

Every which way? On predicting tumor evolution using cancer progression models

Ramon Diaz-Uriarte, Claudia Vasallo
Dept. Biochemistry, Universidad Autónoma de Madrid
Instituto de Investigaciones Biomédicas “Alberto Sols” (UAM-CSIC)
Madrid, Spain*

2018-11-16 (Release: Rev: bea16f1)

Abstract

Successful prediction of the likely paths of tumor progression is valuable for diagnostic, prognostic, and treatment purposes. Cancer progression models (CPMs) use cross-sectional samples to identify restrictions in the order of accumulation of driver mutations and thus encode the paths of tumor progression. Here we examine whether CPMs can be used to predict the true distribution of tumor progression paths and to estimate evolutionary unpredictability. Employing simulations we show that if fitness landscapes are single peaked (have a single fitness maximum), there is good agreement between true and predicted distributions of evolutionary paths when sample sizes are large, but performance is poor with the currently common much smaller sample sizes. Under multi-peaked fitness landscapes (i.e., those with multiple fitness maxima), performance is poor and improves only slightly with sample size. In all cases, detection regime (when tumor samples are taken) is a key determinant of performance. Estimates of evolutionary unpredictability from CPMs tend to overestimate the true unpredictability and the bias is affected by detection regime; CPMs could be useful for estimating upper bounds to the true evolutionary unpredictability. Analysis of twenty two cancer data sets shows estimates of evolutionary unpredictability in regions where useful prediction might be possible for at least some data sets. But most of the evolutionary trajectory predictions themselves are very unreliable, and unreliability increases with numbers of features analyzed. Our results indicate that, currently, obtaining useful predictions of tumor progression paths from CPMs is dubious and emphasize the need for methodological work that can account for the probably multi-peaked fitness landscapes in cancer.

*To whom correspondence should be addressed: ramon.diaz@iib.uam.es, rdiaz02@gmail.com, <http://ligarto.org/rdiaz>

1 Introduction

Improving our ability to predict the paths of tumor progression is helpful for diagnostic, prognostic, and treatment purposes as, for example, it would allow us to identify genes that block the most common paths of disease progression (Greaves, 2015; Lipinski *et al.*, 2016; McPherson *et al.*, 2018; Williams *et al.*, 2018). This interest in predicting paths of progression is, of course, not exclusive to cancer (see e.g., reviews in Lässig *et al.*, 2017; Losos, 2018). For example, in some cases antibiotic resistance shows parallel evolution with mutations being acquired in a similar order (Toprak *et al.*, 2012), and here “Even a modest predictive power might improve therapeutic outcomes by informing the selection of drugs, the preference between monotherapy or combination therapy and the temporal dosing regimen (...)” (Palmer and Kishony, 2013, p. 243i). But detailed information about the paths of tumor evolution and their distribution, obtained from multiple within-patient samples with timing information, is not available.

Cancer progression models (CPMs), such as CBN (Gerstung *et al.*, 2009), CAPRI (Ramazzotti *et al.*, 2015), or OT (Szabo and Boucher, 2008), are a tool that can be used to predict paths of tumor progression. CPMs were originally developed to identify restrictions in the order of accumulation of mutations during tumor progression from cross-sectional data (Beerenwinkel *et al.*, 2015, 2016). But CPMs also implicitly encode all the possible mutational paths or trajectories of tumor progression, from the initial genotype to the genotype with all driver genes mutated, and the identification of these paths or “evolutionary trajectories” is a prominent idea in recent CPM publications (e.g. Caravagna *et al.*, 2016; Ramazzotti *et al.*, 2015). Thus, CPMs could improve our ability to predict disease progression by leveraging on the available, and growing, number of cross-sectional data sets.

The first two questions we address in this paper are whether we can predict the paths of tumor evolution using CPMs and what are the main factors that affect the quality of these predictions. To answer these questions we will examine how close to the truth are the predictions made by CPMs about the distribution of paths of tumor evolution. When addressing this question we need to take into account possible deviations from the models assumed by CPMs. In particular, CPMs assume that the acquisition of a mutation in a driver gene does not decrease the probability of gaining a mutation in another driver gene (Misra *et al.*, 2014): this implies that the fitness landscapes assumed by CPMs cannot have reciprocal sign epistasis (Diaz-Uriarte, 2018). Because acquiring driver mutations cannot decrease fitness, this also implies that the fitness landscapes assumed by CPMs only have a single global fitness maximum (the genotype with all drivers mutated). But reciprocal sign epistasis is likely to be common in cancer (Chiotti *et al.*, 2014), an argument supported by how common synthetic lethality is in both cancer cells (Beijersbergen *et al.*, 2017; O’Neil *et al.*, 2017) and the human genome (Blomen *et al.*, 2015). Moreover, if there are many combinations of a small number of drivers, out of a larger pool of drivers (Tomasetti *et al.*, 2015), that result in the escape genotype, it is likely that cancer landscapes will have several local fitness maxima (i.e., be multi-peaked). As we have shown before (Diaz-Uriarte, 2018), the performance of CPMs for predicting what genotypes can and cannot exist degrades considerably when the assumption of absence of reciprocal sign epistasis is violated. Those results, however, do not provide a direct answer to the question of predictability: if our objective is predicting paths of tumor progression we want to measure directly the quality of the predictions of paths of progression. For example, getting some of the edges of the DAG of restrictions wrong, or predicting some of the genotypes incorrectly, might be of little importance if the main paths of disease are captured. Thus, to answer the question of whether CPMs can be used to predict paths of progression we will need to look directly at the prediction of paths and do that both under scenarios where CPM’s assumptions are met and under scenarios with relevant deviations from the assumptions (see Figure 1).

And relevant scenarios bring us to the third question addressed in this paper: regardless of the performance when predicting the actual paths of tumor progression, can we use CPMs to estimate evolutionary unpredictability? Some tumors seem to follow a few, highly repeatable trajectories, whereas others show a large diversity of trajectories. So that a particular method can reconstruct well the distribution of paths of tumor progression might be of little impor-

tance if that happens in a scenario where the true evolutionary unpredictability itself is very large (where disease progression follows a very large number of possible paths); for practical purposes, forecasting here would be useless. Conversely, a method could be helpful if it suggests few paths are possible, even if its actual predictions are not trustworthy.

To address the above three questions (can we predict the paths of tumor evolution using CPMs?; what factors affect the quality of these predictions?; can we estimate evolutionary unpredictability using CPMs?) we use evolutionary simulations on 1260 fitness landscapes that include from none to severe deviations from the assumptions of CPMs, and analyze the data with six different methods for inferring CPMs. Since the role of evolutionary unpredictability is an important focus of this paper, we simulate evolution under different population sizes and mutation rates, so as to generate different amounts of evolutionary unpredictability. This paper does not attempt to understand the determinants of evolutionary predictability (see, e.g., Bank *et al.*, 2016; de Visser and Krug, 2014; Lässig *et al.*, 2017; Losos, 2018; Szendro *et al.*, 2013) but, instead, we focus on the effects of evolutionary unpredictability for CPMs. This is why we use variation in key determinants of predictability (e.g., variation in population sizes and mutation rates) but these factors, themselves, are only used to generate variability in unpredictability, and not themselves the focus of the study. To better assess the quality of predictions, we use sample sizes that cover the range from what is commonly used to what are much larger sample sizes than currently available. We also include variation in the cancer detection process or detection regime (when cancer samples are taken, or when are patients sampled), since previous studies have shown that it affects the quality of inferences (Diaz-Uriarte, 2015, 2018). Here we find that the agreement between the predicted and true distributions of paths is generally poor, unless sample sizes are very large and fitness landscapes conform to the assumptions of CPMs. Both detection regime and evolutionary unpredictability itself have major effects on performance. But in spite of the unreliability of the predictions of paths of tumor progression, CPMs can be useful for estimating upper bounds to the true evolutionary unpredictability.

What are the implications of our results for the analysis and interpretation of the use of CPMs with cancer data sets? We use 22 real cancer data sets to address these issues. We cannot examine how close predictions are to the truth, since the truth is unknown; thus, we use bootstrap samples to examine the reliability of the inferences. Many of the cancer data sets reflect conditions where useful predictions could be possible, based on the estimates of evolutionary unpredictability from CPMs. But for most data sets these results are thwarted by the unreliability of the predictions themselves, which increases with the number of features analyzed. Our results question the use of CPMs for predicting paths of tumor progression, and suggest the need for methodological work that can account for the probably multi-peaked fitness landscapes in cancer.

1.1 Assumptions

CPMs assume that the different individuals in a data set constitute independent realizations of the same evolutionary process and therefore that the same constraints hold for all tumors (Beerenwinkel *et al.*, 2015, 2016; Gerstung *et al.*, 2011). Thus, a data set can be regarded as a set of replicate evolutionary experiments where all individuals are under the same genetic constraints, though they might later be exposed to different conditions. We also make other common assumptions in this field. Briefly (see details in section 1.1 of Diaz-Uriarte, 2018) we use biallelic loci, and back mutations and crossing valleys in the fitness landscape using a single multi-mutation are not allowed (Beerenwinkel *et al.*, 2007; Bozic *et al.*, 2010; McFarland *et al.*, 2013). All tumors start cancer progression without any of the mutations considered, but other mutations could be present that have caused the initial tumor growth, so we absorb the cancer initiation process in the root node (Attolini *et al.*, 2010); this is necessary to simulate data consistent with cross-sectional sampling. The driver genes are known and there are no observational errors.

2 Materials and methods

2.1 Overview of the simulation study

We have used simulations of tumor evolution on fitness landscapes of different types (see Figure 1), for landscapes of 7 and 10 genes, under different initial population sizes and mutation rates. As explained above, variation in initial population size and mutation rates is used to generate variation in evolutionary predictability, but not of interest *per se*. We have used a total of 1260 fitness landscapes = 35 random fitness landscapes \times 2 conditions of numbers of genes \times 3 types of fitness landscapes \times 3 initial population sizes \times 2 mutation regimes. For each one of the 1260 fitness landscapes, we simulated 20000 evolutionary runs (with the specified parameters for initial population size and mutation rate) using a logistic-like growth model until one of the genotypes at the local fitness maxima (or the single global fitness maximum) reached fixation. Each set of 20000 simulated runs was then sampled under three detection regimes (so that each fitness landscape generated three sets of 20000 simulated genotypes). From each of these sets, we obtained five different splits of the genotypes for each of three sample sizes (50, 200, 4000); thus a total of 56700 (= 1260 \times 3 \times 3 \times 5 combinations of 1260 fitness landscapes, 3 detection regimes, 3 sample sizes, 5 splits) data sets were produced. Each of these 56700 data sets was analyzed with every six one the CPM methods.

2.2 Evolutionary simulations and data sampling

We used three **initial population sizes**, 2000, 50000, and 1×10^6 cells, for the simulations; these cover a range of population sizes at tumor initiation that have previously been used in the literature (e.g. Beerenwinkel *et al.*, 2007; Gerstung *et al.*, 2011; McFarland *et al.*, 2013; Wodarz and Komarova, 2014). We also used two **mutation regimes**; in the first one, all genes had a common mutation rate of 1×10^{-5} ; in the second, genes had different mutation rates, uniformly distributed in the log scale between $(1/5) 1 \times 10^{-5}$ and 5×10^{-5} (i.e., the largest ratio between largest and smallest mutation rates was 25), so that the arithmetic mean of mutation rates was 1.5×10^{-5} and the geometric mean 1×10^{-5} . These mutation rates are within ranges previously used in the literature (Bozic *et al.*, 2010; McFarland *et al.*, 2013; Nowak *et al.*, 2004), with a bias towards larger numbers (since we use only 7 or 10 genes relevant for population growth and we could be modeling pathways, not individual genes). Initial population size and mutation rates are not of intrinsic interest here (since our focus is not the determinants of evolutionary predictability *per se*), but are used to generate variability in evolutionary predictability; see section 3.1.

For each of the combinations of number of genes (7 and 10), initial population size (2000, 50000, 1×10^6), and mutation rate (constant, variable), we generated random fitness landscapes of three kinds (see Figure 1). We generated the DAG-derived **representable** fitness landscapes by generating a random DAG of restrictions and from it the fitness graph. We then assigned birth rates to genotypes using an iterative procedure on the fitness graph where, starting from the genotype without any driver mutation with a birth rate of 1, the birth rate of each descendant genotype was set equal to the maximum fitness of its parent genotypes times a random uniform variate between 1.01 and 1.19 (yielding, therefore, an average multiplicate increase in fitness of 0.1, again within values previously used; Bozic *et al.*, 2010; McFarland *et al.*, 2013; Williams *et al.*, 2018). Birth rate of genotypes without dependencies satisfied was set to 0. (Note that for the growth model used here —see below— birth rates determine fitness at any population size as death rates are identical for all genotypes and depend only on population size. Note also that genotypes with birth rate of 0 are never added to the population; thus, they cannot mutate before dying, so this simulation scheme strictly adheres to the assumptions about accessible and non-accessible genotypes under the CPM model). The DAG-derived **local-maxima** fitness landscapes were obtained by generating a random DAG and from it the fitness graph. Before assigning fitness to genotypes, a random selection of edges of the fitness graph were removed so that all accessible genotypes remained accessible but from a possibly much smaller

set of parents. Fitness was then assigned as above (with the iterative procedure on the fitness graph, where fitness of child = $\max(\text{fitness parents}) U(1.01, 1.19)$). For each DAG we repeated this procedure 50 times, and kept the one that introduced the largest number of local maxima. Creating local maxima almost always resulted in creating reciprocal sign epistasis (but see Supplementary Material, [section 1, "Generating random fitness landscapes"](#)). The local maxima fitness landscapes used in this paper are representable in the weaker sense of Diaz-Uriarte (2018), as all genotypes that should be accessible under the DAG of restrictions are accessible. What the local-maxima landscapes are missing are mutational paths to the genotype with all genes mutated, because we have introduced local fitness maxima (and once we introduce local maxima there is no longer a one-to-one correspondence between DAGs of restrictions and fitness graphs and, thus, there is no longer a one-to-one correspondence between DAGs of restrictions and sets of tumor progression paths). These local maxima landscapes are "easier" than the DAG-derived fitness landscapes used in Diaz-Uriarte (2018), as those also missed some genotypes that should exist under the DAG of restrictions. Our local maxima are easier by design as we want to isolate the effect of multi-peaked landscapes or local maxima (or, equivalently, missing paths), without the additional burden of missing genotypes. The Rough Mount Fuji (RMF) fitness landscapes were obtained from an RMF model, a model that has been useful to model empirical fitness landscapes (de Visser and Krug, 2014; Franke *et al.*, 2011; Neidhart *et al.*, 2014), where the reference genotype and the decrease in birth rate of a genotype per each unit increase in Hamming distance from the reference genotype were randomly chosen (see Supplementary Material, [section 1.2, "Rough Mount Fuji"](#)). These fitness landscapes cannot be represented by DAGs of restrictions with respect to neither paths to the maximum nor accessible genotypes (see also Diaz-Uriarte, 2018).

Once a fitness landscape had been generated, we simulated 20000 evolutionary processes. We used the continuous-time, logistic-like model of McFarland *et al.* (2013), in which death rate depends on total population size, as implemented in OncoSimulR (Diaz-Uriarte, 2017), with the specified parameters of initial population size and mutation rate. Each individual simulation was run until one of the genotypes at the local fitness maxima (or the single global fitness maximum) reached fixation (see details in Supplementary Material, [section 3, "Simulations"](#)). We also verified that all 7 or 10 genes had appeared in at least some genotypes, i.e., were part of the paths of tumor progression. If this condition was not fulfilled, a new fitness landscape was generated and the processes started again. This procedure is independent of the detection process that generates the samples of genotypes (next).

To obtain the samples of genotypes that were analyzed by the CPMs, we used three different **detection regimes** to emulate single-cell sampling at total tumor sizes (number of cells) that are, in the log scale, approximately uniformly distributed (**uniform** detection regime), biased towards large sizes (**large**) or biased towards small sizes (**small**). (Working on the log-scale of tumor size is appropriate as in the model of McFarland *et al.*, 2013, tumor population size increases logarithmically with number of driver mutations). We drew random deviates from beta distributions with parameters $B(1, 1)$, $B(5, 3)$, and $B(3, 5)$ (for uniform, large, and small, respectively), rescaled them to the range of observed sizes, and obtained the sample with population size closest to the target (see details in Supplementary Material, [section 3.3, "Detection regimes: sampling"](#)). For each sample, the genotype returned was the single genotype with the largest frequency (so we did not introduce possible additional noise due to bulk sequencing—i.e., obtaining a single, inexistent, genotype from a sample that could contain multiple different genotypes). Finally, for each of the three **sample sizes** of 50, 200, and 4000, we splitted the 20000 simulations into five sets of non-overlapping data sets. These are the data sets that were analyzed with the six CPMs.

2.3 Inferring Cancer Progression Models (OT, CBN, CAPRI, CAPRESE) and paths of tumor progression

We have used four distinct CPM methods (methods not considered here are either too slow for routine work, have no software available, or have dependencies on non-open source external

libraries —see Supplementary Material, [section 5.2, “CPM software”](#)). Two of the methods used have two variants, yielding a total of six methods. Only a brief overview is provided here; detailed descriptions can be found in Caravagna *et al.* (2016); Desper *et al.* (1999); Gerstung *et al.* (2009, 2011); Montazeri *et al.* (2016); Olde Loohuis *et al.* (2014); Ramazzotti *et al.* (2015); Szabo and Boucher (2008). CPM methods assume that the different individuals in a data set constitute independent realizations of the same evolutionary process —see assumptions. These methods try to identify restrictions in the order of accumulation of mutations from cross-sectional data. The cross-sectional data is a matrix of subjects or samples by driver alteration events, where each entry in the matrix is binary coded as mutated or not-mutated. For the simulations, we will refer to these driver alteration events as “genes”, but they can be individual genes, parts or states of genes, or modules or pathways made from several genes (e.g. Caravagna *et al.*, 2016; Gerstung *et al.*, 2011). When we analyze the 22 cancer data sets (see [section 2.5](#)) we will use the generic term “features” as some of those data sets use genes whereas others use pathway or module information. Both Oncogenetic trees (OT) (Desper *et al.*, 1999; Szabo and Boucher, 2008) and CAPRESE (Olde Loohuis *et al.*, 2014) describe the accumulation of mutations with order constraints that can be represented as trees. Thus, among the “representable” fitness landscapes used in this paper ([section 2.2](#)), OT and CAPRESE can only represent the subset that are trees. A key difference between the two is that CAPRESE reconstructs these models using a probability raising notion of causation in the framework of Suppes’ probabilistic causation, whereas in OT weights along edges can be directly interpreted as probabilities of transition along the edges by the time of observation (Szabo and Boucher, 2008, p. 4). Both CAPRI and CBN allow modeling the dependence of an event on more than one previous event: the output of the model are graphs (DAGs) where some nodes have multiple parents, instead of a single parent (as in trees). CAPRI tries to identify events (alterations) that constitute “selective advantage relationships” again using probability raising in the framework of Suppes’ probabilistic causation. We have used two versions of CAPRI, that we will call CAPRI.AIC and CAPRI.BIC, that differ in the penalization used in the maximum likelihood fit (AIC or BIC, respectively). For CBN we have also used two variants, the one described in Gerstung *et al.* (2009, 2011) that uses simulated annealing with a nested EM algorithm for estimation, and MCCBN, described in Montazeri *et al.* (2016), that uses a Monte-Carlo EM algorithm that allows it to fit data sets with many more genes. See Supplementary Material, [section 5.2, “CPM software”](#), for further details.

Because (the transitive reduction of) a DAG of restrictions determines a fitness graph (see [Figure 1](#) and Diaz-Uriarte, 2018), the set of paths to the maximum encoded by the output from a CPM is obtained from the fitness graph. This we did for all methods. From CBN and MCCBN we can also obtain the estimated probability of each path of tumor progression to the fitness maximum, since both CBN and MCCBN return the parameters of the transition rates between genotypes (see e.g., p. i729 in Montazeri *et al.*, 2016, [section 2.2](#) in Gerstung *et al.*, 2009, or Hosseini, 2018; see details and example in Supplementary Material, [section 5.4, “Computing probabilities of paths”](#)). It is also possible to perform a similar operation with the output of OT, and use the edge weights from the fits of OT to obtain the probabilities of transition to each descendant genotype and, from them, the probabilities of the different paths to the global maximum. It must be noted that these probabilities are not really returned by the model, since the OTs used are untimed oncogenetic trees (Desper *et al.*, 1999; Szabo and Boucher, 2008). We will refer to paths with probabilities assigned in the above way as **probability-weighted paths**. For CAPRESE and CAPRI, it is not possible to map the output to different probabilities of paths of progression (see also Supplementary Material, [section 7, “CAPRI, CAPRESE, and paths of tumor progression”](#)) and in all computations that required probability of paths we assigned the same probability to each path.

2.4 Measures of performance and predictability

We have characterized evolutionary unpredictability using the diversity of Lines of Descent (LODs). LODs were introduced by Szendro *et al.* (2013) and “(...) represent the lineages that arrive at the most populated genotype at the final time” (p. 572). In other words, in our sim-

ulations a LOD is a sequence of parent-child genotypes, from the initial genotype to a local maximum: a LOD is the path that a tumor has taken until fixation. The final genotype in a LOD is a local fitness maximum, but there are no guarantees that any intermediate genotype in the LOD will have been the most common genotype at any time point (especially if clonal interference and stochastic tunneling are present — de Visser and Krug, 2014; Sniegowski and Gerrish, 2010). As in Szendro *et al.* (2013), we can use the entropy of these paths to measure the indeterminism of the paths of evolution, or evolutionary unpredictability, and we will define $S_p = -\sum p_i \ln p_i$, where p_i is the observed probability of each LOD (each path) computed from the 20000 simulations, and the sum is over all paths or LODs. Evolutionary unpredictability, as estimated by the CPMs, will analogously be defined as $S_c = -\sum q_j \ln q_j$, where q_j is the probability of each path to the maximum according to the cancer progression model considered, and the sum is over all paths predicted by the CPMs. (Hosseini, 2018, normalizes predictability by dividing by the maximum entropy, similar to dividing by the prior entropy in the “information gain” statistic in Lässig *et al.*, 2017; but the maximum entropy is a constant for each number of genes, i.e., 7! or 10! for our simulations).

To measure how well CPMs predict tumor progression, we used three different statistics. To compare the overall similarity of the distribution of paths predicted by CPMs with the true observed one (i.e., the distribution of LODs) we used the Jensen-Shannon divergence (**JS**) (Crooks, 2017; Lin, 1991), scaled between 0 and 1 (equivalent to using the logarithm of base 2). JS is a symmetrized Kullback-Leibler divergence between two distributions and is defined even if the two distributions do not have the same sample space, i.e., even if $P(i) \neq 0$ and $Q(i) = 0$ (or $Q(i) \neq 0$ and $P(i) = 0$), as can often be the case for our data. A JS value of 0 means that the distributions are identical, and a value of 1 that they do not overlap. Therefore, predictions of CPMs are closer to the truth the smaller the value of JS. The sum of the probabilities of the paths in the LODs that are not among the paths allowed by the CPMs, $P(\neg DAG|LOD)$, is equivalent to **1 - recall**. Larger values of 1-recall mean that the CPM is not capturing a large fraction of the evolutionary paths to the maximum (or maxima). The sum of the predicted probabilities of paths according to the CPMs that are not used by evolution (i.e., that are not LODs), $P(\neg LOD|DAG)$, is equivalent to **1 - precision**. Larger values of 1-precision mean that the CPMs predict larger numbers of paths to the maximum that are not used by evolution. In the Supplementary Material we also use as statistic the **probability of recovering the most common LOD**; we will rarely refer to this statistic in the main paper since it follows a pattern very similar to recall (see Supplementary Material, Figures S19 and S20). Statistics 1-recall and 1-precision can overestimate performance: they could both have a value of 0, even when JS is very close to 1 (see example in Supplementary Material, [section 5.5, “Example where perfect recall and precision do not guarantee Jensen-Shannon divergence of 0”](#)). Thus, the basic overall performance measure will be JS.

2.4.1 Comparing paths from CPMs with LODs of different lengths

When all paths from the CPM and the LOD have equal length (they end in a genotype with the same number of genes mutated, K) computing the above statistics is straightforward. But paths could differ in length. In fitness landscapes with local maxima, LODs can differ in length; some LODs could have a length (or number mutations of the fixated genotype), K_i , shorter than the length of the paths from the CPM, K_C (all paths from a CPM have the same number of mutations, since all arrive at the genotype with all K_C genes mutated). It is also possible that some or all $K_i > K_C$, i.e., some or all LODs have a length larger than the length of the paths from the CPM. This will happen if the CPM has been built from a data set that contains fewer genes than the number of genes in the landscape (e.g., because one or more genes were absent —see Supplementary material [section 5.3, “Preprocessing of data for CPMs”](#)); if the sampled data set has fewer genes than the landscape in a representable fitness landscape, then all $K_i > K_C$ (as K_i will be equal to either 7 or 10).

To compute JS, 1-recall, and 1-precision that will cover all those cases we used the following procedure (that reduces to the simpler procedure in the above section when all $K_i = K_C$). Let

i and j denote two paths, one from the LOD and the other from the CPM, with corresponding probabilities p_i and q_j ; in contrast to the previous section, and to minimize notation, i, j (and p_i, q_j) could refer to a path from the LOD and a path from the CPM or, alternatively, a path from the CPM and a path from a LOD. Let K_i, K_j denote the length of paths i and j , respectively. At least one set of either K_i s or K_j s has all elements identical (e.g., if j refers to indices of the paths from the CPM, it is necessarily the case that $K_1 = K_2 = \dots = K_m = K_C$, with m the total number of different paths from the CPM).

Now if $K_i > K_j$ and the path i up to K_j mutations (i.e., from the WT genotype to the genotype with K_j mutations) is identical to j , then path j is included in path i : all of q_j is accounted for by i . This also means that path i is partially included in (or accounted for by) path j , but a fraction of it, $(K_i - K_j)/K_i$, is missing or unaccounted for. The above applies directly to calculations of 1-recall and 1-precision. For computing JS, there will be two entries in the vectors with the probability distributions that will be compared: $P = [p_i \frac{K_j}{K_i}, p_i \frac{K_i - K_j}{K_i}]$, $Q = [q_j, 0]$. This procedure can be applied to all elements i, j , summing all unmatched entries: $\sum p_i \frac{K_i - K_j}{K_i}$ is the total flow in the set of paths i that cannot be matched by the j s because they are shorter. To simplify computations, that unmatched term can also include $\sum p_u$, where u denote those paths in i that do not match any j . Conversely, all paths i with $K_i > K_j$ such that the paths become indistinguishable up to K_j can be summed in a single entry so that we obtain $\sum p_i \frac{K_j}{K_i}$ and $\sum p_i \frac{K_i - K_j}{K_i}$ for the matched and unmatched fractions, respectively. All computations have their corresponding counterparts for elements i, j when $K_i < K_j$. This procedure results in unique JS (remember the K are all the same for at least one of the sets of paths) as well as unique 1-precision and 1-recall, and it reduces to the procedure (see above) when all K_i are equal and equal to all K_j . A commented example and further details are provided in the Supplementary Material ([section 5.6.1, "Commented example for paths of unequal length"](#)).

2.4.2 Statistical modeling of performance

We have used generalized linear mixed-effects models, with a beta model for the dependent variable (Ferrari and Cribari-Neto, 2004; Grün *et al.*, 2012; Smithson and Verkuilen, 2006), to model how JS, 1-recall, and 1-precision, are affected by S_p , detection regime, sample size, number of genes, and type of fitness landscape. In all models, the response variable was the average from the five split replicates of each fitness landscape by sample size by detection regime combination, and fitness landscape id (not type) was a random effect. When the dependent variable had values exactly equal to 0 or 1, we used the transformation suggested in Smithson and Verkuilen (2006). Models were fitted using sum-to-zero contrasts (McCullagh and Nelder, 1989) and all regressors were used as discrete regressors, except S_p , which has been scaled (mean 0, variance 1) for easier interpretation; the coefficients of the main effect terms of the discrete regressors are the deviations from the average (see further details in Supplementary Material, [section 5.7, "Coefficients of linear models"](#)). We have used the glmmTMB (Brooks *et al.*, 2017) and car (Fox and Weisberg, 2011) packages for R (R Core Team, 2018) for statistical model fitting and analysis.

2.5 Cancer data sets

We have used 22 cancer data sets (including six different cancer types: breast, glioblastoma, lung, ovarian, colorectal, and pancreatic cancer); some code mutations in terms of genes (somatic mutations and/or copy number alterations) and some in terms of pathways or modules. All of these data, except for the breast cancer data sets BRCA.ba.s and BRCA.he.s (from Cancer Genome Atlas Research Network, 2012b), have been used previously as input for some CPM algorithms in Attolini *et al.* (2010); Caravagna *et al.* (2016); Cheng *et al.* (2012); Gerstung *et al.* (2011); Misra *et al.* (2014); Olde Loohuis *et al.* (2014); Ramazzotti *et al.* (2015), with the original sources of the data being Bamford *et al.* (2004); Brennan *et al.* (2013); Cancer Genome Atlas Research Network (2008, 2011, 2012a); Ding *et al.* (2008); Jones *et al.* (2008); Knutsen *et al.* (2005);

Parsons *et al.* (2008); Piazza *et al.* (2013); Wood *et al.* (2007). Details on sources, names, and how the data were obtained and processed are provided in the Supplementary Material (section 6, “Cancer data sets”).

These data sets vary in sample size (27 to 594 samples), number of features (from 7 to over 100), data types (nonsynonymous somatic mutations and copy number aberrations or both), levels of analysis (genes, modules and pathways, exclusivity groups), patterns of number of mutations per subject and frequency of mutations analyzed, and procedures for driver selection, and restriction of patient subtypes. The data sets, therefore, are a large representative ensemble of data sets to which researchers have previously applied or might apply CPMs.

3 Results

3.1 Simulated fitness landscapes: characteristics, evolutionary predictability, sampled genotypes

In the Supplementary Material (section 2, “Plots of fitness landscapes and inferred DAGs”) we show all the fitness landscapes used. We also show (section 4, “Fitness landscapes: characteristics, evolutionary predictability, clonal interference, and sampled genotypes”) the main characteristics of the fitness landscapes used, the variability in evolutionary predictability, and the characteristics of the samples obtained under the three detection regimes. The three types of fitness landscapes had comparable numbers of accessible genotypes but differed strongly in the number of local fitness maxima and reciprocal sign epistasis (Figures S1 to S3). Simulations resulted in varied amounts of clonal interference, as measured by the average frequency of the most common genotype (Figure S4 or, similarly, the inverse of the average number of clones with frequency > 5%: Figure S5); scenarios where clonal sweeps dominated (i.e., those characterized by the smallest clonal interference) corresponded to initial population sizes of 2000, with clonal interference being much larger at the other population sizes (Figure S4).

Simulations resulted in observed numbers of paths to the maximum (number of distinct LODs) that showed a wide range (Figure S6), from 2 to 3082 (median of 228, 95, and 55, for representable, local maxima, and RMF, respectively), with fitness landscapes with 10 genes with a much larger number than those with 7 genes (105 vs. 1340, 55 vs. 261, 33 vs. 113, for representable, local maxima, and RMF, respectively). LOD diversities (S_p) ranged from 0.3 to 8.7 (Figure S7) with RMF models showing smaller S_p ; RMF landscapes had the largest number and diversity of observed local fitness maxima (Figure S8 and S9) and S_p was strongly associated to the number of accessible genotypes (Figure S10). Of course, the number of mutations of the fitness maxima were 7 and 10 in the representable landscapes, and smaller in the local maxima and RMF landscapes (Figure S11).

The number of different sampled genotypes was comparable between detection regimes (Figure S12), but diversity differed (Figure S13), with the uniform detection regime showing generally larger sampled diversity. The mean and median number of mutations of sampled genotypes (Figures S14 and S15) differed between detection regimes in the expected direction (largest in the large detection regime, and smallest in the small detection regime); the standard deviation and coefficient of variation in the number of mutations (Figures S16 and S17) were largest in the uniform detection regime (thus, the uniform detection regime showed both the largest variation in number of mutations of genotypes and the largest diversity of genotypes). Sample characteristics and the difference in sample characteristics between detection regimes were affected by type of fitness landscape (e.g. Figures S13 and S16).

3.2 Predicting paths of evolution with CPMs

The six methods used can be divided into three groups: methods that return trees (OT and CAPRESE) and two families of methods that return DAGs, CAPRI (CAPRI.AIC and CAPRI.BIC) and CBN (CBN and MCCBN). Comparing within groups with respect to JS, and as seen in the Supplementary Material (section 8, “Overall patterns for the six methods”), one member of the

pair consistently outperformed the other (see Figure S18). OT (using probability-weighted paths, see below) was significantly better than CAPRESE (paired t -test over all non-missing 56595 pairs of results: $t_{56594} = -161.1$, $P < 0.0001$), CBN was significantly better than MCCBN ($t_{56593} = -42.6$, $P < 0.0001$), and CAPRI_AIC was significantly better than CAPRI_BIC ($t_{56594} = -41.9$, $P < 0.0001$).

This ranking within types of methods does not always apply to the other two measures of performance, most notably CAPRESE with respect to 1-recall, where its performance can be one of the best, and often better than that of OT. CAPRESE's better recall, however, is more than offset by its poor precision (often the worst or among the worst). Similar comments apply to other reversals (e.g., MCCBN's slightly better precision in some scenarios being offset by its considerably worse recall). Remember we will assess performance using mainly JS (see "*Measures of performance and predictability*", section 2.4). In what follows, therefore, and for the sake of brevity, we will focus on OT, CBN, and CAPRI_AIC, since the overall performance of their alternatives is worse.

Figure 2 shows how the performance measures for OT, CBN, and CAPRI_AIC change with sample size for all combinations of type of landscape, detection regime, and number of genes (results for the probability of recovering the most common LOD are shown in the Supplementary Material, Figure S19, and the patterns are essentially those of recall, Figure S20). The measures of JS and 1-precision for OT and CBN (and MCCBN) use probability-weighted paths computed as explained in 2.4, because there was strong evidence for all three methods that the probability-weighted paths led to better results (JS, paired t -test over all pairs: OT, $t_{56594} = -195.8$, $P < 0.0001$; CBN: $t_{56594} = -222.3$, $P < 0.0001$; MCCBN: $t_{56593} = -149.0$, $P < 0.0001$; 1-precision: OT: $t_{56594} = -187.6$, $P < 0.0001$; CBN: $t_{56594} = -217.6$, $P < 0.0001$; MCCBN: $t_{56593} = -130.3$, $P < 0.0001$). (See also Supplementary Material, Figures S21, S22, S23).

Overall, CBN was the method with the best performance ($P < 0.0001$ from all pairwise comparisons between the six methods with Tukey's contrasts and single-step multiple testing p-value adjustment—Hothorn *et al.*, 2008—on linear mixed-effects models with landscape by split replicate as random effect). It must be noted, however, that all methods can show large variability in performance, as shown in Figure 3 (also Supplementary Material, Figure S24).

JS differed between type of landscape, number of genes, detection regime, and sample size, but the magnitude and even direction of effects differed between combinations of those factors, as seen in Figure 2 and 4. Generalized linear mixed-effects models fitted to the complete data set and to the different combinations of method and type of landscape (see Supplementary Material, section 17, "*Analysis of deviance tables for fitted models*") also showed highly significant ($P < 0.0001$) two-, three-, and four-way interactions between most of the factors, in particular those involving type of landscape and detection regime. As seen in Figure 3, type of landscape and detection regime also had very strong effects in the variability of the estimates, with relative variabilities that could reach 20% with small sample sizes.

Under representable fitness landscapes, performance improved with increasing sample size and with the uniform detection regime. Performance was worse in fitness landscapes of 10 genes (Figure 2, panel A; Figure 4, top row); the decrease in performance with increasing number of genes is related to methods both missing evolutionary paths (Figure 2B), and allowing paths that are not used by evolution (Figure 2C). Notably, with CAPRI the effect of sample size was much weaker and increases in sample size could even lead to decreases in performance, specially under the uniform detection regime (highly significant, $P < 0.0001$, interactions of detection and sample size—see Supplementary Material, section 17, "*Analysis of deviance tables for fitted models*"). This is attributable to CAPRI excluding many paths taken during evolution (Figure 2B). This behavior of CAPRI can also be seen in Figure 6A, where the ratio of estimated to true number of paths went from slight to very severe underestimation as sample size increased under the uniform detection regime. This was itself caused by CAPRI sometimes allowing only a few or even just one path to the maximum (Supplementary Material, Figure S25).

Under the RMF landscape overall performance was worse. Increasing sample size for OT and CBN led to minor decreases in performance (Figure 2 and Figure 4 bottom row). CPMs failed to capture about 50% of the evolutionary paths (or fractions of paths) to the local maxima (Figure 2B) and included more than 75% of paths (or fractions of paths) that were never taken by evolution (Figure 2C). The behavior under local maxima was similar to that of representable fitness landscapes in terms of the direction of most effects, but effects were generally weaker, with the exception of evolutionary unpredictability (see next).

What about the effect of evolutionary unpredictability itself on performance? There were no marginal effects of evolutionary unpredictability (as measured with S_p) on performance in representable fitness landscapes (Figure 4) for CBN and OT. But the effects of evolutionary unpredictability were, in fact, more complex than depicted in Figure 4, as there were highly significant interactions ($P < 0.0001$) between S_p , detection regime, and sample size, within representable and local maxima landscapes, as well as in the overall models (see Supplementary Material, [section 17, "Analysis of deviance tables for fitted models"](#)). In many cases, the sign of the slope was reverted from its main effect, as shown in Figure 5 (see also Supplementary Material, Figure S26). In most scenarios, performance was worse with larger unpredictability (larger S_p) as seen by the positive slopes of JS on S_p (Figure 5). But under representable landscapes, in the large detection regime and for sample sizes 50 and 200, larger evolutionary unpredictability was associated with better performance; the difference in effects was itself significantly affected by the number of genes (see also Supplementary Material, [section 17, "Analysis of deviance tables for fitted models"](#)). Under RMF fitness landscapes, large evolutionary unpredictability was associated with poorer performance over all sample sizes. Under local maxima, the effect of evolutionary unpredictability depended strongly on sample size and detection regime, with reversal of effects from sample size of 50 compared to 4000 under the large detection regime, similar to those mentioned above for representable landscapes (Figure 5).

3.3 Inferring evolutionary unpredictability from CPMs

Figure 6 shows the relationship between the estimated and true numbers and diversities of paths of tumor progression. Under representable fitness landscapes, and for the two methods with the best behavior, CBN and OT, there were large differences in the ratio of number of paths to the maximum over true number of paths to the maximum, associated to differences in sample size and number of genes, as shown in Figure 6A. Average ratios of estimated paths to the maximum over true paths to the maximum were 1.4 and 6.9 for 7 and 10 genes for CBN (and 0.4 and 1.9 for OT). But values for CBN ranged from 0.5 (7 genes, sample size 50, uniform detection regime), to 33.9 (10 genes, sample size 4000, large detection regime); for OT they ranged from 0.2 (7 genes, sample size 200, uniform, detection) to 5.6 (10 genes, sample size 4000, large detection regime). In section 3.1 (see also Supplementary Material, [section 4, "Fitness landscapes: characteristics, evolutionary predictability, clonal interference, and sampled genotypes"](#)) we saw that the true number of evolutionary paths increased with the number of genes; what we see here is that the inferred number of evolutionary paths to the maximum from CBN and OT often increased even faster, a consequence of worse recall under 10 genes. Detection regime and sample size (again, for OT and CBN) had a large effect: number of paths inferred increased with sample size, specially under the large detection regime (Figure 6).

But for both CBN and OT that disproportionate increase in the number of inferred paths carried only a small penalty in terms of correctly estimating evolutionary unpredictability (the diversity of paths to the maximum, S_p), as can be seen from Figure 6B—and this is a consequence of both using probability-weighted paths and of changes in scale (diversities used logarithms). For example, for CBN the ratio of inferred to observed diversities, S_c/S_p , remained close to 1 over all combinations of detection regime, number of genes, and sample size (averages of $0.81\times$ and $0.93\times$ for 7 and 10 genes); the values were closest to one with sample size 4000 and under the uniform detection regime.

In contrast to CBN and OT, patterns for CAPRI seemed dominated by the tendency of CAPRI to only allow very few paths as the sample size grows large, and mainly under the

uniform detection regime (see also Supplementary Material, Figure S25). Under representable landscapes, CAPRI underestimates, sometimes severely, the true diversity of paths to the maximum (Figure 6B) and can lead to very large variability of the estimates (Supplementary Material, Figure S27, for coefficient of variation of S_c).

Type of landscape affected the quality of estimates. Under RMF, the number of paths tended to be overestimated by very large factors (averages over 7 and 10 genes: paths: CBN $2.9\times$ and $55.6\times$; OT: $5.1\times$ and $128\times$; CAPRI: $3.5\times$ and $61\times$), especially with 10 genes and sample sizes of 4000 (CBN: $112\times$; OT: $236\times$; CAPRI: $61\times$). Diversity was also overestimated but, as was the case for representable landscapes, by smaller factors (averages over 7 and 10 genes: CBN $1.1\times$ and $1.6\times$; OT: $2.1\times$ and $2.8\times$; CAPRI: $3.0\times$ and $3.5\times$; values for 10 genes and sample sizes of 4000: CBN: $2.2\times$; OT: $3.5\times$; CAPRI: $4.0\times$).

And how does the estimated evolutionary unpredictability change with the true evolutionary unpredictability? Figure 6C shows that the slopes of regressions of estimated unpredictability from CPMs (S_c) on true unpredictability (S_p) changed depending on fitness landscape, detection regime, and sample size, including slopes over and under 1, and even inversion of signs (ranges of slopes over all combinations of type of landscape by detection regime by number of genes by sample size: CBN: 0.47 to 1.27; OT: 0.43 to 1.50; CAPRI: -1.04 to 1.19).

3.4 Cancer data sets

We will use CPMs on 22 cancer data sets to examine their usefulness for predicting tumor evolution. As explained in section 2.5, these data sets are a large representative set to which CPMs have been applied or might be applied in the future; they include six different cancer types, and show wide variation in sample size, number of features, data type, levels of analysis (genes, modules, pathways), methods for driver and patient subtype selection, and distribution of number of mutations per subject and frequency of mutations (see Supplementary Material, section 6, “Cancer data sets” and Figures S31 and S32).

We have analyzed all the data sets with CBN (the best performing method —see sections 3.2 and 3.3). We have run the analysis three times per data set, limiting the number of features analyzed to the 7, 10, and 12 most common ones, so as to examine how our assessments depend on the number of features analyzed; the first two thresholds use the same number of features as the simulations. (Of course, for data sets with 7 or fewer features, there are no differences in the data sets used under the 7, 10, and 12 thresholds, so the values shown below reflect variability between runs; ditto for data sets with 8 to 10 features with respect to thresholds 10 and 12).

We do not know the true paths of tumor progression, but we can use the bootstrap to assess the robustness of the inferences. To do so, we repeated the process above with 100 bootstrap samples (see Supplementary Material, section 6.2, “Bootstrapping on the cancer data sets”). We computed $JS_{o,b}$, the JS between the distribution of paths to the maximum from the original data set and each of the bootstrapped samples. Large differences in the distribution of paths between the analyses with the bootstrap samples and the analysis with the original sample suggests that the inferences are unreliable and cannot be trusted (but small differences do not indicate that the inferred paths match the distribution of the true ones).

The results are shown in Figure 7; summary patterns are shown in Figure 8. Unreliability ($JS_{o,b}$) was large for most data sets, and very large for some of them. These results would be expected, even if the true fitness landscapes were representable ones, as most of the data sets have small sample sizes (less than 1000), and we have seen that performance is poor (large JS) for that range of sample sizes (Figure 2A). For these data sets, as can be seen from Figure 8, there was no relationship between $JS_{o,b}$ and sample size, and when the same data set was analyzed using pathways/modules and genes, performance was generally better using pathways or modules (Pan_pa vs. Pan_ge, Col_pa vs. Col_g, GBM_ge vs. GBM_pa, GBM_mo vs. GBM_CNA). Within data sets, and for all data sets, as the number of features analyzed increased performance either decreased or stayed the same (i.e., for data sets with more than 7 features, unreliability at the 10 feature threshold, $JS_{o,b}^{10}$, was larger or equal to unreliability at the 7 feature threshold, $JS_{o,b}^7$; for data sets with more than 10 features, $JS_{o,b}^7 \leq JS_{o,b}^{10} \leq JS_{o,b}^{12}$).

Figure 7A).

Smaller numbers of features and smaller S_c should be associated with smaller $JS_{o,b}$, and there were mild trends for these patterns (Figure 8B, C), with notable exceptions: the Pancreas Pathways (Pan_pa) data set had very small $JS_{o,b}$ even for moderate number of features, and the All Pathways (all_pa) data set had a relatively small $JS_{o,b}$ even though it used 12 features and had a large S_c ; the GBM CNA modules (GBM_mo) data set also showed moderate $JS_{o,b}$ in spite of having 9 features and relatively large S_c . Conversely, some data sets with small S_c had extremely unreliable path predictions (e.g., BRCA_ba_s, Col_mss.co, Col_msi.co, GBM_ge).

Values for S_c were well within the ranges of S_c estimated by CBN for the simulated data (see Supplementary Material, Figure S28). Of course, S_c increased with number of features analyzed (see also Supplementary Material, Figure S33). Given the results from section 3.3, where generally $S_p < S_c$, this suggests that the true evolutionary unpredictability (when analyzing up to 12 features) for 13 of the data sets should be less than that corresponding to about 100 equiprobable paths to the maximum, but only eight are below the much more manageable, and useful, 20 equiprobable paths. The Pan_pa, GBM_coo, and BRCA_he_s show outstanding patterns in Figures 7 and 8. Examination of the output showed that there was one single path with estimated probability > 0.97 for Pan_pa, and two paths to the maximum of about equal probability that together added > 0.95 for GBM_coo. BRCA_he_s had only four features but mutations in SRPRA and PIK3R1 were present each in only four individuals (different individuals for the two mutations); repeated runs of CBN led to different sets of restrictions being inferred which, because there are few paths to the maximum, and some had large probabilities (> 0.5) which resulted in large differences in JS statistic between runs (and bootstrap runs will exacerbate these differences).

4 Discussion

Can we predict the likely course of tumor progression using CPMs? CBN was the best performing CPM method in our study. Using CBN under the representable fitness landscapes (the easiest scenario, as it fits the underlying model) returned estimates of the probability of paths of tumor evolution that were not far from the true distribution of paths of evolution (Figure 2A) when sample size was very large. But we find that, even under representable fitness landscapes, performance with moderate (and more realistic) sample sizes was considerably worse and was affected by detection regime. The analysis of the 22 cancer data sets revealed that performance (as measured by $JS_{o,b}$, an indicator of unreliability of inferences) was poor or very poor for most data sets. Even data sets with few features and small diversity of paths to the maximum, S_c , showed very unreliable predictions.

What factors, and how, affect performance? Under representable fitness landscapes, performance on simulated data was of course affected by the number of features, the dimension of the fitness landscape: JS was worse with 10 than with 7 genes (Figures 2, 4). Increasing sample size improved performance (Figures 2 and 4). Detection regime and evolutionary unpredictability, as measured by LOD diversity (S_p), affected individually and jointly all performance measures (Figures 2, 4, 5). Increased evolutionary unpredictability hurt performance under most conditions (Figure 5). Detection regime was a key determinant of performance, as already found in previous work (Diaz-Uriarte, 2015, 2018); performance was better under the uniform detection regime and, more importantly, it affected how the rest of the factors (evolutionary unpredictability, sample size, and number of features) impacted on performance (Figures 2 to 5).

The analysis of the 22 cancer data sets also indicated number of features as major determinant of performance. Across data sets, unreliability of inferences ($JS_{o,b}$) increased with number of features. More importantly, within data set unreliability increased as the number of features increased; note that an increase in the number of features analyzed leads to an increase in the number of features with low frequency events. Interestingly, the driver-selected data sets (Col_mss, Col_msi, BRCA_he_s, BRCA_ba_s) did not perform much better

than data sets with a simple frequency-based selection of features (e.g., Lu, Ov, or comparison Ov with Ov_drv). Even data sets with very careful, manually-curated selection of drivers and “exclusivity groups” and where variability due to subtypes has been minimized (Col_msi, Col_msi_co, Col_mss, Col_mss_co, ACML_co, BRCA_he_s and BRCA_ba_s) show very large $JS_{o,b}$. And BRCA_he_s, with only four features, showed much larger $JS_{o,b}$ than GBM_coo and Pan_pa (with 3 and 7 features, respectively), due to the presence of two low frequency alterations. These results bring forth the problem of the selection of the relevant features for analysis (Caravagna *et al.*, 2016; Cristea *et al.*, 2016; Gerstung *et al.*, 2011) and whether sample size is large enough relative to the number and frequency of features considered. We have previously shown that feature selection can have a very detrimental impact on the performance of CPM methods (Diaz-Uriarte, 2015). Using pathways instead of genes in the analyses (see, e.g., Cristea *et al.*, 2016; Raphael and Vandin, 2015) can alleviate some of the problems of feature selection. For example, data sets coded as pathways or modules generally reduce the presence of low-frequency alterations (see Supplementary Material, Figures S31 and S32). Pathways can also improve predictability and how close the estimates of path distributions are to the truth because they are more similar to heritable phenotypes, which often have smoother phenotype-fitness maps and tend to show more repeatable evolution (Lässig *et al.*, 2017; see also Wang *et al.*, 2015, but also Chebib and Guillaume, 2017; Sailer and Harms, 2017). Gerstung *et al.* (2011) found that analysis using pathways gave stronger evidence for order constraints than analysis using genes, and we also see in Figure 7 that both S_c and $JS_{o,b}$ tend to decrease if we use pathways or modules (Pan_pa vs. Pan_ge, Col_pa vs. Col_g, GBM_ge vs. GBM_pa, GBM_mo vs. GBM_CNA). “All Pathways” constitutes a promising case because it has large S_c but moderate $JS_{o,b}$. Using so-called “exclusivity groups” (sensu Caravagna *et al.*, 2016) to identify “fitness equivalent alterations” is a similar, though not identical, procedure that in this paper showed only modest improvements in $JS_{o,b}$ (Col_mss_co vs. Col_mss, Col_msi_co vs. Col_msi, ACML_co vs. ACML); this can of course be due to particularities of these data sets (e.g., large number of features relative to number of subjects) or the intrinsic difficulties of identifying true fitness equivalent groups via “hard/soft exclusivities”. However, note that although analysis using pathways/modules/exclusivity groups might lead to more reliable results from the predictability point of view, the identification of paths at the gene level is still the ultimate goal for therapeutic interventions (see Ashworth *et al.*, 2011). Regardless of the details of the procedure for collapsing and reducing features, our results suggest that further work on feature selection should consider reduction of variability of estimates of evolutionary paths as a key component.

Hosseini (2018) has reanalyzed the DAG-derived representable and a subset (those where the fully mutated genotype has largest fitness) of the DAG-derived non-representable fitness landscapes in Diaz-Uriarte (2018). He finds good agreement between the distributions of paths to the maximum from CBN and the fitness landscape-based probability distribution of paths to the maximum. Our results for CBN under the best conditions are not as optimistic. Two differences in the studies explain the differences. First, Hosseini (2018) computes the fitness landscape-based probability of paths assuming a strong selection weak mutation regime, not by directly examining the distribution of the paths to the maximum in each simulation (i.e., he does not use the LODs) and, second, he uses CBN with the very large sample size of 20000 (the full data sets in Diaz-Uriarte, 2018).

Even very good performance, though, needs to be interpreted with care. Very good performance simply tells us that the true and estimated probability distributions of the paths to the maximum agree closely. If the true evolutionary unpredictability is large, then for practical purposes our capacity to predict what will happen (in the sense of providing a small set of likely outcomes) is very limited. Ranges of diversities of 3.2 to 6.0, equivalent to 25 to 400 equiprobable paths, were common in the simulated data (see Supplementary Material, Figure S28) and are comparable to the ranges in most cancer data sets with 7 and 10 feature thresholds (see Figure 7). The inability to narrow down the likely paths to a small set of paths in these cases is, of course, not a limitation of the methods, but a problem inherent to the unpredictability of the evolutionary process in many scenarios, which could severely limit the usefulness of

even perfect predictions.

The discussion above has centered on representable fitness landscapes. As argued before, fitness landscapes with local fitness maxima are probably common in cancer. Interestingly, for small sample sizes, recall was sometimes better in local-maxima and RMF than under representable landscapes (Figure 2): with local fitness maxima, achieving good recall involves the relatively easier task of getting right the first part of short paths to the maximum (see Supplementary Material, Figure S29 and S30, where 1-recall increases with the average number of mutations of local fitness maxima). But good recall was more than offset by low precision: overall predictability was very poor. The decrease in precision is the consequence of local fitness maxima: CPMs are fitting models with paths of tumor progression that extend beyond the true end point of the progression. In addition, RMF fitness landscapes strongly violate the CPM assumption that acquiring a mutation in one gene does not decrease the probability of acquiring a mutation in another gene (see Diaz-Uriarte, 2018). The violations of assumptions in RMF and local fitness maxima explain the decreases in the relevance of sampling regime and why increasing sample size has negligible (or even detrimental) effects in these regimes (Figures 2 and 4). Remarkably, regardless of type of fitness landscape (i.e., even under violation of assumptions), and for all tasks considered (prediction of paths and estimating unpredictability) performance of methods that could return probability-weighted paths (CBN, MCCBN, OT) was better when using probability-weighted paths; thus, further improvement in these methods, even under violations of assumptions, might be possible by recalibrating their output.

And we return to our third original question, as even if achieving good performance in predicting the paths of tumor progression is unlikely, inferring evolutionary unpredictability could be an easier task. Can we use inferences of evolutionary unpredictability from CPMs as estimates of the true evolutionary unpredictability? Under representable fitness landscapes, CBN, the best performing method also for this task (Figure 6B), returned values of S_c very similar to S_p , the evolutionary unpredictability estimated from the diversity of paths, and this held over detection regimes and sample sizes. Hosseini (2018) also finds that the estimates of predictability from CBN correlate well with the true evolutionary predictability, with slopes of the regression of CPM-based on landscape-based predictability generally slightly below 1, similar to our Figure 6C (left-most column). These good results do not hold under the other two fitness landscapes: evolutionary unpredictability is overestimated, and increasing sample sizes made the problems worse and, as shown in Figure 6C, different evolutionary scenarios, sample sizes, and detection regimes have different relationships of estimated and true unpredictability. But our results indicate that we can use CBN to set upper bounds on the true S_p ; obtaining tighter estimates is an objective for further research to explore. And here our analysis of 22 cancer data sets suggests that the true evolutionary unpredictability of at least some cancer scenarios might be reasonably small, specially if S_c is overestimating the true unpredictability.

4.1 Conclusion

The answer to the question “can we predict the likely course of tumor progression using CPMs?” is, unfortunately, “only with moderate success and only under representable fitness landscapes and with very large sample sizes; but even perfect predictions might be of little use if evolutionary unpredictability is large”. Estimating upper bounds to evolutionary unpredictability is a more modest, though more likely to succeed, use of CPMs. There are three key difficulties for successful prediction: the sheer size of the problem even for moderate numbers of genes, the intrinsic evolutionary unpredictability in many scenarios, and the deviations from the assumptions of CPMs that are likely to hold in most cancer data. In addition to the caveat about using these methods under scenarios where performance is very poor, this paper raises the general question of what can we really predict about likely paths of tumor progression from cross-sectional data, for instance to guide therapeutic interventions. At a minimum, measures such as $JS_{o,b}$ and S_c on methods that return probability-weighted paths should probably become routine as ways of providing a sense of the reliability of predictions and for assessing whether the predictions could be of any practical use.

5 Acknowledgements

N. Beerenwinkel, S. Posada-Céspedes, and G. Caravagna for discussion about progression models or software; S.-R. Hosseini for providing a preprint of his MSc. thesis and for comments that helped us clarify our methods. C. Lázaro-Perea for comments on the ms.

6 Funding

Supported by BFU2015-67302-R (MINECO/FEDER, EU) to RDU. CV supported by PEJD-2016/MED-2116 from Comunidad de Madrid.

7 Data and code availability

All data for this article, as well as source code, is available from the supplementary material.

8 Author Contributions

Conceptualization: RDU, CV; **Data curation:** CV, RDU; **Methodology:** RDU; **Software:** RDU; **Formal analysis:** RDU; **Funding Acquisition:** RDU; **Investigation:** RDU, CV; **Visualization:** RDU, CV; **Writing - original draft:** RDU; **Writing - review & editing:** RDU, CV.

References

- Ashworth, A., Lord, C.J. and Reis-Filho, J.S. (2011). Genetic interactions in cancer progression and treatment. *Cell*, **145**(1), 30–38.
- Attolini, C. et al (2010). A mathematical framework to determine the temporal sequence of somatic genetic events in cancer. *Proceedings of the National Academy of Sciences*, **107**(41), 17604–17609.
- Bamford, S. et al (2004). The COSMIC (Catalogue of Somatic Mutations in Cancer) database and website. *Br J Cancer*, **91**(2), 355–358.
- Bank, C. et al (2016). On the (un)predictability of a large intragenic fitness landscape. *PNAS*, **113**(49), 14085–14090.
- Beerenwinkel, N. et al (2007). Genetic progression and the waiting time to cancer. *PLoS computational biology*, **3**(11), e225.
- Beerenwinkel, N. et al (2015). Cancer evolution: Mathematical models and computational inference. *Systematic Biology*, **64**(1), e1–e25.
- Beerenwinkel, N., Greenman, C.D. and Lagergren, J. (2016). Computational Cancer Biology: An Evolutionary Perspective. *PLoS Comput. Biol.*, **12**(2), e1004717.
- Beijersbergen, R.L., Wessels, L.F.A. and Bernards, R. (2017). Synthetic Lethality in Cancer Therapeutics. *Annual Review of Cancer Biology*, **1**(1), 141–161.
- Blomen, V.A. et al (2015). Gene essentiality and synthetic lethality in haploid human cells. *Science*, **350**(6264), 1092–1096.
- Bozic, I. et al (2010). Accumulation of driver and passenger mutations during tumor progression. *Proceedings of the National Academy of Sciences of the United States of America*, **107**, 18545–18550.
- Brennan, C.W. et al (2013). The Somatic Genomic Landscape of Glioblastoma. *Cell*, **155**(2), 462–477.

- Brooks, M.E. et al (2017). glmmTMB Balances Speed and Flexibility Among Packages for Zero-inflated Generalized Linear Mixed Modeling. *The R Journal*, **9**(2), 378–400.
- Brouillet, S. et al (2015). MAGELLAN: A tool to explore small fitness landscapes. *bioRxiv*, page 031583.
- Cancer Genome Atlas Research Network (2008). Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature*, **455**(7216), 1061–1068.
- Cancer Genome Atlas Research Network (2011). Integrated genomic analyses of ovarian carcinoma. *Nature*, **474**(7353), 609–615.
- Cancer Genome Atlas Research Network (2012a). Comprehensive molecular characterization of human colon and rectal cancer. *Nature*, **487**(7407), 330–337.
- Cancer Genome Atlas Research Network (2012b). Comprehensive molecular portraits of human breast tumours. *Nature*, **490**(7418), 61–70.
- Caravagna, G. et al (2016). Algorithmic methods to infer the evolutionary trajectories in cancer progression. *PNAS*, **113**(28), E4025–E4034.
- Chebib, J. and Guillaume, F. (2017). What affects the predictability of evolutionary constraints using a G-matrix? The relative effects of modular pleiotropy and mutational correlation. *Evolution*, **71**(10), 2298–2312.
- Cheng, Y.K. et al (2012). A mathematical methodology for determining the temporal order of pathway alterations arising during gliomagenesis. *PLoS computational biology*, **8**(1), e1002337.
- Chiotti, K.E. et al (2014). The Valley-of-Death: Reciprocal sign epistasis constrains adaptive trajectories in a constant, nutrient limiting environment. *Genomics*, **104**(6, Part A), 431–437.
- Cristea, S., Kuipers, J. and Beerenwinkel, N. (2016). pathTiMEx: Joint Inference of Mutually Exclusive Cancer Pathways and Their Progression Dynamics. *Journal of Computational Biology*.
- Crona, K., Greene, D. and Barlow, M. (2013). The peaks and geometry of fitness landscapes. *Journal of Theoretical Biology*, **317**, 1–10.
- Crooks, G.E. (2017). On measures of entropy and information. Technical report.
- de Visser, J.A.G.M. and Krug, J. (2014). Empirical fitness landscapes and the predictability of evolution. *Nat Rev Genet*, **15**(7), 480–490.
- Desper, R. et al (1999). Inferring tree models for oncogenesis from comparative genome hybridization data. *J Comput Biol*, **6**(1), 37–51.
- Diaz-Uriarte, R. (2015). Identifying restrictions in the order of accumulation of mutations during tumor progression: Effects of passengers, evolutionary models, and sampling. *BMC Bioinformatics*, **16**(41).
- Diaz-Uriarte, R. (2017). OncoSimulR: Genetic simulation with arbitrary epistasis and mutator genes in asexual populations. *Bioinformatics*, **33**(12), 1898–1899.
- Diaz-Uriarte, R. (2018). Cancer progression models and fitness landscapes: A many-to-many relationship. *Bioinformatics*, **34**(5), 836–844.
- Ding, L. et al (2008). Somatic mutations affect key pathways in lung adenocarcinoma. *Nature*, **455**(7216), 1069–1075.
- Ferrari, S. and Cribari-Neto, F. (2004). Beta Regression for Modelling Rates and Proportions. *Journal of Applied Statistics*, **31**(7), 799–815.

- Fox, J. and Weisberg, S. (2011). *An R Companion to Applied Regression, 2nd Ed.* Sage, Thousand Oaks, CA.
- Franke, J. et al (2011). Evolutionary Accessibility of Mutational Pathways. *PLoS Comput Biol*, **7**(8), e1002134.
- Gerstung, M. et al (2009). Quantifying cancer progression with conjunctive Bayesian networks. *Bioinformatics (Oxford, England)*, **25**(21), 2809–2815.
- Gerstung, M. et al (2011). The Temporal Order of Genetic and Pathway Alterations in Tumorigenesis. *PLoS ONE*, **6**(11), e27136.
- Greaves, M. (2015). Evolutionary Determinants of Cancer. *Cancer Discovery*, **5**(8), 806–820.
- Grün, B., Kosmidis, I. and Zeileis, A. (2012). Extended Beta Regression in R : Shaken, Stirred, Mixed, and Partitioned. *Journal of Statistical Software*, **48**(11).
- Hosseini, S.R. (2018). Quantifying the predictability of cancer progression using Conjunctive Bayesian Networks. M.Sc. Thesis, Swiss Federal Institute of Technology, Zürich.
- Hothorn, T., Bretz, F. and Westfall, P. (2008). Simultaneous inference in general parametric models. *Biom J*, **50**(3), 346–363.
- Jones, S. et al (2008). Core signaling pathways in human pancreatic cancers revealed by global genomic analyses. *Science (New York, N.Y.)*, **321**(5897), 1801–6.
- Knutsen, T. et al (2005). The Interactive Online SKY/M-FISH & CGH Database and the Entrez Cancer Chromosomes Search Database: Linkage of Chromosomal Aberrations with the Genome Sequence. *Genes, Chromosomes and Cancer*, **44**(1), 52–64.
- Lässig, M., Mustonen, V. and Walczak, A.M. (2017). Predicting evolution. *Nature Ecology & Evolution*, **1**(3), s41559–017–0077–017.
- Lin, J. (1991). Divergence measures based on the Shannon entropy. *IEEE Transactions on Information theory*, **37**(1), 145–151.
- Lipinski, K.A. et al (2016). Cancer Evolution and the Limits of Predictability in Precision Cancer Medicine. *Trends in Cancer*, **2**(1), 49–63.
- Losos, J.B. (2018). *Improbable Destinies: Fate, Chance, and the Future of Evolution*. Riverhead Books, S.I.
- McCullagh, P. and Nelder, J. (1989). *Generalized Linear Models, 2nd Ed.* Chapman and Hall/CRC, London.
- McFarland, C.D. et al (2013). Impact of deleterious passenger mutations on cancer progression. *Proceedings of the National Academy of Sciences of the United States of America*, **110**(8), 2910–5.
- McPherson, A.W., Chan, F.C. and Shah, S.P. (2018). Observing Clonal Dynamics across Spatiotemporal Axes: A Prelude to Quantitative Fitness Models for Cancer. *Cold Spring Harb Perspect Med*, **8**(2).
- Misra, N., Szczurek, E. and Vingron, M. (2014). Inferring the paths of somatic evolution in cancer. *Bioinformatics (Oxford, England)*, **30**(17), 2456–2463.
- Montazeri, H. et al (2016). Large-scale inference of conjunctive Bayesian networks. *Bioinformatics*, **32**(17), i727–i735.
- Neidhart, J., Szendro, I.G. and Krug, J. (2014). Adaptation in Tunably Rugged Fitness Landscapes: The Rough Mount Fuji Model. *Genetics*, **198**(2), 699–721.

- Nowak, M.A. et al (2004). Evolutionary dynamics of tumor suppressor gene inactivation. *PNAS*, **101**(29), 10635–10638.
- Olde Loohuis, L. et al (2014). Inferring Tree Causal Models of Cancer Progression with Probability Raising. *PLOS ONE*, **9**(10), e108358.
- O’Neil, N.J., Bailey, M.L. and Hieter, P. (2017). Synthetic lethality and cancer. *Nat Rev Genet*, **18**(10), 613–623.
- Palmer, A.C. and Kishony, R. (2013). Understanding, predicting and manipulating the genotypic evolution of antibiotic resistance. *Nature Reviews Genetics*, **14**(4), 243–248.
- Parsons, D.W. et al (2008). An Integrated Genomic Analysis of Human Glioblastoma Multiforme. *Science*, **321**(5897), 1807–1812.
- Piazza, R. et al (2013). Recurrent SETBP1 mutations in atypical chronic myeloid leukemia. *Nature Genetics*, **45**(1), 18–24.
- R Core Team (2018). *R: A Language and Environment for Statistical Computing*. Vienna, Austria.
- Ramazzotti, D. et al (2015). CAPRI: Efficient inference of cancer progression models from cross-sectional data. *Bioinformatics*, **31**(18), 3016–3026.
- Raphael, B.J. and Vandin, F. (2015). Simultaneous Inference of Cancer Pathways and Tumor Progression from CrossSectional Mutation Data. *Journal of Computational Biology*, **22**(00), 250–264.
- Sailer, Z.R. and Harms, M.J. (2017). Molecular ensembles make evolution unpredictable. *PNAS*, **114**(45), 11938–11943.
- Smithson, M. and Verkuilen, J. (2006). A better lemon squeezer? Maximum-likelihood regression with beta-distributed dependent variables. *Psychological methods*, **11**(1), 54–71.
- Sniegowski, P.D. and Gerrish, P.J. (2010). Beneficial mutations and the dynamics of adaptation in asexual populations. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, **365**(1544), 1255–1263.
- Szabo, A. and Boucher, K.M. (2008). Oncogenetic trees. In W.-Y. Tan and L. Hanin, editors, *Handbook of Cancer Models with Applications*, pages 1–24. World Scientific.
- Szendro, I.G. et al (2013). Predictability of evolution depends nonmonotonically on population size. *PNAS*, **110**(2), 571–576.
- Tomasetti, C. et al (2015). Only three driver gene mutations are required for the development of lung and colorectal cancers. *PNAS*, **112**(1), 118–123.
- Toprak, E. et al (2012). Evolutionary paths to antibiotic resistance under dynamically sustained drug selection. *Nature Genetics*, **44**(1), 101–105.
- Wang, E. et al (2015). Predictive genomics: A cancer hallmark network framework for predicting tumor clinical phenotypes using genome sequencing data. *Seminars in Cancer Biology*, **30**, 4–12.
- Williams, M.J. et al (2018). Quantification of subclonal selection in cancer from bulk sequencing data. *Nature Genetics*, **50**(6), 895–903.
- Wodarz, D. and Komarova, N.L. (2014). *Dynamics of Cancer: Mathematical Foundations of Oncology*.
- Wood, L.D. et al (2007). The Genomic Landscapes of Human Breast and Colorectal Cancers. *Science*, **318**(5853), 1108–1113.

9 Figures

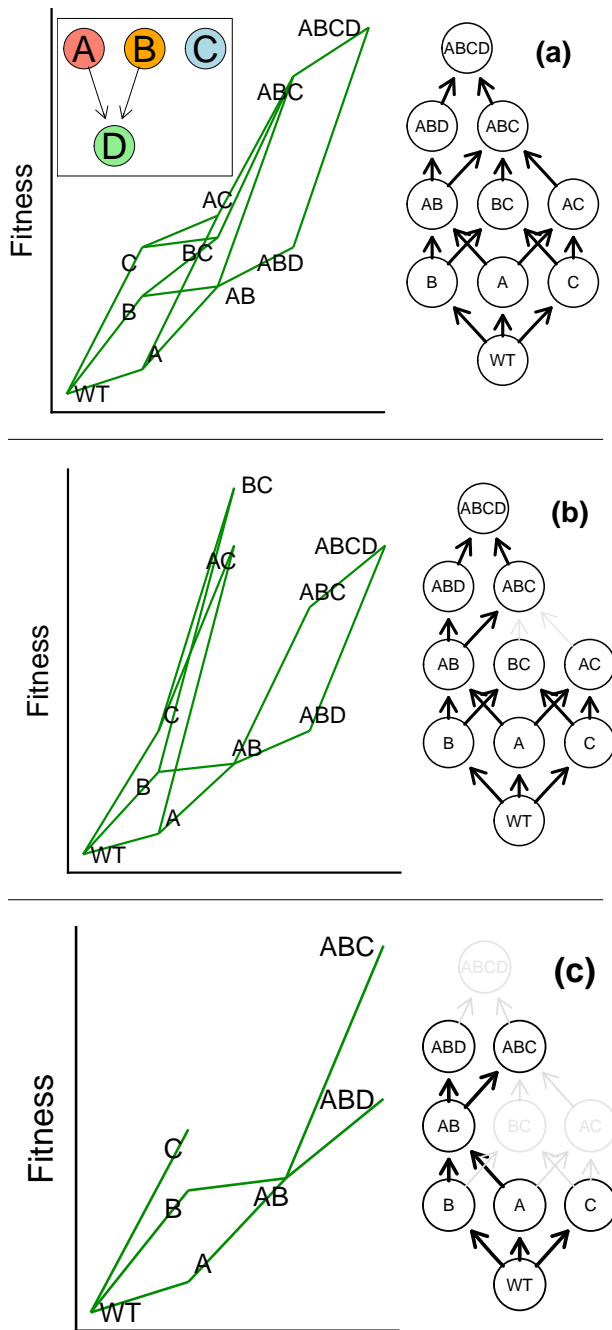


Figure 1: Fitness landscapes, paths of tumor progression, and DAGs of restrictions in the order of accumulation of mutations for the three types of landscapes used. (a) Representable; (b) local maxima; (c) RMF. In each row, on the left, the fitness landscape (representation based on Brouillet *et al.*, 2015) that shows the accessible genotypes (where the notation “AB” means a genotype with both genes A and B mutated) and on the right the fitness graphs or graphs of mutational paths (Crona *et al.*, 2013; de Visser and Krug, 2014; Franke *et al.*, 2011), where nodes are genotypes and arrows point toward mutational neighbors of higher fitness. These fitness graphs show all the paths of tumor progression, the set of accessible mutational paths and adaptive walks that, under the restriction that there can be no back mutations, start from the “wild type” (WT) genotype—where we absorb all cancer initiation events—and end in the local fitness maxima (or single global fitness maximum). Each path from WT to a maximum corresponds to a different Line of Descent (LOD). For (b) and (c), gray edges and nodes denote those that are present in (a) but missing in (b) or (c). The inset in the first row shows the DAG of restrictions in the order of accumulation of mutations that applies to (a) and (b). A DAG of restrictions shows genes in the nodes; an arrow (directed edge) from gene i to gene j indicates a direct dependency of a mutation in j on a mutation in i ; a mutation in j cannot be observed unless i is mutated. In the example, a mutation in gene D can only be observed if both A and B are mutated; note that, among the methods considered in this paper, CAPRESE and OT can only represent trees (so they can not account for D having two, or more, incoming arrows). The absence of an arrow between two genes indicates a lack of direct dependencies between the two genes. The set of genotypes that can exist under both (a) and (b) is the same, and all of them satisfy the restrictions in the DAG of restrictions. But the fitness landscape in (b) has three maxima; there are fewer paths to “ABCD” and several paths end in the other two maxima (“AC”, “BC”). Thus, the fitness graph of (b) does not fulfil the assumptions of CPMs. The defining features of (b) are that the set of accessible genotypes can be represented by a DAG of restrictions, but there are missing paths. The fitness landscape in (c) cannot be represented by any DAG of restrictions; e.g., no DAG of restrictions can account at the same time for the presence of genotypes “A”, “B”, “C”, and the absence of every double mutant with “C”. Relative to (a), (c) is missing both paths and genotypes (relative to other DAGs of restrictions it could either be missing and/or adding genotypes and paths).

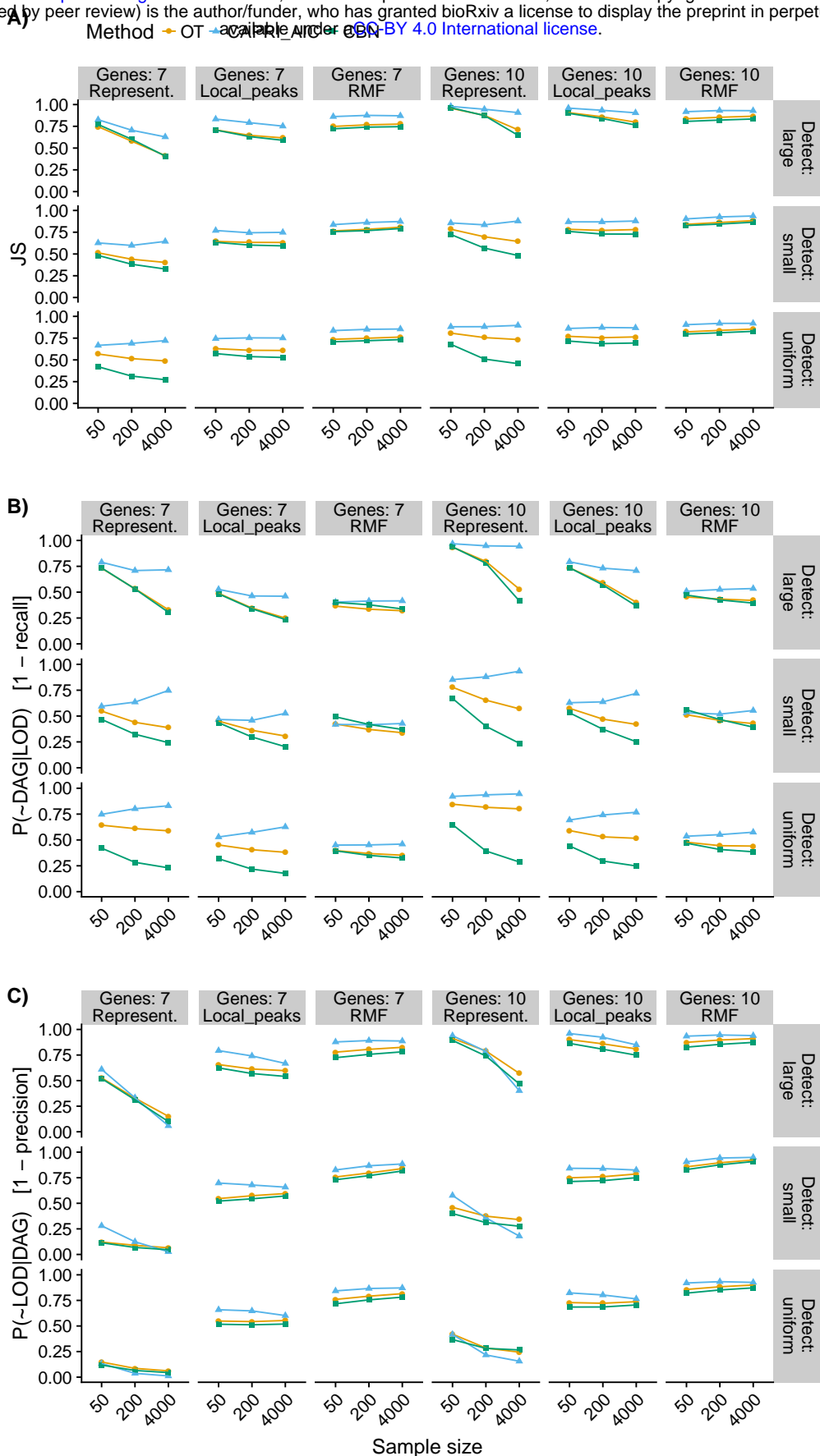


Figure 2: Summary performance measures (see definitions in 2.4) for OT, CAPRI (with AIC penalty) and CBN for all combinations of sample size, type of landscape, detection regime, and number of genes. For all measures, smaller is better. For OT and CBN, Jensen-Shannon divergence (JS) and 1-precision use probability-weighted paths (see text). Each point represented is the average of 210 points (35 replicates of each one of the six combinations of 3 initial size by 2 mutation rate regimes — see 2.1); we are thus marginalizing over mutation rate by initial simulation size combinations. Each one of the 210 points is, itself, the average of five runs on different partitions of the simulated data. See Supplementary Material, Figure S18, for results for all six methods used.

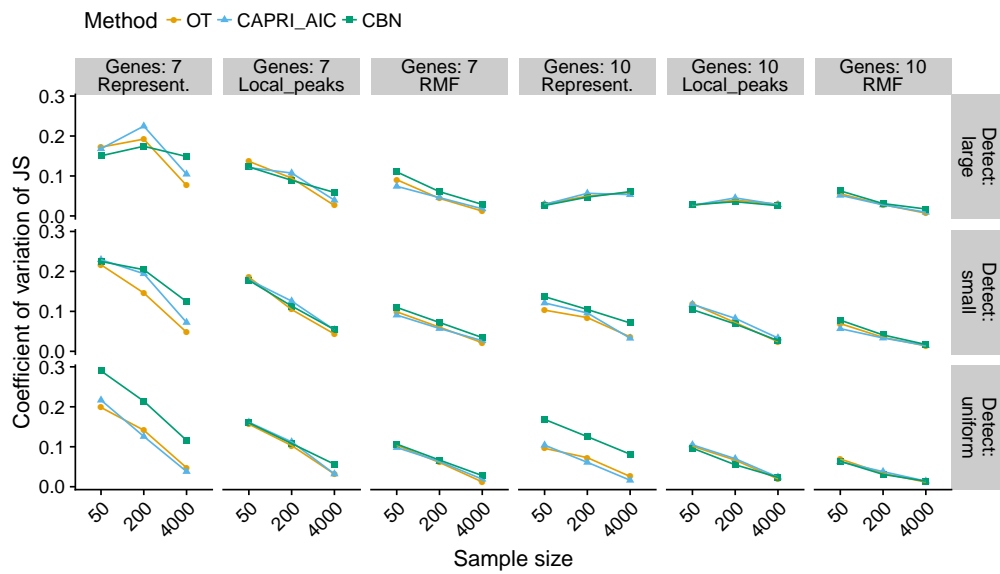


Figure 3: Coefficient of variation (standard deviation/mean) of JS for each method for all combinations of sample size, type of landscape, detection regime, and number of genes. The coefficient of variation has been computed from the five runs for each landscape on each combination of sample size and detection regime. For OT and CBN, JS is computed using the probability-weighted paths (see text). Each point plotted is the average of 210 points (see Figure 2).

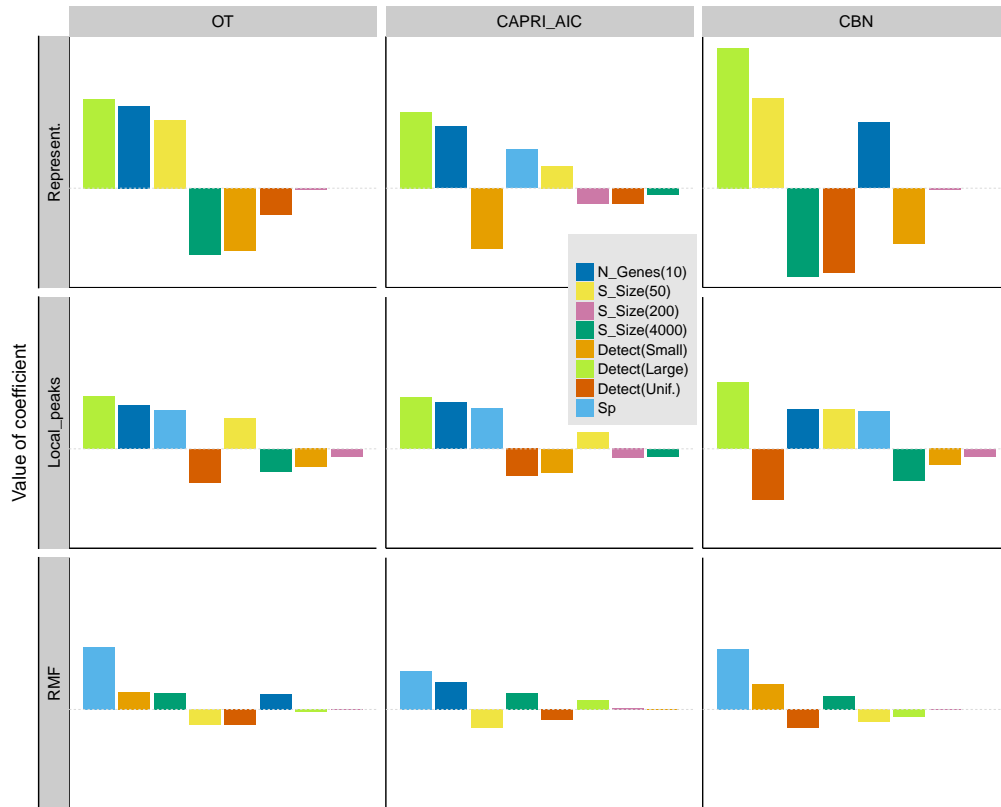


Figure 4: Coefficients from generalized linear mixed-effects models, for JS as dependent variable, with separate models fitted to each combination of method and type of fitness landscape. Coefficients are from models with sum-to-zero contrasts (see text and Supplementary Material, [section 5.7, “Coefficients of linear models”](#)). Within each panel, coefficients have been ordered from left to right according to decreasing absolute value of coefficient. The dotted horizontal gray line indicates 0 (i.e. no effect). Coefficients with a large positive value indicate factors that lead to a large decrease in performance (increase in JS). Only coefficients that correspond to a term with a P-value < 0.05 in Type II Wald chi-square tests are shown. The coefficient that corresponds to Number of genes 7 is not shown (as it is minus the coefficient for 10 genes—from using sum-to-zero contrasts). “N_Genes”: number of genes; “S_Size”: sample size; “Detect”: detection regime; “Sp”: LOD diversity (S_p).

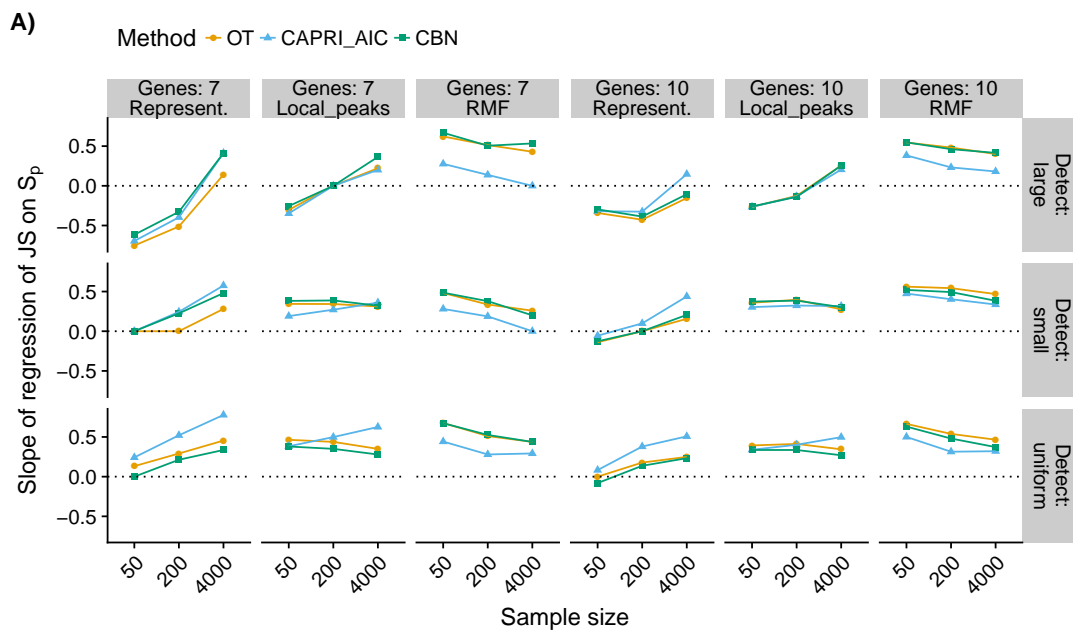


Figure 5: Estimated slopes of the regression of Jensen-Shannon divergence (JS) on LOD diversity (S_p) for all combinations of sample size by type of landscape by detection regime by number of genes. A beta regression was fitted to each subset of data. Slopes not significantly different from 0 ($P > 0.05$) shown as 0. Each regression was fitted to 210 points, each of which is itself the average of five replicates, one for each of the five runs on different partitions of the simulated data.

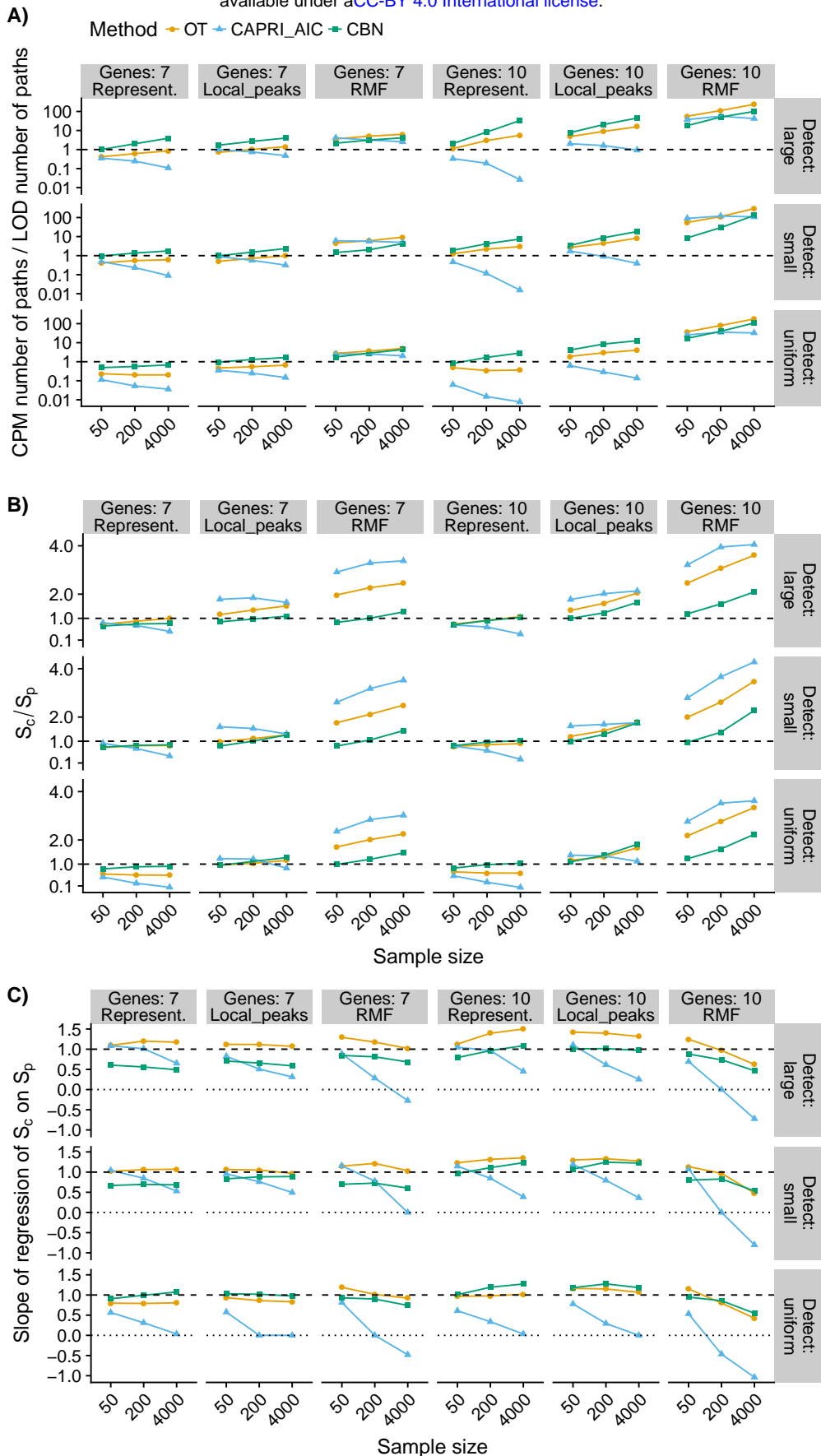


Figure 6: Number of paths and path diversities inferred from CPMs relative to the values from LODs. Panel A: average of the ratio of number of paths to the maximum from the CPMs relative to the observed number of distinct LODs for all combinations of type of landscape by detection regime by number of genes by sample size. Panel B like panel A, but for diversities of paths to the maxima. As in Figure 2, each point is the average of 210 points. Panel C shows the slope of the regression S_c on S_p ; each point is thus a slope from a regression of 210 points, each of which is itself the average of 5 replicates (see Figure 5). Panels B and C show different features of the data: panel B shows whether evolutionary unpredictability (S_p) tends to be over- or under-estimated by S_c ; panel C shows how S_c changes with S_p —see Supplementary Material, [section 19, “LOD and CPM diversity: ratios and slopes”](#), for an example of positive ratios with negative slopes.

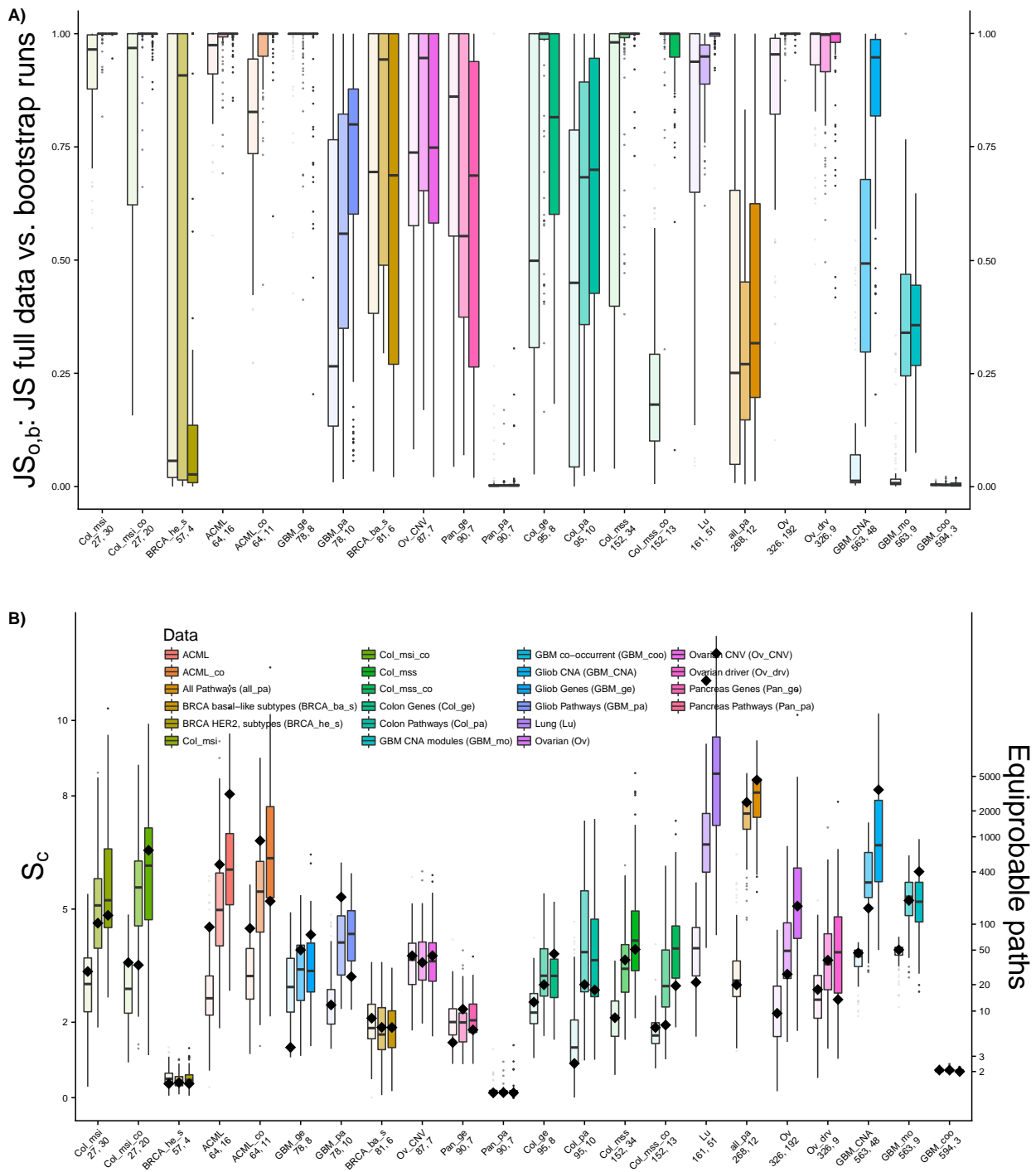


Figure 7: Results from the cancer data sets analyzed with CBN. Data sets have been ordered by increasing sample size, and the x-axis labels provide the acronym (shown in full in the inset legend). Below the data set acronym are the number of subjects and the total number of features, respectively. Analysis were run three times, limiting the number of features analyzed to the 7, 10, and 12 most common ones; the boxplots for each data set are shown in increasing order of number of features. For data sets such as, say, Pancreas genes (PG), using 7, 10, or 12 maximum features makes no difference in the number of features analyzed; the three replicate runs show run-to-run variability. A) $JS_{0,b}$: JS statistic for the comparison of the distribution of paths from running CBN on the original data set against the distribution of paths from running CBN on each one of the bootstrap runs. B) Diamonds show the S_c from the full data, and boxplots the S_c from the bootstrap runs. Right axis labeled by number of equiprobable paths equivalent to the S_c .

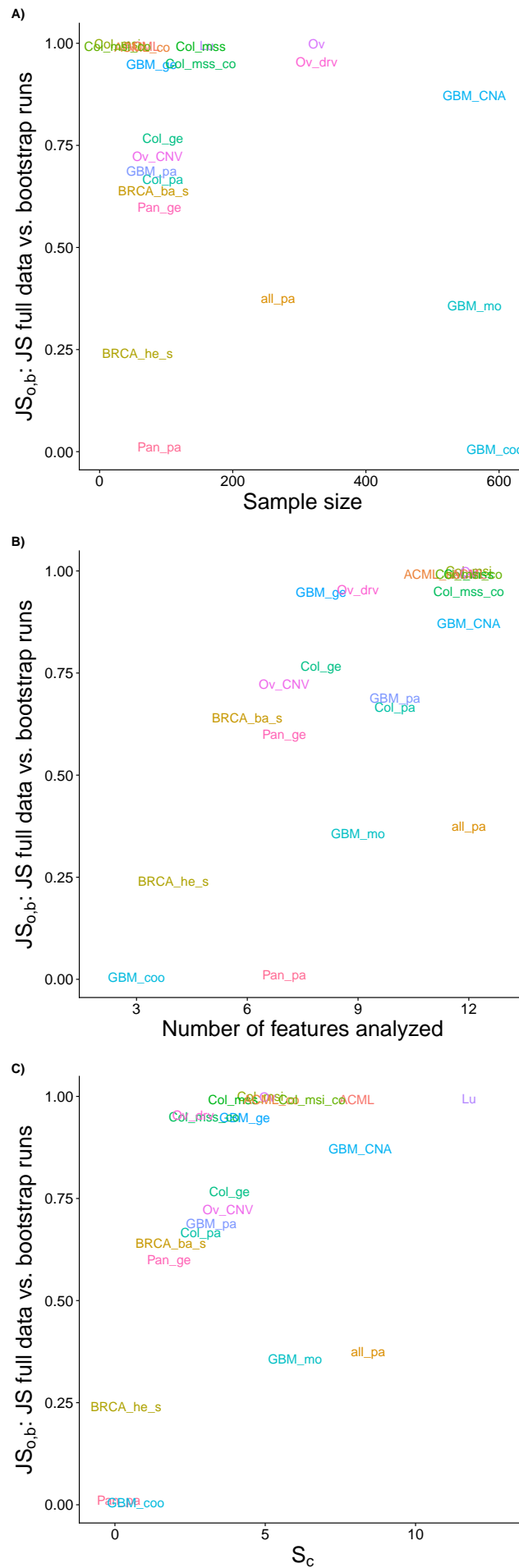


Figure 8: Summary patterns for average $JS_{0,b}$, JS for the full data vs. bootstrap runs as a function of sample size, number of features analyzed, and S_c (from the full original data set) for the cancer data sets using the statistics from the analysis with 12 features. See legend of Figure 7 for labels.