Projection-Displacement based Query Performance Prediction for Embedded Space of Dense Retrievers

SUCHANA DATTA, University College Dublin, Ireland

GUGLIELMO FAGGIOLI, University of Padua, Italy

NICOLA FERRO, University of Padua, Italy

DEBASIS GANGULY, University of Glasgow, UK

CRISTINA IOANA MUNTEAN, ISTI-CNR, Italy

RAFFAELE PEREGO, ISTI-CNR, Italy

NICOLA TONELLOTTO, University of Pisa, Italy

Recent advances in representation learning allow neural Information Retrieval (IR) systems to use learned dense representations for queries and documents to effectively handle semantics, language nuances, and vocabulary mismatch problems. In contrast to traditional IR systems that rely on word matching, dense IR models exploit query/document similarities in dense latent spaces but need substantial training data and come with increased computational demands. Thus, it would be beneficial to predict how a system will perform for a given query to decide whether a dense IR model is the best option or alternatives should be used. Traditional Query Performance Predictors (QPP) are designed for lexical IR approaches and hence they perform sub-optimally when applied to (dense) neural IR systems. Therefore, there has been a renewed interest in QPP to make it more effective for (dense) neural IR models. While the results of the new QPP methods are generally encouraging, there is ample room for improvement in terms of absolute performance and stability. We argue that by using features that are more aligned with the inner rationale underneath dense IR models, we can improve the performance of QPP. In this respect, we propose the Projection-Displacement based QPP (PDQPP) that, exploiting the geometric properties of dense IR models, projects queries and retrieved documents onto sub-spaces defined by pseudo-relevant documents and considers the changes in retrieval scores in such sub-spaces as proxy for retrieval incoherence. Minor score changes suggest coherent retrieval, while significant alterations indicate semantic divergence and potentially poor performance. Results over a wide range of experiment settings on both traditional (TREC Robust) and neural-oriented (TREC Deep Learning) test collections show that PDQPP mostly outperforms the state-of-the-art QPP baselines.

ACM Reference Format:

1 INTRODUCTION

The advent of pretrained Large Language Models (LLMs) has accelerated the development of supervised Information Retrieval (IR) models that use these LLMs as foundation models, the parameters of which are then fine-tuned on examples of relevant and non-relevant documents for queries [33, 35, 36, 38, 62, 67]. The parameters of a fine-tuned bi-encoder

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.

© 2024 Association for Computing Machinery.

Manuscript submitted to ACM

model represent the dense vector representations (embeddings) of documents and queries [38, 46]. A bi-encoder model encodes the query and the documents in the latent embedding space separately. This has the advantage that it is possible to embed documents beforehand and only compute the query representation at run-time. Conversely, cross-encoders embed queries and documents jointly, thus requiring reindexing of the documents and the query at runtime. Dense end-to-end IR models operate by conducting an approximate nearest neighbor search on an indexed embedding space of document and query vectors [33, 35, 36, 46, 62, 67].

While dense representations are more effective in bridging the semantic gap between queries and documents, they are also more computationally expensive. Recognizing which queries benefit from using dense models and which ones can be managed with traditional lexical approaches would allow us to reduce query latency and save computational resources [15, 40]. A second major drawback of dense IR models is the need for vast training data. In particular, if the training set does not contain enough examples for a specific query type, we might observe low performance for such queries. In this sense, recognizing which queries are likely to fail may help collect aimed annotations to improve the performance of the dense models on such queries [9, 27]. Consequently, developing effective Query Performance Prediction (QPP) approaches can potentially help develop adaptive pipelines for IR systems, where only a subset of queries on which lexical models do not perform well may be routed to more computationally expensive rankers [15, 40] (e.g., dense end-to-end models). Additionally, QPP estimates may also be used to select queries for deeper relevance assessments to help develop more effective rankers [28].

Most classical QPP approaches leverage discrete term statistics and hence operate on sparse retrieval pipelines [53, 55, 70]. Off-the-shelf application of these classical QPP approaches on neural ranking models (NRMs) have been shown not to produce sufficiently effective results mainly because these approaches do not factor in term semantics [16, 21, 22].

This paper focuses on improving the QPP effectiveness for end-to-end dense rankers. It has been recently shown that the use of query variants (i.e., alternative formulations of the information need of a query) plays an important role in improving QPP effectiveness [16, 66], mainly because the ranked list of documents retrieved with these variants provide additional sources of information about the retrieval quality of the original top-retrieved list itself. Existing works on generating query variants operate in the discrete term space, e.g., reformulating a query 'five stages of grief' to a more specific version 'five stages of grief in sports' by adding terms. While such variants can be used for QPP estimates via methods such as [16, 66], the variant generation process does not take into account the topology of the embedded space itself, which is somewhat limiting in nature. In this paper, we try to address this limitation towards devising an effective QPP approach for dense NRMs.

The main idea of our method is conceptually similar to the idea of aggregating the relative changes in QPP estimates measured across the variants [16]. However, instead of generating variants in the discrete term space and then embedding them as dense vectors, we rather measure these relative changes across the embedded vector representations of the top-retrieved documents. In other words, a top-retrieved document in our proposed method acts as a proxy for a query variant vector. More specifically speaking, our proposed QPP estimator Projection Displacement Query Performance Predictor (PDQPP), projects both the query and the retrieved documents on the subspaces defined by a set of pivot vectors constituted of top-k ranked documents. Our method then aggregates the relative changes in the similarities between the projected vectors and the original ones for each query document pair.

Indeed, the pseudo-relevant documents provide us with an unsupervised way to describe different facets of the topic underlying a query. Suppose there are no major changes in the retrieval scores when we project the query and documents to the subspaces identified by each pseudo-relevant document. In that case, we can hypothesize the retrieval

 is of good quality. In contrast, significant changes in retrieval scores suggest that different pseudo-relevant documents define semantically quite different spaces, which, in turn, indicates a possibly incoherent retrieval, potentially indicating low performance.

Our research question can be formalized as follows: can we employ the projection displacement to exploit topological properties of the latent embedding space of a dense IR model, to devise a query performance predictor that achieves state-of-the-art effectiveness?

To answer this question, we conducted extensive experimentation, relying on both traditional experimental collections (TREC Robust) and neural-oriented ones (TREC Deep Learning), considering several dense IR retrieval models (ANCE, Contriever, and TAS-B) and a range of state-of-the-art QPP approaches. Our experiments show that our proposed predictor – PDQPP – is very often the top-performing approach or, at least, in the top-performing group. Moreover, it delivers very stable performance across experimental collections and IR models, different from current state-of-the-art approaches, which suffer from performance variability under various operating conditions.

The paper is organized as follows: Section 2 summarizes the relevant literature; Section 3 introduces the projection displacement operator based QPP (PDQPP); Sections 4 and 5 present the experiment setup and results; Section 6 concludes the paper with future directions.

2 RELATED WORK

Dense IR. Traditionally, IR systems relied on lexical signals, such as the presence of the query terms within the documents. However, the emergence of neural models transformed how we represent and match queries and documents. Dense IR approaches are traditionally divided into three main categories: bi-encoders, cross-encoders, and late-interaction models [68].

Bi-encoders (a.k.a dual-encoders) are models that use two separate (but possibly identical) neural networks to represent documents and queries [68]. In the most typical scenario, a placeholder token, such as the [CLS] token [45], is appended to the text (i.e., the query or the document) and the string is fed to a transformer architecture. The latent representation of the placeholder token is then used as a representation of the text. To compute the similarity between the query and a document, the inner product between the representation of the two is used. This has the major advantage of allowing to precompute the representation of all documents. At runtime, it is sufficient to compute the representation of the query and compute its inner product with the representation of all the documents. In recent years, several such models have been released, e.g., STAR [67], ANCE [62], Contriever [35], TAS-B [33]. In this work, we focus on this category of models as they allow us to represent in the same latent space both queries and documents separately. More in detail, we focus on symmetric bi-encoders using the same neural network to encode queries and the document.

Traditional cross-encoders jointly represent documents and queries [68]. To do so, such models concatenate a placeholder token to the query and the document, obtaining a final string with the format "[CLS] \(query \) [SEP] \(\document \)", where the special token [SEP] indicates where the query finished and the document begins. The string is fed to a transformer architecture that produces a contextual representation of each token. Then, the representation of the [CLS] token is further fed to a fully connected layer that outputs the probability that the query is relevant to the token. To obtain the aforementioned representation, it is necessary to have access to the query. This would require computing a new representation for the documents every time a new query is received. Therefore, cross-encoders are mostly used to operate on a small set of documents, such as for reranking.

Late-interaction models, such as ColBERT [36], require computing and storing a contextual representation of each term of the query and documents. For what concerns documents, such representation can be computed beforehand and

 stored in an efficient index structure. At runtime, the contextual vector representation of each query term is matched with the most similar document term representation for each document. The QPP proposed in this paper focuses on approaches that employ a single representation for the query or the documents, while late-interaction models employ multiple vectors (i.e., on for each word) to represent queries and documents. How to adapt PDQPP for late-interaction models is left as a future work.

Dense IR systems produce smaller but denser representations than those produced by the traditional lexical IR approaches. Indeed, traditional approaches are based on representations whose dimensionality ranges from the tens of thousands to hundreds of thousands of dimensions: the number of terms in the vocabulary considered by the IR. Vice-versa, dense IR systems learn representations whose dimensionality falls within the range of hundreds to thousands.

Traditional QPP. Depending on the features they rely upon, traditional QPPs are divided into pre- and post-retrieval predictors [6, 31, 32]. The former relies on signals that can be derived without considering the ranked list of documents produced in response to the query. Such signals are, for example, the collection frequency of terms appearing in the query [42, 69]. On the other hand, post-retrieval predictors infer their predictions by taking as input also the ranked list of documents in response to the query. Depending on which aspects are considered to compute the prediction, there are three main classes of post-retrieval predictors: coherence-, score-, and robustness-based. Coherence-based predictors rely on measuring how strongly documents retrieved are clustered together: the most well-known representative of this class of approaches is Clarity [12]. PDQPP, the predictor proposed in this paper, is a representative of this class. Score-based predictors employ heuristics computed on the retrieval score of the retrieved documents, some examples include Weighted Information Gain (WIG) [70], Normalized Query Commitment (NQC) [55], and Score Magnitude and Variance (SMV) [57]. Finally, robustness-based predictors compare the original ranking of documents with one produced by introducing noise in the query, the index, or documents, e.g., the Utility Estimation Framework (UEF) [53], the Reference Lists framework [49, 54], and Robust Standard Deviation (RSD) [51]. Given their similarity with our approach, we include the UEF framework and RSD as baselines.

Traditional QPPs were meant and designed to operate on lexical IR methods, such as BM25 [48] or the traditional language models, that relied on the presence of the same terms in both queries and documents to determine the relevance of a document. With the advent of Neural IR and semantic matching-based IR systems, it was highlighted the need for novel QPPs explicitly designed to cooperate with such novel IR systems [22]. In this regard, we recognize two novel classes of QPPs: those that employ semantic signals but are aimed at predicting the performance of lexical IR systems, and those explicitly designed to cooperate with Neural IR models.

Semantic QPPs for lexical IR systems. The advent of word embeddings fostered the development of QPP models that exploit them to compute their predictions. NeuralQPP, proposed by Zamani et al. [65], uses Deep Learning to integrate three diverse signals: the query text, the retrieval scores, and aspects related to the distribution of the terms. On the same line, Roy et al. [52] show that, by utilizing the semantic similarity aspect of word embedding, it is possible to estimate the local neighborhood of a query using Gaussian Mixture Models. Roy et al. observe that the spatial properties of such neighborhood correlate with system performance. In a similar manner, Arabzadeh et al. [4, 5] propose a set of measures derived from neural embeddings that allow for quantifying the term specificity. They observe that the presence of highly specific terms in a query is an indicator of more effective retrieval. Khodabakhsh and Bagheri [37] propose three neural features based on dense word representations: Neural Matching, Neural Aggregated Matching, and Neural Distance. These features combine the embeddings of query and document tokens to capture the semantic relationships

 occurring between them. The authors use the matching signals provided by such features to encode semantic aspects within classic predictors. Different from this, Datta et al. [14] proposes to make use of interaction between query and document terms as signals for QPP. More specifically, they employ 3D convolutional neural networks with shared parameters to train an end-to-end pairwise predictor, called Deep-QPP.

Arabzadeh et al. [1] introduce BERT-QPP, one of the first methods harnessing LLMs for QPP. Specifically, they fine-tune BERT [18] by utilizing BM25's performance on each training query and the BERT representation of the first retrieved document as supervision to train a QPP. Several subsequent works build upon BERT-QPP. Similarly, Chen et al. [8] extend BERT-QPP, introducing a groupwise approach enabling the prediction of query performance using signals from multiple queries simultaneously.

Arabzadeh et al. [3] also utilize LLMs to create a predictor for conversational search. They leverage BERT to construct a document graph and cluster documents. If, for a document, multiple clusters exist, they identify the user's information need by posing clarifying questions to determine the cluster containing relevant documents. Subsequently, they test this approach using BM25. While these methods lean towards Neural IR models, their primary application remains associated with lexical IR approaches. This leads to a discrepancy between the query/document representations utilized for ranking and prediction phases, with the former relying more on lexical aspects and the latter emphasizing semantic information. All the aforementioned predictors are evaluated on IR systems that rely on lexical matching and therefore are hindered when used to predict the performance of IR systems that exploit semantic matching [22]. Given that most of these predictors are designed for and tested on lexical IR models, they do not align with the focus of this paper, which instead addresses OPP predictors tailored for dense IR models.

QPP for Neural IR. Among QPPs explicitly designed to work the best with neural IR systems, Hashemi et al. [30] introduce Non-Factoid Question Answering QPP (NQAQPP), a methodology incorporating retrieval scores, query lexical features, and both query and answer lexical features within a deep neural network framework for addressing Non-Factoid Question Answering. Hashemi et al. is also one of the early works evaluating the effectiveness of QPP on neural IR models. They specifically evaluate it on BM25, aNMM [64], and Conv-KNRM [13], noting a substantial gap in predictive accuracy between BM25 and neural IR models, attributed to distinct score distributions generated by neural models. In a recent investigation, Faggioli et al. [22] scrutinize the capability of traditional QPP techniques in predicting the performance of neural IR systems and, through a series of experiments, they find a significant decline in the performance of current QPP models when applied to neural IR systems. This trend persists even when employing BERT-QPP as a predictive model for neural IR. Similarly, Datta et al. [16, 17] observe the diminished effectiveness of prior QPP methods when employed for neural IR compared to lexical IR. In response, they propose Weighted Relative Information Gain-based model (WRIG), a statistical approach employing probabilistic combinations of retrieval scores for multiple query formulations. To demonstrate the efficacy of their approach, they utilize WRIG to predict performance in BM25, four variations of DRMM [29], and the initial stage neural IR model, ColBERT [36]. Singh et al. [56] propose a novel QPP that employs an auxiliary pairwise ranker (DuoT5) as an unsupervised QPP model to measure how often the ranking produced by the IR system agrees with the pairwise comparison of the auxiliary model. Similarly to [16, 17], Singh et al. test the performance of the proposed model on multiple neural IR models, both considered end-to-end retrieval as well as reranking. Faggioli et al. [20] utilize the geometric characteristics of dense representations for performance prediction in conversational search, by devising the Hypervolume (HV) predictor which consists of computing the volume on the axes-aligned bounding box containing the top-k retrieved documents and the query. More recently, Arabzadeh et al. [2] proposed a strategy explicitly designed to be applied for dense IR systems. The predictor

311 312 proposed by Arabzadeh et al., called DenseQPP (DQPP), is based on measuring the similarity between the original ranked list and the ranked list obtained after perturbing the query with appositely crafted Gaussian noise. Faggioli et al. [19] propose a novel framework, called Dense-Centroid (DC) framework, to adapt traditional predictors to the dense IR systems. They start by noticing that classical predictors require regularizing predictions by the retrieval score that the corpus would achieve in response to the query. This score cannot be computed for dense models, as it would require feeding the entire corpus to the dense IR system and obtaining its representation. Therefore, they propose to use, as a proxy representation of the corpus, the centroid of the documents. More concretely, in their approach, the dot product between the original query and the centroid is used as a regularization factor within the classical QPPs. As these approaches share similar characteristics with that of our proposed approach PDQPP, in Section 3.6, we provide a detailed comparison between PDQPP and the above-mentioned state-of-the-art predictors.

3 PROPOSED METHODOLOGY

In this section, we describe our proposed methodology of projection-based QPP estimation.

3.1 Notations and Core concepts

In this section, we introduce the notations for embedded query and document vectors and outline the concept of vector projection, an essential component of our proposed predictor.

Embedded documents and queries. Since we aim to predict QPP for dense neural ranking models, we introduce the notations that will be useful to understand how our methodology works on the space of embedded vectors obtained via a bi-encoder-based neural representation model [35, 63]. Let ϕ be a bi-encoder-based supervised neural representation model, which has learned the parameterised representations of queries and documents from a training dataset. Embeddings of the textual representation of a query Q and that of a document $D \in \mathcal{D}$ (\mathcal{D} denotes a document corpus) are then denoted, respectively, as **q** and **d**, where both **q** and $\mathbf{d} \in \mathbb{R}^p$. The retrieval score of a document for a query Q is then usually obtained by computing a dot product between the embedded representations of the query and a document from a candidate set, i.e., $\pi(Q, D) \stackrel{\text{def}}{=} \mathbf{q} \cdot \mathbf{d}$, where $D \in \mathcal{D}_k$ denoting a candidate set of k documents obtained via approximate nearest neighbour search on the embedded space.

Vector projections. We now introduce the concept of projection and discuss how it plays an important role in our proposed predictor. Informally speaking, a projection of a vector d onto another vector v leads to aligning d in the direction of \mathbf{v} and also changes its magnitude. The standard notation to denote the projected vector is $\mathbf{d}_{\mathbf{v}}$ (read as \mathbf{d} projected onto v), and is defined as

$$\mathbf{d_v} = \left(\frac{\mathbf{d} \cdot \mathbf{v}}{||\mathbf{v}||}\right) \hat{\mathbf{v}},\tag{1}$$

where $||\mathbf{v}||$ denotes any norm (e.g., L^2) of the vector \mathbf{v} , and $\hat{\mathbf{v}}$ denotes the unit vector along \mathbf{v} , i.e., $\hat{\mathbf{v}} = \mathbf{v}/||\mathbf{v}||$. Note that the quantity within parenthesis is a scalar, and hence the projected vector $\mathbf{d}_{\mathbf{v}}$ is a scaled version of $\hat{\mathbf{v}}$.

Projection displacement. We now introduce the concept of *projection displacement*, which represents how much the similarity between a pair of vectors (in terms of their dot product) or the angular distance between them (in terms of the cosine inverse of their dot product) changes when both are projected onto a different vector. Formally, we define the projection displacement of a pair of vectors (\mathbf{q}, \mathbf{d}) given a third vector \mathbf{v} as

$$\delta_{\mathbf{v}}(\mathbf{q}, \mathbf{d}) = \mathbf{q} \cdot \mathbf{d} - \mathbf{q}_{\mathbf{v}} \cdot \mathbf{d}_{\mathbf{v}}, \tag{2}$$

where the notation $\mathbf{x}_{\mathbf{v}}$, as per Equation 1, denotes the projection of \mathbf{x} onto \mathbf{v} . The projection displacement of Equation 2 represents the relative gain (or loss) of the estimated similarity between two vectors q and d when a different frame of reference (v) is used to estimate this similarity.

3.2 Relative Changes in Retrieval Scores

In this section, we discuss the idea of projection displacement under the specific context of embedded documents and query vectors. Revisiting Equation 2 with an assumption that q refers to the embedding of a query Q and d refers to that of a document D, projection displacement can be interpreted as the relative change in the similarity between the query and the document when a different frame of reference is used. In terms of retrieval using an NRM, this affects the relative rank of the document D.

To better understand the characteristics of projection displacement within the specific context of dense retrievers, let us revisit Equation 2 and express it in terms of the angles between the vectors. Substituting the identity $\mathbf{x} \cdot \mathbf{y} =$ $||\mathbf{x}|| ||\mathbf{y}|| \cos(\mathbf{x}, \mathbf{y})$ into Equation 2, we see that the dot product between a pair of vectors \mathbf{q} and \mathbf{d} when projected on an arbitrary vector v can be expressed as

$$\begin{split} q_{v} \cdot d_{v} &= \frac{||q|| ||v|| \cos(q, v)}{||v||} \hat{v} \cdot \frac{||d|| ||v|| \cos(d, v)}{||v||} \hat{v} \\ &= ||q|| \cos(q, v) \hat{v} \cdot ||d|| \cos(d, v) \hat{v} \\ &= ||q|| \cos(q, v)||d|| \cos(d, v) \cos(\hat{v}, \hat{v}) \\ &= ||q|| ||d|| \cos(q, v) \cos(d, v). \end{split} \tag{3}$$

The last step is derived from the fact that both q_v and d_v are vectors along the same direction, and hence $\cos(\hat{v},\hat{v})=1$. Equation 3 expresses the similarity between a query and a document vector projected along the same direction as a product of their norms and their angles with the axis of projection, which when substituted into Equation 2 yields the expression for projection displacement as

$$\begin{split} \delta_{\mathbf{v}}(\mathbf{q}, \mathbf{d}) &= \mathbf{q} \cdot \mathbf{d} - \mathbf{q}_{\mathbf{v}} \cdot \mathbf{d}_{\mathbf{v}} \\ &= ||\mathbf{q}|| ||\mathbf{d}|| \cos(\mathbf{q}, \mathbf{d}) - ||\mathbf{q}|| ||\mathbf{d}|| \cos(\mathbf{q}, \mathbf{v}) \cos(\mathbf{d}, \mathbf{v}) \\ &= ||\mathbf{q}|| ||\mathbf{d}|| \Big(\cos(\mathbf{q}, \mathbf{d}) - \cos(\mathbf{q}, \mathbf{v}) \cos(\mathbf{d}, \mathbf{v}) \Big). \end{split} \tag{4}$$

The formulation of projection displacement in Equation 4 allows relating it to QPP estimation. As a boundary case realise that $\delta_{\mathbf{v}}(\mathbf{q},\mathbf{d})=0$ if $\mathbf{v}=\mathbf{q}$ or $\mathbf{v}=\mathbf{d}$, e.g., if $\mathbf{v}=\mathbf{q}$ then $\cos(\mathbf{q},\mathbf{v})=1$ and $\cos(\mathbf{q},\mathbf{v})=\cos(\mathbf{q},\mathbf{d})$, as a result $\delta_{\mathbf{v}}(\mathbf{q}, \mathbf{d}) = ||\mathbf{q}|| ||\mathbf{d}|| (\cos(\mathbf{q}, \mathbf{d}) - \cos(\mathbf{q}, \mathbf{d})) = 0.$

By a similar argument, if the projection axis v is close to either the query or the document, i.e., $|1 - \cos(q, v)| < \epsilon$ for a sufficiently small $\epsilon \in \mathbb{R}^+$, it is easy to see that $\delta_{\mathbf{v}}(\mathbf{q},\mathbf{d}) \to 0$. In other words, projection axes \mathbf{v} close to either the query or the document induces small projection displacements.

Choosing the Projection Vectors

Till now, we have defined the projection displacement (Equations 2 and 4) in a generic way for an arbitrary vector v. We now consider the situation when this vector v corresponds to an alternative formulation of the same information need as expressed by the embedding q of a query Q. In such a case, $\mathbf{q}_{\mathbf{v}} \cdot \mathbf{d}_{\mathbf{v}}$ can be interpreted as the similarity between the

 query and a document D (embedded as \mathbf{d}) in this transformed space of an alternative representation of the information need.

According to the *Clusering Hypothesis* [59], if a document D is relevant to the query Q, then we expect their representation to be similar. Furthermore, assume V represents a piece of information highly related to Q, such as a reformulation or the response to Q. If D is relevant to a query Q, it is also likely to be relevant (and hence likely to yield a high similarity score) to a query variant V [7, 16, 66].

In the context of QPP, this means that for a query and a relevant document pair (Q, D), the projection displacements or the relative changes in the retrieval scores for a different way of expressing the information need (i.e., V) should be small. This is the key idea of our proposed QPP estimator which measures the relative stability of the retrieval scores of top-retrieved documents along different projection vectors.

While previous QPP approaches, such as [16, 66], have leveraged manually created or automatically generated query variants for discrete text, it is inconvenient to generate such variants in the embedded space of vectors. For QPP on dense vectors, we propose to make use of the top-retrieved documents themselves as the different axes for computing the projection displacements.

Now we define the fundamental component of our proposed QPP predictor, that is the *projection displacement deviation* (PDD) for a pivot document (say D) over a set of top-ranked k documents. Formally, given a query embedding \mathbf{q} , a set of top-retrieved documents \mathcal{D}_k for the query and the embedding \mathbf{d} of a pivot document $D \in \mathcal{D}_k$, we define PDD as the standard deviation of the projection displacement values for each top-ranked document when projected along the pivot, i.e.,

PDD(
$$\mathbf{q}, \mathbf{d}, \mathcal{D}_k$$
) = $\sqrt{\frac{\sum_{i=1}^k (\delta_{\mathbf{d}}(\mathbf{q}, \mathbf{d}_i) - \mu)^2}{k}}$, where $\mu = \sum_{i=1}^k \delta_{\mathbf{d}}(\mathbf{q}, \mathbf{d}_i)$. (5)

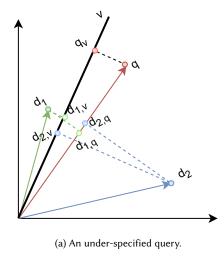
Intuitively, we expect the function $PDD(\mathbf{q}, \mathbf{d}, \mathcal{D}_k)$ to yield a small value if the pivot document is topically aligned with the query and every other document in the top-retrieved set. This is likely to happen if the pivot document is relevant to the query.

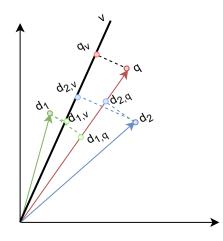
For under-specified queries, but also in the case of unsuccessful retrieval, the top-retrieved set of documents likely corresponds to different aspects of information need. In such a situation, selecting a pivot document that corresponds to a particular aspect of information need may lead to larger PDD values due to the presence of other documents corresponding to a different aspect. This also means that PDD concerning a pivot top-ranked document (Equation 5) can potentially act as a component to define an effective query performance estimator for dense vector spaces of queries and documents because a small value of this quantity is indicative of a likely well-specified query and vice-versa.

3.4 PDD-based QPP predictor

With the PDD definition (Section 3.3) and its geometric illustration (Section 3.5) we now formulate the predictor in terms of the PDD values. The key idea behind the proposed predictor is to aggregate the evidence for PDD values along several top-retrieved documents, which is similar to the idea of aggregating QPP estimates over multiple query variants [16, 66].

While the PDD values indicate the standard deviation of the topical alignment of the top-retrieved documents, it is potentially useful to scale these values relative to the similarities between the query and the document vectors in the embedded space, i.e., the retrieval scores. This scaling is likely to help calibrate these values over a range of different queries and potentially leads to an effective comparison between the QPP estimates.





(b) A well-specified query.

Fig. 1. A 2D visualisation of the local geometry of the embedding spaces of two queries and their top-retrieved documents. a) It is easy to find a pivot vector for which some of the document embeddings will be well-aligned whereas others will not. This will lead to a large PDD value (Equation 5). b) document vectors are well-aligned to the pivot vector, which means that the PDD values will be smaller. While v can be any possible direction, we observe empirically that the best results are achieved when v is aligned with pseudo-relevant documents vectors

Since our predictor, which we call **PDQPP**, aggregates the scaled PDD values over multiple pivots, we introduce an additional parameter to allow provision for how many documents to consider for this aggregation. Formally speaking, we call \bar{l} the mean of the retrieval scores of the top-l retrieved document, i.e., $\bar{l} = \frac{1}{l} \sum_{D \in \mathcal{D}_l} \mathbf{q} \cdot \mathbf{d}$. Then, PDQPP is defined as:

$$PDQPP(Q) = -\frac{\sum_{j=1}^{k} PDD(\mathbf{q}, \mathbf{d}_{j}, \mathcal{D}_{h})}{k \cdot \sqrt{\frac{1}{l} \sum_{i=1}^{l} (\mathbf{q} \cdot \mathbf{d}_{i} - \bar{l})^{2}}},$$
(6)

where the numerator represents the PDD values (Equation 5) computed for k pivots over a set of h top-ranked documents (h and k being two different parameters), whereas the denominator corresponds to the scaling factor of average similarity values between the query and a set of top-l ranked documents (again the parameter l is different from k and h).

The predictor is an additive inverse of these aggregated displacement values (minus sign at the front of Equation 6) because the higher the displacements the higher is the likelihood that the query itself is under-specified and the top documents potentially correspond to different aspects of information need, some of which could be non-relevant thus degrading the retrieval effectiveness of such queries.

The generic form of our proposed predictor has three hyper-parameters to control the sizes of the top-retrieved sets for three different computation purposes - i) a top-set of h documents to compute the PDD values with respect to a particular pivot document (Equation 5), ii) k, which specifies how many pivot documents to consider for aggregating the PDD values, and iii) l, the number of documents considered to compute the standard deviation of the retrieval scores.

3.5 A Geometric Illustration

Dense vector representations of the top-retrieved documents addressing different aspects of the information need are likely to be aligned along different subspaces while all of these are still similar to the query subspace. However, this means that a document addressing a specific aspect of the information need is likely to less similar to another on a different aspect.

We provide here an illustrative example using a query that contains the polysemous word "bank", such as "Where is the closest bank?". The word "bank" might refer to the "financial institution" or "the land alongside a river", among many other meanings. Therefore, in response to the query, the retrieval system has retrieved documents concerning both financial institutions and geographic structures. Assume now we can somehow disambiguate the meanings of bank by transforming the projection space. If we could move to the "financial" projection space, we would observe documents concerning financial institutions to be close to the query, as in this new space the word "bank" refers to the financial institution, while documents regarding river sides would be demoted. Vice-versa, if we were to move on the "geography" projection space, we would observe documents concerning river banks to be close to the query. Depending on which meaning we attribute to the query, the ranking of documents dramatically changes. This is a clear indication of a complex query for which the IR system is likely to fail – not even a human being would be able to answer the query "Where is the closest bank?" without asking further questions!

Consider now a more specific query, such as "Where is the closest financial institution?". In this case, we could assume that our IR system will retrieve almost exclusively documents where the word bank refers to the financial meaning. Thus, regardless of the space we consider, we will observe the documents to be close to the query.

Our example assumes we are capable of doing two types of geometric operations: i) we are able to define subspaces that represent the different semantic meanings of the query ii) we can change our projection space to reflect different semantic aspects. The second operation is handled using the projection operator defined in Eq. 1.

Conversely, to address the first operation, Equation 6 employs the first top k documents as pivot documents. We assume that each of these documents conveys a specific semantic meaning and defines a subspace characterized by a latent semantic. These subspaces might be very similar if the pivot documents have similar semantics (e.g., they refer to closely related topics), or might be very different if different documents refer to completely unrelated subjects. Considering our example again, when it comes to the query "Where is the closest bank", the top-k documents could for example focus on different meanings of the word bank. Therefore, depending on which document is used as the pivot, we will observe differences in ranking when things are projected onto such pivot. This hints at a weak retrieval. Vice-versa, when we consider the second query, "Where is the closest financial institution?", if the top-k retrieved documents have similar meanings, then, when using each of them as a pivot document, we will observe a relatively stable ranking.

Figure 1 visualises the idea in two dimensions. Figure 1b shows the embeddings of the top-retrieved documents for an under-specified query, where the angles between the top-retrieved documents can be large if they represent different topics, whereas, for a well-specified query (Figure 1a), it is likely that all the top-retrieved documents are likely to be similar to each other (and also to the query).

In terms of the project displacement deviation (as defined in Equation 5), choosing any direction as the pivot document for computing PDD values over the embeddings of Figure 1b is likely to lead to a large value because there potentially will be documents that are not well aligned with the pivot direction v. On the other hand, the PDD values for the

embeddings in Figure 1a are likely to be small because each top-retrieved document will potentially be aligned well with any pivot vector.

Considering the bank example, Figure 1a could represent the situation where the query is "Where is the closest bank?". In line with our example, d_1 is a document about financial institutions, while d_2 regards river banks. If our pivot document v concerns financial aspects, then when we project the query and the documents on the subspace defined by v, we observe d_1 being closer to the query in the subspace than d_2 – i.e., when projected on v, d_1 is closer to q than d_2 . This is exactly the opposite of what happens when we consider the default situation (i.e., d_1 and d_2 projected on the query). Thus, our change or reference space induces a switch between d_1 and d_2 in the ranking. Vice-versa, Figure 1b represents the scenario where all the retrieved documents are closely related. Then, when we observe the projection on v of d_1 and d_2 , we do not notice any switch in their ranking.

3.6 PDQPP vs. other existing predictors

After presenting our predictor, we now discuss how PDQPP is different from existing predictors, while still resembling them in certain ways. This will be useful to see how PDQPP generalises some predictors seeking to mitigate their limitations.

PDQPP vs. Score Variance (SV)-based predictors. Several classical predictors [44, 50, 55], as well as DC predictors [21], employ the retrieval score variance to produce predictions. The rationale is that a high variance indicates that the IR system scored much higher on the documents retrieved in the top positions within a ranked list as compared to lower positions. This utilises the hypothesis that such high scores are reflective of the relevance of the top-ranked documents. Our predictor PDQPP uses the same signal (denominator of Equation 6) with a different underlying objective - which is to normalise the projection displacement values to produce its predictions.

The major improvement obtained by our predictor PDQPP over score variance-based ones (results later in Section 5) can most likely be attributed to the additional factor incorporated as the projection displacement deviation. While an existing score variance-based predictor can only compute how topically distinct is a set of a very-top list of documents from the ones that follow it, such predictors cannot predict the topical coherence of the top-retrieved set - a coherent set potentially indicating better quality retrieval.

PDQPP vs. the UEF estimator. The UEF framework for QPP estimation [53] relies on using pseudo-relevant documents to expand a query, retrieve a new set of documents, and compare the original ranked list with the one obtained from the expanded query. Conceptually, the UEF framework and PDQPP share several common characteristics. The UEF framework, in a sense, transforms a query into the reference space induced by the pseudo-relevant documents and then estimates how this transformed representation affects the ranking of the documents. Our proposed PDQPP operates in a similar but more explicit manner in that the query (and the documents) are explicitly projected within the pseudo-relevant space. Furthermore, similar to UEF, the projection displacement measures the (dis-)similarity between the results of the query in the original space as against the ones induced by the pseudo-relevant documents (Equation 5). The major difference is that by explicitly relying on the geometrical representation of the various elements, PDQPP better suits the end-to-end dense IR models.

PDQPP and **DQPP**. DQPP [2] projects the top-ranked documents on a subspace obtained by a perturbed version of the query, and then computes the robustness of the ranking of documents relative to this change. It can therefore be argued that both DQPP and PDQPP models project information on a different space, and hence estimate the robustness

Table 1. Evaluation (nDCG@10) of the dense IR models on the respective test collections subsequently used for our QPP experiments.

Topic set	ANCE	Contriever	TAS-B
DL '19	0.645	0.676	0.716
DL '20	0.646	0.671	0.684
DL Hard	0.328	0.376	0.376
Robust '04	0.362	0.499	0.453

of an IR model relative to this transformed representation. While DQPP obtains this directly by comparing the two retrieved lists, PDQPP on the other hand, achieves this using the projection displacement operator. Moreover, PDQPP offers an advantage over DQPP in the sense that the new subspace where documents and queries are projected is not random. Instead, this reference subspace aligns with the pseudo-relevant documents, thus allowing provision to leverage the latent semantics of these documents for query performance estimation.

PDQPP and **WRIG**. WRIG computes the relative changes in the QPP estimates with reference to a set of query variants with the idea that a large increase potentially indicates that the original query itself was under-specified (poor retrieval quality), whereas a large decrease suggests that the original query itself was well-specified (effective retrieval quality) [16]. The idea of transforming a query via projection onto a reference subspace relates to that of leveraging information from variants in WRIG. While WRIG uses the relative gains computed via a base QPP estimator, our predictor PDQPP, instead, uses deviations of projection displacements.

4 EXPERIMENT SETTINGS

4.1 Datasets and Models

Dense Neural Models. In our experimental analysis, we consider three dense retrieval models: ANCE¹ [62], Contriever² [35], and TAS-B³ [33]. We use the model weights fine-tuned on the MS MARCO collection and publicly available on the huggingface repository. All the models that we experimented with use 768 dimensional embeddings for documents and queries.

Dataset. As benchmark datasets, we employ the following four collections: TREC Deep Learning '19 (DL '19) [11], TREC Deep Learning '20 (DL '20) [10], Deep Learning Hard (DL Hard) [39], and TREC Robust '04 (Robust '04) [60]. DL '19, DL '20, and DL Hard datasets constitute 43, 54, and 50 queries, respectively, with depth pooled relevance assessments (depth 10). The underlying task is ad-hoc passage retrieval on MS MARCO corpus, which contains over 8M passages [43]. As a part of the experiment setup, all the dense IR systems were fine-tuned on the MS MARCO training set of pairs of queries and relevant passages. The respective topic sets of DL '19, DL '20 and DL Hard, the predictions are in-domain in nature.

Additionally, to evaluate the QPP effectiveness for the neural models for out-domain ranking predictions, we employ Robust '04, constituted of disks 4 and 5 (minus congressional records) of the Tipster collection. The Robust '04 collection uses a deeper pool (depth 100) for relevance assessments, as a result of which recall plays a crucial role in determining a query's performance. It thus offers a different evaluation setting as compared to MS MARCO passage collection.

 $^{^{1}} https://hugging face.co/sentence-transformers/msmarco-roberta-base-ance-firstp\\$

²https://huggingface.co/facebook/contriever

³https://huggingface.co/sentence-transformers/msmarco-distilbert-base-tas-b

4.2 Baselines and Evaluation Measures

Since our proposed QPP estimator is an unsupervised approach, for a fair comparison we employ a wide range of existing unsupervised predictors as baselines. More specifically, we consider two different categories of QPP models - i) those that are agnostic of an IR model, and ii) those that are explicitly designed to operate on embedding spaces of dense IR models.

IR Model agnostic QPP approaches. As QPP baselines that can work on both sparse and dense retrievers (i.e., agnostic QPP) we employ the following:

- SV [44] is an approach that predicts the variance of the retrieval scores of the first top-k retrieved as the QPP estimate.
- Clarity [12] computes the Kullback–Leibler (KL) divergence between the language model of the entire corpus and the one of the top-*k* retrieved documents. Clarity operates under the assumption that observing a large KL divergence indicates a well-characterized and coherent set of top-*k* documents, which hints at a good retrieval.
- NQC [55] is the standard deviation of the retrieval scores of the first top-k retrieved documents, regularized by the retrieval score of the entire corpus.
- **RSD** [51] iterates over the retrieved list of documents, computing at each position the unbiased standard deviation of the scores reweighed by the WIG score of the ranking list up to that position and sums all these values.
- **SMV** [57] Combines NQC and WIG by taking into consideration both the magnitude and the variance of the retrieval scores of the top-*k* documents.
- WIG [70] is the average retrieval score of the first top-k retrieved documents, regularized by the retrieval score of the entire corpus.
- **UEF Framework** [53] The UEF framework operates by reweighing any of the aforementioned predictors (Clarity, NQC, SMV and WIG) by the similarity between the original retrieved list of documents and the list of documents retrieved after rewriting the query via Pseudo-Relevance Feedback (PRF).
- WRIG [16] is a variant-based predictor that computes the changes in the QPP estimates as obtained from a base predictor on a set of query variants relative to the original query. As suggested in [16], we employed NQC as the baseline predictor of WRIG. Additionally, we worked with a set of query variants automatically generated by skipgram embeddings [41] as suggested in [16]. Notice that WRIG was observed to supersed reference lists based methods [66].

Dense IR-based approaches. This class of QPP models are explicitly formulated to operate with dense IR models. As baselines we use the following:

- DC framework [21] instantiates traditional predictors (e.g., WIG [70], NQC [55], SMV [57]) by considering the centroid of all documents as an approximated corpus representation. Among the DC class of predictors, we consider DCWIG, DCNQC, and DCSMV as suggested in [21].
- HV [20] predictor correlates the IR system performance with the volume of the n-parallelepiped encompassing the top-k documents retrieved.
- **DQPP** [2] introduces a small calibrated noise to a query's dense representation, and then as the QPP score, measures the similarity between the original ranked list of documents and the one obtained with the perturbed query.

QPP Evaluation Metrics. As evaluation metrics for QPP, we follow the standard protocol of reporting the correlation of predicted QPP estimates and a target metric (measured with Pearson's ρ), and also the rank correlation between the

ideal ordering of query performance as obtained by a target metric vs. the predicted ordering obtained via QPP scores (measured with Kendall's τ) [25]. As the target metric, we employ nDCG@10 following previous work [2, 21, 22], and being the official evaluation metric of TREC DL [10, 11]. In addition, we also employ a recently proposed error-based metric - scaled Mean Absolute Rank Error (sMARE) [23, 24], smaller values of which indicate better performance. To provide a consistent interpretation across the metrics, we report the values of one minus the sMARE scores (the range of sMARE values is in [0, 1]), which we call $\overline{\text{sMARE}}$.

4.3 Hyper-parameter tuning

For each predictor, we validate the hyper-parameters using the commonly adopted 2-fold validation strategy [16, 20, 55, 65, 66]. Specifically, this commonly used validation strategy involves randomly splitting a set of queries into two partitions, one used as a 'training set' for tuning parameters (for supervised approaches) or hyper-parameters (for unsupervised approaches), and the other partition is used as a 'test set' to evaluate the model performance. The roles of the two partitions are then switched, and the average performance over the two folds is then used as an evaluation measure. Evaluation measures collected this way are then aggregated over 30 random 2-fold splits of the data.

Recall that the hyper-parameters of our proposed method are the three cut-off values k, h and l, denoting the number of top documents to used to aggregate PDD values, the number of ones used as pivots for computing PDD values and the number of documents used to compute the scaling factor based on retrieval scores, respectively (see Equation 6). For a tractable choice of the number of experiments, we set the value of k to 5, which means that PDD values are aggregated over 5 documents. Later, in Section 5.3, we analyze the sensitivity of PDQPP to the number of documents used as pivots. The other two cut-offs in PDQPP, namely h and l, were optimised via grid search over the training splits from the set $\{5, 10, 50, 100, 250, 500\}$.

For a fair comparison, the hyper-parameter k (the number of top-documents used for estimation cut-off) of all the other baseline predictors were also optimised over the training folds. The baseline DQPP involves an additional parameter - the standard deviation of the Gaussian noise used to perturb query vectors. This parameter was validated in the range [0.01, 0.09] with a step of 0.01 following the implementation in the repository provided by Arabzadeh et al.⁴.

5 RESULTS

5.1 Comparison with other predictors

As a sanity checking step, we first report in Table 1 the nDCG@10 (our QPP target metric) values for the various datasets used for each IR model considered in our experiments. It can be seen that the results are consistent with existing numbers reported in the literature [33, 35, 62], which, in turn, shows that our retrieval setup is at par with previous findings.

Tables 2, 3, and 4 report the effectiveness of the proposed PDQPP model in comparison to the different baseline models. The best results obtained for a particular collection is shown bold-faced, whereas the second-best ones are shown underlined. To denote whether the best performing approach (bold-faced) is statistically indistinguishable from other methods, we append an asterisk to both. In particular, for the significance testing we employed ANOVA with Tukey's honestly significant difference (HSD) test with significance level $\alpha = 0.05$ [58].

In addition, to provide further insights on the relative performance of the QPP models, we report the number of times a particular method turns out to be a winner model either by outperforming other approaches or being statistically

⁴https://github.com/Narabzad/Dense-QPP

Table 2. Performance of PDQPP compared to the baselines in predicting ANCE's nDCG@10 in terms of Kendall's τ , Pearson's ρ and $\overline{\text{sMARE}}$ (1 – sMARE). For each dataset and IR model, we highlight in bold the best method and underline the runner-up. The postfix '*' indicates QPP models that are statistically at par with the best. The last column reports the effectiveness index (EI), which is the number of times a QPP model either is the winner or ends up being a 'star' competitor (i.e., statistically indistinguishable from the best method).

	DL '19				DL '20			DL Hai	:d		EI		
	τ	ρ	sMARE	τ	ρ	sMARE	τ	ρ	sMARE	τ	ρ	sMARE	EI
ANCE													
SV	.355	.497	.781	.254	.318	.743	.348*	.330	.783*	.407	.430	.797*	3
Clarity	.100	.151	.697	.044	.013	.678	.306	$.486^*$.758	.111	.167	.702	1
NQC	.285	.363	.760	.137	.201	.697	.243	.256	.735	.251	.410	.745	0
RSD	.267	.445	.744	.317	.443	.765	.378*	.450	.786*	.380	.498	.785	2
SMV	.185	.189	.732	.007	051	.662	.170	.100	.710	.118	.220	.700	0
WIG	.323	.483	.759	.286	.456	.746	.184	.279	.713	.414*	.546*	.794	2
UEFClarity	.136	.222	.692	.121	.139	.701	.089	.119	.689	.164	.214	.723	0
UEFNQC	.190	.312	.714	.148	.215	.704	.177	.169	.714	.231	.277	.742	0
UEFSMV	.172	.256	.717	.047	.123	.673	.177	.048	.731	.185	.229	.726	0
UEFWIG	.177	.281	.705	.173	.254	.708	.080	.101	.691	.272	.299	.754	0
WRIG	.292	.451	.768	.196	$.498^{*}$.733	.141	.223	.719	.099	.161	.690	1
DCNQC	.361	.528	.784*	.255	.340	.744	.353*	.405	.783*	.407	.509	.796*	4
DCSMV	.334	.492	.769	.267	.364	.744	.276	.351	.762	.403	.483	.796*	2
DCWIG	.411*	.537	.798*	.269	.389	.740	.158	.215	.710	.276	.400	.747	2
HV	.224	.357	.735	.252	.321	.735	182	187	.616	.235	.347	.733	0
DQPP	<u>.369</u> *	<u>.600</u> *	.787*	.205	.275	.728	.123	.183	.700	.262	.377	.749	3
PDQPP	.367	.611*	.787*	.403*	.502*	.790*	.309	.394	.764	.405	.547*	.793	6

indistinguishable from the best performing model. We call this count the *effectiveness index* (EI) of a model and report its values in the last column of Tables 2 to 4 (higher values of this number indicating better effectiveness). Intuitively speaking, it does not over-penalise a model for not yielding the best results. Instead, it rewards the runner-up model for being statistically indistinguishable from the best one thus factoring in the variational effects of random 2-fold splits - the commonly used setup of QPP experiments [16, 20, 25, 55, 65, 66], as well as the well-known problem of the intrinsic variability of the OPP measurements (see Section 5.2).

In Tables 2, 3, and 4 we see that with a few exceptions, the proposed PDQPP model outperforms the current state of the art models, or is at par with the best approach. With a few exceptions, PDQPP can easily outperform agnostic QPPs (upper part of the Tables). This phenomenon is unsurprising as previous work showed the diminished effectiveness of classic QPP models in dealing with dense and semantic-based IR systems [21, 22]. Furthermore, PDQPP makes use of topological characteristics of the embedded space in an explicit manner via leveraging subspace projections, which is the likely reason for its superior performance. Indeed, dense-IR based approaches are a more effective comparison with DCNQC, DCWIG or DQPP being particularly effective, depending on the collection/predicted system. IR system-wise, PDQPP is the most effective in predicting the retrieval performance for Contriever (Table 3) and TAS-B (Table 4). As can be seen from the Tables, for both these models PDQPP turns out to be the best or indistinguishable from the best in 11 out of 12 setups. No other baseline approach exhibits this high consistency in predicting the retrieval performance for Contriever and TAS-B.

Table 3. Performance of PDQPP compared to the baselines in predicting Contriever's nDCG@10 in terms of Kendall's τ , Pearson's ρ and $\overline{\text{sMARE}}$ (1 – sMARE). For each dataset and IR model, we highlight in bold the best method and underline the runner-up. The postfix '*' indicates QPP models that are statistically at par with the best. The last column reports the effectiveness index (EI), which is the number of times a QPP model either is the winner or ends up being a 'star' competitor (i.e., statistically indistinguishable from the best method).

	DL '19				DL '20)		DL Hai	rd		EI		
	τ	ρ	sMARE	τ	ρ	$\overline{\text{sMARE}}$	τ	ρ	$\overline{\text{sMARE}}$	τ	ρ	$\overline{\text{sMARE}}$	EI
Contriever													
SV	.233	.411	.744	.114	.276	.708	.251*	.197	.760*	.298	.336	.767	2
Clarity	.198	.297	.725	015	019	.657	.261*	.344*	.750	.100	.137	.699	2
NQC	.241	.197	.743	.103	016	.701	.172	.282	.729	.227	.339	.739	0
RSD	.163	.283	.709	.264	.304	.742	.272*	s.291	$.752^{*}$.212	.330	.733	2
SMV	.214	.148	.732	.039	188	.676	.210	.200	.716	.131	.187	.707	0
WIG	.198	.372	.728	.104	.239	.696	.139	.278	.711	.248	.355	.745	0
UEFClarity	.255	.298	.743	045	098	.655	105	171	.633	.154	.239	.719	0
UEFNQC	.250	.210	.743	023	117	.669	018	068	.670	.197	.331	.733	0
UEFSMV	.204	.103	.730	001	226	.670	.030	120	.671	.168	.282	.724	0
UEFWIG	.254	.315	.739	015	075	.668	107	211	.627	.191	.292	.730	0
WRIG	.214	.273	.725	.113	.252	.702	011	.193	.663	.048	.151	.681	0
DCNQC	.263	.439	.757*	.125	.250	.712	.252*	.248	.760*	.279	.376	.764	3
DCSMV	.238	.411	.744	.162	.279	.715	.229*	.245	.751*	.264	.379	.756	2
DCWIG	.309*	.506*	.752*	.259*	<u>.415</u> *	<u>.747</u> *	.124	.170	.689	.219	.345	.734	6
HV	.129	.259	.724	.241	.322	.733	139	194	.625	.252	.384	.742	0
DQPP	.328*	.538*	.763*	078	011	.652	.128	.228	.714	.188	.297	.726	3
PDQPP	.268*	.479*	.744	.288*	.438*	.754*	.225*	<u>.340</u> *	<u>.757</u> *	.310*	.415*	$\boldsymbol{.774}^*$	11

PDQPP appears to be slightly less effective on ANCE (Table 2), where it belongs to the top-tier of predictors only 8 times out of 12. Notice that it is still the predictor with the highest EI. Furthermore, as per our observations on ANCE, the best baseline predictor depends heavily on the collection considered: for DL '19, we observe good high performance for DCWIG, while for DL Hard and Robust '04 the most effective baselines are SV and DCNQC. If we inspect the results collection-wise, we notice that PDQPP is particularly effective on DL '19,DL '20, and Robust '04.

5.2 On the improved QPP stability of DPQPP

Overall, the high variance in terms of the quality of the predictions is a well-known problem in the QPP domain [6, 22, 24, 31]. The variability is influenced by a number of different factors, such as which queries are considered, collections, retrieval models and evaluation measures. For example Hauff [31, p. 83-84] considers different subsets of queries of three collections, TREC Vol. 4 and 5, WT10g, and GOV2, observing how the best QPP heavily depends on which subset of queries is considered. On a different line, but with similar conclusions, Carmel and Yom-Tov [6, p. 23-24,35-36] apply several predictors on different collections, observing high volatility in terms of which QPP can be considered the most effective, depending on the collection. [47] employed 9 different corpora, observing again variability in which system performs the best. More recently, Ganguly et al. [26] explore the impact that several factors have on the QPP effectiveness, observing important consequences linked to the chosen metric as well as the IR system. Finally, Faggioli et al. [22] investigate several predictors applied on both lexical and neural IR systems, observing a strong variability on what is the best predictor, depending on which IR system we are trying to predict the performance for.

Table 4. Performance of PDQPP compared to the baselines in predicting TAS-B's nDCG@10 in terms of Kendall's τ , Pearson's ρ and $\overline{\text{sMARE}}$ (1 – sMARE). For each dataset and IR model, we highlight in bold the best method and underline the runner-up. The postfix '*' indicates QPP models that are statistically at par with the best. The last column reports the effectiveness index (EI), which is the number of times a QPP model either is the winner or ends up being a 'star' competitor (i.e., statistically indistinguishable from the best method).

	DL '19				DL '20)	DL Hard			Robust '04			EI
	τ	ρ	sMARE	τ	ρ	$\overline{\text{sMARE}}$	τ	ρ	$\overline{\text{sMARE}}$	τ	ρ	$\overline{\text{sMARE}}$	EI
	TAS-B												
SV	.167	.241	.709	.214	.394*	.727	.360*	.411*	.779*	.402*	.464	.792*	6
Clarity	.171	.268	.727	045	014	.653	.238	.334	.745	.209	.287	.733	0
NQC	.131	.196	.714	.101	.139	.691	.213	.389	.722	.282	.417	.751	0
RSD	.151	.244	.701	.275	.406*	.758*	.289	$.432^{*}$.763	.362	.507	.781	3
SMV	.124	.145	.705	.035	164	.681	.162	.264	.719	.161	.226	.719	0
WIG	.195	.353	.727	.174	.279	.711	.192	.311	.716	.318	.468	.765	0
UEFClarity	.205	.281	.724	032	105	.669	063	080	.644	.198	.304	.728	0
UEFNQC	.217	.244	.740*	.004	072	.667	.048	.134	.689	.269	.371	.750	1
UEFSMV	.228	.249	.735	009	206	.667	004	.027	.669	.236	.299	.739	0
UEFWIG	.223	.332	.726	.004	.003	.676	023	031	.660	.241	.354	.740	0
WRIG	.228	<u>.353</u> *	<u>.744</u> *	.151	.096	.708	.235	.179	.739	.166	.230	.717	2
DCNQC	.170	.284	.712	.204	.346	.724	.355*	.433*	.775*	.399*	.530*	.791*	6
DCSMV	.164	.259	.711	.185	.360	.721	.284	.405*	.750	.395*	.516	.785	2
DCWIG	.164	.182	.716	.178	.253	.726	190	203	.609	.303	.432	.764	0
HV	.086	.145	.694	.251	.327	.743	088	100	.637	.183	.299	.718	0
DQPP	.209	.190	.738*	024	041	.658	.066	.210	.699	.268	.382	.746	1
PDQPP	.293*	.401*	.749*	.303*	.437*	.761*	.339*	.454*	.759	.405*	.543*	.792*	11

The very same behaviour can be observed in our results reported in Tables 2, 3, and 4). Depending on what collection is considered and which retrieval model is the target of our predictions, we observe most of the baselines exhibit a high variance in evaluation metric values. For example, consider Table 2, where we observe that when predicting the performance of ANCE on Robust '04, WIG is the best system. If we apply WIG on Contriever and DL '20 (Table 3), WIG performance is 63% worse than PDQPP, the most effective predictor in those cases. Similarly, Clarity, which is the best for predicting the performance of Contriever on the DL Hard collection, perfroms quite poorly on DL '20 for all IR systems sometimes leading to negative correlation. Generally speaking, this pattern is more severe for agnostic predictors than for dense ones (with few exceptions, such as HV and DQPP which also exhibit instability).

The major advantage of PDQPP is indeed able to provide a more stable performance than the current baseline predictors (as can be seen from the consistency in the EI values from Tables 2 to 4). Even in scenarios where PDQPP fails to outperform all the baselines, it is either statistically at par with the best, or reasonably close to the best.

To better exemplify this, we report the *critical difference diagram* of the evaluated QPP in Figure 2. The critical difference diagram reports on the x axis (on top) the rank, indicating what is the average rank for a QPP over the various experimental settings (i.e., retrieval model, collection, and correlation measure considered). Furthermore, the thick horizontal lines represent groups of statistically equivalent approaches, according to the Wilcoxon test [61] corrected according to the Holm correction procedure [34]. For example, in Figure 2, we observe that the average rank of PDQPP is 2.14. Furthermore, the second-best approach is DCNQC with an average rank of 4.22. Furthermore, PDQPP is statistically the best according to the multiple-comparison adjusted Wilcoxon test, while the second-best, DCNQC,

Fig. 2. Critical difference diagram across all experimental settings (IR system, collection, correlation measure). The average rank for PDQPP is 2.14, and it is statistically better than the average rank of the second best (DCNQC, with average rank of 4.22).

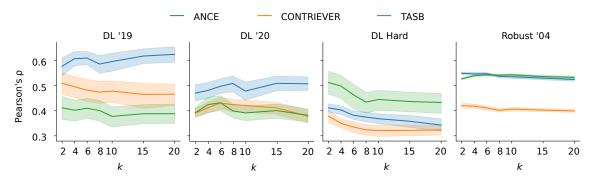


Fig. 3. The performance of the PDQPP when varying the number of pseudo-relevant documents. The general trend suggests that choosing 2-6 documents as pseudo-relevant is the most effective strategy, but large confidence intervals (or the almost flat lines for Robust '04), indicate a relatively small impact on the performance due to picking a wrong amount of pseudo-relevant documents.

given the high variance in its rank across different scenarios, is statistically at par with SV, DCSMV, RSD, WIG and DCWIG.

While designing a QPP that performs the best on all possible situations – predicted ir system, measure, collection – is a very complex task, we argue that the QPP systems should be reasonably reliable, without major drops in performance which render them untrustworthy. Our choice of including in our analyses multiple IR systems and collections is aimed at showing the overall stability of the proposed PDQPP. Indeed, where PDQPP is not the best, it still provides reasonable guarantees of effectiveness, even if compared against an always different most effective predictor. This observation, in fact, leads to a justification of our choice of considering three different IR systems and multiple collections.

5.3 Sensitivity to the pivot documents

To keep the number of experiments to tractable limits, we used the set of top 5 documents as pivots for computing the PDD values, i.e., we set k = 5, for all the results reported in Table 2 to 4. We now analyse the sensitivity of our predictor on this parameter. Figure 3 reports the effect of modifying the number of pivot documents from which we can make some interesting observations.

Firstly, we observe that since DL '19, DL '20, and DL Hard contain a much smaller number of queries than Robust '04, as a result of which, the performance of PDQPP on such collections is affected by a larger variance. Secondly, as a general trend, we observe that the performance tends to decrease with an increase in the number of pivot documents. This is in

Table 5. A comparison between the original PDQPP and its three variants where directions to project the embedded document and query vectors are sampled from different distributions. Similar to the results of Tables 2, 3, and 4, the target IR metric to compute QPP effectiveness is nDCG@10.

	DL '19				DL '2	20		DL Ha	ırd	Robust '04			
	τ	ρ	sMARE	τ	ρ	sMARE	τ	ρ	sMARE	τ	ρ	sMARE	
	ANCE												
R-PDQPP	.155	.236	.723	.171	.143	.718	.137	.148	.712	.291	.416	.753	
Q-PDQPP	.387	.558	.792	.300	.353	.765	.302	.418	.769	.382	.507	.787	
D-PDQPP	.380	.562	.788	.415	.564	.788	.282	.384	.755	.386	.524	.787	
PDQPP	.360	.611	.782	.399	.504	.791	.304	.394	.763	.405	.547	.793	
						CONTI	RIEVER	}					
R-PDQPP	.090	.136	.694	.161	.237	.709	.053	.054	.696	.227	.288	.741	
Q-PDQPP	.208	.432	.724	.230	.269	.732	.270	.374	.746	.298	.396	.760	
D-PDQPP	.272	.493	.745	.278	.399	.748	.239	.336	.755	.312	.416	.772	
PDQPP	.274	.475	.740	.275	.437	.753	.254	.340	.757	.310	.415	.774	
						TA	.SB						
R-PDQPP	.212	.284	.729	.080	.131	.688	.221	.220	.725	.273	.373	.748	
Q-PDQPP	.211	.352	.726	.267	.298	.746	.385	.462	.773	.418	.564	.798	
D-PDQPP	.368	.456	.782	.268	.357	.751	.304	.427	.756	.393	.514	.789	
PDQPP	.286	.401	.743	.297	.436	.759	.344	.456	.762	.405	.543	.792	

line with our hypothesis that such documents provide a way to disambiguate the meanings of the query in a latent space. The more documents we use further down a ranked list to define the reference spaces for computing the projection displacements, the more the chances are that such documents are not relevant to the query, thus incorporating noise in the prediction.

We also observe that the QPP effectiveness mostly decreases (often monotonically with a small number of exceptions) with an increase in k, e.g., see the results for the DL Hard collection. For some collections, we observe that the QPP effectiveness peaks at a value close to the range of about 2 to 4 documents beyond which it decreases almost steadily, e.g., see the plots for DL '19 and DL '20 collections. The ANCE model on DL '19 and DL '20 collections shows a reverse trend of improved QPP effectiveness with a larger number of pivots.

5.4 Using arbitrary subspaces for projections

As mentioned in Section 3, PDQPP relies on pseudo-relevant documents to identify the axes on which to project the query and the retrieved documents. While we argue that this approach allows leveraging the information within the ranked list itself, there might be alternative approaches to sample directions that might also be effective. Therefore, we conduct additional experiments with three different variations of PDQPP, each with its own way of obtaining the directions on which to compute the projection displacements, as detailed below.

• Random PDQPP (R-PDQPP) samples the directions from a Normal distribution centered at zero with standard deviation tuned in [0.1, 0.9] with steps of 0.1.

998

1000

1001

1002

1008 1009

1007

1010 1011

1012 1013

1014

1015

1016

1017 1018

1019

1020

1021

1022 1023

1024

1025

1026

5.5 PDQPP Limitations

1031

1032

1033 1034 1035

1036

1038 1039

1040

- Query PDQPP (Q-PDQPP) samples directions by perturbing the query with random noise drawn from a Normal distribution centred at zero, with standard deviation tuned in [0.1, 0.9] with steps of 0.1. DQPP uses the same method for generating perturbed queries. However, DQPP does not involve computing projection displacements as is the case for the variant Q-PDQPP.
- Documents PDOPP (D-PDOPP) samples directions by from an isotropic multivariate Normal distribution with parameters estimated from the top-5 document vector samples.

Table 5 reports a comparison of these variants with the originally proposed predictor (Equation 6). As a general trend, we observe that R-PDQPP is the worst-performing solution, the likely reason for which can be attributed to the fact that the projection axes being randomly sampled do not contain enough semantic information to differentiate between the different aspects of the information need inherent in a query. Q-PDQPP and D-PDQPP, on the other hand, yield much better results than R-PDQPP; in some cases, they even outperform the original predictor PDQPP based on pseudo-relevant documents. Interestingly, Q-PDQPP yields the best results on DL Hard, which constitutes a set of manually chosen 'hard' queries, i.e., queries for which retrieval performance is low.

While compared to the current state of the art PDQPP appears more robust and capable of achieving good performance across all scenarios, three major limitations that affect PDQPP need to be discussed. Limitation 1: PDQPP is not a model agnostic QPP. Indeed, PDQPP can be applied to predict the performance only of dense IR models. Two mitigating conditions should be taken into consideration. First, dense models are more and more popular in the IR community. It is usually common for IR pipelines to include a dense component for the purposes of first stage retrieval (as in this work), for reranking, or for both. Even though PDQPP, in principle, can also be applied for sparse vectors, the method is particularly suitable for dense vectors. This is because the projection of a sparse vector over another sparse one can lead to an abrupt effect of removing term weights from the former thus making it more sparse. Whereas for dense vectors the projections over subspaces retain more information. The popularity of dense IR models motivates its importance. Secondly, a model agnostic QPP cannot take into consideration specific additional information available to the IR model that might lead to an improvement in performance. In this case, the predictor exploits the geometric properties of the embedding space to better identify queries whose documents are affected by high variability in their semantics, suggesting possibly weak retrieval.

Limitation 2: PDOPP may not be suited for scenarios where the diversity in results is particularly important. PDOPP operates under the assumption that a stable and coherent retrieval list is likely more effective than a highly diversified one. These assumptions underly many QPPs such as Clarity [12], the UEF framework [53] or the reference lists framework [49]. This might not be the case of a fairness-oriented IR system which aims at maximizing the diversity of the results. Nevertheless, as future research direction, PDQPP should be tested for fairness-oriented IR tasks.

Limitation 3: PDOPP is not always the best performing OPP. This limitation has been extensively discussed in 5.2. To summarize such discussion, PDQPP is the most stable predictor compared to all other baselines, making it reliable even when it is not the most effective predictor. Conversely, most of the other approaches exhibit both gains and losses of high magnitudes in effectiveness depending on the experimental setup considered.

1050

1051 1052 1053 1054 1055

1056 1057 1058

1059 1060

1067 1068 1069 1070

1066

1072 1073 1074

1071

1076 1077 1078

1079

1080

1081

1082 1083 1084

1085

1087 1088 1089

1090 1091 1092

6 CONCLUSIONS AND FUTURE WORK

In this work, we proposed PDQPP, a novel QPP model capable of exploiting geometric properties in a dense embedding space to predict IR performance. The proposed predictor is based on the concept of *projection displacement*: we project the query and the retrieved documents on a reference subspace induced by the pseudo-relevant documents. The change of retrieval scores observed in the novel space represents a measure of the incoherence of the IR system. If, in the novel subspace, the query and the documents remain closely related, then we can assume the dense IR system to be successful. On the other hand, if we observe major changes in the novel subs-pace, then it is possible that the retrieval was unsuccessful and the performance will be low. In terms of effectiveness, the proposed QPP model can overcome several state-of-the-art baselines under a wide range of settings. Additionally, we also show that using pseudo-documents as subspaces yield better solutions than to use randomly selected ones.

In future directions, we plan to extend our predictor to other types of representation-learning based IR systems, including distillation models of late-interaction systems and sparse IR systems. We also plan to investigate other strategies to devise projection spaces, such as the space defined by previous utterances for a conversational search system or clustering of documents.

REFERENCES

- [1] Negar Arabzadeh, Maryam Khodabakhsh, and Ebrahim Bagheri. 2021. BERT-QPP: Contextualized Pre-trained transformers for Query Performance Prediction. In CIKM '21: The 30th ACM International Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021, Gianluca Demartini, Guido Zuccon, J. Shane Culpepper, Zi Huang, and Hanghang Tong (Eds.). ACM, 2857–2861. https://doi.org/10.1145/3459637.3482063
- [2] Negar Arabzadeh, Radin Hamidi Rad, Maryam Khodabakhsh, and Ebrahim Bagheri. 2023. Noisy Perturbations for Estimating Query Difficulty in Dense Retrievers. In Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM 2023, Birmingham, United Kingdom, October 21-25, 2023, Ingo Frommholz, Frank Hopfgartner, Mark Lee, Michael Oakes, Mounia Lalmas, Min Zhang, and Rodrygo L. T. Santos (Eds.). ACM, 3722–3727. https://doi.org/10.1145/3583780.3615270
- [3] Negar Arabzadeh, Mahsa Seifikar, and Charles L. A. Clarke. 2022. Unsupervised Question Clarity Prediction through Retrieved Item Coherency. In Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022, Mohammad Al Hasan and Li Xiong (Eds.). ACM, 3811–3816. https://doi.org/10.1145/3511808.3557719
- [4] Negar Arabzadeh, Fattane Zarrinkalam, Jelena Jovanovic, Feras N. Al-Obeidat, and Ebrahim Bagheri. 2020. Neural embedding-based specificity metrics for pre-retrieval query performance prediction. *Inf. Process. Manag.* 57, 4 (2020), 102248. https://doi.org/10.1016/j.ipm.2020.102248
- [5] Negar Arabzadeh, Fattane Zarrinkalam, Jelena Jovanovic, and Ebrahim Bagheri. 2020. Neural Embedding-Based Metrics for Pre-retrieval Query Performance Prediction. In Advances in Information Retrieval - 42nd European Conference on IR Research, ECIR 2020, Lisbon, Portugal, April 14-17, 2020, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 12036), Joemon M. Jose, Emine Yilmaz, João Magalhães, Pablo Castells, Nicola Ferro, Mário J. Silva, and Flávio Martins (Eds.). Springer, 78–85. https://doi.org/10.1007/978-3-030-45442-5_10
- [6] David Carmel and Elad Yom-Tov. 2010. Estimating the Query Difficulty for Information Retrieval. Morgan & Claypool Publishers. https://doi.org/10.2200/S00235ED1V01Y201004ICR015
- [7] Anirban Chakraborty, Debasis Ganguly, and Owen Conlan. 2020. Retrievability based Document Selection for Relevance Feedback with Automatically Generated Query Variants. In CIKM. ACM, 125–134.
- [8] Xiaoyang Chen, Ben He, and Le Sun. 2022. Groupwise Query Performance Prediction with BERT. In Advances in Information Retrieval 44th European Conference on IR Research, ECIR 2022, Stavanger, Norway, April 10-14, 2022, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 13186), Matthias Hagen, Suzan Verberne, Craig Macdonald, Christin Seifert, Krisztian Balog, Kjetil Nørvåg, and Vinay Setty (Eds.). Springer, 64-74. https://doi.org/10.1007/978-3-030-99739-7
- [9] Gordon V. Cormack, Christopher R. Palmer, and Charles L. A. Clarke. 1998. Efficient Construction of Large Test Collections. In SIGIR '98: Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, August 24-28 1998, Melbourne, Australia, W. Bruce Croft, Alistair Moffat, C. J. van Rijsbergen, Ross Wilkinson, and Justin Zobel (Eds.). ACM, 282–289. https://doi.org/10.1145/290941.291009
- [10] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, and Daniel Campos. 2021. Overview of the TREC 2020 deep learning track. CoRR abs/2102.07662 (2021). arXiv:2102.07662 https://arxiv.org/abs/2102.07662
- [11] Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. Overview of the TREC 2019 deep learning track. CoRR abs/2003.07820 (2020). arXiv:2003.07820 https://arxiv.org/abs/2003.07820
- [12] Stephen Cronen-Townsend, Yun Zhou, and W. Bruce Croft. 2002. Predicting query performance. In SIGIR 2002: Proceedings of the 25th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, August 11-15, 2002, Tampere, Finland, Kalervo Järvelin,

1097

1098

1099

1106

1107

1108

1109

1110

1111

1112

1124

1125

1130

1131

1132

1133

1134

1135

1136

1139

1141

- Micheline Beaulieu, Ricardo A. Baeza-Yates, and Sung-Hyon Myaeng (Eds.). ACM, 299-306. https://doi.org/10.1145/564376.564429
- [13] Zhuyun Dai, Chenyan Xiong, Jamie Callan, and Zhiyuan Liu. 2018. Convolutional Neural Networks for Soft-Matching N-Grams in Ad-hoc Search.
 In Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM 2018, Marina Del Rey, CA, USA, February 5-9,
 2018, Yi Chang, Chengxiang Zhai, Yan Liu, and Yoelle Maarek (Eds.). ACM, 126-134. https://doi.org/10.1145/3159652.3159659
 - [14] Suchana Datta, Debasis Ganguly, Derek Greene, and Mandar Mitra. 2022. Deep-QPP: A Pairwise Interaction-based Deep Learning Model for Supervised Query Performance Prediction. In WSDM '22: The Fifteenth ACM International Conference on Web Search and Data Mining, Virtual Event / Tempe, AZ, USA, February 21 - 25, 2022, K. Selcuk Candan, Huan Liu, Leman Akoglu, Xin Luna Dong, and Jiliang Tang (Eds.). ACM, 201–209. https://doi.org/10.1145/3488560.3498491
- [15] Suchana Datta, Debasis Ganguly, Sean MacAvaney, and Derek Greene. 2024. A Deep Learning Approach for Selective Relevance Feedback. In
 Advances in Information Retrieval 46th European Conference on Information Retrieval, ECIR 2024, Glasgow, UK, March 24-28, 2024, Proceedings, Part II
 (Lecture Notes in Computer Science, Vol. 14609), Nazli Goharian, Nicola Tonellotto, Yulan He, Aldo Lipani, Graham McDonald, Craig Macdonald, and
 Iadh Ounis (Eds.). Springer, 189-204. https://doi.org/10.1007/978-3-031-56060-6
- [16] Suchana Datta, Debasis Ganguly, Mandar Mitra, and Derek Greene. 2023. A Relative Information Gain-based Query Performance Prediction
 Framework with Generated Query Variants. ACM Trans. Inf. Syst. 41, 2 (2023), 38:1–38:31. https://doi.org/10.1145/3545112
 - [17] Suchana Datta, Sean MacAvaney, Debasis Ganguly, and Derek Greene. 2022. A 'Pointwise-Query, Listwise-Document' Based Query Performance Prediction Approach. In Proceedings of 45th international ACM SIGIR conference research development in information retrieval. 2148—-2153. https://doi.org/10.1145/3477495.3531821
 - [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proc. of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers). ACL, 4171–4186. https://doi.org/10.18653/ v1/n19-1423
 - [19] Guglielmo Faggioli, Nicola Ferro, Josiane Mothe, and Fiana Raiber. 2023. QPP++ 2023: Query-Performance Prediction and Its Evaluation in New Tasks. In Advances in Information Retrieval 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2-6, 2023, Proceedings, Part III (Lecture Notes in Computer Science, Vol. 13982). Springer, 388-391. https://doi.org/10.1007/978-3-031-28241-6_42
- [115] [20] Guglielmo Faggioli, Nicola Ferro, Cristina Muntean, Raffaele Perego, and Nicola Tonellotto. 2023. A Geometric Framework for Query Performance
 Prediction in Conversational Search. In Proceedings of 46th international ACM SIGIR Conference on Research & Development in Information Retrieval,
 SIGIR 2023 July 23–27, 2023, Taipei, Taiwan. ACM. https://doi.org/10.1145/3539618.3591625
- 1118 [21] Guglielmo Faggioli, Thibault Formal, Simon Lupart, Stefano Marchesin, Stéphane Clinchant, Nicola Ferro, and Benjamin Piwowarski. 2023. Towards

 Query Performance Prediction for Neural Information Retrieval: Challenges and Opportunities. In Proceedings of the 2023 ACM SIGIR International

 Conference on Theory of Information Retrieval, ICTIR 2023, Taipei, Taiwan, 23 July 2023, Masaharu Yoshioka, Julia Kiseleva, and Mohammad Aliannejadi

 (Eds.). ACM, 51–63. https://doi.org/10.1145/3578337.3605142
- [22] Guglielmo Faggioli, Thibault Formal, Stefano Marchesin, Stéphane Clinchant, Nicola Ferro, and Benjamin Piwowarski. 2023. Query Performance
 Prediction for Neural IR: Are We There Yet?. In Advances in Information Retrieval 45th European Conference on IR Research, ECIR 2023, Dublin,
 Ireland, April 2-6, 2023. 1-18. https://doi.org/10.48550/ARXIV.2302.09947
 - [23] Guglielmo Faggioli, Oleg Zendel, J. Shane Culpepper, Nicola Ferro, and Falk Scholer. 2021. An Enhanced Evaluation Framework for Query Performance Prediction. In Advances in Information Retrieval 43rd European Conference on IR Research, ECIR 2021, Virtual Event, March 28 April 1, 2021, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 12656), Djoerd Hiemstra, Marie-Francine Moens, Josiane Mothe, Raffaele Perego, Martin Potthast, and Fabrizio Sebastiani (Eds.). Springer, 115–129. https://doi.org/10.1007/978-3-030-72113-8_8
- [24] Guglielmo Faggioli, Oleg Zendel, J. Shane Culpepper, Nicola Ferro, and Falk Scholer. 2022. sMARE: a new paradigm to evaluate and understand query performance prediction methods. *Inf. Retr. J.* 25, 2 (2022), 94–122. https://doi.org/10.1007/s10791-022-09407-w
 - [25] Debasis Ganguly, Suchana Datta, Mandar Mitra, and Derek Greene. 2022. An Analysis of Variations in the Effectiveness of Query Performance Prediction. In ECIR (1) (Lecture Notes in Computer Science, Vol. 13185). Springer, 215–229.
 - [26] Debasis Ganguly, Suchana Datta, Mandar Mitra, and Derek Greene. 2022. An Analysis of Variations in the Effectiveness of Query Performance Prediction. In Advances in Information Retrieval - 44th European Conference on IR Research, ECIR 2022, Stavanger, Norway, April 10-14, 2022, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 13185), Matthias Hagen, Suzan Verberne, Craig Macdonald, Christin Seifert, Krisztian Balog, Kjetil Nørvåg, and Vinay Setty (Eds.). Springer, 215–229. https://doi.org/10.1007/978-3-030-99736-6_15
 - [27] Debasis Ganguly and Emine Yilmaz. 2023. Query-specific Variable Depth Pooling via Query Performance Prediction towards Reducing Relevance Assessment Effort. CoRR abs/2304.11752 (2023). https://doi.org/10.48550/arXiv.2304.11752 arXiv:2304.11752
- 1137 [28] Debasis Ganguly and Emine Yilmaz. 2023. Query-specific Variable Depth Pooling via Query Performance Prediction towards Reducing Relevance
 1138 Assessment Effort. CoRR abs/2304.11752 (2023).
 - [29] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce Croft. 2016. A Deep Relevance Matching Model for Ad-hoc Retrieval. In Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, USA, October 24-28, 2016, Snehasis Mukhopadhyay, ChengXiang Zhai, Elisa Bertino, Fabio Crestani, Javed Mostafa, Jie Tang, Luo Si, Xiaofang Zhou, Yi Chang, Yunyao Li, and Parikshit Sondhi (Eds.). ACM, 55-64. https://doi.org/10.1145/2983323.2983769
- [10] Helia Hashemi, Hamed Zamani, and W. Bruce Croft. 2019. Performance Prediction for Non-Factoid Question Answering. In Proceedings of the 2019
 ACM SIGIR International Conference on Theory of Information Retrieval, ICTIR 2019, Santa Clara, CA, USA, October 2-5, 2019, Yi Fang, Yi Zhang, James

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1193

1194

1195

- Allan, Krisztian Balog, Ben Carterette, and Jiafeng Guo (Eds.). ACM, 55–58. https://doi.org/10.1145/3341981.3344249
 - [31] Claudia Hauff. 2010. Predicting the effectiveness of queries and retrieval systems. SIGIR Forum 44, 1 (2010), 88. https://doi.org/10.1145/1842890.1842906
 - [32] Claudia Hauff, Djoerd Hiemstra, and Franciska de Jong. 2008. A survey of pre-retrieval query performance predictors. In Proceedings of the 17th ACM Conference on Information and Knowledge Management, CIKM 2008, Napa Valley, California, USA, October 26-30, 2008, James G. Shanahan, Sihem Amer-Yahia, Ioana Manolescu, Yi Zhang, David A. Evans, Aleksander Kolcz, Key-Sun Choi, and Abdur Chowdhury (Eds.). ACM, 1419–1420. https://doi.org/10.1145/1458082.1458311
 - [33] Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently Teaching an Effective Dense Retriever with Balanced Topic Aware Sampling. In SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021, Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai (Eds.). ACM, 113-122. https://doi.org/10.1145/3404835.3462891
 - [34] Sture Holm. 1979. A Simple Sequentially Rejective Multiple Test Procedure. Scandinavian Journal of Statistics 6, 2 (1979), 65–70. http://www.jstor.org/stable/4615733
 - [35] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Towards Unsupervised Dense Information Retrieval with Contrastive Learning. CoRR abs/2112.09118 (2021). arXiv:2112.09118 https://arxiv.org/abs/2112.09118
 - [36] Omar Khattab and Matei Zaharia. 2020. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020, Jimmy X. Huang, Yi Chang, Xueqi Cheng, Jaap Kamps, Vanessa Murdock, Ji-Rong Wen, and Yiqun Liu (Eds.). ACM, 39–48. https://doi.org/10.1145/3397271.3401075
 - [37] Maryam Khodabakhsh and Ebrahim Bagheri. 2021. Semantics-enabled query performance prediction for ad hoc table retrieval. Inf. Process. Manag. 58, 1 (2021), 102399. https://doi.org/10.1016/j.ipm.2020.102399
 - [38] Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. 2021. Sparse, Dense, and Attentional Representations for Text Retrieval. Trans. Assoc. Comput. Linguistics 9 (2021), 329–345. https://doi.org/10.1162/tacl a 00369
 - [39] Iain Mackie, Jeffrey Dalton, and Andrew Yates. 2021. How Deep is your Learning: the DL-HARD Annotated Deep Learning Dataset. In SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021, Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai (Eds.). ACM, 2335–2341. https://doi.org/10.1145/3404835.3463262
 - [40] Stefano Marchesin, Alberto Purpura, and Gianmaria Silvello. 2020. Focal elements of neural information retrieval models. An outlook through a reproducibility study. Inf. Process. Manag. 57, 6 (2020), 102109. https://doi.org/10.1016/J.IPM.2019.102109
 - [41] Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. In 1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings. http://arxiv.org/abs/1301.3781
 - [42] Josiane Mothe and Ludovic Tanguy. 2005. Linguistic features to predict query difficulty. In ACM Conference on research and Development in Information Retrieval, SIGIR, Predicting query difficulty methods and applications workshop. Salvador de Bahia, Brazil, 7–10.
 - [43] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A Human Generated MAchine Reading COmprehension Dataset. In Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, December 9, 2016 (CEUR Workshop Proceedings, Vol. 1773), Tarek Richard Besold, Antoine Bordes, Artur S. d'Avila Garcez, and Greg Wayne (Eds.). CEUR-WS.org. https://ceur-ws.org/Vol-1773/CoCoNIPS_2016_paper9.pdf
 - [44] Joaquín Pérez-Iglesias and Lourdes Araujo. 2010. Standard Deviation as a Query Hardness Estimator. In String Processing and Information Retrieval 17th International Symposium, SPIRE 2010, Los Cabos, Mexico, October 11-13, 2010. Proceedings (Lecture Notes in Computer Science, Vol. 6393), Edgar Chávez and Stefano Lonardi (Eds.). Springer, 207–212. https://doi.org/10.1007/978-3-642-16321-0_21
 - [45] Yifan Qiao, Chenyan Xiong, Zhenghao Liu, and Zhiyuan Liu. 2019. Understanding the Behaviors of BERT in Ranking. CoRR abs/1904.07531 (2019). arXiv:1904.07531 http://arxiv.org/abs/1904.07531
 - [46] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 139), Marina Meila and Tong Zhang (Eds.). PMLR, 8748–8763. http://proceedings.mlr.press/v139/radford21a.html
 - [47] Fiana Raiber and Oren Kurland. 2014. Query-performance prediction: setting the expectations straight. In The 37th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '14, Gold Coast, QLD, Australia - July 06 - 11, 2014, Shlomo Geva, Andrew Trotman, Peter Bruza, Charles L. A. Clarke, and Kalervo Järvelin (Eds.). ACM, 13-22. https://doi.org/10.1145/2600428.2609581
 - [48] Stephen E. Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. Found. Trends Inf. Retr. 3, 4 (2009), 333–389. https://doi.org/10.1561/1500000019
 - [49] Haggai Roitman. 2017. An Enhanced Approach to Query Performance Prediction Using Reference Lists. In Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval, Shinjuku, Tokyo, Japan, August 7-11, 2017, Noriko Kando, Tetsuya Sakai, Hideo Joho, Hang Li, Arjen P. de Vries, and Ryen W. White (Eds.). ACM, 869-872. https://doi.org/10.1145/3077136.3080665
 - [50] Haggai Roitman. 2019. Normalized Query Commitment Revisited. In Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019, Benjamin Piwowarski, Max Chevalier, Éric Gaussier, Yoelle Maarek, Jian-Yun Nie, and Falk Scholer (Eds.). ACM, 1085–1088. https://doi.org/10.1145/3331184.3331334

1203

1207

1208

1213

1214

1215

1233

1234

1235

- [51] Haggai Roitman, Shai Erera, and Bar Weiner. 2017. Robust Standard Deviation Estimation for Query Performance Prediction. In Proceedings of the
 ACM SIGIR International Conference on Theory of Information Retrieval, ICTIR 2017, Amsterdam, The Netherlands, October 1-4, 2017, Jaap Kamps,
 Evangelos Kanoulas, Maarten de Rijke, Hui Fang, and Emine Yilmaz (Eds.). ACM, 245–248. https://doi.org/10.1145/3121050.3121087
- 1200 [52] Dwaipayan Roy, Debasis Ganguly, Mandar Mitra, and Gareth J. F. Jones. 2019. Estimating Gaussian mixture models in the local neighbourhood of embedded word vectors for query performance prediction. *Inf. Process. Manag.* 56, 3 (2019), 1026–1045. https://doi.org/10.1016/j.ipm.2018.10.009
 - [53] Anna Shtok, Oren Kurland, and David Carmel. 2010. Using statistical decision theory and relevance models for query-performance prediction. In Proceeding of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2010, Geneva, Switzerland, July 19-23, 2010, Fabio Crestani, Stéphane Marchand-Maillet, Hsin-Hsi Chen, Efthimis N. Efthimiadis, and Jacques Savoy (Eds.). ACM, 259–266. https://doi.org/10.1145/1835449.1835494
- [54] Anna Shtok, Oren Kurland, and David Carmel. 2016. Query Performance Prediction Using Reference Lists. ACM Trans. Inf. Syst. 34, 4 (2016),
 19:1–19:34. https://doi.org/10.1145/2926790
 - [55] Anna Shtok, Oren Kurland, David Carmel, Fiana Raiber, and Gad Markovits. 2012. Predicting Query Performance by Query-Drift Estimation. ACM Trans. Inf. Syst. 30, 2 (2012), 11:1–11:35. https://doi.org/10.1145/2180868.2180873
- [56] Ashutosh Singh, Debasis Ganguly, Suchana Datta, and Craig MacDonald. 2023. Unsupervised Query Performance Prediction for Neural Models
 with Pairwise Rank Preferences. In Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval,
 SIGIR 2023, Taipei, Taiwan, July 23-27, 2023, Hsin-Hsi Chen, Wei-Jou (Edward) Duh, Hen-Hsen Huang, Makoto P. Kato, Josiane Mothe, and Barbara
 Poblete (Eds.). ACM, 2486-2490. https://doi.org/10.1145/3539618.3592082
 - [57] Yongquan Tao and Shengli Wu. 2014. Query Performance Prediction By Considering Score Magnitude and Variance Together. In Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management, CIKM 2014, Shanghai, China, November 3-7, 2014, Jianzhong Li, Xiaoyang Sean Wang, Minos N. Garofalakis, Ian Soboroff, Torsten Suel, and Min Wang (Eds.). ACM, 1891–1894. https://doi.org/10.1145/2661829.2661906
- 1216 [58] John W. Tukey. 1949. Comparing Individual Means in the Analysis of Variance. Biometrics 5, 2 (1949), 99-114. http://www.jstor.org/stable/3001913
- [59] C. J. van Rijsbergen, 1979, Information Retrieval, Butterworth.
- 1218 [60] Ellen Voorhees. 2005. Overview of the TREC 2004 Robust Retrieval Track. https://doi.org/10.6028/NIST.SP.500-261
- [61] Frank Wilcoxon. 1945. Individual Comparisons by Ranking Methods. Biometrics Bulletin 1, 6 (1945), 80–83. http://www.jstor.org/stable/3001968
- [62] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. 2021. Approximate Nearest
 Neighbor Negative Contrastive Learning for Dense Text Retrieval. In 9th International Conference on Learning Representations, ICLR 2021, Virtual
 Event, Austria, May 3-7, 2021. OpenReview.net. https://openreview.net/forum?id=zeFrfgyZln
- [63] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. 2021. Approximate Nearest
 Neighbor Negative Contrastive Learning for Dense Text Retrieval. In 9th International Conference on Learning Representations, ICLR 2021, Virtual

 Event, Austria, May 3-7, 2021. OpenReview.net. https://openreview.net/forum?id=zeFrfgyZln
- [64] Liu Yang, Qingyao Ai, Jiafeng Guo, and W. Bruce Croft. 2016. aNMM: Ranking Short Answer Texts with Attention-Based Neural Matching Model.
 In Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, USA, October 24-28,
 2016, Snehasis Mukhopadhyay, ChengXiang Zhai, Elisa Bertino, Fabio Crestani, Javed Mostafa, Jie Tang, Luo Si, Xiaofang Zhou, Yi Chang, Yunyao
 Li, and Parikshit Sondhi (Eds.). ACM, 287-296. https://doi.org/10.1145/2983323.2983818
- [65] Hamed Zamani, W. Bruce Croft, and J. Shane Culpepper. 2018. Neural Query Performance Prediction using Weak Supervision from Multiple Signals.
 In The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08-12, 2018,
 Kevyn Collins-Thompson, Qiaozhu Mei, Brian D. Davison, Yiqun Liu, and Emine Yilmaz (Eds.). ACM, 105-114. https://doi.org/10.1145/3209978.
 3210041
 - [66] Oleg Zendel, Anna Shtok, Fiana Raiber, Oren Kurland, and J. Shane Culpepper. 2019. Information Needs, Queries, and Query Performance Prediction. In Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019, Benjamin Piwowarski, Max Chevalier, Éric Gaussier, Yoelle Maarek, Jian-Yun Nie, and Falk Scholer (Eds.). ACM, 395–404. https://doi.org/10.1145/3331184.3331253
- [67] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. Optimizing Dense Retrieval Model Training with Hard
 Negatives. In Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (<conf-loc>),
 city>Virtual Event</city>, <country>Canada</country>, </conf-loc>) (SIGIR '21). Association for Computing Machinery, New York, NY, USA,
 1503–1512. https://doi.org/10.1145/3404835.3462880
- 1240 [68] Wayne Xin Zhao, Jing Liu, Ruiyang Ren, and Ji-Rong Wen. 2022. Dense Text Retrieval based on Pretrained Language Models: A Survey. CoRR
 1241 abs/2211.14876 (2022). https://doi.org/10.48550/ARXIV.2211.14876 arXiv:2211.14876
- [69] Ying Zhao, Falk Scholer, and Yohannes Tsegay. 2008. Effective Pre-retrieval Query Performance Prediction Using Similarity and Variability Evidence.
 In Advances in Information Retrieval, 30th European Conference on IR Research, ECIR 2008, Glasgow, UK, March 30-April 3, 2008. Proceedings (Lecture
 Notes in Computer Science, Vol. 4956), Craig Macdonald, Iadh Ounis, Vassilis Plachouras, Ian Ruthven, and Ryen W. White (Eds.). Springer, 52–64.
 https://doi.org/10.1007/978-3-540-78646-7_8
- [70] Yun Zhou and W. Bruce Croft. 2007. Query performance prediction in web search environments. In SIGIR 2007: Proceedings of the 30th Annual
 International ACM SIGIR Conference on Research and Development in Information Retrieval, Amsterdam, The Netherlands, July 23-27, 2007, Wessel
 Kraaij, Arjen P. de Vries, Charles L. A. Clarke, Norbert Fuhr, and Noriko Kando (Eds.). ACM, 543-550. https://doi.org/10.1145/1277741.1277835