

Finding Driver Pathways in Cancer: Models and Algorithms

Fabio Vandin^{1,2}, Eli Upfal^{1,2}, and Benjamin J. Raphael^{1,2}

¹ Department of Computer Science, and

² Center for Computational Molecular Biology, Brown University, Providence, RI.
{vandinfa,eli,braphael}@cs.brown.edu

Abstract Cancer sequencing projects are now measuring somatic mutations in large numbers of cancer genomes. A key challenge in interpreting these data is to distinguish *driver mutations*, mutations important for cancer development, from *passenger mutations* that have accumulated in somatic cells but without functional consequences. A common approach to identify genes harboring driver mutations is a *single gene test* that identifies individual genes that are mutated in a significant number of cancer genomes. However, the power of this test is reduced by the mutational heterogeneity in most cancer genomes and by the necessity of estimating the *background mutation rate* (BMR). We investigate the problem of discovering *driver pathways*, groups of genes containing driver mutations, directly from cancer mutation data and without prior knowledge of pathways or other interactions between genes. We introduce two generative models of somatic mutations in cancer and study the algorithmic complexity of discovering driver pathways in both models. We show that a single gene test for driver genes is highly sensitive to the estimate of the BMR. In contrast, we show that an algorithmic approach that maximizes a straightforward measure of the mutational properties of a driver pathway successfully discovers these groups of genes without an estimate of the BMR. Moreover, this approach is also successful in the case when the observed frequencies of passenger and driver mutations are indistinguishable, a situation where single gene tests fail.

1 Introduction

Cancer is a disease driven in part by somatic mutations that accumulate during the lifetime of an individual. These mutations include single nucleotide substitutions, small indels, and larger copy number aberrations and structural aberrations. A key challenge in cancer genomics is to distinguish *driver mutations*, mutations important for cancer development, from random *passenger mutations* that have accumulated in somatic cells but do not have functional consequences. Recent advances in DNA sequencing technologies allow the measurement of somatic mutations in large numbers of cancer genomes. Thus, a common approach to identify driver mutations, and the driver genes in which they reside, is to identify genes with recurrent mutations in a large cohort of cancer patients. The standard approach to identify such recurrently mutated genes is to perform a

single gene test, in which individual genes are tested to determine if their observed frequency of mutation is significantly higher than expected [11,5,12]. This approach has identified a number of important cancer genes, but has not revealed all of the driver mutations and driver genes in individual cancers.

There are two difficulties with the identification of driver genes by a single gene test of recurrent mutation. First, the test requires a reasonable estimate of the *background mutation rate* (BMR), which quantifies the accumulation of passenger mutations. Obtaining such an estimate is not a straightforward task, as the BMR is not just the rate of somatic mutation per nucleotide per cell generation, but also must account for selection and clonal amplification in the somatic evolution of a tumor [7,11]. Second, it is widely observed that there is extensive mutational heterogeneity in cancer, with mutations occurring in different genes in different patients. This mutational heterogeneity is a consequence of both the presence of passenger mutations in each cancer genome, and the fact that driver mutations typically target genes in cellular signaling and regulatory pathways [8,16]. Since each of these pathways contains multiple genes, there are numerous combinations of driver mutations that can perturb a pathway important for cancer. This mutational heterogeneity inflates the number of patients required to distinguish passenger from driver mutations, as rare driver mutations may not be observed at frequencies above the background. Thus, a common alternative to single gene tests is to test the recurrence of mutations in groups of genes derived from known pathways [6,2] or genome-scale gene interaction networks [3,13]. However, these approaches require prior knowledge of the interactions between genes/proteins, and this knowledge is presently far from complete. Moreover, pathway/network based approaches typically also require an estimate of the BMR.

The availability of somatic mutation data from increasing numbers of cancer patients motivates the question of whether is it possible to identify *driver pathways*, groups of genes with recurrent driver mutations, *de novo*; i.e. without prior knowledge of interactions between genes/proteins. At first glance, this seems implausible because there are an enormous number of possible sets of genes to consider. For example, there are more than 10^{26} sets of 7 human genes. However, we previously showed that mild additional constraints on the expected patterns of somatic mutations considerably reduce the number of gene sets to examine, and make *de novo* discovery of driver pathways possible [14]. These constraints are consistent with the current understanding of the somatic mutational process of cancer [9,16]. In particular, we assume that an important cancer pathway should be perturbed in a large number of patients. Thus, given genome-wide measurements of somatic mutations, we expect that a driver pathway will have high *coverage*: i.e. most patients will have a mutation in some gene in the pathway. Second, a driver mutation in a single gene of the pathway is often assumed to be sufficient to perturb the pathway. Combined with the fact that driver mutations are relatively rare, most patients exhibit only a single driver mutation in a pathway. Thus, we expect that the genes in a pathway ex-

hibit a pattern of *mutually exclusive* driver mutations, where driver mutations are observed in exactly one gene in the pathway in each patient [17].

Note that the exclusivity constraint is assumed only for driver mutations in the *same* pathway. As a cancer genome likely has multiple driver pathways, the exclusivity assumption does not preclude the presence of co-occurring, and possibly cooperative, mutations, examples of which are known [15,4]. It is also possible that co-occurring mutations are necessary to perturb a pathway. In this case, there will likely remain a large subset of genes in the pathway whose mutations are exclusive, e.g. a subset obtained by removing one gene from each co-occurring pair. The identification of these subsets of genes by the approaches described here can be a starting point to later identify the other genes with co-occurring mutations.

1.1 Our Contribution

This work proposes a mathematical framework to study the problem of *de novo* discovery of driver genes and pathways. We define two generative models of driver mutations in cancer and study the algorithmic complexity of the discovery problem in each of the models, both analytically and in simulations. The two generative models differ in how conditioning on a sample being from a cancer patient affects the ratio between the driver and passenger mutation probabilities in that sample. While the difference is relatively small, it has a major implication on the practicality of the standard single gene test for identifying the driver genes. In the first model we prove a bound on the number of patients required to detect all driver genes with high probability using a single gene test, while in the second model it is not possible to identify the driver genes using such a test for *any* number of patients.

Next, we study a weight function on sets of genes that quantifies the coverage and exclusivity properties of a driver pathway. We introduced this function in [14], and showed that finding sets with high weight provides an alternative approach for identifying driver mutations. Here, we prove that for both generative models, when mutation data from enough patients is available, the weight function is monotone in the number of discovered driver genes and is maximized by the driver pathway. Based on this observation we prove that a simple greedy algorithm identifies the driver pathways with high probability. This improves the result in [14], where we showed that the discovery problem is NP-hard for arbitrary mutation data and that a greedy algorithm performs well under different conditions that did not arise from a generative model of the data. We also show that our earlier Markov Chain Monte Carlo (MCMC) approach for identifying the driver pathways rapidly converges to the driver pathway in both generative models, thus improving the convergence result of [14] for arbitrary mutation data. These results show that we can identify driver pathways *without* an estimate of the background mutation rate (BMR), giving a more reliable and robust solution for the problem.

We complement our analytical results with experiments on simulated data from the first model. We compare the number of patients required to identify

driver genes using the single gene test with the number required using the greedy algorithm that maximizes the weight function. We show that the number of patients is similar when a perfect estimate of the BMR is available, but that the greedy algorithm requires a smaller number of patients when the estimate of the BMR deviates from its real value. Our analytical and experimental results help characterize the limitations of detecting driver genes and pathways under reasonable models of somatic mutation.

2 Stochastic Models for Somatic Mutations in Cancer

In this section we introduce two stochastic models for somatic mutations in cancer. In both models driver mutations occur in *sets* of genes, which we refer to as *driver pathways*. Passenger mutations occur randomly across all genes. We assume that mutations have been measured in n genes in a collection of m cancer patients, and represent the somatic mutations as a $m \times n$ binary mutation matrix A . The entry A_{ig} in row i and column g is equal to 1 if gene g is mutated in patient i , and it is 0 otherwise. Let \mathcal{G} be the set of all columns (genes). In both models, we assume that the mutation matrix contains a *driver pathway*: a subset $\mathcal{D} \subseteq \mathcal{G}$ of genes, with $|\mathcal{D}| = k$, such that in each patient *exactly one* of the genes of \mathcal{D} contains a driver mutation. Thus, a driver pathway \mathcal{D} exhibits the properties of high *coverage* – every patient has a mutation in a gene in \mathcal{D} – and *mutual exclusivity* – no patient has a driver mutation in more than one gene in \mathcal{D} . In both models, random *passenger* mutations occur at random in all genes, including genes in \mathcal{D} . The difference between the two models is in the relative mutation rates in driver and passenger genes. In the following we consider the case in which the mutation matrix contains only one driver pathway. However, our results can be generalized to the case of multiple disjoint driver pathways. In particular the following iterative procedure identifies all driver pathways using our algorithms: once we identify a driver pathway, we remove its genes from the mutation matrix, and look for driver pathways in the reduced mutation matrix.

Following the hypothesis that cancer is triggered by a mutation in a driver gene, the sample of cancer patients can be viewed as a subset of a larger initial population. The genome of each member of the initial population was subject to random mutations, where each gene was mutated independently, and our sample is the subset of the initial population with a driver mutation in a gene of \mathcal{D} .

The first stochastic model captures the above intuition by modeling the distribution of mutations in patients as independent with fixed probability q , conditioning on having a driver mutation. The mutation matrix A is generated by the following process: in each row (patient) we choose one gene $d \in \mathcal{D}$ uniformly at random to contain the driver mutation, and set the corresponding entry A_{id} to 1. All other entries at that row are set to 1 with probability $q < 1$ and to 0 otherwise, and all events are independent. We call the parameter q the *passenger*

*mutation probability*³, as it is the probability that a gene contains a passenger mutation. We denote the model above as the D>P model.

A possible limitation of the D>P model is that it implies a conditional distribution in which driver genes have higher expected frequency of mutation than the passenger genes (thus the name D>P model) in a cohort of patients. In practice the driver pathway could contain dozens of genes, and some of them may have rare driver mutations. Thus the expected frequency of mutation of some genes in \mathcal{D} may be indistinguishable from the expected frequency of mutation of some passenger genes. To examine this situation we introduce a second model, which we call the D=P model, in which all genes in \mathcal{G} are mutated with the same probability in the patients, regardless of whether they are driver or passenger genes. Of course, this is a “worst case” model, as any cancer cohort with a reasonable number of patients will have some driver genes mutated at appreciable frequency. Nevertheless, we study the D=P model to consider the limits of driver pathway identification. The mutation matrix A in the D=P model is generated by the following process: in row (patient) i an entry A_{id} is chosen uniformly at random for $d \in \mathcal{D}$ and is set to 1. All other entries $A_{id'}$ for $d' \in \mathcal{D}$ are set to 1 with probability $r = \frac{qk-1}{k-1}$, and all entries A_{ig} , for $g \in \mathcal{G} \setminus \mathcal{D}$ are set to 1 with probability q . All events are independent. We require $q \geq 1/k$ so that r is a proper probability. Note that for any $g \in \mathcal{G}$ the probability that g is mutated is the same since for $d \in \mathcal{D}$, $\frac{1}{k} + (1 - \frac{1}{k})r = q$.

Note that both models differ from a simple *binomial* model, where each entry of A is mutated independently with a fixed probability. Since we condition on each patient having at least one mutation in \mathcal{D} , the entries of A corresponding to genes in \mathcal{D} are not independent. In what follows, we let $\Gamma(g) = \{i : A_{ig} = 1\}$ denote the set of patients in which a gene g is mutated. Similarly, for a set M of genes, let $\Gamma(M)$ denote the set of patients in which at least one of the genes in M is mutated: $\Gamma(M) = \cup_{g \in M} \Gamma(g)$.

3 Finding Recurrently Mutated Genes

The standard approach to identify the driver genes is to identify recurrently mutated genes, i.e. those genes whose observed frequency of mutations is significantly higher than the expected *passenger mutation probability*[11,5,12]. This approach assumes a prior knowledge or a good estimate of the passenger mutation probability, the parameter q in our models. This approach is combined with a multi-hypothesis test to identify a list of genes, each mutated in significantly more patients than expected. The pseudocode for such a test is given in Algorithm RMG (Figure 1). (In Algorithm RMG we use Bonferroni correction for multiple hypothesis testing. Other corrections, like Benjamini-Hochberg [1] to control the *False Discovery Rate*, are possible. The results of this section also apply to those other corrections.)

³ Note that q is greater than the BMR, since it is the probability that a *gene* has a passenger mutation. For example, estimates of the BMR are typically $\approx 10^{-5}$, and since the length of most genes is around 10^4 , we have that $q \approx 10^{-1}$.

Algorithm RMG

Input: An $m \times n$ mutation matrix A , a probability q that a gene contains a passenger mutation in a patient, a significance level α .

Output: Set \mathcal{O} of recurrently mutated genes.

```
1  $\mathcal{O} \leftarrow \emptyset$ ;  
2 for  $g \in \mathcal{G}$  do  
3    $\Gamma(g) \leftarrow \{i : A_{ig} = 1\}$ ;  
4    $p_g \leftarrow \Pr[B(m, q) \geq |\Gamma(g)|]$ ;  
5   if  $p_g \leq \frac{\alpha}{n}$  then  $\mathcal{O} \leftarrow \mathcal{O} \cup \{g\}$ ;  
6 return  $\mathcal{O}$ ;
```

Figure 1: Pseudocode of the algorithm for finding recurrently mutated genes, based on a single-gene test.

We first analyze the D>P model of Section 2. We start by showing that if q is known and the number of patients is sufficiently large, then Algorithm RMG outputs all the driver genes with high probability.

Theorem 1. *Suppose an $m \times n$ mutation matrix A is generated by the the D>P model, the family wise error rate of the test is $\alpha = \frac{1}{2n^\varepsilon}$ and Algorithm RMG outputs \mathcal{O} . If $m \geq \frac{2k^2(1+\varepsilon)}{(1-q)^2} \ln 2n$ for a constant $\varepsilon > 0$, then $\Pr[\mathcal{O} \neq \mathcal{D}] \leq \frac{1}{n^\varepsilon}$.*

Proof. The p -value calculations and the Bonferroni correction in Algorithm RMG guarantee that the probability that any gene $g \notin \mathcal{D}$ is included in the output set \mathcal{O} is bounded by $\alpha = \frac{1}{2n^\varepsilon}$. It remains to prove that if $m \geq \frac{2k^2(1+\varepsilon)}{(1-q)^2} \ln 2n$ the probability that any $d \in \mathcal{D}$ is not included in \mathcal{O} is bounded by $\frac{1}{2n^\varepsilon}$.

Consider a gene $d \in \mathcal{D}$. Let $X_i = 1$ if gene d is mutated in patient i , and $X_i = 0$ otherwise. Note that for $i \neq j$, X_i and X_j are independent. Let X be the number of patients in which d is mutated. We have $X = \sum_{i=1}^m X_i$. To compute $\mathbf{E}[X_i]$ we observe that a driver gene is mutated with probability 1 when it contains the driver mutation, and with probability q otherwise. Since the gene d containing the driver mutation is chosen uniformly at random among all the k genes in \mathcal{D} , we have $\mathbf{E}[X_i] = \frac{1}{k} + (1 - \frac{1}{k})q$. Thus $\mathbf{E}[X] = \sum_{i=1}^m \mathbf{E}[X_i] = m(\frac{1}{k} + (1 - \frac{1}{k})q) > mq$. Let $t = \frac{1}{k}(\frac{1-q}{2})$. By the Chernoff-Hoeffding bound:

$$\Pr[X \leq \mathbf{E}[X] - tm] = \Pr[X \leq m\mathbf{E}[X_i] - tm] \leq e^{-\frac{2m^2t^2}{m}} \leq \frac{1}{2n^{1+\varepsilon}}.$$

Since $|\mathcal{D}| < n$, by union bound we have:

$$\Pr[\exists d \in \mathcal{D} \text{ mutated in } \leq (\mathbf{E}[X] - tm) \text{ patients}] \leq n \frac{1}{2n^{1+\varepsilon}} = \frac{1}{2n^\varepsilon}.$$

Thus with probability at least $1 - \frac{1}{2n^\varepsilon}$ all genes in \mathcal{D} are mutated in at least $\mathbf{E}[X] - tm$ patients. Let $B(m, q)$ be a binomial random variable with parameters

m, q . Using the Chernoff-Hoeffding bound we can upper bound the p -value p_d that Algorithm RMG derives for $d \in D$:

$$p_d \leq \Pr[|B(m, q) - mq| \geq tm] \leq e^{-2\frac{t^2 m^2}{m}} \leq \frac{1}{2n^{1+\varepsilon}}.$$

Thus, with probability at least $1 - \frac{1}{2n^\varepsilon}$ for any $d \in \mathcal{D}$ the number of patients with a mutation in d is such that its p -value satisfies $p_d < \alpha/n$ and thus it is included in the output set \mathcal{O} . \square

Theorem 1 shows that in the D>P model an estimate of the passenger mutation probability q and a sufficient number of patients are enough to identify the driver genes. This is not the case in the D=P model. It is easy to see that in D=P model the expected number of rows in which a column g is mutated is the same for all $g \in \mathcal{G}$, that is for all $g \in \mathcal{G}$ we have $\mathbf{E}[|\Gamma(g)|] = qm$. In fact, the number $|\Gamma(d)|$ of patients in which a gene $d \in \mathcal{D}$ is mutated and the number $|\Gamma(g)|$ of patients in which gene $g \notin \mathcal{D}$ is mutated are both binomial random variables $B(m, q)$. We thus have the following.

Fact 1 *Under the D=P model, the probability distribution of $|\Gamma(d)|$ for $d \in \mathcal{D}$ and $|\Gamma(g)|$ for $g \notin \mathcal{D}$ are the same. Thus Algorithm RMG cannot identify the genes in \mathcal{D} for any number of patients m .*

4 A Weight Function to Identify Driver Pathways

In this section we analyze a method that identifies the set \mathcal{D} of driver genes with no prior information on the passenger mutation probability q , and works for both the D>P and D=P models. The method relies on a weight function $W(M)$, defined on sets of genes, first introduced in [14]. The measure W quantifies the extent to which a set simultaneously exhibits both: (i) high *coverage*: most patients have at least one mutation in the set; (ii) high *exclusivity*: nearly all patients have no more than one mutation in the set. (For lack of space, some proofs of the results in this section are omitted. They will be included in the full version of this work.)

For a set of genes, M , we define the coverage overlap $\omega(M) = \sum_{g \in M} |\Gamma(g)| - |\Gamma(M)|$. Note that $\omega(M) \geq 0$, with equality if and only if the mutations in M are mutually exclusive. To account for both the coverage, $\Gamma(M)$, and the coverage overlap, $\omega(M)$, we define the weight function of M :

$$W(M) = |\Gamma(M)| - \omega(M) = 2|\Gamma(M)| - \sum_{g \in M} |\Gamma(g)|.$$

Finding a set M of genes with maximum weight is in general a computationally challenging problem (it is NP-hard in the worst case). Nonetheless, we showed in [14] that under some assumptions on the distribution of mutations in patients, a greedy algorithm will identify the maximum weight set. We also proposed a

Markov Chain Monte Carlo (MCMC) approach that samples sets of genes with probability proportional to their weight.

Based on the coverage and exclusivity properties of a driver pathway we expect it has the highest weight among all sets of size k . In this section we formalize this intuition for our generative models and show that under the two models the maximum weight set is easy to compute. We use M_k^* to denote the set of size k with maximum weight (M_k^* may not be unique).

We start with the D>P model. Note that the parameter q controls the expected number of passenger mutations in a set of k passenger genes. Since passenger mutations are relatively rare and k (the number of genes in a driver pathway) is relatively small, we expect that a set of k passenger genes will not have a mutation in the majority of the patients. Thus we assume that the probability $1 - (1 - q)^k$ that a set of k passenger genes contains a mutation is less than a constant $a < \frac{1}{2}$. Since $1 - (1 - q)^k \approx qk$ we have $q \leq \frac{a}{k}$. For ease of exposition in what follows we use $a = \frac{1}{4}$, so that $q \leq \frac{1}{4k}$.

Let $M_{k,\ell} \subset \mathcal{G}$ be a set of k genes with exactly ℓ genes of \mathcal{D} , that is $M_{k,\ell} = \{d_1, d_2, \dots, d_\ell\} \cup \{g_1, \dots, g_{k-\ell}\}$ with $d_j \in \mathcal{D}$ for $1 \leq j \leq \ell$, and $g_j \in \mathcal{G} \setminus \mathcal{D}$ for $1 \leq j \leq k - \ell$. We first prove that $\mathbf{E}[W(M_{k,\ell})]$ is monotone in ℓ .

Lemma 1. *Let $q \leq \frac{1}{4k}$. For $0 \leq \ell \leq k - 1$: $\mathbf{E}[W(M_{k,\ell+1})] \geq \mathbf{E}[W(M_{k,\ell})] + \frac{m}{2k}$.*

Next we show that for sufficiently large number of patients m , the random value $W(M_{k,\ell})$ is concentrated near its expectation.

Theorem 2. *Suppose A is generated by the D>P model with $q \leq \frac{1}{4k}$. For $m \geq 8k^3(k + \varepsilon) \ln n$, and for $0 \leq \ell \leq k - 1$, $\Pr[\exists M_{k,\ell} \text{ s.t. } |W(M_{k,\ell}) - \mathbf{E}[W(M_{k,\ell})]| \geq \frac{m}{4k}] \leq \frac{1}{n^\varepsilon}$.*

Combining the results of Lemma 1 and Theorem 2 we have

Corollary 1. *If $m \geq 8k^3(k + \varepsilon) \ln n$, then $\Pr[M_k^* \neq \mathcal{D}] \leq \frac{1}{n^\varepsilon}$.*

Corollary 1 shows that with sufficient number of patients the set \mathcal{D} can be identified by finding the set of maximum weight, without an estimate of the probability q that a gene is mutated as a passenger. It was shown in [14] that with an arbitrary mutation distribution identifying the set of maximum weight is NP-Hard. However, a simple corollary of Theorem 2 shows that in our generative model computing a set of maximum weight is easy.

Corollary 2. *If $m \geq 8k^3(k + \varepsilon) \ln n$ and $q \leq \frac{1}{4k}$, a greedy algorithm that computes the weight function of up to $O(nk)$ sets finds M_k^* with failure probability $\leq \frac{1}{n^\varepsilon}$.*

Proof. Start with an arbitrary set M of k genes. Now consider the elements of $M = \{g_1, \dots, g_k\}$ one after the other in a greedy process: for $g_j \in M$, find $w = \arg \max_{g \in \mathcal{G} \setminus M} W(M \setminus \{g_j\} \cup \{g\})$. If $W(M) < W(M \setminus \{g_j\} \cup \{w\})$, substitute g_j with w in M ; then move to g_{j+1} . Theorem 2 guarantees that if w is inserted in M , it is in \mathcal{D} , and that when a gene $g_j \in M \setminus \mathcal{D}$ is considered, it will be switched with a gene $d \in \mathcal{D} \setminus M$. \square

We now consider the D=P model. Analogously to what we proved under the D>P model, we prove that maximizing the weight function W identifies the driver pathway \mathcal{D} when mutation data from enough patients is available.

Theorem 3. *Suppose A is generated by the D=P model. If $m \geq \frac{k^3(k+\varepsilon)}{2(1-q)^{2k+2}} \left(\frac{k-1}{k}\right)^{2k} \ln n$, then $\Pr[M_k^* \neq \mathcal{D}] \leq \frac{1}{n^\varepsilon}$.*

We prove that a simple greedy algorithm, similar to the one proposed for the D>P model, identifies the set M_k^* of maximum weight under the D=P model.

Corollary 3. *If $m \geq \frac{k^3(k+\varepsilon)}{2(1-q)^{2k+2}} \left(\frac{k-1}{k}\right)^{2k} \ln n$, a greedy algorithm that computes the weight function of up to $O(n^2)$ sets finds M_k^* with failure probability $\leq \frac{1}{n^\varepsilon}$.*

Thus under the D=P model we identify the driver pathway \mathcal{D} by maximizing $W(M)$. Recall that Algorithm RMG cannot find driver genes under this model (Section 3, Fact 1). Also note that when $q \leq 1/2$ and the probability $(1-q)^k$ that a set of k genes in $\mathcal{G} \setminus \mathcal{D}$ is not mutated in a patient is greater than $\frac{1}{2} \left(\frac{k-1}{k}\right)^k$ (this occurs when passenger mutations are relatively rare, for example when $q \approx 1/k$) the bound on m in Corollary 3 is the same as the bound in Corollary 2. That is, the weight W identifies the set \mathcal{D} under both models with the same number of patients.

For completeness, we also analyze the Monte-Carlo Markov Chain approach proposed in [14] to sample sets of genes with distribution exponentially proportional to their weight. The states of the Markov chain are the subsets of \mathcal{G} of size k . If $M^{(t)}$ is the state at time t , $M^{(t+1)}$ is computed choosing uniformly at random a gene $w \in \mathcal{G}$ and a gene $v \in M^{(t)}$, and setting $M^{(t+1)} = M^{(t)} \setminus \{v\} \cup \{w\}$ with probability $\min[1, e^{cW(M^{(t)} \setminus \{v\} \cup \{w\}) - cW(M^{(t)})}]$, and $M^{(t+1)} = M^{(t)}$ otherwise. It is easy to verify that the chain is ergodic with a unique stationary distribution $\pi(M) = \frac{e^{cW(M)}}{\sum_{R \in \mathcal{M}_k} e^{cW(R)}}$, where $\mathcal{M}_k = \{M \subset \mathcal{G} \mid |M| = k\}$. The efficiency of this algorithm depends on the speed of convergence of the Markov chain to its stationary distribution.

In [14], we show that there is a non-trivial interval of values for c for which the chain is rapidly mixing without assuming any generative model for the mutation matrix. The analysis of [14] applied to D>P and D=P models requires $c < 1/k$. However, applying Lemma 1 and 2 under the D>P model, and Theorem 3 under the D=P model we show that for any $c > 0$ the process rapidly converges to the set \mathcal{D} .

Theorem 4. *Suppose that A is generated by the D>P model with $q \leq \frac{1}{4k}$, or the D=P model with $q \leq 1/2$ and $(1-q)^k \geq \frac{1}{2} \left(\frac{k-1}{k}\right)^k$. For $m \geq 8k^3(k+\varepsilon) \ln n$ and any $c > 0$, the MCMC converges to the set \mathcal{D} in $O(nk \log k)$ iterations with probability $\geq 1 - \frac{1}{n^\varepsilon}$.*

5 Experimental Results

In this section we compare the single gene test provided in Algorithm RMG and the weight function $W(M)$ to detect the set of driver genes using mutation data simulated using the D>P model. In particular, we use the greedy algorithm of Section 4 (see Corollary 2) to identify the set M_k^* of maximum weight, where $k = |\mathcal{D}|$.

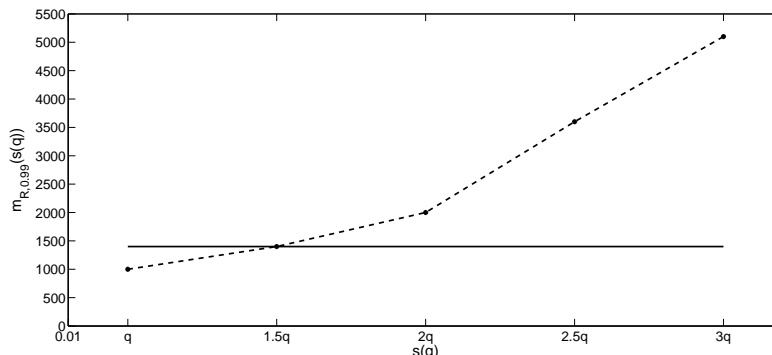


Figure 2: Number of patients $m_{R,0.99}(s(q))$ required to identify the driver pathway \mathcal{D} with Algorithm RMG, for different estimates $s(q)$ of the probability q (dashed). Number of patients $m_{G,0.99}$ required to identify \mathcal{D} with the greedy algorithm (solid).

We generated mutation data according to the D>P model with $k = |\mathcal{D}| = 20$, $q = 0.0125$, $n = 10000$. We set $\alpha = 0.005$ for Algorithm RMG which corresponds to $\varepsilon = 0.5$. To compare the performance of the two approaches, we measured the minimum number of patients required to detect the driver pathway \mathcal{D} over a range of estimates of the passenger mutation probability q . Specifically, let $E_{s(q)} = \text{“estimate } s(q) \text{ of } q \text{ is used by Algorithm RMG”}$. Let $m_{R,x}(s(q)) = \min_m \{\Pr[\mathcal{O} = \mathcal{D} | E_{s(q)}] > x\}$ be the minimum number of patients required for Algorithm RMG to output $\mathcal{O} = \mathcal{D}$ with probability $> x$ over all $m \times n$ mutation matrices generated by the model when the estimate $s(q)$ is used. Similarly, let \mathcal{P} be the output of the greedy algorithm of Corollary 2. Let $m_{G,x} = \min_m \{\Pr[\mathcal{P} = \mathcal{D}] > x\}$ be the minimum number of patients required for the greedy algorithm to output \mathcal{D} with probability $> x$ over all $m \times n$ mutation matrices generated by the model. Recall that $m_{G,x}$ does not depend on $s(q)$ by Corollary 2. Figure 2 shows the values of $m_{R,0.99}(s(q))$ and $m_{G,0.99}$ as a function of $s(q)$. We varied $s(q)$ starting from $s(q) = q$ (i.e., q is perfectly estimated) and gradually increased $s(q)$ while maintaining $s(q) < 1/k$. The latter condition assures that $s(q)$ is strictly smaller than the expected probability of mutation of any gene in \mathcal{D} , a necessary condition for Algorithm RMG to be able to identify \mathcal{D} . To estimate $m_{R,0.99}$ and $m_{G,0.99}$ we generated 100 mutation matrices for each $m_i = i * 100$ patients for $1 \leq i \leq 52$. Figure 2 shows that $m_{R,0.99}(s(q))$ is monotonically increasing with $s(q)$. When the estimate of q is perfect, the greedy algorithm requires more patients than Algorithm RMG to correctly identify the set \mathcal{D} , but when the estimate $s(q)$ is larger than the true value of q , $m_{R,0.99}(s(q))$ increases and

becomes much larger than $m_{G,0.99}$. (Typically, an overestimate of q is used so that the test for recurrent genes in conservative [10]). Note that even when $s(q) = q$, $m_{G,0.99}$ is close to $m_{R,0.99}(q)$, while the bounds in Theorem 1 and Corollary 2 give $\frac{m_{G,0.99}}{m_{R,0.99}(s(q))} \geq 1000$. Similar results were obtained when comparing $m_{R,0.95}(s(q))$ and $m_{G,0.95}$; i.e. the minimum number of patients for which Algorithm RMG and the greedy algorithm report the driver set \mathcal{D} at least 95% of the time (data not shown).

Finally, we consider the case $s(q) < q$ where the estimate of q is smaller than its true value. In this case, some genes not in \mathcal{D} (false positives) are eventually reported by Algorithm RMG. For example, with $s(q) = 0.8q$ and $m = 1000$ (for which the correct result is always reported when $s(q) = q$), Algorithm RMG reports false positives in approximately 16% of the datasets.

6 Conclusions

We investigate the problem of detecting recurrently mutated genes and pathways using two simple generative models of driver mutations in cancer. In the first $D > P$ model, where the driver mutation probability is larger than the passenger mutation probability, we prove a bound on the number of patients required to detect all driver genes with high probability using a single gene test of recurrence. In the second $D = P$ model, where the driver mutation probability and passenger mutation probability cannot be distinguished, it is impossible to identify driver genes using the single gene test for *any* number of patients. We prove that under either model, the weight function that we defined in [14] is maximized by a driver pathway. Thus, with mutation data from enough patients, it is possible to identify driver pathways *without* an estimate of the passenger mutation probability q . In particular, we show that a simple greedy algorithm finds driver pathways with high probability. We also show that an MCMC approach converges rapidly. Finally, we present results on simulated data showing that the greedy algorithm successfully identifies the driver pathway with fewer patients than the single gene test when the estimate of q deviates from its real value.

In practice, any test that identifies driver genes by recurrent mutations requires a good estimate of q . An underestimate of q leads to false positive predictions of driver genes, while an over estimate (i.e. a conservative estimate to minimize false positives) increases the number of patients required to find driver genes. The passenger mutation probability is derived from the background mutation rate (BMR), which is difficult to measure as it depends on a number of parameters whose values are not easily determined. There has been extensive discussion in the community about appropriate ways to estimate the BMR and find recurrently mutated genes [7,11]. Therefore, methods that do not require an estimate of the BMR, as the ones we provide here, can give increased power in the discovery of driver genes. However, further study of more sophisticated mutation models is necessary. For example, we assume a constant passenger mutation probability q across all genes, but models that allow q to vary by gene would be useful in applications and warrant further investigation.

Acknowledgements

We thank the anonymous reviewers for helpful suggestions that improved the manuscript. This work is supported by NSF grant IIS-1016648, the Department of Defense Breast Cancer Research Program, the Alfred P. Sloan Foundation, and the Susan G. Komen Foundation. BJR is also supported by a Career Award at the Scientific Interface from the Burroughs Wellcome Fund.

References

1. Y. Benjamini and Y. Hochberg. Controlling the false discovery rate. *J. Royal Statistical Society*, 57:289–300, 1995.
2. S. M. Boca, K. W. Kinzler, V. E. Velculescu, B. Vogelstein, and G. Parmigiani. Patient-oriented gene set analysis for cancer mutation data. *Genome Biol.*, 11:R112, 2010.
3. E. Cerami, E. Demir, N. Schultz, B. S. Taylor, and C. Sander. Automated network analysis identifies core pathways in glioblastoma. *PLoS ONE*, 5:e8918, 2010.
4. K. Deguchi and D. G. Gilliland. Cooperativity between mutations in tyrosine kinases and in hematopoietic transcription factors in AML. *Leukemia*, 16:740–744, 2002.
5. L. Ding et al. Somatic mutations affect key pathways in lung adenocarcinoma. *Nature*, 455:1069–1075, 2008.
6. S. Efroni, R. Ben-Hamo, M. Edmonson, S. Greenblum, C. F. Schaefer, and K. H. Buetow. Detecting cancer gene networks characterized by recurrent genomic alterations in a population. *PLoS ONE*, 6:e14437, 2011.
7. G. Getz, H. Hoffing, J. P. Mesirov, T. R. Golub, M. Meyerson, R. Tibshirani, and E. S. Lander. Comment on "The consensus coding sequences of human breast and colorectal cancers". *Science*, 317:1500, 2007.
8. W. C. Hahn and R. A. Weinberg. Modelling the molecular circuitry of cancer. *Nat. Rev. Cancer*, 2:331–341, 2002.
9. F. McCormick. Signalling networks that cause cancer. *Trends Cell Biol.*, 9:M53–56, 1999.
10. G. Parmigiani et al. Response to comments on "the consensus coding sequences of human breast and colorectal cancers". *Science*, 317(5844):1500, 2007.
11. T. Sjoblom et al. The consensus coding sequences of human breast and colorectal cancers. *Science*, 314:268–274, 2006.
12. The Cancer Genome Atlas Research Network. Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature*, 455(7216):1061–8, 2008.
13. F. Vandin, E. Upfal, and B. J. Raphael. Algorithms for detecting significantly mutated pathways in cancer. *J. Comput. Biol.*, 18:507–522, 2011.
14. F. Vandin, E. Upfal, and B. J. Raphael. *De novo* Discovery of Mutated Driver Pathways in Cancer. *Genome Research*, in press, 2011.
15. I. Varela et al. Exome sequencing identifies frequent mutation of the SWI/SNF complex gene PBRM1 in renal carcinoma. *Nature*, 469:539–542, 2011.
16. B. Vogelstein and K. W. Kinzler. Cancer genes and the pathways they control. *Nat. Med.*, 10:789–799, 2004.
17. C.H. Yeang, F. McCormick, and A. Levine. Combinatorial patterns of somatic gene mutations in cancer. *The FASEB Journal*, 2008.