- **Link:** <https://github.com/rudinger/winogender-schemas>
- **Further info:** Rudinger et al. [440]
- **Variants:** Winogender-NER [368] is a modified version of the template appropriate for named entity recognition.

A.221 Word Embedding Association Test (WEAT)

- **Description:** this resource was created to audit biases in English WEs. Following the Implicit Association Test (IAT) from social psychology [202], this dataset defines two groups of target words, relating e.g. to flowers and insects, and two groups of attribute words, relating e.g. to pleasantness and unpleasantness. The dataset can be used to measure biased associations between the target words and the attribute words represented by a set of WEs. WEAT comprises ten tests across different word categories. The most salient for the purposes of algorithmic fairness support tests of associations between race and pleasantness, age and pleasantness, gender and career (vs family), gender and propensity to math (vs arts). Race-related words are first names predominantly associated with African American or European American individuals. Gender is encoded in a similar fashion, or with intrinsically gendered words (e.g. mother).
- **Affiliation of creators:** Princeton University; University of Bath.
- **Domain:** linguistics.
- **Tasks in fairness literature:** bias evaluation in WEs [62, 208].
- **Data spec:** text.
- **Sample size:** ~10 groups of words, with ~10-60 words in each group.
- **Year:** 2017.
- **Sensitive features:** race, gender.
- **Link:** <https://arxiv.org/pdf/1608.07187.pdf>
- **Further info:** Caliskan et al. [71]

A.222 Yahoo! A1 Search Marketing

- **Description:** this dataset contains bids from all advertisers who participated in Yahoo! Search Marketing auctions for the top 1000 search queries from June 15, 2002, to June 14, 2003. The identities of advertisers and the queries they target are anonymized for confidentiality reasons.
- **Affiliation of creators:** Yahoo! Labs.
- **Domain:** marketing.
- **Tasks in fairness literature:** fair advertising [79, 380].
- **Data spec:** advertiser-keyword pairs.
- **Sample size:** ~ 20M bids by ~ 10K advertisers over ~ 1K search queries.
- **Year:** after 2003.
- **Sensitive features:** none.
- **Link:** <https://webscope.sandbox.yahoo.com/catalog.php?datatype=a>
- **Further info:**

A.223 Yahoo! c14B Learning to Rank

- **Description:** this resource consists of 2 datasets which encode the interactions of Yahoo! users with the search engine in the US and an unknown Asian country. This data is a subset of the entire training set used internally to train the ranking functions of the Yahoo! search engine. Textual features are deliberately obfuscated and the final data consists of numerical features which encode query-document pairs. Query-document pairs are assigned multigraded relevance judgements by a professional editor.
- **Affiliation of creators:** Yahoo! Labs.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair ranking [463].
- **Data spec:** query-document pairs.
- **Sample size:** ~ 40K queries, ~ 900K documents.
- **Year:** 2011.
- **Sensitive features:** none.
- **Link:** <https://webscope.sandbox.yahoo.com/catalog.php?datatype=c>
- **Further info:** Chapelle and Chang [88]

A.224 YouTube Dialect Accuracy

- **Description:** this dataset was curated to audit the accuracy of YouTube’s automated captioning system across two genders and five dialects of English. Eighty speakers were sampled from videos matching the query “accent challenge <region>” or “accent tag <region>”, where <region> is one of five areas selected for geographic separation and distinct local dialects: California, Georgia, New England, New Zealand and Scotland. This curation choice targets a popular internet phenomenon (called “accent tag”, “dialect meme” or “accent challenge”) consisting of videos of people from different areas presenting themselves and their linguistic background, subsequently reading a list of words designed to elicit pronunciation differences dependent on dialect. This resource focuses only on the word portion of these videos, with a “phonetically-trained listener familiar with the dialects” performing the annotation for word caption accuracy.
- **Affiliation of creators:** University of Washington.
- **Domain:** social media.
- **Tasks in fairness literature:** fairness evaluation of speech recognition [491].
- **Data spec:** tabular data.
- **Sample size:** ~ 100 speakers.
- **Year:** 2016.
- **Sensitive features:** gender, geography.
- **Link:** <https://github.com/rctatman/youtubeDialectAccuracy>
- **Further info:** Tatman [491]

A.225 Yow news

- **Description:** this dataset was collected to support research on personalized information integration and retrieval. The data, consisting of implicit and explicit user feedback stored in interaction logs, was gathered in a user study via a special browser accessing a web-based news story filtering system. The task associated with this resource is personalized news recommendation.
- **Affiliation of creators:** Carnegie Mellon University.
- **Domain:** news, information systems.
- **Tasks in fairness literature:** fair ranking [462].
- **Data spec:** user-story pairs.

- **Sample size:** ~ 10K interaction logs.
- **Year:** 2009.
- **Sensitive features:** news provider.
- **Link:** <https://users.soe.ucsc.edu/~yiz/papers/data/YOWStudy/>
- **Further info:** Zhang [573]; <https://users.soe.ucsc.edu/~yiz/piir/>

A.226 Zillow Searches

- **Description:** this is a proprietary dataset from Zillow, a famous real estate marketplace. It consists of a random sample of over 13,000 search sessions covering more than 36,000 property listings. Each listing consists of several features, some of which are considered salient by the creators and a sensible target for fair ranking algorithms. Among these are the ownership of the house (Zillow, independent realtor, new construction listed by builders) and the availability of 3D/video tours of the property. This dataset was collected internally to study the problem of fair recommendation and ranking on Zillow data.
- **Affiliation of creators:** Boston University; Zillow Group.
- **Domain:** information systems.
- **Tasks in fairness literature:** fair ranking [89].
- **Data spec:** unknown.
- **Sample size:** ~ 10K search sessions featuring ~ 40K property listings.
- **Year:** 2020.
- **Sensitive features:** ownership, tour availability.
- **Link:** not available
- **Further info:** Chaudhari et al. [89]

B ADULT

Key references include Cohany et al. [112], Ding et al. [141], Kohavi [285], McKenna [353, 354], UCI Machine Learning Repository [502], U.S. Dept. of Commerce Bureau of the Census [505].

B.1 Datasheet

B.1.1 Motivation.

- **For what purpose was the dataset created?**
The Adult dataset was created as a resource to benchmark the performance of machine learning algorithms. Rather than powering a specific task or application, the dataset was likely chosen as a real-world source of socially relevant data [285].
- **Who created the dataset?**
Barry Becker extracted this dataset from the 1994 Census database. Ronny Kohavi and Barry Becker donated it to UCI Machine Learning Repository in 1996. At that time, both were working for Silicon Graphics Inc [502]
- **Who funded the creation of the dataset?**
The underlying database is a product of the Current Population Survey (CPS) of March 1994, a joint effort by the US Census Bureau and the US Bureau of Labor Statistics (BLS), funded by the US federal government. The extraction of Adult from the larger database was plausibly part of work remunerated by Silicon Graphics.

B.1.2 Composition.

- **What do the instances that comprise the dataset represent?**
Each instance is a **March 1994 CPS respondent**, represented along demographic and socio-economic dimensions.
- **How many instances are there in total?**
The dataset consists of **48,842 instances**.
- **Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?**
Adult contains individuals from a **sample** of US households, extracted from the 1994 Annual Social and Economic Supplement (ASEC) of the CPS with the following query:

$$(AAGE > 16) \&\& (AGI > 100) \&\& (AFNLWGT > 1) \&\& (HRSWK > 0).$$

This means Adult focuses on a subset of ASEC respondents aged 17 or older, whose income is above \$100, working at least 1 hour per week. While these were conceived as conditions to filter out noisy records [502], they may introduce sampling effects. Moreover, the 1994 CPS data was itself a sample, selected according to Census Bureau best practices, reaching over 70,000 households in nearly 2,000 US counties. The March 1994 CPS sample aimed at obtaining more reliable information on the Hispanic population, and was hence extended to an additional 2,500 eligible housing units.

- **What data does each instance consist of?**

Each instance consists of a combination of nominal, ordinal and continuous attributes, denominated age, workclass, fnlwgt, education, education-num, marital-status, occupation, relationship, race, sex, capital-gain, capital-loss, hours-per-week, native-country. See Table 7 for a detailed explanation of features and their values.

- **Is there a label or target associated with each instance?**
Yes. Each person instance comes with a binary label encoding whether their income is above a 50,000 threshold.
- **Is any information missing from individual instances?**
Yes. Over 7% of the instances have missing values. This is likely due to issues with data recording and coding or respondents' inability to recall information.
- **Are relationships between individual instances made explicit e.g., users' movie ratings, social network links?**
No. Some instances are related persons from the same household [505] but this information is not reported in the dataset.
- **Are there recommended data splits?**
Yes. The dataset comes with a specified train/test split made using MLC++ GenCVFiles, resulting in a 2/3–1/3 random split [502]. The training set consists of 32561 instances, the test set of 16281 instances.
- **Are there any errors, sources of noise, or redundancies in the dataset?**
Yes. Sources of error include definitional difficulties, differences in interpretation of questions, respondents inability or unwillingness to provide correct information, errors made during data collection, data processing or missing value imputation. The tendency in household surveys for respondents to under-report their income was an explicit concern. Finally, noise infusion such as topcoding (saturation to \$99,999) was applied to avoid re-identification of certain individuals [505].
- **Is the dataset self-contained, or does it link to or otherwise rely on external resources?**
The dataset is **self-contained**.
- **Does the dataset contain data that might be considered confidential?**
Yes. The data is protected by Title 13 of the United States Code, protecting individuals against identification from Census data.²³
- **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?**
No, not strictly. Interpreting the question more broadly, however, the envisioned racial and sexual categories may be deemed inadequate.
- **Does the dataset identify any subpopulations (e.g., by age, gender)?**
Yes. The dataset provides information on sex, age and race of respondents. These were self-reported, although self-identification was bounded by envisioned categories. These are (female, male) for sex and (White, Black, American Indian/Aleut Eskimo, Asian or Pacific Islander, Other) for race. Table 4 summarizes the marginal distribution of the Adult dataset across these subpopulations.
- **Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset?**
Unknown. Important variables for data re-identification, such as birth date or ZIP code, are absent from the Adult dataset. However, instances in this dataset may be linked to the original CPS 1994 data [141]. Moreover, re-identification studies internal to the Census Bureau pointed to combinations of variables that could potentially be used to re-identify respondents from Census microdata [354].
- **Does the dataset contain data that might be considered sensitive in any way?**
Yes. This dataset contains sensitive data, such as sex, race, native country and financial situation of respondents.
- **Any other comments?**
A precise definition for the variable called fnlwgt is unknown. It was used by Census Bureau statisticians to obtain population-level estimates from the CPS sample. For this reason, its use in classification tasks would be unusual.

B.1.3 Collection process.

- **How was the data associated with each instance acquired?**
Trained interviewers asked questions directly to respondents [505]. The data was made available through US Census data products which were used by Barry Becker to extract the Adult dataset.
- **What mechanisms or procedures were used to collect the data?**
Interviewers conducted the survey either in person at the respondent's home or by phone. They used laptop computers with ad-hoc software to prompt questions and record answers. At the end of each day, interviewers transmitted the collected data via modem to the Bureau headquarters [505].
- **If the dataset is a sample from a larger set, what was the sampling strategy?**
A probabilistic sample was selected according to US Census Bureau best practice, with a multi-stage stratified design. The US territory was divided into strata, from which one county (or group of counties) was selected. From each selected county a sample of addresses was later obtained and added to the sample [504]. Barry Becker extracted a "set of reasonably clean records" using the following conditions:

$$(AAGE > 16) \&\& (AGI > 100) \&\& (AFNLWGT > 1) \&\& (HRSWK > 0).$$

²³https://www.census.gov/about/policies/privacy/data_stewardship/title_13_-_protection_of_confidential_information.html

Demographic Characteristic	Values
Percentage of male subjects	66.85%
Percentage of female subjects	33.15%
Percentage of White subjects	85.50%
Percentage of Black subjects	9.60%
Percentage of Asian-Pac-Islander subjects	3.11%
Percentage of Amer-Indian-Eskimo subjects	0.96%
Percentage of people belonging to other races	0.83%
Percentage of people between 16-19 years old	5.14%
Percentage of people between 20-29 years old	24.58%
Percentage of people between 30-39 years old	26.47%
Percentage of people between 40-49 years old	21.95%
Percentage of people between 50-59 years old	13.55%
Percentage of people between 60-69 years old	6.25%
Percentage of people between 70-79 years old	1.67%
Percentage of people between 80-89 years old	0.27%
Percentage of people between 90-99 years old	0.11%

Table 4: Demographic Characteristics of the Adult dataset.

- **Who was involved in the data collection process and how were they compensated?**
Interviewers trained by the US Census Bureau were involved in the data collection process. Data extraction was later performed by Barry Becker while affiliated with Silicon Graphics. Their compensation is unknown.
- **Over what timeframe was the data collected?**
Respondents were interviewed in March 1994, while the Adult dataset was donated to UCI ML Repository in May 1996.
- **Were any ethical review processes conducted?**
The Microdata Review Panel likely reviewed this data for compliance with Title 13 [354] and authorized its publication.
- **Was the data collected from the individuals in question directly, or obtain it via third parties or other sources?**
Directly. US Census Bureau interviewers collected the data through interviews, conducted in person or over the phone. Danny Kohavi and Barry Becker later processed this data, obtaining it from the Census Bureau website.
- **Were the individuals in question notified about the data collection?**
Yes. Individuals knew they were part of a sample chosen by the Census Bureau chosen for statistical analysis. They were not notified about their data being included in the Adult dataset.
- **Did the individuals in question consent to the collection and use of their data?**
Yes. For the CPS, participation is voluntary. A recent version of the information provided to respondents before interviews is available on the US Census Website.²⁴
- **If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses?**
Unknown.
- **Has an analysis of the potential impact of the dataset and its use on data subjects been conducted?**
Yes. Re-identification studies have been conducted both internally [354] and externally [431] on Census Bureau data. McKenna [354] mention finding combinations of variables on Census files that can lead to successful re-identification, which were subsequently removed or protected with noise injection. Rocher et al. [431] demonstrate on the Adult dataset that the likelihood of a specific individual to have been correctly re-identified can be estimated with high accuracy. We are unaware of studies about the potential impact of successful re-identification on respondents.

B.1.4 Preprocessing/cleaning/labelling.

- **Was any preprocessing/cleaning/labeling of the data done?**

²⁴https://www2.census.gov/programs-surveys/cps/advance_letter.pdf

Yes. Preprocessing operations by the Census Bureau include missing value imputation and topcoding. Furthermore, Barry Becker and Ron Kohavi binarized the income variable (> \$50K) and discarded several CPS respondents who are not included in the Adult dataset.

- **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data?**

Unknown.

- **Is the software used to preprocess/clean/label the instances available?**

Likely no. It seems unlikely for the code to be available 25 years after its last known use.

B.1.5 *Uses.*

- **For what tasks has the dataset been used?**

This dataset probably owes its status in the ML community to an early position of publicly-available and interesting resource based on real-world data. For this reason, rather than powering specific applications, Adult is used as a benchmark for classifiers in many fields of machine learning. Due to its encoding of sensitive attributes, it has also become the most used dataset in the fair ML literature.

- **Is there a repository that links to any or all papers or systems that use the dataset?**

Yes. A selection of early works (pre-2005) using this dataset can be found in UCI Machine Learning Repository [502]. A more recent list is available under the beta version of the UCI ML Repository.²⁵ See Appendix A.7 for a (non-exhaustive) list of algorithmic fairness works using this resource.

- **What (other) tasks could the dataset be used for?**

The Adult dataset is used in tasks where data of social significance is deemed important, for example privacy-preserving ML.

- **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?**

Yes. The threshold used to quantize income for a binary classification task is very high (\$50K). As a result a trivial rejector achieves very large accuracy on the black subpopulation (93%). For the same reason, models are often more accurate for the female subpopulation than for the male one [141]. Some numerical results on Adult may be an artifact of this threshold choice.

- **Are there tasks for which the dataset should not be used?**

Based on the previous answer, we caution against drawing overarching conclusions based on experimental results obtained on this dataset alone.

B.1.6 *Distribution.*

- **Is the dataset distributed to third parties outside of the entity on behalf of which the dataset was created?**

Yes. The dataset is publicly available [502].

- **How is the dataset distributed?**

The dataset is available as a **csv file**.

- **When was the dataset distributed?**

The dataset was released on the UCI ML Repository in **May 1996**.

- **Is the dataset distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?**

Yes. The UCI ML repository has a citation policy. Terms of Use concerning the privacy of CPS respondents are likely to apply.

- **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?**

Likely no. We are unaware of any IP-based restrictions.

- **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?**

Likely no.

B.1.7 *Maintenance.*

- **Who is supporting/hosting/maintaining the dataset?**

The dataset is hosted and maintained by the **UCI Machine Learning Repository** [502].

- **How can the owner/curator/manager of the dataset be contacted?**

Comments and inquiries may be directed at ml-repository@ics.uci.edu. Ronny Kohavi is the primary contact for this specific resource, available at ronnyk@live.com.

- **Is there an erratum?**

Likely no. We are unaware of any erratum.

- **Will the dataset be updated?**

A superset of the dataset without quantization of the target income variable is available [141].

- **If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances?**

Unknown.

- **Will older versions of the dataset continue to be supported/hosted/maintained?**

Unless otherwise indicated, the Adult dataset will remain hosted on the UCI ML Repository in its current version.

²⁵<https://archive-beta.ics.uci.edu/ml/datasets/2>

- **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?**
Unknown.

B.2 Data Nutrition Label

METADATA	
Filenames	adult
Format	csv
Url	https://archive.ics.uci.edu/ml/datasets/adult
Domain	Economics
Keywords	US census, income
Type	Tabular
Rows	48842
Columns	14
% of missing cells	0.9%
Rows with missing cells	7%
License	UCI Repository citation policy
Released	May 1996
Range	1994
Description	A benchmark for classifiers tasked with predicting whether individual income exceeds \$50K/yr based on demographic and socio-economic information. Also known as "Census Income" dataset.

Table 5: Metadata of the Adult dataset

PROVENANCE	
Source	
Name	U. S. Census Bureau
Url	https://www.census.gov/en.html
email	//
Authors	
Names	Ronny Kohavi and Barry Becker
Url	https://archive.ics.uci.edu/ml/datasets/
email	ronnyk@live.com

Table 6: Provenance of the Adult dataset

VARIABLES	
age	Respondent's age.
workclass	Broad classification of employment, with following envisioned classes. Private Self-emp-not-inc (Self employed not-incorporated) Self-emp-inc (Self employed incorporated) Federal-gov Local-gov State-gov Without-pay (Without pay in family business) Never-worked
fnlwtg	Variable used to produce population estimates from the CPS sample.
education	Educational attainment of respondent. Preschool 1st-4th 5th-6th 7th-8th 9th 10th 11th 12th (no diploma) HS-grad (High school graduation) Some-college (no degree) Assoc-voc (associate degree in college, vocation program) Assoc-acdm (associate degree in college, academic program) Bachelors Masters Prof-school (professional school) Doctorate
education-num	Ordinal encoding of previous variable.

Table 7: Variables of the Adult dataset (1/3).

VARIABLES	
marital-status	Respondent's marital status, with following envisioned classes. Married-civ-spouse (married, civilian spouse present) Divorced Never-married Separated Widowed Married-spouse-absent Married-AF-spouse (married, armed force spouse)
occupation	Job of respondent. Tech-support (Technical, sales, and administrative support) Craft-repair (Precision production, craft, and repair) Other-service Sales Exec-managerial (Managerial and professional speciality) Prof-specialty (Professional speciality) Handlers-cleaners (Handlers, equipment cleaners, helpers, and laborers) Machine-op-inspct (Operators, fabricators, and laborers) Adm-clerical (Administrative support occupations, including clerical) Farming-fishing (Farming, forestry, and fishing) Transport-moving (Transportation and material moving) Priv-house-serv (Private household service, e.g. cooks, cleaners) Protective-serv (Protective service, e.g. firefighters, police) Armed-Forces
relationship	Familial role within household. Wife Own-child Husband Not-in-family Other-relative Unmarried

Table 8: Variables of the Adult dataset (2/3).

VARIABLES	
race	Respondent's race. Amer-Indian-Eskimo Asian-Pac-Islander Black White Other
sex	Respondent's sex. Female Male
capital-gain	Profits from sale of assets.
capital-loss	Losses from sale of assets.
hours-per-week	Average hours of work per week.
native-country	Native Country of respondent
target variable	Does respondent's income exceed \$50,000?

Table 9: Variables of the Adult dataset (3/3).

STATISTICS						
Ordinal						
name	type	count	uniqueEntries	mostFrequent	leastFrequent	missing
education-num	int	48842	16	9	1	0

Table 10: Ordinal variables statistics of the Adult dataset

Categorical						
name	type	count	uniqueEntries	mostFrequent	leastFrequent	missing
workclass	string	48842	8	Private	Never-worked	2799
education	string	48842	16	HS-grad	Preschool	0
marital-status	string	48842	7	Married-civ-spouse	Married-AF-spouse	0
occupation	string	48842	14	Prof-specialty	Armed-Forces	2809
relationship	string	48842	6	Husband	Other-relative	0
race	string	48842	5	White	Other	0
sex	string	48842	2	Male	Female	0
native-country	string	48842	41	United-States	Holland-Netherlands	857
target variable	string	48842	2	<= 50K	> 50K	0

Table 11: Categorical variables statistics of the Adult dataset

Quantitative									
name	type	count	min	median	max	mean	stdDev	miss	zeros
age	int	48842	17	37	90	38.64	13.71	0	0
fnlwgt	int	48842	12285	178144.5	1490400	189664.13	105604.03	0	0
capital-gain	int	48842	0	0	99999	1079.07	7452.02	0	44807
capital-loss	int	48842	0	0	4356	87.50	403	0	46560
hours-per-week	int	48842	1	40	99	40.42	12.39	0	0

Table 12: Quantitative variables statistics of the Adult dataset.

C COMPAS

Key references include Angwin et al. [15], Bao et al. [28], Barenstein [31], Brennan et al. [56], Dieterich et al. [140], Equivant [156], Larson et al. [301], ProPublica [408].

C.1 Datasheet

C.1.1 Motivation.

- **For what purpose was the dataset created?**

This dataset was created for an external audit of racial biases in the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) risk assessment tool developed by Northpointe (now Equivant), which estimates the likelihood of a defendant becoming a recidivist.

- **Who created the dataset and on behalf of which entity?**

The dataset was created by Julia Angwin (senior reporter), Jeff Larson (data editor), Surya Mattu (contributing researcher), Lauren Kirchner (senior reporting fellow). All four contributors were affiliated with ProPublica at the time.

- **Who funded the creation of the dataset?**

The dataset curation work was likely remunerated by ProPublica.

C.1.2 Composition.

- **What do the instances that comprise the dataset represent?**

Each instance is a person that was scored for risk of recidivism by the COMPAS system in Broward County, Florida, between 2013–2014. In other words, instances are **defendants**.

- **How many instances are there in total?**

The COMPAS dataset [408] consists of **11,757** defendants assessed at the pretrial stage (`compas-scores.csv`). A separate dataset is released for a subset of 7,214 defendants that were observed for two years after screening (`compas-scores-two-years.csv`). Finally a smaller subset of 4,743 defendants focuses on violent recidivism (`compas-scores-two-years-violent.csv`).

- **Does the dataset contain all possible instances or is it a sample of instances from a larger set?**

The dataset represents a **convenience sample** of all individuals that were scored by the COMPAS tool. It concentrates on defendants in Broward County, as it is a large jurisdiction in a state with strong open-records laws [301]. Moreover, due to Broward County using COMPAS primarily in release/detain decisions prior to a defendant's trial, scores assessed at parole, probation or other stages were discarded. A notable anomaly in the sample is the low amount of defendants screened between June and July 2013 compared to the remaining time span of the COMPAS dataset [31].

- **What data does each instance consist of?**

Instances represent Broward County defendants scored with COMPAS for risk of recidivism. For each defendant the data provided by ProPublica includes tens of variables (~ 50) summarizing their demographics, criminal record, custody and COMPAS scores.

- **Is there a label or target associated with each instance?**

Yes. Instances are associated with two target variables (`is_recid` and `is_violent_recid`), indicating whether defendants were booked in jail with a criminal offense (potentially violent) that took place after their COMPAS screening but within two years. The definition of recidivism and the two-year cutoff were selected by ProPublica staff to align their audit with definitions by Northpointe [15, 56].

- **Is any information missing from individual instances?**

Yes. There are several columns where data is missing for one or more instances, including dates when defendants committed the offense (`c_offense_date`) were incarcerated (`c_jail_in`) or released (`c_jail_out`). Missingness in this dataset is not surprising as its curation was a complex endeavour that required cross-referencing information from three separate sources, namely Broward County Sheriff's Office, Broward County Clerk's Office and Florida Department of Corrections. Moreover, Northpointe's response to the ProPublica's study points out important risk factors considered by the COMPAS algorithm that are not present in the dataset, among which the criminal involvement scale, drug problems sub-scale, age at first adjudication, arrest rate and vocational educational scale [140]. Finally, a clear indication of whether defendants were released or detained pretrial seems to be missing.

- **Are relationships between individual instances made explicit?**

No. While it is plausible for some Broward County defendants to be connected, this information is not available.

- **Are there recommended data splits?**

No.

- **Are there any errors, sources of noise, or redundancies in the dataset?**

Yes. Clerical errors in records caused incorrect matches between individuals' COMPAS scores and their criminal records, leading to an error rate close to 4% [301]. Moreover, an important temporal trend was spuriously introduced by ProPublica's preprocessing in `compas-scores-two-years.csv` and `compas-scores-two-years-violent.csv`, due to which defendants with a screening date after April 2014 are all recidivists [31]. In terms of redundancies, `compas-scores.csv` contains two identical columns (called `decile_score` and `decile_score.1`).

- **Is the dataset self-contained, or does it link to or otherwise rely on external resources?**

The dataset is **self-contained**.

- **Does the dataset contain data that might be considered confidential?**

No. However it does contain first names and last names of defendants, connecting them to their criminal history.

- **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?**

Yes. The column `vr_charge_desc` describing violent recidivism charges is one such example.

- **Does the dataset identify any subpopulations (e.g., by age, gender)?**

Yes. The dataset identifies population by age, sex and race. The curators of the COMPAS dataset maintained the race classifications used by the Broward County Sheriff's Office, identifying individuals as Asian, Black, Hispanic, Native American and White [301]. Age is reported as an integer, sex as either Male or Female. A distribution along these dimensions is reported in Table 13 which summarizes data in `compas-scores-two-years.csv`. Distributions in remaining files are similar.

compas-scores-two-years	
Demographic Characteristic	Values
Percentage of male subjects	80.83%
Percentage of female subjects	19.17%
Percentage of African-American subjects	51.46%
Percentage of Caucasian subjects	33.63%
Percentage of Hispanic subjects	8.67%
Percentage of Asian subjects	0.48%
Percentage of Native American subjects	0.20%
Percentage of people belonging to other races	5.56%
Percentage of people under-19 years old	0.42%
Percentage of people between 20-29 years old	42.41%
Percentage of people between 30-39 years old	28.04%
Percentage of people between 40-49 years old	14.60%
Percentage of people between 50-59 years old	11.00%
Percentage of people between 60-69 years old	3.01%
Percentage of people over-70 years old	0.51%

Table 13: Demographic Characteristics of compas-scores-two-years.

- **Is it possible to identify individuals, either directly or indirectly from the dataset?**

Yes. The dataset reports defendants' first name, last name and date of birth.

- **Does the dataset contain data that might be considered sensitive in any way?**

Yes. The COMPAS dataset reports individuals' race, criminal history, full name and date of birth.

C.1.3 Collection process.

- **How was the data associated with each instance acquired?**

The data was obtained cross-referencing three sources. From the Broward County Sheriff's Office in Florida, ProPublica obtained COMPAS scores associated with all 18,610 people scored in 2013 and 2014. Defendants' public criminal records were obtained from the Broward County Clerk's Office website matching them based on date of birth, first and last names. The dataset was augmented with jail records provided by the Broward County Sheriff's Office. Finally public incarceration records were downloaded from the Florida Department of Corrections website.

- **What mechanisms or procedures were used to collect the data?**

The original data was plausibly recorded by employees of the Broward County Sheriff's Office, Broward County Clerk's Office, and Florida Department of Corrections. The curators of the COMPAS dataset obtained records from the County Sheriff's Office through a public records request, while data from the County Clerk's Office and the Florida Department of Correction was downloaded from their official website, matching the methodology of a COMPAS validation study [301].

- **If the dataset is a sample from a larger set, what was the sampling strategy?**

In terms of auditing the COMPAS risk assessment tool, this dataset represents a **convenience sample**, focused on a single county and scoring period 2013–2014. Considering a single county in a state with strong open-records laws reduced the data cross-referencing overhead. Concentrating on recent scores predating the study by 2–3 years kept the study timely and permitted a measurement of recidivism aligned with the one by Northpointe. The fact that Northpointe’s response to the ProPublica study only contains minor criticism of the sample (concerning the definition of pretrial defendants [140]) may be interpreted as testimony to its overall quality. More broadly and beyond the COMPAS audit, arrest data as a proxy for crime brings about specific sampling effects, inevitably mediated by law enforcement practices [229, 547].

- **Who was involved in the data collection process and how were they compensated?**

The original data was plausibly recorded by Broward County and Florida Department of Corrections employees. On ProPublica’s side, we assume that key curation choices were made and implemented by four employees credited in the article [15] and accompanying technical report [301], namely Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner. Given the focus on arrest data, the Broward County law enforcement community is also important in the data sampling process.

- **Over what timeframe was the data collected?**

COMPAS scores are from 2013 and 2014, while jail records cover the period from January 2013 to April 2016. The dataset was first released by ProPublica in May 2016 [408].

- **Were any ethical review processes conducted?**

Unknown.

- **Was the data collected from the individuals in question directly, or obtained via third parties or other sources?**

The data was obtained **via third parties**, namely the Broward County Sheriff’s Office in Florida through a public records request, from the Broward County Clerk’s Office through the official website and through the Florida Department of Corrections through the official website. Collection from interested individuals would not have been viable.

- **Were the individuals in question notified about the data collection?**

Likely no. Most of the COMPAS data was publicly available and downloaded from the official websites of Broward County Clerk’s Office and the Florida Department of Corrections.

- **Did the individuals in question consent to the collection and use of their data?**

Likely no. Public availability of arrest/conviction records is associated with collateral consequences that typically damage subjects socially and financially [15, 403].

- **If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses?**

Likely no.

- **Has an analysis of the potential impact of the dataset and its use on data subjects been conducted?**

Likely no. We are unaware of analyses specifically focused on the COMPAS dataset. More broadly, public availability of criminal records is related to studies on the employability of offenders [200].

C.1.4 Preprocessing/cleaning/labelling.

- **Was any preprocessing/cleaning/labeling of the data done?**

Yes. Instances were discarded if assessed with COMPAS at parole, probation or other stages in the criminal justice system. This data is unavailable. Moreover, ProPublica published its datasets with accompanying preprocessing code which has become standard [408]. The standard preprocessing removes instances for which (1) arrest dates or charge dates are not within 30 days of the COMPAS assessment, (2) true recidivism cannot be decided, (3) charge degree is not defined as misdemeanor or felony, (4) the COMPAS score is not clearly defined. The remaining COMPAS scores were bucketed into low, medium and high risk.

- **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data?**

Yes. The data is available in the official ProPublica github repository [408]. This is an intermediate data artifact, already cross-referenced by ProPublica across three separate sources.

- **Is the software used to preprocess/clean/label the instances available?**

Yes. The standard preprocessing software can be found in the official ProPublica github repository [408]. The software used to cross-reference data from separate sources is not publicly available.

C.1.5 Uses.

- **For what tasks has the dataset been used?**

The creators used this dataset to audit the COMPAS tool for racial bias. In the literature it has also been used to evaluate the fairness and accuracy of different algorithms and, more broadly, to study definitions of algorithmic fairness.

- **Is there a repository that links to any or all papers or systems that use the dataset?**

See Appendix A.41 for a (non-exhaustive) list of algorithmic fairness works using this resource.

- **What (other) tasks could the dataset be used for?**

In terms of immediate applications, the dataset could be used to train novel recidivism risk assessment tools. From a methodological perspective, COMPAS may be used in high-stakes domains connected with decision-making about human subjects, including explainable and privacy-preserving ML.

- **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?**

From a very narrow perspective, the fact that all defendants with a screening date after April 2014 are recidivists introduces artificially inflated recidivism base rates [31], which would likely be inherited by tools trained on the COMPAS dataset. Moreover, the dataset contains no clear indication concerning pretrial detention or release of defendants. Therefore, researchers must come up with subjective criteria to label individuals as detained or released if they are interested in studying pretrial detention as an intervention deviating from a default course of action [367]. From a broader perspective, the data is likely influenced by historical biases in criminal justice, with differential impact on different communities [15, 229, 547]. Zooming out further, the use of automated risk assessment tools in pretrial decisions is the subject of controversial debate [29] which cannot be overlooked.

- **Are there tasks for which the dataset should not be used?**

Given the above considerations and the narrow geographical scope of the dataset, COMPAS should not be used to train and deploy risk assessment tools for the judicial system. In research settings, users should exercise care in selecting both rows and columns. Bao et al. [28] suggest avoiding the use of COMPAS to demonstrate novel approaches in algorithmic fairness, as considering data without proper context may bring to misleading conclusions which could misguidedly enter the broader debate on criminal justice.

C.1.6 *Distribution.*

- **Is the dataset distributed to third parties outside of the entity on behalf of which the dataset was created?**

Yes. The COMPAS dataset is publicly available.

- **How is the dataset distributed?**

The dataset is hosted on ProPublica's official github repository [408].

- **When was the dataset distributed?**

Since **May 2016**.

- **Is the dataset distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?**

As of June 2021 the COMPAS dataset is freely distributed under ProPublica's standard ToU [409]. The dataset cannot be republished in its entirety, it cannot be sold, and can only be used for publication if ProPublica's work is properly referenced.

- **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?**

Likely no.

- **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?**

Unknown.

C.1.7 *Maintenance.*

- **Who is supporting/hosting/maintaining the dataset?**

The dataset is currently hosted and maintained by **ProPublica** on github.

- **How can the owner/curator/manager of the dataset be contacted?**

The contact for ProPublica's data store is data.store@propublica.org.

- **Is there an erratum?**

No. There is no official erratum. An external report highlighting anomalies in the data is available [31].

- **Will the dataset be updated?**

Likely no. In the event of an update, ProPublica's data store ToU specifies users are solely responsible for checking their sites for updates [409]

- **If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances?**

Unknown.

- **Will older versions of the dataset continue to be supported/hosted/maintained?**

Unknown.

- **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?**

Likely no.

C.2 Data Nutrition Label

The following analysis refers to `compas-scores-two-years.csv` after applying the standard COMPAS preprocessing [408].

METADATA	
Filename	compas-scores-two-years
Format	csv
Url	https://www.propublica.org/datastore/dataset/compas-recidivism-risk-score-data-and-analysis
Domain	Law
Keywords	risk assessment, pretrial, recidivism
Type	Tabular
Rows	6,172
Columns	57
% missing cells	5%
Rows with missing cells	100%
License	ProPublica's ToU [409]
Released	May 2016
Range	2013-2014 for COMPAS scores, 2013-2016 for arrest and detention history.
Description	Dataset curated by ProPublica to audit COMPAS software for racial biases, focusing on Broward County 2013-2014.

Table 14: Metadata of COMPAS dataset.

PROVENANCE	
Source	
Name	Broward County Sheriff's Office
Url	http://www.sheriff.org/
email	//
Name	Broward County Clerk's Office
Url	https://www.browardclerk.org
email	Eclerk@browardclerk.org
Name	Florida Department of Corrections
Url	http://www.dc.state.fl.us/
email	FDCCitizenServices@fdc.myflorida.com
Authors	
Names	Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner
Url	https://www.propublica.org/datastore/dataset/compas-recidivism-risk-score-data-and-analysis
email	data.store@propublica.org

Table 15: Provenance of COMPAS dataset.

VARIABLES	
id	Unique identifier assigned by the authors
name	Defendant's first and last name
first	Defendant's first name
last	Defendant's last name
compas_screening_date	Day defendant was scored by COMPAS
sex	Defendant's sex
dob	Defendant's date of birth
age	Defendant's age
age_cat	Age quantization: less than 25 25-45 greater than 45
race	Defendant's race: African-American Asian Caucasian Hispanic Native American Other
juv_fel_count	Number of juvenile felonies
decile_score	COMPAS recidivism score (10-point scale)
juv_misd_count	Number of juvenile misdemeanors
juv_other_count	Number of other juvenile convictions (not considering misdemeanor and felonies)
priors_count	Number of prior crimes
days_b_screening_arrest	Days between imprisonment (c_jail_in) and COMPAS screening (compas_screening_date)

Table 16: Variables of COMPAS dataset (1/3).

VARIABLES	
c_jail_in	Date of imprisonment
c_jail_out	Date of release
c_case_number	Alpha-numeric case identifier
c_offense_date	Date on which the offense was committed
c_arrest_date	Date on which defendant was arrested
c_days_from_compas	Days elapsed between offense/arrest and the date of COMPAS screening
c_charge_degree	Degree of charge: F (felony) M (misdemeanor)
c_charge_desc	Textual description of charge
is_recid	Binary indication of recidivism.
r_case_number	Alpha-numeric case identifier for recidivist offense
r_charge_degree	Degree of recidivist charge
r_days_from_arrest	Days elapsed between date of recidivist offense (r_offense_date) and date of recidivist incarceration (r_jail_in)
r_offense_date	Date of recidivist offense
r_charge_desc	Textual description of recidivist charge
r_jail_in	Date of incarceration for recidivist offense
r_jail_out	Date of release for recidivist offense

Table 17: Variables of COMPAS dataset (2/3).

VARIABLES	
violent_recid	Unknown; all nan
is_violent_recid	Binary indication of violent recidivism. If true, then is_recid is true.
vr_case_number	Alpha-numeric case identifier for violent recidivist offense
vr_charge_degree	Degree of violent recidivist offense
vr_offense_date	Date of violent recidivist offense
vr_charge_desc	Textual description of the violent recidivist charge
type_of_assessment	Type of COMPAS assessment - all 'Risk of Recidivism'.
decile_score_1	Identical to decile_score
score_text	Quantization of decile_score: LOW (1-4) MEDIUM (5-7) HIGH (8-10).
screening_date	Identical to compas_screening_date
v_type_of_assessment	Type of COMPAS violent assessment - all 'Risk of Violence'.
v_decile_score	COMPAS violent recidivism score (10-point scale)
v_score_text	Quantization of v_decile_score: LOW (1-4) MEDIUM (5-7) HIGH (8-10).
v_screening_date	Identical to compas_screening_date.
in_custody	Unknown
out_custody	Unknown
priors_count.1	Identical to priors_count.
start	Unknown
end	Unknown
event	Unknown
two_year_recid	Unknown

Table 18: Variables of COMPAS dataset (3/3).

STATISTICS							
Ordinal							
name	type	count	uniqueEntries	mostFrequent	leastFrequent	missing	
id	int	6,172	6,172	multiple	multiple	0	
compas_screening_date	date	6,172	685	2013-04-20	multiple	0	
dob	date	6,172	4,830	multiple	multiple	0	
age_cat	string	6,172	3	25 - 45	Greater than 45	0	
c_jail_in	date	6,172	6,172	multiple	multiple	433	
c_jail_out	date	6,172	6,161	2013-09-14 05:58:00	multiple	433	
c_offense_date	date	6,172	737	multiple	multiple	1388	
c_arrest_date	date	6,172	417	2013-02-06	multiple	8425	
r_offense_date	date	6,172	1,041	2014-12-08	multiple	3,182	
r_jail_in	date	6,172	928	multiple	multiple	4,175	
r_jail_out	date	6,172	893	multiple	multiple	4,175	
vr_offense_date	date	6,172	505	2015-08-15	multiple	5,480	
v_score_text	string	6,172	3	Low	High	0	
v_screening_date	date	6,172	685	2013-04-20	multiple	0	
score_text	string	6,172	3	Low	High	0	
screening_date	date	6,172	685	2013-04-20	multiple	0	
in_custody	date	6,172	1,087	multiple	multiple	0	
out_custody	date	6,172	1,097	2020-01-01	multiple	0	

Table 19: Ordinal variables statistics of COMPAS dataset

Categorical							
name	type	count	uniqueEntries	mostFrequent	leastFrequent	missing	
name	string	6,172	9,128	mutiple	multiple	0	
first	string	6,172	2,493	michael	multiple	0	
last	string	6,172	3,465	williams	multiple	0	
sex	string	6,172	2	Male	Female	0	
race	string	6,172	6	African-American	Native American	0	
c_case_number	string	6,172	6,172	multiple	multiple	0	
c_charge_desc	string	6,172	390	Battery	multiple	5	
c_charge_degree	string	6,172	2	F	M	0	
r_case_number	string	6,172	2,991	multiple	multiple	3,182	
r_charge_desc	string	6,172	319	Possess Cannabis/ 20 Grams Or Less	multiple	3,228	
r_charge_degree	string	6,172	11	(M1)	(F5)	0	
vr_case_number	string	6,172	693	multiple	multiple	5,480	
vr_charge_desc	string	6,172	82	Battery	multiple	5,480	
vr_charge_degree	string	6,172	10	(M1)	(F5)	5,480	
type_of_assessment	string	6,172	1	Risk of Recidivism	Risk of Recidivism	0	
v_type_of_assessment	string	6,172	1	Risk of Violence	Risk of Violence	0	
is_recid	binary	6,172	2	0	1	0	
is_violent_recid	binary	6,172	2	0	1	0	
event	binary	6,172	2	0	1	0	
two_year_recid	binary	6,172	2	0	1	0	

Table 20: Categorical variables statistics of COMPAS dataset

Quantitative										
name	type	count	min	median	max	mean	stdDev	miss	zeros	
age	int	6,172	18	31	96	34.53	11.73	0	0	
juv_fel_count	int	6,172	0	0	20	0.06	0.46	0	5,964	
juv_misd_count	int	6,172	0	0	13	0.09	0.50	0	5,820	
juv_other_count	int	6,172	0	0	9	0.11	0.47	0	5,711	
priors_count	int	6,172	0	1	38	3.25	4.74	0	2,085	
days_b_screening_arrest	int	6,172	-30.0	-1	30.0	-1.74	5.08	0	1,379	
c_days_from_compas	int	6,172	0	1	9,485	24.90	276.81	0	869	
r_days_from_arrest	int	6,172	-1	0	993	20.10	76.54	4,175	1,452	
decile_score	int	6,172	1	4	10	4.42	2.84	0	0	
v_decile_score	int	6,172	1	3	10	3.64	2.49	0	0	
start	int	6,172	0	0	937	13.32	50.14	0	3,485	
end	int	6,172	0	539	1,186	555.05	400.26	0	1	

Table 21: Quantitative variables statistics of COMPAS dataset.

D GERMAN CREDIT

Key references include Grömping [204], Häußler [242], UCI Machine Learning Repository [501, 503].

D.1 Datasheet

D.1.1 Motivation.

- **For what purpose was the dataset created?**

This dataset was created to study the problem of automated credit decisions at a regional Bank in southern Germany.

- **Who created the dataset and on behalf of which entity?**

The dataset was created at a regional Bank of southern Germany (most likely Hypo Bank) and first used by Walter Häußler in the late 1970s as part of his PhD thesis. Hans Hofmann, affiliated with Universität Hamburg at the time, is credited as dataset source [501]. Presumably, he donated the dataset to the European Statlog project and a representative of Strathclyde University donated it to UCI [204].

- **Who funded the creation of the dataset?**

The first known work using the dataset describes it as originating from a regional Bank of southern Germany [242]. Given the affiliation of the author is Hypo Bank, which fit the description at the time, we assume the dataset was collected, curated and funded at Hypo Bank.

D.1.2 Composition.

- **What do the instances that comprise the dataset represent?**

Instances represent Hypo bank **loan recipients** from 1973–1975.

- **How many instances are there in total?**

The dataset consists of **1,000** instances.

- **Does the dataset contain all possible instances or is it a sample of instances from a larger set?**

In principle this is a **convenience sample**, consisting of people who were deemed creditworthy by a bank clerk. A representative sample stemming from indiscriminate credit grants would not have been viable [242]. However, if the envisioned application was *post-screening* credit decisions, the influence of this selection bias would be reduced. Finally loan recipients associated with delayed payment or loan default (“bad credit”) are oversampled (30%).

- **What data does each instance consist of?**

For each instance, 13 categorical and 7 quantitative variables are provided, summarizing their financial situation, credit history, and personal situation, including housing, number of liable people, and a mixed variable encoding marital status and sex. A more through description is deferred to Tables 25-27.

- **Is there a label or target associated with each instance?**

Yes. A binary label encodes whether loan recipients punctually payed each installment (“good credit”) or not (“bad credit”). The latter label includes a range of situations from delayed payment up to loan default.

- **Is any information missing from individual instances?**

No. No cell is missing, however the variable “property” has a level jointly encoding the conditions “no property” and “unknown”. A similar joint encoding exists for “savings”, so some values may actually be deemed missing for these variables.

- **Are relationships between individual instances made explicit?**

No. There are no known relationships between instances.

- **Are there recommended data splits?**

No.

- **Are there any errors, sources of noise, or redundancies in the dataset?**

Yes. The dataset documentation is filled with errors, so that several levels of categorical variables do not correspond to what they should according to the official documentation from UCI Machine Learning Repository [501]. This is not necessarily an issue if one is purely interested in the evaluation of a method. For example, according to the official documentation, a majority of loan recipients are foreign workers, while in reality this should appear rather strange and indeed is not true [204]. Computationally, this will make no difference, as the input to a machine learning method will remain the same. However if one is interested to the context surrounding the data, as should be the case with fairness research, the wrong encoding poses several problems. The most significant problem is the impression that one can retrieve people’s sex from the joint sex-marital-status encoding, which is simply false as a single level corresponds to both single males and divorced/separated/married females [204]. Despite this information being available since 2019, the fairness community does not seem to have taken notice. Several experiments of algorithmic fairness on this dataset consider the protected attribute “sex” (sometimes even called “gender”). These experiments are part of work recently published in the most reputable venues for fairness research (Appendix A.73). More mistakes in the documentation of eight variables and the relative errata are outlined in Grömping [204]. A clean version of the dataset is available at UCI Machine Learning Repository [503].

- **Is the dataset self-contained, or does it link to or otherwise rely on external resources?**

The dataset is **self-contained**.

- **Does the dataset contain data that might be considered confidential?**

Yes. The dataset summarizes customers' financial and personal situation, including past credit history.

- **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?**

No.

- **Does the dataset identify any subpopulations?**

Yes. The dataset identifies subpopulation by age and sex. Sex is jointly encoded with marital status and cannot be retrieved, contrary to documentation accompanying the dataset [501]. A summary based on amended documentation [204] is presented in Table 22.

Demographic Characteristic	Values
Percentage of people under-19 years old	0.20%
Percentage of people between 20-29 years old	36.70%
Percentage of people between 30-39 years old	33.20%
Percentage of people between 40-49 years old	17.60%
Percentage of people between 50-59 years old	7.20%
Percentage of people between 60-69 years old	4.40%
Percentage of people over-70 years old	0.70%
Percentage of people who are male : divorced/separated	5.00%
Percentage of people who are female : non-single or male : single	31.00%
Percentage of people who are male : married/widowed	54.80%
Percentage of people who are female : single	9.20%

Table 22: Demographic characteristics of the German credit dataset.

- **Is it possible to identify individuals, either directly or indirectly, from the dataset?**

Likely no, especially given the fact that these records date back to almost 50 years ago. Also, important variables for re-identification, such as ZIP code and date of birth are missing and many other variables are bucketed.

- **Does the dataset contain data that might be considered sensitive in any way?**

Yes. For each instance, the dataset encodes sex, marital status and financial situation.

D.1.3 Collection process.

- **How was the data associated with each instance acquired?**

The data was collected by Hypo bank clerks. Some variables were observable (e.g. credit history with the bank), other variables were reported by subjects (e.g. loan purpose).

- **What mechanisms or procedures were used to collect the data?**

Unknown.

- **If the dataset is a sample from a larger set, what was the sampling strategy?**

The so-called "bad credits" are heavily oversampled to make the classification problem more balanced. A natural selection bias is present in the data, as it only consists of applicants who were deemed creditworthy and were thus granted a loan.

- **Who was involved in the data collection process and how were they compensated?**

The data was likely collected by Hypo bank clerks. Walter Häußler was likely involved in sample selection.

- **Over what timeframe was the data collected?**

The dataset covers loans granted in the period 1973–1975. Its first publicly-known use dates back to 1979 [242]. It became publicly available in November 1994 [501].

- **Were any ethical review processes conducted?**

Unknown.

- **Was the data collected from the individuals in question directly, or obtained via third parties or other sources?**

Likely both. Some variables were necessarily collected from loan applicants (e.g. loan purpose), while other variables were likely available from bank records (e.g. credit history with the bank).

- **Were the individuals in question notified about the data collection?**

Individuals provided some of this data as part of a loan application. Collection and notification practices for variables like credit history are unclear.

- **Did the individuals in question consent to the collection and use of their data?**
Likely yes, for the purposes of the immediate credit decision. However it seems implausible they agreed to their data becoming publicly available in an anonymized fashion.
- **If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses?**
Likely no.
- **Has an analysis of the potential impact of the dataset and its use on data subjects been conducted?**
Unknown.

D.1.4 Preprocessing/cleaning/labelling.

- **Was any preprocessing/cleaning/labeling of the data done?**
Yes. Some instances were discarded. Remaining instances were associated with a binary label according to compliance with the contract. Bucketing took place on several variables, including balance on checking and savings account (A1, A6) and duration of current employment (A7). Sex and marital status were jointly coded (A9).
- **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data?**
Unknown.
- **Is the software used to preprocess/clean/label the instances available?**
Likely no.

D.1.5 Uses.

- **For what tasks has the dataset been used?**
The dataset was originally used to study the problem of automated credit scoring [242]. Similarly to the Adult dataset, since becoming publicly available it has been used as a benchmark in various machine learning fields.
- **Is there a repository that links to any or all papers or systems that use the dataset?**
Yes. A selection of early works (pre-2005) using this dataset can be found in UCI Machine Learning Repository [501]. A more recent list is available under the beta version of the UCI ML Repository.²⁶ See Appendix A.73 for a (non-exhaustive) list of algorithmic fairness works using this resource.
- **What (other) tasks could the dataset be used for?**
The German Credit could be used in fields that concentrate on socially relevant goals and require socially relevant data, such as privacy and explainability. The task at hand is always credit scoring.
- **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?**
Contrary to documentation accompanying the dataset [501], the sex of loan recipients cannot be reliably retrieved. Works of algorithmic fairness should not use this feature.
- **Are there tasks for which the dataset should not be used?**
In its most common version [501] the German Credit dataset should not be used in works of explainability/interpretability as the incorrect documentation would result in counter-intuitive explanations. The 2019 version [503] associated with the erratum [204] is recommended.

D.1.6 Distribution.

- **Is the dataset distributed to third parties outside of the entity on behalf of which the dataset was created?**
Yes. The dataset is publicly available [501]
- **How is the dataset distributed?**
The dataset is available as a **csv file**.
- **When was the dataset distributed?**
The dataset was released to the UCI ML Repository in **November 1994**.
- **Is the dataset distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?**
Yes. The UCI ML repository has a citation policy.
- **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?**
Likely no. We are unaware of any IP-based restrictions.
- **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?**
Unknown.

D.1.7 Maintenance.

- **Who is supporting/hosting/maintaining the dataset?**

²⁶<https://archive-beta.ics.uci.edu/ml/datasets/144>

The dataset is hosted and maintained by the **UCI Machine Learning Repository** [501]. A clean and well-documented version of the same dataset donated by Ulrike Gromping [503] is also available on the same repository.

- **How can the owner/curator/manager of the dataset be contacted?**

The dataset donor, Hans Hofmann retired in 2008. Comments and inquiries for UCI may be sent to ml-repository@ics.uci.edu.

- **Is there an erratum?**

Yes. A clean data release [503] and accompanying report [204] are available online.

- **Will the dataset be updated?**

Likely no. The recently released South German Credit Data Set [503] may be considered an update.

- **If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances?**

Unknown.

- **Will older versions of the dataset continue to be supported/hosted/maintained?**

Unless otherwise indicated, both the new [503] and the old version [501] of the German Credit dataset will remain hosted on the UCI ML Repository in its current version.

- **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?**

Unknown.

D.2 Data Nutrition Label

For the sake of correctness, we report redacted information based on the new South German Credit Data Set [503] and accompanying documentation [204].

METADATA	
Filenames	SouthGermanCredit
Format	.asc
Url	https://archive.ics.uci.edu/ml/datasets/South+German+Credit
Domain	Economics
Keywords	credit scoring, Germany, loan, classification
Type	Tabular
Rows	1000
Columns	21
% missing cells	0%
Rows with missing cells	0%
License	UCI Repository citation policy
Released	November 2019
Range	1973-1975
Description	This dataset encodes socio-economical features of loan recipients from a bank in southern Germany, along with binary variable encoding whether they punctually payed every installment, which is the target of a classification task.

Table 23: Metadata of South German Credit dataset.

PROVENANCE	
Source	
Name	Walter Häußler
Url	https://archive.ics.uci.edu/ml/datasets/Statlog+%28German+Credit+Data%29
email	//
Authors	
Names	Ulrike Grömping
Url	https://archive.ics.uci.edu/ml/datasets/South+German+Credit
email	groemping@bht-berlin.de

Table 24: Provenance of South German Credit dataset

VARIABLES	
status	Checking account balance (in Deutsche Mark) 1 (no checking account) 2 (< 0 DM) 3 ($0 \leq \dots < 200$ DM) 4 (≥ 200 DM)
duration	Credit duration (in months)
credit_history	Applicant's credit history 0 (delay in past payments) 1 (critical account/other credits elsewhere) 2 (no credits taken/all credits paid back duly) 3 (existing credits paid back duly till now) 4 (all credits at this bank paid back duly)
purpose	Purpose of loan 0 (other) 1 (new car) 2 (used car) 3 (furniture/equipment) 4 (radio/television) 5 (domestic appliances) 6 (repairs) 7 (education) 8 (vacation) 9 (retraining) 10 (business)
amount	Credit amount (result of unknown monotonic transformation)

Table 25: Variables of South German Credit dataset (1/3).

VARIABLES	
savings	Savings account balance (in Deutsche Mark) 1 (unknown/ no savings account) 2 (< 100 DM) 3 (100 ≤ ... < 500 DM) 4 (500 ≥ ... < 1000 DM) 5 (≥ 1000 DM)
employment_duration	Duration of applicant's current employment 1 (unemployed) 2 (< 1 year) 3 (1 ≤ ... < 4 years) 4 (4 ≤ ... < 7 years) 5 (≥ 7 years)
installment_rate	Installment amount to disposable income ratio [%] 1 (≥35) 2 (25 ≤ ... < 35) 3 (20 ≤ ... < 25) 4 (< 20)
personal_status_sex	Joint encoding of sex and marital status of applicant 1 (male - divorced/separated) 2 (female - non single or male - single) 3 (male - married/widowed) 4 (female - single)
other_debtors	Presence of co-debtor or guarantor 1 (none) 2 (co-applicant) 3 (guarantor)
present_residence	Years living at current address 1 (< 1 year) 2 (1 ≤ ... < 4 years) 3 (4 ≤ ... < 7 years) 4 (≥ 7 years)
property	Applicant's most valuable property 1 (unknown / no property) 2 (car or other) 3 (building soc. savings agr / life insurance) 4 (real estate)

Table 26: Variables of South German Credit dataset (2/3).

VARIABLES	
age	Applicant's age (years)
other_installment_plans	Installment plans with other banks 1 (bank) 2 (stores) 3 (none)
housing	Type of housing 1 (for free) 2 (rent) 3 (own)
number_credits	Number of credits (ongoing or past, including current) with this bank 1 (1) 2 (2-3) 3 (4-5) 4 (≥ 6)
job	Applicant's job and employability 1 (unemployed/ unskilled - non-resident) 2 (unskilled - resident) 3 (skilled employee / official) 4 (manager / self-empl. / highly qualif. employee)
people_liable	Number of people who financially depend on the applicant 1 (3 or more) 2 (0 to 2)
telephone	Presence of telephone landline registered under applicant's name (2) or not (1)
foreign_worker	Foreign worker (1) or not (2)
credit_risk	Punctually payed back every installment (1) or not (2)

Table 27: Variables of South German Credit dataset (3/3).

STATISTICS							
Ordinal							
name	type	count	unique	mostFrequent	leastFrequent	missing	
status	string	1000	4	4 (≥ 200)	3 ($0 \leq \dots < 200$)	0	
savings	string	1000	5	1 (unknown/no savings)	4 ($500 \leq \dots < 1000$)	0	
employment_duration	string	1000	5	3 ($1 \leq \dots < 4$)	1 (unemployed)	0	
installment_rate	string	1000	4	4 (< 20)	1 ≥ 35	0	
present_residence	string	1000	4	4 (≥ 7 yrs)	1 (< 1 yr)	0	
number_credits	string	1000	4	1 (1)	4 (≥ 6)	0	
people liable	string	1000	2	2 (0 to 2)	1 (3 or more)	0	

Table 28: Ordinal variables statistics of South German Credit dataset

Categorical							
name	type	count	uniqueEntries	mostFrequent	leastFrequent	missing	
credit_history	string	1000	5	2 (no credits taken)	0 (delay in paying off)	0	
purpose	string	1000	11	3 (furniture/equipment)	8 (vacation)	0	
status_sex	string	1000	4	3 (male-marr/widow)	1 (male-divorc/separ)	0	
other_debtors	string	1000	3	1 (none)	2 (co-appliant)	0	
property	string	1000	4	3 (building soc. savings)	4 (real estate)	0	
other_plans	string	1000	3	3 (none)	2 (stores)	0	
housing	string	1000	3	2 (rent)	3 (own)	0	
job	string	1000	4	3 (skilled empl/office)	1 (unempl/unsk non-res)	0	
telephone	string	1000	2	1 (no)	2 (yes)	0	
foreign_worker	string	1000	2	2 (no)	1 (yes)	0	
credit_risk	string	1000	2	1 (good)	0 (bad)	0	

Table 29: Categorical variables statistics of South German Credit dataset

Quantitative									
name	type	count	min	median	max	mean	stdDev	miss	zeros
duration	number	1000	4	18	72	20.90	12.06	0	0
amount	number	1000	250	2319.50	18424	3271.25	2822.75	0	0
age	number	1000	19	33	75	35.54	11.35	0	0

Table 30: Quantitative variables statistics of South German Credit dataset.