

# Modeling Cumulative Biological Phenomena with Suppes-Bayes Causal Networks

Daniele Ramazzotti<sup>1</sup>, Alex Graudenzi<sup>2</sup>, Giulio Caravagna<sup>3</sup> and Marco Antoniotti<sup>2</sup>

<sup>1</sup>Department of Pathology, Stanford University, Stanford, CA, USA. <sup>2</sup>Department of Informatics, Systems and Communication, University of Milano-Bicocca, Milan, Italy. <sup>3</sup>School of Informatics, University of Edinburgh, Edinburgh, UK.

Evolutionary Bioinformatics  
Volume 14: 1–10  
© The Author(s) 2018  
Reprints and permissions:  
sagepub.co.uk/journalsPermissions.nav  
DOI: 10.1177/1176934318785167



**ABSTRACT:** Several diseases related to cell proliferation are characterized by the accumulation of somatic DNA changes, with respect to wild-type conditions. Cancer and HIV are 2 common examples of such diseases, where the mutational load in the cancerous/viral population increases over time. In these cases, selective pressures are often observed along with competition, co-operation, and parasitism among distinct cellular clones. Recently, we presented a mathematical framework to model these phenomena, based on a combination of Bayesian inference and Suppes' theory of probabilistic causation, depicted in graphical structures dubbed Suppes-Bayes Causal Networks (SBCNs). The SBCNs are generative probabilistic graphical models that recapitulate the potential ordering of accumulation of such DNA changes during the progression of the disease. Such models can be inferred from data by exploiting likelihood-based model selection strategies with regularization. In this article, we discuss the theoretical foundations of our approach and we investigate in depth the influence on the model selection task of (1) the poset based on Suppes' theory and (2) different regularization strategies. Furthermore, we provide an example of application of our framework to HIV genetic data highlighting the valuable insights provided by the inferred SBCN

**KEYWORDS:** cumulative, phenomena, Bayesian graphical models, probabilistic causality

**RECEIVED:** January 11, 2018. **ACCEPTED:** May 27, 2018.

**TYPE:** Algorithm Development for Evolutionary Biological Computation – Review

**FUNDING:** The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: This work has been partially supported by grants from the SysBioNet project, an MIUR initiative for the Italian Roadmap of European Strategy Forum on Research Infrastructures (ESFRI).

**DECLARATION OF CONFLICTING INTERESTS:** The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

**CORRESPONDING AUTHOR:** Daniele Ramazzotti, Department of Pathology, Stanford University, Stanford, CA 94305, USA. Email: daniele.ramazzotti@stanford.edu

## Introduction

A number of diseases are characterized by the accumulation of genomic lesions in the DNA of a population of cells. Such lesions are often classified as *mutations*, if they involve one or few nucleotides, or *chromosomal alterations*, if they involve wider regions of a chromosome. The effect of these lesions, occurring randomly and inherited through cell divisions (ie, they are *somatic*), is that of inducing a phenotypic change in the cells. If the change is advantageous, then the clonal population might enjoy a *fitness* advantage over competing clones. In some cases, a natural selection process will tend to select the clones with more advantageous and inheritable traits. This particular picture can be framed in terms of Darwinian evolution as a scenario of survival of the fittest where, however, the prevalence of multiple heterogeneous populations is often observed.<sup>1</sup>

Cancer and HIV are 2 diseases where the mutational (from now on, we will use the term mutation to refer to the types of genomic lesions mentioned above) load in the cancerous/viral population of cells increases over time and drives phenotypic switches and disease progression. In this article, we specifically focus on these diseases, but many biological systems present similar characteristics.<sup>2–4</sup>

The emergence and development of cancer can be characterized as an *evolutionary process* involving a large population of cells, heterogeneous both in their genomes and in their epigenomes. In fact, genetic and epigenetic random alterations commonly occurring in any cell can occasionally be beneficial to the

neoplastic cells and confer to these clones a *functional selective advantage*. During clonal evolution, clones are generally selected for increased proliferation and survival, which may eventually allow the cancer clones to outgrow the competing cells and, in turn, may lead to invasion and metastasis.<sup>5,6</sup> By means of such a multistep stochastic process, cancer cells acquire over time a set of biological capabilities, ie, *hallmarks*.<sup>7,8</sup> However, not all the alterations are involved in this acquisition; as a matter of fact, in solid tumors, we observe an average of 33 to 66 genes displaying somatic mutations.<sup>9</sup> But only some of them are involved in the hallmark acquisition, ie, *drivers*, whereas the remaining ones are present in the cancer clones without increasing their *fitness*, ie, *passengers*.<sup>9</sup>

The onset of AIDS is characterized by the collapse of the immune system after a prolonged asymptomatic period, but its progression's mechanistic basis is still unknown. It was recently hypothesized that the elevated turnover of lymphocytes throughout the asymptomatic period results in the accumulation of deleterious mutations, which impairs immunological function, replicative ability, and viability of lymphocytes.<sup>10,11</sup> The failure of the modern combination therapies (ie, highly active antiretroviral therapy) of the disease is mostly due to the capability of the virus to escape from drug pressure by developing drug resistance. This mechanism is determined by HIV's high rates of replication and mutation. In fact, under fixed drug pressure, these mutations are



virtually nonreversible because they confer a strong selective advantage to viral populations.<sup>12,13</sup>

In the past decades, huge technological advancements led to the development of *next-generation sequencing* (NGS) techniques. These allow, in different forms and with different technological characteristics, to read out genomes from single cells or bulk.<sup>14–17</sup> Thus, we can use these technologies to quantify the presence of mutations in a sample. With these data at hand, we can therefore investigate the problem of inferring a *progression model* (PM) that recapitulates the ordering of accumulation of mutations during disease origination and development.<sup>18</sup> This problem allows different formulations according to the type of diseases that we are considering, the type of NGS data that we are processing, and other factors. We point the reader to the works by Caravagna et al and Beerenwinkel et al<sup>18,19</sup> for a review on PMs.

This work is focused on a particular class of mathematical models that are becoming successful to represent such mutational ordering. These are called *SBCNs* (the first use of these networks appears in the work by Ramazzotti et al,<sup>20</sup> and its earliest formal definition in the work by Bonchi et al<sup>21</sup>; SBCN<sup>21</sup>), and derived from a more general class of models, Bayesian Networks (BN<sup>22</sup>), that has been successfully exploited to model cancer and HIV progressions.<sup>23–25</sup> The SBCNs are probabilistic graphical models that are derived within a statistical framework based on Patrick Suppes<sup>26</sup> theory of *probabilistic causation*. Thus, the main difference between standard BNs and SBCNs is the encoding in the model of a set of causal axioms that describe the accumulation process. Both SBCNs and BNs are generative probabilistic models that induce a distribution of observing a particular mutational signature in a sample. But, the distribution induced by an SBCN is also consistent with the causal axioms and, in general, is different from the distribution induced by a standard BN.<sup>20</sup>

Informally, SBCNs are BNs depicting a set of well-defined statistical relations between pairs of events. In fact, when a first event precedes a second event in the network (ie, there is an arrow starting from the first event and pointing toward the second), this implies (1) a temporal relation where the first event happens *invariably before* the second, (2) statistical positive correlation between the 2 events, and (3) *relevance* of the first event in terms of being statistically informative in explaining the occurrences of the second event.

The motivation for adopting a causal framework on top of standard BNs is that, in the particular case of cumulative biological phenomena, SBCNs allow better inferential algorithms and data analysis pipelines to be developed.<sup>18,20,27</sup> Extensive studies in the inference of cancer progression have indeed shown that model selection strategies to extract SBCNs from NGS data achieve better performance than algorithms that infer BNs. In fact, SBCN's inferential algorithms have higher rate of detection of true-positive ordering relations and higher rate of filtering out false-positive ones. In general, these

algorithms also show better *scalability*, *resistance to noise in the data*, and *ability to work with datasets with few samples*.<sup>20,27</sup>

In this article, we give a formal definition of SBCNs, and we assess their relevance in modeling cumulative phenomena and investigate the influence of (1) Suppes' poset and (2) distinct maximum likelihood regularization strategies for model selection. We do this by performing extensive synthetic tests in operational settings that are representative of different possible types of progressions and data-harboring signals from heterogeneous populations.

## Suppes-Bayes Causal Networks

Theories of causality enjoy an old and prolific literature comprising contributions from many fields. Among them, some of the most prominent results are due to the efforts by Judea Pearl,<sup>28</sup> whose theories have gained a huge impact over the computational community. However, algorithms derived from this theory may sometimes lead to computational intractability. For this reason, in this work, we follow a different approach based on the theory of probabilistic causation by Patrick Suppes<sup>26</sup> that is particularly effective in modeling cumulative phenomena, yet still being computationally tractable.

Suppes<sup>26</sup> introduced the notion of *prima facie causation*. A *prima facie* relation between a cause  $u$  and its effect  $v$  is verified when the following 2 conditions are true: (1) *temporal priority* (TP), ie, any cause happens before its effect and (2) *probability raising* (PR), ie, the presence of the cause raises the probability of observing its effect.

*Definition 1. Probabilistic causation.*<sup>26</sup> For any 2 events  $u$  and  $v$ , occurring, respectively, at times  $t_u$  and  $t_v$ , under the mild assumptions that  $0 < P(u), P(v) < 1$ , the event  $u$  is called a *prima facie cause* of  $v$  if it occurs *before* and *raises the probability of*  $v$ , ie,

$$TP: t_u < t_v \quad (1)$$

$$PR: P(v|u) > P(v|\bar{u}) \quad (2)$$

Although the notion of *prima facie causation* has known limitations in the context of the general theories of causality,<sup>29</sup> this formulation seems to intuitively characterize the dynamics of phenomena driven by the *monotonic accumulation of events*. In these cases, in fact, a temporal order among the events is implied and, furthermore, the occurrence of an early event *positively correlates* to the subsequent occurrence of a later one. Thus, this approach seems appropriate to capture the notion of selective advantage emerging from somatic mutations that accumulate during, eg, cancer or HIV progression.

Let us now consider a graphical representation of the aforementioned dynamics in terms of a Bayesian graphical model.

*Definition 2. Bayesian network.*<sup>22</sup> The pair  $\mathcal{B} = \langle G, P \rangle$  is a BN, where  $G$  is a *directed acyclic graph* (DAG)  $G = (V, E)$  of  $V$

nodes and  $E$  arcs, and  $P$  is a distribution induced over the nodes by the graph. Let  $V = \{v_1, \dots, v_n\}$  be *random variables* and the *edges/arcs*  $E \subseteq V \times V$  encode the conditional dependencies among the variables. Define, for any  $v_i \in V$ , the *parent set*  $\pi(v_i) = \{x \mid x \rightarrow v_i \in E\}$ , then  $P$  defines the *joint probability distribution* induced by the BN as follows:

$$P(v_1, \dots, v_n) = \prod_{v_i \in V} P(v_i \mid \pi(v_i)) \quad (3)$$

All in all, a BN is a statistical model which succinctly represents the conditional dependencies among a set of *random variables*  $V$  through a DAG. More precisely, a BN represents a factorization of the *joint distribution*  $P(v_1, \dots, v_n)$  in terms of marginal (when  $\pi(v) = \emptyset$ ) and conditional probabilities  $P(\cdot \mid \cdot)$ .

We now consider a common situation when we deal with data (ie, observations) obtained at one (or a few) points in time, rather than through a time line. In this case, we are restricted to work with cross-sectional data, where no explicit information of time is provided. Therefore, we can model the dynamics of cumulative phenomena by means of a specific set of the general BNs where the nodes  $V$  represent the accumulating events as *Bernoulli random variables* taking values in  $\{0, 1\}$  based on their occurrence: the value of the variable is 1 if the event is observed and 0 otherwise. We then define a data set  $D$  of  $s$  cross-sectional samples over  $n$  Bernoulli random variables as follows:

$$\begin{matrix} v_1 & v_2 & \dots & v_n \\ \begin{matrix} s_1 \\ s_2 \\ \vdots \\ s_m \end{matrix} \begin{pmatrix} d_{1,1} & d_{1,2} & \dots & d_{1,n} \\ d_{2,1} & d_{2,2} & \dots & d_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ d_{m,1} & d_{m,2} & \dots & d_{m,n} \end{pmatrix} \end{matrix} = D \quad (4)$$

To extend BNs to account for Suppes' theory of probabilistic causation, we need to estimate for any variable  $v \in V$  its timing  $t_v$ . Because we are dealing with cumulative phenomena and, in the most general case, data that do not harbor any evident temporal information, we can use the *marginal probability*  $P(v)$  as a proxy for  $t_v$  (see also the commentary at the end of this section). (In many cases, the data that we can access are *cross-sectional*, meaning that the samples are collected at independent and unknown time points. For this reason, we have to resort on the simplest possible approach to estimate timings. However, in the case we were provided with explicit observations of time, the TP would be directly and, yet, more efficiently assessable.) In cancer and HIV, for instance, this makes sense because mutations are inherited through cells divisions and thus will fixate in the clonal populations during disease progression, ie, they are persistent.

**Definition 3.** *SBCN.*<sup>21</sup> A BN  $\mathcal{B}$  is an SBCN if and only if, for any edge  $v_i \rightarrow v_j \in E$ , Suppes' conditions (Definition 1) hold, that is,

$$P(v_i) > P(v_j) \text{ and } P(v_j \mid v_i) > P(v_j \mid \neg v_i) \quad (5)$$

It should be noted that SBCNs and BNs *have the same likelihood function*. Thus, SBCNs *do not embed any constraint of the cumulative process in the likelihood computation*, whereas approaches based on cumulative BNs do.<sup>25</sup> Instead, the structure of the model,  $E$ , is consistent with the causal model à-la-Suppes and, of course, this in turn reflects in the induced distribution. Even though this difference seem subtle, this is arguably the most interesting advantage of SBCNs over ad hoc BNs for cumulative phenomena.

*Model selection to infer a network from data.* The structure  $G$  of a BN (or of a SBCN) can be inferred from a data matrix  $D$ , as well as the parameters of the conditional distributions that define  $P$ . The model selection task is that of inferring such information from data; in general, we expect different models (ie, edges) if we infer a SBCN or a BN, as SBCNs encode Suppes' additional constraints.

The general structural learning, ie, *the model selection problem*, for BNs is NP-HARD<sup>22</sup>; hence, one needs to resort on approximate strategies. For each BN  $\mathcal{B}$ , a *log-likelihood function*  $\mathcal{L}(D \mid E)$  can be used to search in the space of structures (ie, the set of edges  $E$ ), together with a regularization function  $\mathcal{R}(\cdot)$  that penalizes overly complicated models. The network's structure is then inferred by solving the following optimization problem:

$$E_* = \underset{E}{\operatorname{argmax}} [\mathcal{L}(D \mid E) - \mathcal{R}(E)] \quad (6)$$

Moreover, the parameters of the conditional distributions can be computed by *maximum-likelihood estimation* for the set of edges  $E_*$ ; the overall solution is *locally optimal*.<sup>22</sup>

Model selection for SBCNs works in this very same way but constrains the search for valid solutions.<sup>20</sup> In particular, it scans only the subset of edges that are consistent with Definition 1, whereas a BN search will look for the full  $V \times V$  space. To filter pairs of edges, Suppes' conditions can be estimated from the data with solutions based, for instance, on bootstrap estimates.<sup>20</sup> The resulting model will satisfy, by construction, the conditions of probabilistic causation. It has been shown that if the underlying phenomenon that produced our data is characterised by an accumulation, then the inference of an SBCN, rather than a BN, leads to *much better precision and recall*.<sup>20,27</sup>

We conclude this section by discussing in detail the characteristics of the SBCNs and, in particular, to which extent they are capable of modeling cumulative phenomena.

*Temporal priority.* Suppes' first constraint ("event  $u$  is temporally preceding event  $v$ ," Definition 1) assumes an underlying *temporal (partial) order*  $\sqsubseteq_{\text{TP}}$  among the events/variables of the SBCN that we need to compute.

Cross-sectional data, unfortunately, are not provided with an explicit measure of time and hence  $\sqsubseteq_{\text{TP}}$  needs to be estimated from data  $D$  (we notice that in the case we were provided with explicit observations of time,  $\sqsubseteq_{\text{TP}}$  would be directly



and, yet, more efficiently assessable). The cumulative nature of the phenomenon that we want to model leads to a simple estimation of  $\sqsubseteq_{TP}$ : *the temporal priority TP* is assessed in terms of marginal frequencies<sup>20</sup>:

$$v_j \sqsubseteq_{TP} v_i \Leftrightarrow P(v_i) > P(v_j) \quad (7)$$

Thus, more frequent events, ie,  $v_i$ , are assumed to occur earlier, which is sound when we assume the accumulating events to be irreversible.

Temporal priority is combined with PR to complete Suppes' conditions for *prima facie* (see below). Its contribution is fundamental for model selection, as we now elaborate.

First of all, recall that the model selection problem for BNs is in general NP-HARD,<sup>22</sup> and that, as a result of the assessment of Suppes' conditions (TP and PR), we constrain our search space to the networks with a given order. Because of time irreversibility, marginal distributions induce a total ordering  $\sqsubseteq_{TP}$  on the  $v_i$ , ie, reflexing  $\leq$ . Learning BNs given a fixed order  $\sqsubseteq$ —even partial<sup>22</sup>—of the variables bounds the cardinality of the parent set as follows:

$$|\pi(v_x)| \leq |\{v_j \mid v_x \sqsubseteq v_j\}| \quad (8)$$

and, in general, it make inference easier than the general case.<sup>22</sup> Thus, by constraining Suppes' conditions via  $\sqsubseteq_{TP}$ 's total ordering, we drop down the model selection complexity. It should be noted that, after model selection, the ordering among the variables that we practically have in the selected arcs set  $E$  is in general partial; in the BN literature, this is sometimes called *poset*.

*Probability raising*. Besides TP, as a second constraint we further require that the arcs are consistent with the condition of PR: this leads to another relation  $\sqsubseteq_{PR}$ . Probability raising is equivalent to constraining for *positive statistical dependence*<sup>27</sup>:

$$\begin{aligned} v_j \sqsubseteq_{PR} v_i \\ \Leftrightarrow P(v_j \mid v_i) > P(v_j \mid \tilde{v}_i) \\ \Leftrightarrow P(v_i, v_j) > P(v_i)P(v_j) \end{aligned} \quad (9)$$

Thus, we model all and only the positive dependant relations. Definition 1 is thus obtained by selecting those PR relations that are consistent with TP

$$\sqsubseteq_{TP} \cap \sqsubseteq_{PR} \quad (10)$$

as the core of Suppes' characterization of causation is relevant.<sup>26</sup>

If  $\sqsubseteq_{TP}$  reduces the search space of the possible valid structures for the network by setting a specific total order to the nodes,  $\sqsubseteq_{PR}$  instead reduces the search space of the possible valid parameters of the network by requiring that the related conditional probability tables, ie,  $P(\cdot)$ , account only for positive

statistical dependencies. It should be noted that these constraints affect the structure and the parameters of the model, but the likelihood function is the same for BNs and SBCNs.

*Network simplification, regularization, and spurious causality*. Suppes'<sup>26</sup> criteria are known to be necessary but not sufficient to evaluate general causal claims. Even if we restrict to causal cumulative phenomena, the expressivity of probabilistic causality needs to be taken into account.

When dealing with small sample sized data sets (ie, small  $m$ ), many pairs of variables that satisfy Suppes' condition may be *spurious causes*, ie, *false positive*. (An edge is spurious when it satisfies Definition 1, but it is not actually the true model edge. For instance, for a linear model  $u \rightarrow v \rightarrow w$ , transitive edge  $u \rightarrow w$  is spurious. Small  $m$  induces further spurious associations in the data, not necessarily related to particular topological structures.). *False negatives* should be few and mostly due to noise in the data. Thus, it follows the following:

- We expect all the “statistically relevant” relations to be also *prima facie*<sup>20</sup>;
- We need to filter out spurious causality instances (a detailed account of these topics, the particular types of spurious structures, and their interpretation for different types of models are available in Ramazzotti and colleagues<sup>20,27,30</sup>), as we know that *prima facie* overfits.

A model selection strategy which exploits a regularization schema seems thus the best approach to the task. Practically, this strategy will select simpler (ie, *sparse*) models according to a *penalized likelihood fit criterion*—for this reason, it will filter out edges proportionally to how much the regularization is stringent. Also, it will rank spurious association according to a criterion that is consistent with Suppes' intuition of causality, as likelihood relates to statistical (in)dependence. Alternatives based on likelihood itself, ie, without regularization, do not seem viable to minimize the effect of likelihood's overfit, that happens unless  $m \rightarrow \infty$ .<sup>22</sup> In fact, one must recall that due to statistical noise and sample size, exact statistical (in)dependence between pair of variables is never (or very unlikely) observed.

### Modeling heterogeneous populations

Complex biological processes, eg, *proliferation, nutrition, apoptosis*, are orchestrated by multiple cooperative networks of proteins and molecules. Therefore, different “mutants” can evade such control mechanisms in different ways. Mutations happen as a random process that is unrelated to the relative fitness advantage that they confer to a cell. As such, different cells will deviate from wild type by exhibiting different mutational signatures during disease progression. This has an implication for many cumulative diseases that arise from populations that are *heterogeneous*, both at the level of the single patient (inpatient heterogeneity) and in the population of patients (interpatient

**Table 1.** Definitions for CMPN, DMPN, and XMPN.

$\text{CMPN: } P\left(v \left  \sum \pi(v) = 1 \pi(v) \right. \right) = \theta P\left(v \left  \sum \pi(v) < 1 \pi(v) \right. \right) \leq \varepsilon$	(11)
$\text{DMPN: } P\left(v \left  \sum \pi(v) > 0 \right. \right) = \theta P\left(v \left  \sum \pi(v) = 0 \right. \right) \leq \varepsilon$	(12)
$\text{XMPN: } P\left(v \left  \sum \pi(v) = 1 \right. \right) = \theta P\left(v \left  \sum \pi(v) \neq 1 \right. \right) \leq \varepsilon$	(13)

heterogeneity). Heterogeneity introduces significant challenges in designing effective treatment strategies, and major efforts are ongoing at deciphering its extent for many diseases such as cancer and HIV.<sup>18,20,28</sup>

We now introduce a class of mathematical models that are suitable at modeling heterogenous progressions. These models are derived by augmenting BNs with logical formulas and are called *monotonic progression networks* (MPNs).<sup>31,32</sup> The MPNs represent the progression of events that accumulate *monotonically* (the events accumulate over time and when later events occur earlier events are observed as well) over time, where the conditions for any event to happen is described by a probabilistic version of the canonical boolean operators, ie, conjunction ( $\wedge$ ), inclusive disjunction ( $\vee$ ), and exclusive disjunction ( $\oplus$ ).

Following Farahani and Lagergren<sup>31</sup> and Korsunsky et al,<sup>32</sup> we define 1 type of MPNs for each boolean operator: the *conjunctive* (CMPN), the *disjunctive semimonotonic* (DMPN), and the *exclusive disjunction* (XMPN). The operator associated with each network type refers to the logical relation among the parents that eventually lead to the common effect to occur.

*Definition 4. Monotonic Progression Networks.*<sup>31,32</sup> The MPNs are BNs that, for  $\theta, \varepsilon \in [0, 1]$  and  $\theta \gg \varepsilon$ , satisfy the conditions shown in Table 1 for each  $v \in V$ .

Here,  $\theta$  represents the conditional probability of any “effect” to follow its preceding “cause” and  $\varepsilon$  models the probability of any noisy observation—that is the observation of a sample where the effects are observed without their causes. Note that the above inequalities define, for each type of MPN, specific constraints to the induced distributions. These are sometimes termed, according to the probabilistic logical relations, *noisy-AND*, *noisy-OR*, and *noisy-XOR* networks.<sup>28,32</sup>

*Model selection with heterogeneous populations.* When dealing with heterogeneous populations, the task of model selection, and, more in general, any statistical analysis, are non-trivial. One of the main reasons for this state of affairs is the emergence of statistical paradoxes such as Simpson’s paradox.<sup>33,34</sup> This phenomenon refers to the fact that sometimes, associations among dichotomous variables, which are similar within subgroups of a population, eg, women and men, change their statistical trend if the individuals of the subgroups are pooled together. Let us now recall a famous example to this regard. The admissions of the University of Berkeley for the fall of

1973 showed that men applying were much more likely than women to be admitted with a difference that was unlikely to be due to chance. But, when looking at the individual departments separately, it emerged that 6 out of 85 were indeed biased in favor of women, whereas only 4 presented a slightly bias against them. The reason for this inconsistency was due to the fact that women tended to apply to competitive departments which had low rates of admissions, whereas men tended to apply to less-competitive departments with high rates of admissions, leading to an apparent bias toward them in the overall population.<sup>35</sup>

Similar situations may arise in cancer when different populations of cancer samples are mixed. As an example, let us consider an hypothetical progression leading to the alteration of gene  $e$ . Let us now assume that the alterations of this gene may be due to the previous alterations of either gene  $c_1$  or gene  $c_2$  exclusively. If this was the case, then we would expect a significant pattern of selective advantage from any of its causes to  $e$  if we were able to stratify the patients accordingly to either alteration  $c_1$  or  $c_2$ , but we may lose these associations when looking at all the patients together.

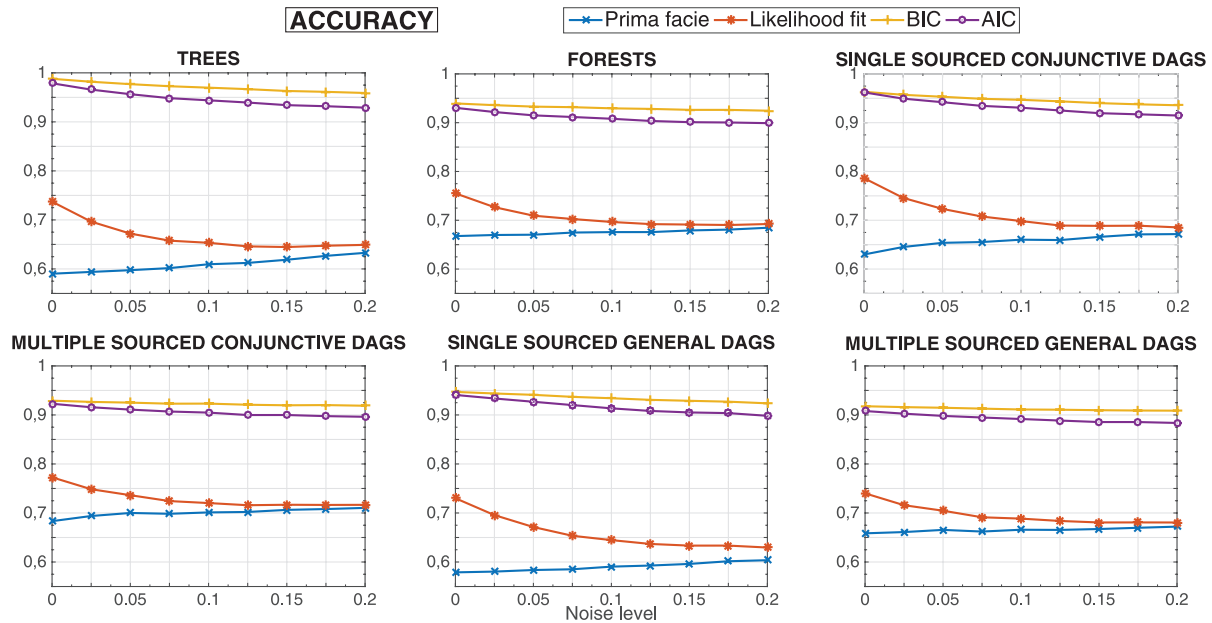
In the work by Ramazzotti et al,<sup>20</sup> the notion of *progression pattern* is introduced to describe this situation, defined as a boolean relation among all the genes, members of the parent set of any node as the ones defined by MPNs. To this extent, the authors extend Suppes’ definition of prima facie causality to account for such patterns rather than for relations among atomic events as for Definition 1. Also, they claim that general MPNs can be learned in polynomial time provided that the data set given as input is *lifted*<sup>20</sup> with a Bernoulli variable per causal relation representing the logical formula involving any parent set.

Following Ramazzotti and colleagues,<sup>20,30</sup> we now consider any formula in *conjunctive normal form* (CNF):

$$\varphi = c_1 \wedge \dots \wedge c_n$$

where each  $c_i$  is a *disjunctive clause*  $c_i = c_{i,1} \vee \dots \vee c_{i,k}$  over a set of literals and each literal represents an event (a Boolean variable) or its negation. By following analogous arguments as the ones used before, we can extend Definition 1 as follows.

*Definition 5. CNF probabilistic causation.*<sup>20,30</sup> For any CNF formula  $\varphi$  and  $e$ , occurring, respectively, at times  $t_\varphi$  and  $t_e$ , under the mild assumptions that  $0 < P(\varphi), P(e) < 1$ ,  $\varphi$  is a *prima facie cause* of  $e$  if



**Figure 1.** Performance of the inference on simulated data sets of 100 samples and networks of 15 nodes in terms of accuracy for the 6 considered topological structures. The y-axis refers to the performance, whereas the x-axis represents the different noise levels.

$$t_\varphi < t_e \text{ and } P(e|\varphi) > P(e|\bar{\varphi}) \quad (14)$$

Given these premises, we can now define the extended SBCNs, an extension of SBCNs which allows to model heterogeneity as defined probabilistically by MPNs.

*Definition 6. Extended SBCN.* A BN  $\mathcal{B}$  is an extended SBCN if and only if, for any edge  $\varphi_i \rightarrow v_j \in E$ , Suppes' generalized conditions (Definition 5) hold, that is,

$$P(\varphi_i) > P(v_j) \text{ and } P(v_j | v_\varphi) > P(v_j | \neg v_\varphi) \quad (15)$$

### Evaluation on Simulated Data

We now evaluate the performance of the inference of SBCN on simulated data, with specific attention on the impact of the constraints based on Suppes' probabilistic causation on the overall performance. All the simulations are performed with the following settings.

We consider 6 different topological structures: the first 2 where any node has at the most one predecessor, ie, (1) trees, (2) forests, and the others where we set a limit of 3 predecessors and, hence, we consider (3) DAGs with a single source and conjunctive parents, (4) DAGs with multiple sources and conjunctive parents, (5) DAGs with a single source and disjunctive parents, and (6) DAGs with multiple sources and disjunctive parents. For each of these configurations, we generate 100 random structures.

Moreover, we consider 4 different sample sizes (50, 100, 150, and 200 samples) and 9 noise levels (ie, probability of a random entry for the observation of any node in a sample) from 0% to 20% with step 2.5%. Furthermore, we repeat the above settings for networks of 10 and 15 nodes. Any

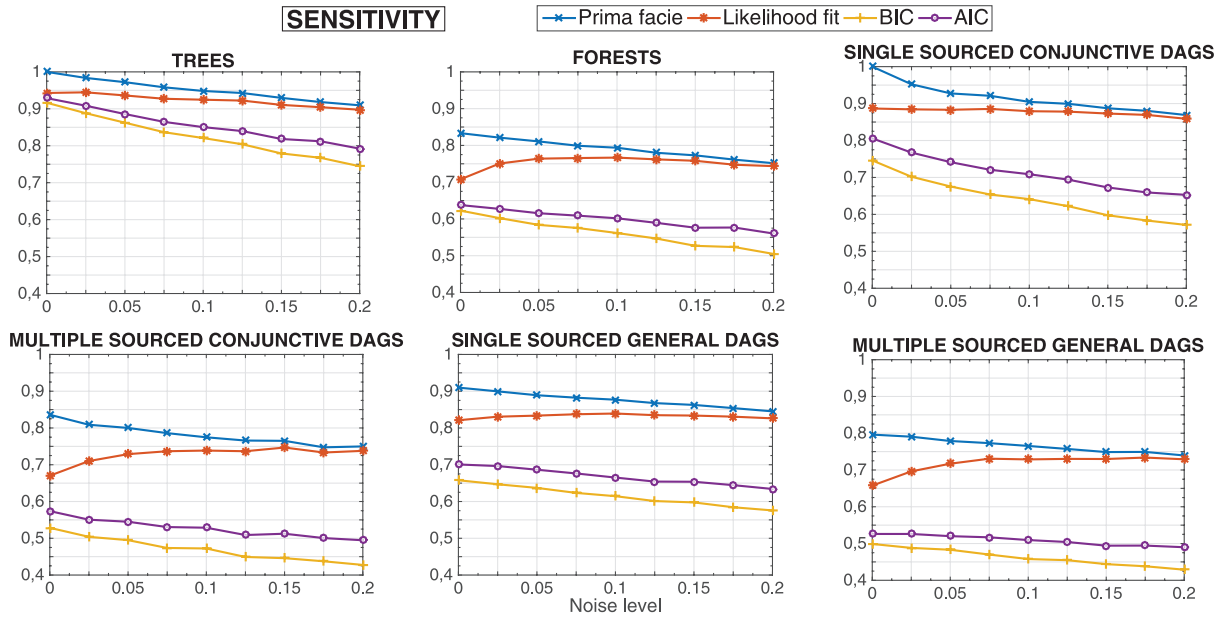
configuration is then sampled 10 times independently, for a total of more than 4 million distinct simulated data sets.

The sequencing quality of mutation profiling for diseases such as cancer and HIV depends on multiple factors such as, but not limited to, depth and coverage of the sequencing. In this work, we introduced errors in the data by means of a random model of noise. A detailed analysis of how more sophisticated models of noise can affect the inference is out of the scope of this study and left for future works.

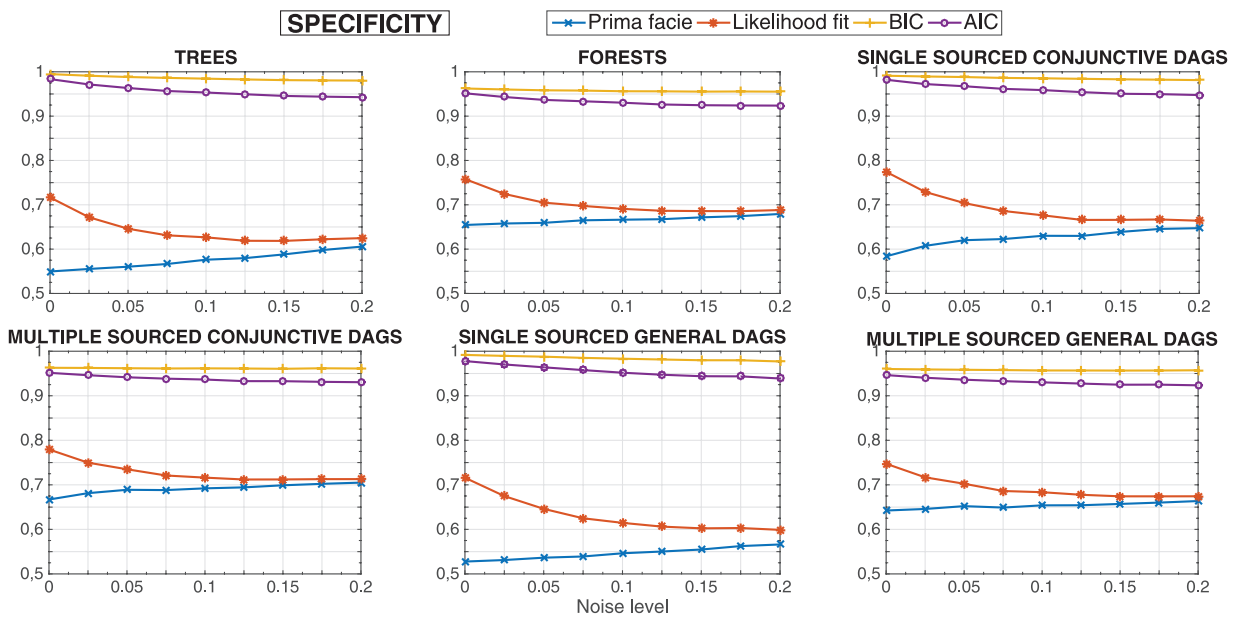
Finally, the inference of the structure of the SBCN is performed using the algorithm proposed in the work by Ramazzotti et al<sup>20</sup> and the performance is assessed in terms of  $accuracy = (TP + TN) / (TP + TN + FP + FN)$ ,  $sensitivity = TP / (TP + FN)$ , and  $specificity = TN / (FP + TN)$  with  $TP$  and  $FP$  being the true and false positive (we define as positive any arc that is present in the network) and  $TN$  and  $FN$  being the true and false negative (we define negative any arc that is not present in the network). All these measures are values in  $[0,1]$  with results close to 1 indicators of good performance.

In Figures 1 to 3, we show the performance of the inference on simulated data sets of 100 samples and networks of 15 nodes in terms of accuracy, sensitivity, and specificity for different settings which we discuss in detail in the next paragraphs.

*Suppes' prima facie conditions are necessary but not sufficient.* We first discuss the performance by applying *only* the prima facie criteria and we evaluate the obtained prima facie network in terms of accuracy, sensitivity, and specificity on simulated data sets of 100 samples and networks of 15 nodes (see Figures 1 to 3). As expected, the sensitivity is much higher than the specificity implying the significant impact of false positives rather than false negatives for the networks of the prima facie arcs. This result is



**Figure 2.** Performance of the inference on simulated data sets of 100 samples and networks of 15 nodes in terms of sensitivity for the 6 considered topological structures. The y-axis refers to the performance, whereas the x-axis represents the different noise levels.



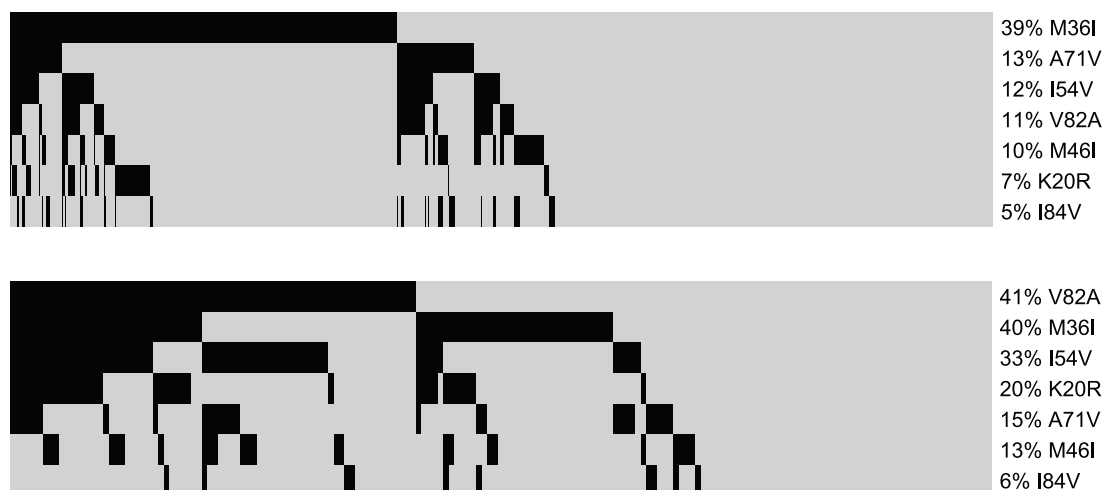
**Figure 3.** Performance of the inference on simulated data sets of 100 samples and networks of 15 nodes in terms of specificity for the 6 considered topological structures. The y-axis refers to the performance while the x-axis represents the different noise levels.

indeed expected being Suppes' criteria mostly capable of removing some of the arcs which do not represent valid causal relations rather than assess the exact set of valid arcs. Interestingly, the false negatives are still limited even when we consider DMPN, ie, when we do not have guarantees for the algorithm of Ramazzotti et al<sup>20</sup> to converge. The same simulations with different sample sizes (50, 150, and 200 samples) and on networks of 10 nodes present a similar trend (results not shown here).

*The likelihood score overfits the data.* In Figures 1 to 3, we also show the performance of the inference by likelihood fit (without any regularizator) on the prima facie network in terms of accuracy, sensitivity, and specificity on simulated data sets of 100

samples and networks of 15 nodes. Once again, in general, sensitivity is much higher than specificity implying also in this case a significant impact of false positives rather than false negatives for the inferred networks. These results make explicit the need for a regularization heuristic when dealing with real (not infinite) sample sized data sets as discussed in the next paragraph. Another interesting consideration comes from the observation that the prima facie networks and the networks inferred via likelihood fit without regularization seem to converge to the same performance as the noise level increases. This is due to the fact that, in general, the prima facie constraints are very conservative in the sense that false positives are admitted as long as false





**Figure 4.** Mutations detected in the genome for 179 patients with HIV under ritonavir (top) and 1035 under indinavir (bottom). Each black rectangle denotes the presence of a mutation in the gene annotated to the right of the plot; percentages correspond to marginal probabilities.

negatives are limited. When the noise level increases, the positive dependencies among nodes are generally reduced and, hence, less arcs pass the prima facie cut for positive dependency. Also in this case, the same simulations with different sample sizes (50, 150, and 200 samples) and on networks of 10 nodes present a similar trend (results not shown here).

*Model selection with different regularization strategies.* We now investigate the role of different regularizations on the performance. In particular, we consider 2 commonly used regularizations: (1) the *Bayesian information criterion* (BIC)<sup>36</sup> and (2) the *Akaike information criterion* (AIC).<sup>37</sup>

Although BIC and AIC are both scores based on maximum likelihood estimation and a penalization term to reduce overfitting, yet with distinct approaches, they produce significantly different behaviors. More specifically, BIC assumes the existence of one *true* statistical model which is generating the data, whereas AIC aims at finding the best approximating model to the unknown data-generating process. As such, BIC may likely underfit, whereas, conversely, AIC might overfit. (Thus, BIC tends to make a trade-off between the likelihood and model complexity with the aim of inferring the statistical model which generates the data. This makes it useful when the purpose is to detect the best model describing the data. Instead, asymptotically, minimizing AIC is equivalent to minimizing the cross validation value.<sup>38</sup> It is this property that makes the AIC score useful in model selection when the purpose is prediction. Overall, the choice of the regularizator tunes the level of sparsity of the retrieved SBCN and, yet, the confidence of the inferred arcs.)

The performance on simulated data sets are shown in Figures 1 to 3. In general, the performance is improved in all the settings with both regularizators, as they succeed in shrinking toward sparse networks.

Furthermore, we observe that the performance obtained by SBCNs is still good even when we consider simulated data generated by DMPN. Although in this case we do not have any guarantee of convergence, in practice, the algorithm seems

efficient in approximating the generative model. In conclusion, without any further input, SBCNs can model CMPNs and, yet, depict the more significant arcs of DMPNs. To infer XMPN, the data set needs to be lifted.<sup>20</sup>

The same simulations with different sample sizes (50, 150, and 200 samples) and on networks of 10 nodes present a similar trend (results not shown here).

### Application to HIV Genetic Data

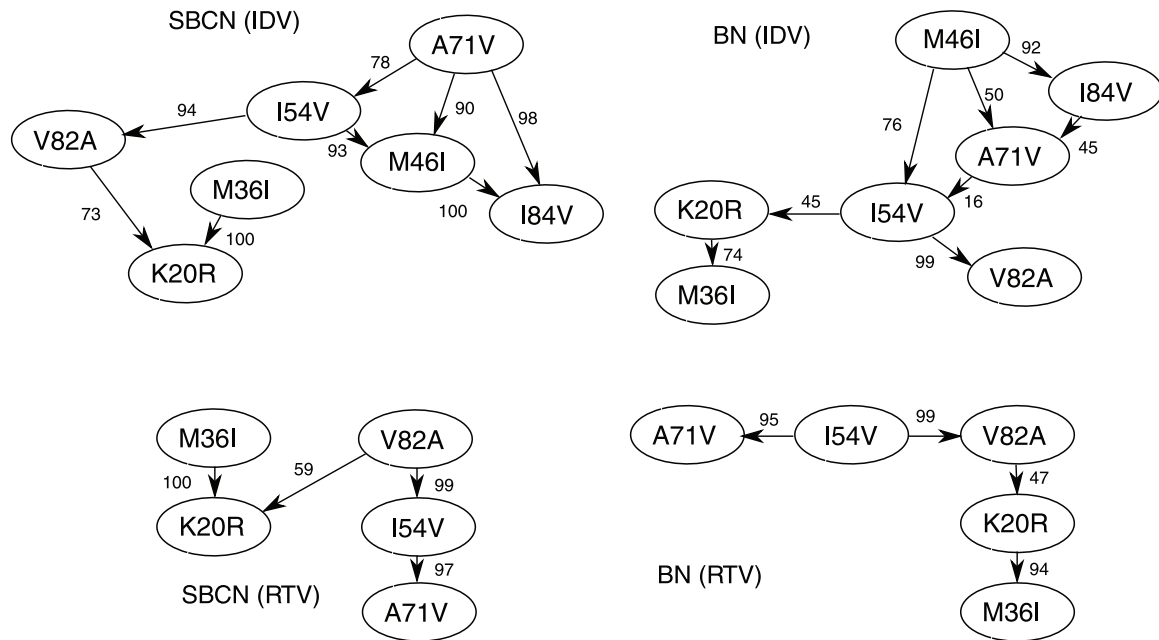
We now present an example of application of our framework on HIV genomic data. In particular, we study drug resistance in patients under antiretroviral therapy and we select a set of 7 amino acid alterations in the HIV genome to be depicted in the resulting graphical model, namely, *K20R*, *M36I*, *M46I*, *I54V*, *A71V*, *V82A*, *I84V*, where, as an example, the genomic event *K20R* describes a mutation from lysine (*K*) to arginine (*R*) at position 20 of the HIV protease.

In this study, we consider data sets from the Stanford HIV Drug Resistance Database<sup>39</sup> for 2 protease inhibitors, ritonavir (RTV) and indinavir (IDV). The first data set consists of 179 samples (see Figure 4) and the second of 1035 samples (see Figure 4).

We then infer a model on these data sets by both BN and SBCN. We show the results in Figures 5 where each node represents a mutation and the scores on the arcs measure the confidence in the found relation by nonparametric bootstrap.

In this case, it is interesting to observe that the set of dependency relations (ie, any pair of nodes connected by an arc, without considering its direction) depicted both by SBCNs and BNs is very similar, with the main difference being the direction of some connection. This difference is expected and can be attributed to the constrain of TP adopted in the SBCNs. Furthermore, we also observe that most of the found relations in the SBCN are more confident (ie, higher bootstrap score) than the one depicted in the related BN, leading us to observe a higher statistical confidence in the models inferred by SBCNs.





**Figure 5.** HIV progression of patients under ritonavir or indinavir (Figure 4) described as a Bayesian Network or as a Suppes-Bayes Causal Network. Edges are annotated with nonparametric bootstrap scores.

### Conclusions

In this work, we investigated the properties of a constrained version of BN, named SBCN, which is particularly sound in modeling the dynamics of system driven by the monotonic accumulation of events, thanks to encoded poset based on Suppes’ theory of probabilistic causation. In particular, we showed how SBCNs can, in general, describe different types of MPN, which makes them capable of characterizing a broad range of cumulative phenomena not limited to cancer evolution and HIV drug resistance.

Besides, we investigated the influence of Suppes’ poset on the inference performance with cross-sectional synthetic data sets. In particular, we showed that Suppes’ constraints are effective in defining a partially order set accounting for accumulating events, with very few false negatives, yet many false positives. To overcome this limitation, we explored the role of 2 maximum likelihood regularization parameters, ie, BIC and AIC, the former being more suitable to test previously conjectured hypotheses and the latter to predict novel hypotheses.

Finally, we showed on a data set of HIV genomic data how SBCN can be effectively adopted to model cumulative phenomena, with results presenting a higher statistical significance compared with standard BNs.

### Author Contributions

All authors performed the analysis and wrote the manuscript.

### REFERENCES

- Burrell RA, McGranahan N, Bartek J, Swanton C. The causes and consequences of genetic heterogeneity in cancer evolution. *Nature*. 2013;501:338–345.
- Weinreich DM, Delaney NF, DePristo MA, Hartl DL. Darwinian evolution can follow only very few mutational paths to fitter proteins. *Science*. 2006;312:111–114.

- Poelwijk FJ, Kiviet DJ, Weinreich DM, Tans SJ. Empirical fitness landscapes reveal accessible evolutionary paths. *Nature*. 2007;445:383–386.
- Lozovsky ER, Chookajorn T, Brown KM, et al. Stepwise acquisition of pyrimethamine resistance in the malaria parasite. *Proc Natl Acad Sci U S A*. 2009;106:12025–12030.
- Nowell PC. The clonal evolution of tumor cell populations. *Science*. 1976;194:23–28.
- Merlo LM, Pepper JW, Reid BJ, Maley CC. Cancer as an evolutionary and ecological process. *Nat Rev Cancer*. 2006;6:924–935.
- Hanahan D, Weinberg RA. The hallmarks of cancer. *Cell*. 2000;100:57–70.
- Hanahan D, Weinberg RA. Hallmarks of cancer: the next generation. *Cell*. 2011;144:646–674.
- Vogelstein B, Papadopoulos N, Velculescu VE, Zhou S, Diaz LA, Kinzler KW. Cancer genome landscapes. *Science*. 2013;339:1546–1558.
- Galvani AP. The role of mutation accumulation in HIV progression. *Proc Biol Sci*. 2005;272:1851–1858.
- Seifert D, Di Giallonardo F, Metzner KJ, Günthard HF, Beerenwinkel N. A framework for inferring fitness landscapes of patient-derived viruses using quasi-species theory. *Genetics*. 2015;199:191–203.
- Perrin L, Telenti A. HIV treatment failure testing for HIV resistance in clinical practice. *Science*. 1998;280:1871–1873.
- Vandamme AM, Van Laethem K, De Clercq E. Managing resistance to anti-HIV drugs. *Drugs*. 1999;57:337–361.
- Navin NE. Cancer genomics: one cell at a time. *Genome Biol*. 2014;15:452.
- Wang Y, Waters J, Leung ML, et al. Clonal evolution in breast cancer revealed by single nucleus genome sequencing. *Nature*. 2014;512:155–160.
- Gerlinger M, Rowan AJ, Horswell S, et al. Intratumor heterogeneity and branched evolution revealed by multiregion sequencing. *N Engl J Med*. 2012;366:883–892.
- Gerlinger M, Horswell S, Larkin J, et al. Genomic architecture and evolution of clear cell renal cell carcinomas defined by multiregion sequencing. *Nat Genet*. 2014;46:225–233.
- Caravagna G, Graudenzi A, Ramazzotti D, et al. Algorithmic methods to infer the evolutionary trajectories in cancer progression. *Proc Natl Acad Sci U S A*. 2016;113:E4025–E4034.
- Beerenwinkel N, Schwarz RF, Gerstung M, Markowitz F. Cancer evolution: mathematical models and computational inference. *Syst Biol*. 2015;64:e1–e25.
- Ramazzotti D, Caravagna G, Olde Loohuis L, et al. CAPRI efficient inference of cancer progression models from cross-sectional data. *Bioinformatics*. 2015;31:3016–3026. doi:10.1093/bioinformatics/btv296.
- Bonchi F, Hajian S, Mishra B, Ramazzotti D. Exposing the probabilistic causal structure of discrimination. *Int J Data Sci Anal*. 2017;3:1–21.
- Koller D, Friedman N. *Probabilistic Graphical Models Principles and Techniques*. Cambridge, MA: MIT Press; 2009.
- Desper R, Jiang F, Kallioniemi OP, Moch H, Papadimitriou CH, Schäffer AA. Inferring tree models for oncogenesis from comparative genome hybridization data. *J Comput Biol*. 1999;6:37–51.

24. Beerenwinkel N, Eriksson N, Sturmfels B. Conjunctive Bayesian networks. *Bernoulli*. 2007;13:893–909.
25. Gerstung M, Baudis M, Moch H, Beerenwinkel N. Quantifying cancer progression with conjunctive Bayesian networks. *Bioinformatics*. 2009;25:2809–2815.
26. Suppes P. *A Probabilistic Theory of Causality*. Amsterdam: North-Holland Publishing Company; 1970.
27. Loohuis LO, Caravagna G, Graudenzi A, et al. Inferring tree causal models of cancer progression with probability raising. *PLoS ONE*. 2014;9:e108358.
28. Pearl J. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Francisco, CA: Morgan Kaufmann Publishers; 2014.
29. Hitchcock C. Probabilistic causation. In: EN Zalta, ed. *Stanford Encyclopedia of Philosophy*. Winter ed; 2012.
30. Ramazzotti D. *A Model of Selective Advantage for the Efficient Inference of Cancer Clonal Evolution* [PhD thesis]. Milan: University of Milan; 2016.
31. Farahani HS, Lagergren J. Learning oncogenetic networks by reducing to mixed integer linear programming. *PLoS ONE*. 2013;8:e65773.
32. Korsunsky I, Ramazzotti D, Caravagna G, Mishra B. Inference of cancer progression models with biological noise. arXiv:1408.6032; 2014.
33. Yule GU. Notes on the theory of association of attributes in statistics. *Biometrika*. 1903;2:121–134.
34. Simpson EH. The interpretation of interaction in contingency tables. *J Roy Stat Soc B Met*. 1951;13:238–241.
35. Bickel PJ, Hammel EA, O'Connell JW, et al. Sex bias in graduate admissions: data from Berkeley. *Science*. 1975;187:398–404.
36. Schwarz G, et al. Estimating the dimension of a model. *Ann Stat*. 1978;6:461–464.
37. Akaike H. Information theory and an extension of the maximum likelihood principle. In: Parzen E, Tanabe K, Kitagawa G, eds. *Selected Papers of Hirotugu Akaike*. New York: Springer; 1998:199–213.
38. Stone M. An asymptotic equivalence of choice of model by cross-validation and Akaike's criterion. *J Roy Stat Soc B Met*. 1977;39:44–47.
39. Rhee SY, Gonzales MJ, Kantor R, Betts BJ, Ravela J, Shafer RW. Human immunodeficiency virus reverse transcriptase and protease sequence database. *Nucleic Acids Res*. 2003;31:298–303.