

Bayesian MAP Model Selection of Chain Event Graphs

G. Freeman^a, J.Q. Smith^a

^a*Department of Statistics, University of Warwick, Coventry, CV4 7AL*

Abstract

Chain event graphs are graphical models that while retaining most of the structural advantages of Bayesian networks for model interrogation, propagation and learning, more naturally encode asymmetric state spaces and the order in which events happen than Bayesian networks do. In addition, the class of models that can be represented by chain event graphs for a finite set of discrete variables is a strict superset of the class that can be described by Bayesian networks. In this paper we demonstrate how with complete sampling, conjugate closed form model selection based on product Dirichlet priors is possible, and prove that suitable homogeneity assumptions characterise the product Dirichlet prior on this class of models. We demonstrate our techniques using two educational examples.

Keywords: chain event graphs, Bayesian model selection, Dirichlet distribution

1. Introduction

Bayesian networks (BNs) are currently one of the most widely used graphical models for representing and analysing finite discrete multivariate distributions, with their explicit coding of conditional independence relationships between a system's variables [1, 2]. However, despite their power and usefulness, it has long been known that BNs cannot fully or efficiently represent certain common scenarios [3]. These include situations where the state space of a variable is known to depend on other variables, or where the conditional independence between variables is itself dependent on the values of other variables, called CONTEXT-SPECIFIC INDEPENDENCE in the literature [4]. Some examples of the latter phenomenon are given by Poole and Zhang [5]. In order to overcome such deficiencies, enhancements have been proposed to the canonical Bayesian network. Poole and Zhang [5], for example, define CONTEXTUAL BELIEF NETWORKS. These, however, don't represent the context-specific independence relationships graphically, thus undermining the rationale for using a graphical model in the first place. Boutilier et al [4], meanwhile, keeps the BN in place but additionally uses trees to describe the structures of the conditional probability distributions.

A new graphical model — the Chain Event Graph (CEG) — first propounded in [6], aims to represent the context-specific independences and asymmetric sample spaces of a model in a single graph. To this end, CEGs are based not on Bayesian networks, but on event trees (ETs) [7]. Event trees are trees where nodes represent situations — i.e. scenarios in which a unit might find itself — and each node's extending edges represent possible future situations that can develop from the

Email addresses: g.freeman@warwick.ac.uk (G. Freeman), j.q.smith@warwick.ac.uk (J.Q. Smith)

current one. It follows that every atom of the event space is encoded by exactly one root-to-leaf path, and each root-to-leaf path corresponds to exactly one atomic event. It has been argued that ETs are expressive frameworks to directly and accurately represent beliefs about a process, particularly when the model is described most naturally, as in the example below, through how situations might unfold [7]. However, as explained in [6], ETs can contain excessive redundancy in their structure, with subtrees describing probabilistically isomorphic unfoldings of situations being represented separately. They are also unable to explicitly express a model's non-trivial conditional independences. The CEG deals with these shortcomings by combining the subtrees that describe identical subprocesses (see [6] for further details), so that the CEG derived from a particular ET has a simpler topology while in turn expressing more conditional independence statements than is possible through an ET.

We illustrate the construction and the types of symmetries that can be coded using a CEG with the following running example, which exemplifies the types of hypotheses we plan to search over in our model selection.

Example 1. *Successful students on a one year programme study components A and B , but not everyone will study the components in the same order: each student will be allocated to study either module A or B for the first 6 months and then the other component for the final 6 months. After the first 6 months each student will be examined on their allocated module and be awarded a distinction (denoted with D), a pass (P) or a fail (F), with an automatic opportunity to resit the module in the last case. If they resit then they can pass and be allowed to proceed to the other component of their course, or fail again and be permanently withdrawn from the programme. Students who have succeeded in proceeding to the second module can again either fail, pass or be awarded a distinction. On this second round, however, there is no possibility of resitting if the component is failed. With an obvious extension of the labelling, this system can be depicted by the event tree given in Figure 1.*

To specify a full probability distribution for this model it is sufficient to only specify the distributions associated with the unfolding of each situation a student might reach. However, in many applications it is often natural to hypothesise a model where the distribution associated with the unfolding from one situation is assumed identical to another. Situations that are thus hypothesised to have the same transition probabilities to their children are said to be in the same *stage*. Thus in Example 1 suppose that as well as subscribing to the ET of Figure 1 we want to consider the plausibility of the following three hypotheses:

1. The chances of doing well in the second component are the same whether the student passed first time or after a resit.
2. The components A and B are equally hard.
3. The distribution of marks for the second component is unaffected by whether students passed or got a distinction for the first component.

Each of these hypotheses can be identified with a partitioning of the non-leaf nodes (situations). In Figure 1 the set of situations is

$$S = \{V_0, A, B, P_{1,A}, P_{1,B}, D_{1,A}, D_{1,B}, F_{1,A}, F_{1,B}, P_{R,A}, P_{R,B}\}.$$

The partition C of S that encodes the above three hypotheses consists of the stages $u_1 = \{A, B\}$, $u_2 = \{F_{1,A}, F_{1,B}\}$, and $u_3 = \{P_{1,A}, P_{1,B}, P_{R,A}, P_{R,B}, D_{1,A}, D_{1,B}\}$ together with the singleton $u_0 =$

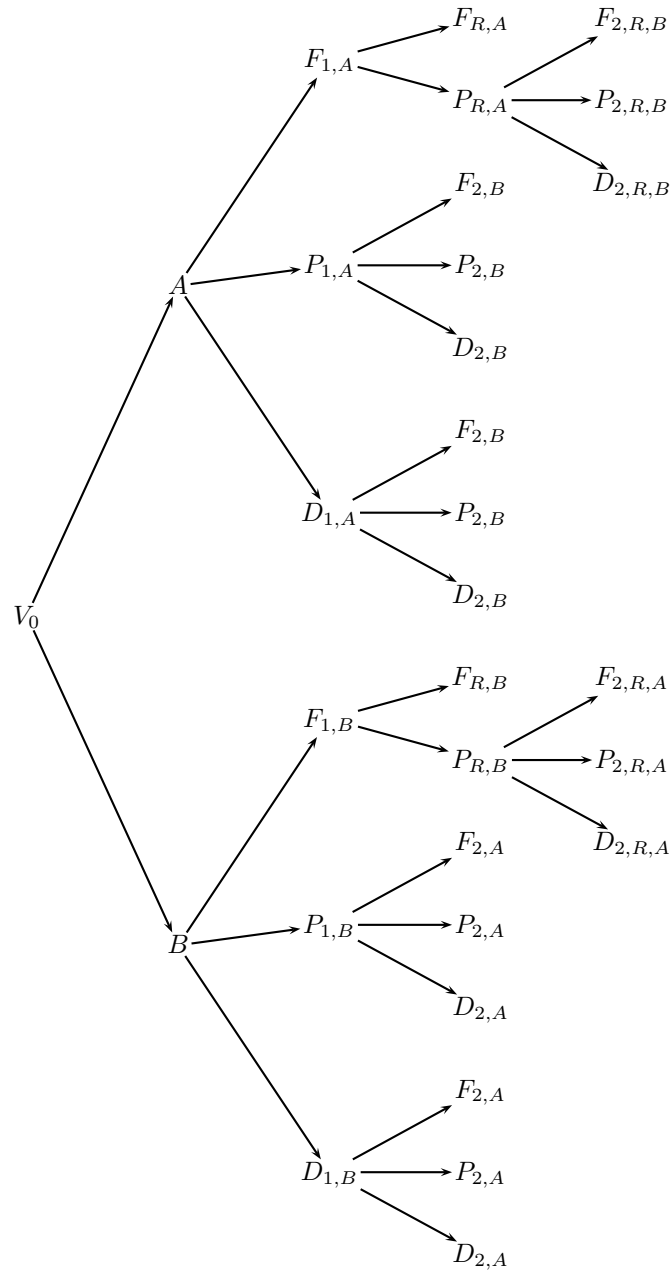


Figure 1: Event tree of a student’s potential progress through a hypothetical course described in Example 1. Each non-leaf node represents a juncture at which a random event will take place, with the selection of possible outcomes represented by the edges emanating from that node. Each edge distribution is defined conditional on the path passed through earlier in the tree to reach the specific node.

$\{V_0\}$. Thus the second stage u_2 , for example, implies that the probabilities on the edges $(F_{1,B}, F_{R,B})$ and $(F_{1,A}, F_{R,A})$ are equal, as are the probabilities on $(F_{1,B}, P_{R,B})$ and $(F_{1,A}, P_{R,A})$. Clearly the joint probability distribution of the model – whose atoms are the root to leaf paths of the tree – is determined by the conditional probabilities associated with the stages. A CEG is the graph that is constructed to encode a model that can be specified through an event tree combined with a partitioning of its situations into stages.

In this paper we suppose that we are in a context similar to that of Example 1, where, for any possible model, with a selection of these types of hypotheses, the sample space of the problem must be consistent with a single event tree. On the basis of a sample of students’ records we want to select one of a number of these different possible CEG models, i.e. we want to find the “best” partitioning of the situations into stages. We take a Bayesian approach to this problem and choose the model with the highest posterior probability — the Maximum A Posteriori (MAP) model. This is the simplest and possibly most common Bayesian model selection method, advocated by, for example, Dennison et al [8], Castelo [9], and Heckerman [10], the latter two specifically for Bayesian network selection.

The paper is structured as follows. In the next section we review the definitions of event trees and CEGs. In Section 3 we develop the theory of how conjugate learning of CEGs is performed. In Section 4 we apply this theory by using the posterior probability of a CEG as its score in a model search algorithm that is derived using an analogous procedure to the model selection of BNs. We characterise the product Dirichlet distribution as a prior distribution for the CEGs’ parameters under particular homogeneity conditions. In Section 5 the algorithm is used to discover a good explanatory model for real students’ exam results. We finish with a discussion.

2. Definitions

In this section we define the CEG and any necessary prerequisites.

2.1. Basic graph theory concepts

There are many good sources of further information about these terms, including [11].

Definition 2. A GRAPH G is a pair $(V(G), E(G))$ where $V(G)$ is its set of vertices (or nodes), $E(G)$ is its set of edges. The set of edges can be thought of as a relation on $V(G)$.

When a graph is drawn, the vertices are displayed as points and the edges as curves between the appropriate points.

Definition 3. A DIRECTED GRAPH (or digraph) is a graph G where the edges are ordered pairs of vertices. Thus the edges $e_1 = (v_1, v_2)$ and $e_2 = (v_2, v_1)$ (where $v_1, v_2 \in V(G)$) are distinct elements of $E(G)$.

All edges in a directed graph are drawn as arrows from the first vertex to the second vertex in the ordered pair.

All graphs in this paper are directed graphs, and the following definitions assume this.

Definition 4. The CHILD of the edge $e = (v_1, v_2) \in E(G)$, written $ch(e)$, is v_2 . Its PARENT $pa(e)$ is v_1 .

By abuse of notation, the children of a vertex $v \in V(G)$, written $ch(v)$, is defined as

$$ch(v) = \{v' : v' \in V(G), (v, v') \in E(G)\} \quad (1)$$

and $pa(v)$ is defined similarly.

Definition 5. A PATH λ between two vertices $v_1, v_2 \in V(G)$ is an ordered sequence of edges $\lambda(v_1, v_2) = (e_1, \dots, e_n)$ where $e_1, \dots, e_n \in E(G)$, $pa(e_1) = v_1$, $ch(e_n) = v_2$ and $ch(e_k) = pa(e_{k+1})$ for $k = 1, \dots, n-1$.

The LENGTH of a path is the number of edges it contains.

By an abuse of notation, we say $v \in \lambda$ when $v = ch(e)$ for some $e \in \lambda$.

Definition 6. A CYCLE is a path $\lambda(v_1, v_2)$ where $v_1 = v_2$.

An ACYCLIC GRAPH contains no cycles.

Definition 7. A graph is CONNECTED if there exists a path in the graph between every pair of vertices, if direction of edges can be changed.

Definition 8. A TREE is a connected acyclic graph where one vertex (denoted v_0) has no parents and all other vertices have exactly one parent.

Definition 9. A LEAF NODE in a tree is a vertex with no children. The set of leaf nodes of a tree T is denoted $L(T)$.

2.2. Event Trees

We now define event trees, the graphical models that form a basis for CEGs. Further details about them can be found in [7].

Let $T = (V(T), E(T))$ be a directed tree where $V(T)$ is its node set and $E(T)$ its edge set.

Definition 10. The set of SITUATIONS of T , $S(T)$, is the set of non-leaf nodes $\{v : v \in V(T) \setminus L(T)\}$.

Let \mathbb{X} be the set of root-to-leaf paths of T , so that $\mathbb{X} = \{\lambda(v_0, v) : v \in L(T)\}$. \mathbb{X} represents the event space of the model, with every root-to-leaf path an atom of the event space.

Let $\mathbb{X}(v)$ denote the set of children of $v \in V(T)$. In an event tree, each situation $v \in S(T)$ has an associated random variable $X(v)$ with state space $\mathbb{X}(v)$, defined conditional on having reached v .

Definition 11. The distribution of $X(v)$ is determined by the PRIMITIVE PROBABILITIES $\{p(v'|v) = P(X(v) = v') : v' \in \mathbb{X}(v)\}$.

The probability of an event $\lambda \in \mathbb{X}$ can therefore be calculated by multiplying the primitive probabilities along the path.

Definition 12. The FLORET of $v \in S(T)$ is

$$\mathcal{F}(v) = (V(\mathcal{F}(v)), E(\mathcal{F}(v)))$$

where $V(\mathcal{F}(v)) = \{v\} \cup \{v' \in V(T) : (v, v') \in E(T)\}$ and $E(\mathcal{F}(v)) = \{e \in E(T) : e = (v, v')\}$.

The floret of a vertex v is thus a sub-tree consisting of v , its children, and the edges connecting v and its children, as shown in Figure 2. The floret represents the situation v , the associated random variable $X(v)$ and its sample space $\mathbb{X}(v)$.

Example 13. Figure 3 shows a tree for two Bernoulli random variables, A and B , with A occurring before B . In an education setting A could be the indicator variable of a student passing one module, and B the indicator variable for a subsequent module.

Here we have random variables $X(v_0) = A$, $X(v_1) = B|(A = 0)$ and $X(v_2) = B|(A = 1)$, and primitive probabilities $\pi(v_1|v_0) = p(A = 0)$, $\pi(v_3|v_1) = p(B = 0|A = 0)$ and so on for every other edge. Path probabilities can be found by multiplying primitive probabilities along a path, e.g. $p(A = 0, B = 0) = p(A = 0)p(B = 0|A = 0) = \pi(v_1|v_0)\pi(v_3|v_1)$ as v_0 and v_1 are on a path.

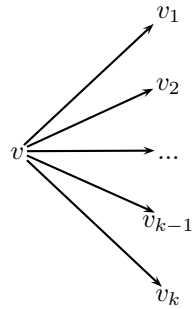


Figure 2: Floret of v . This subtree represents both the random variable $X(v)$ and its state space $\mathbb{X}(v)$.

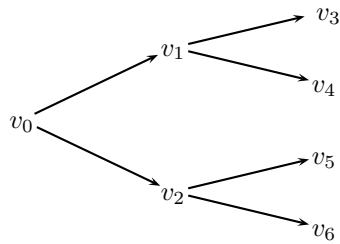


Figure 3: Simple event tree. The non-zero-probability events in the joint probability distribution of two Bernoulli random variables, A and B , with A observed before B , can be represented by this tree. Here, all four joint states are possible, because there are four root-to-leaf paths through the nodes.

2.3. Chain Event Graphs

Starting with an event tree T , we extend the definition with three new concepts to form the CEG: STAGES, EDGE COLOURS and POSITIONS, along the lines of [6] and [13].

One of the redundancies that can be eliminated from an ET is that of two situations, v and v' say, which have identical associated edge probabilities despite being defined by different conditioning paths. We say these two situations are at the same STAGE. This concept is formally defined below.

Definition 14. *Two situations $v, v' \in S(T)$ are in the same stage u if and only if $X(v)$ and $X(v')$ have the same distribution under a bijection*

$$\psi_u(v, v') : \mathbb{X}(v) \rightarrow \mathbb{X}(v') \quad (2)$$

The set of stages of an event tree T is written $J(T)$. This set partitions the set of situations $S(T)$, due to the associated set of bijections $\{\psi_u(v, v') : v, v' \in u, u \in J(T)\}$ forming an equivalence relation on $S(T)$.

Definition 15. *Any two edges $(v, v^*), (v', v'^*) \in E(T)$ have the same colour if and only if $v, v' \in u \in J(T)$ and $\psi_u(v, v')(v^*) = v'^*$.*

The edge colours make it clear, when drawn, which edges represent the same primitive probabilities and hence which situations are in the same stage.

We can construct a STAGED TREE $U(T)$ of an event tree T with $V(\mathcal{G}) = V(T)$, $E(\mathcal{G}) = E(T)$, and edge colours as given above.

Sometimes two situations have even more in common than the distribution over their respective variables: the entire subtrees with the two situations as roots share the same distribution over their paths. We define this concept formally.

Definition 16. *Two situations $v, v' \in S(T)$ are in the same POSITION w if and only if there exists a bijection*

$$\phi_w(v, v') : \Lambda(v, T) \rightarrow \Lambda(v', T)$$

where $\Lambda(v, T)$ is the set of paths in T from v to a leaf node of T , such that for every path $\lambda(v) \in \Lambda(v, T)$, the ordered sequence of colours in $\lambda(v)$ equals the ordered sequence of colours in $\lambda(v') := \phi_w(v, T)(\lambda(v)) \in \Lambda(v', T)$

We denote the set of positions as $K(T)$. It is clear that $J(T)$ is a partition of $K(T)$, as situations in the same position will always be in the same stage, and that therefore $K(T)$ is a finer partition of $S(T)$ than $J(T)$.

Now the CEG can finally be constructed by taking the staged tree $U(T)$ of an event tree and merging situations that are in the same position.

Definition 17. *The CHAIN EVENT GRAPH (CEG) $C(T)$ of an event tree T is the coloured directed graph with vertex set $V(C)$ and edge set $E(C)$ defined as follows:*

- $V(C) = K(T) \cup w_\infty$, so that each non-leaf node in the CEG represents one position and w_∞ represents the set of leaf nodes.
- Each edge in $E(C)$ exists for one of the following two reasons.
 - For $w, w' \in V(C) \setminus w_\infty$, there is an edge $(w, w') \in E(C)$ if and only if there exist situations $v, v' \in S(T)$ such that $v \in w$, $v' \in w'$ and $(v, v') \in E(T)$.

– For $w \in V(C) \setminus w_\infty$, there is an edge $(w, w_\infty) \in E(C)$ if and only if there exist situations $v \in S(T)$ and $v' \in L(T)$ such that $v \in w$ and $(v, v') \in E(T)$.

- The edge $(w, w') \in E(C)$ has the same colour as $(v, v') \in E(T)$ where $v \in w$, $v' \in w'$.

An example of a CEG that could be constructed from the event tree in Figure 1 is shown in Figure 5.

It is worth noting that for a finite number of discrete variables that the set of possible CEG models over those variables is a strict superset of the set of possible BN models. While a probability model that can be described by a BN will look different when described by a CEG it will still be the same model. The conditional independence statements described by a BN can always be represented by a CEG through stages and positions. This is shown in [6].

3. Conjugate learning of CEGs

One convenient property of CEGs is that conjugate updating of the model parameters proceeds in a closely analogous fashion to that on a BN. Conjugacy is a crucial part of the model selection algorithm that will be described in Section 4, because it leads to closed form expressions for the posterior probabilities of candidate CEGs. This in turn makes it possible to search the often very large model space quickly to find optimal models. The CEG model class will in general be bigger than the BN class for the same random variables, so that a model search will generally take longer but with the benefit that a richer model class is being considered. We demonstrate here how a conjugate analysis on a CEG proceeds.

Let a CEG C have set of stages $J(C) = \{u_1, \dots, u_k\}$, and let each stage u_i have k_i emanating edges (labelled e_1, \dots, e_{k_i}) with associated probability vector $\boldsymbol{\pi}_i = (\pi_{i1}, \pi_{i2}, \dots, \pi_{ik_i})'$ (where $\sum_{j=1}^{k_i} \pi_{ij} = 1$ and $\pi_{ij} > 0$ for $j \in \{1, \dots, k_i\}$). Then, under random sampling, the likelihood of the CEG can be decomposed into a product of the likelihood of each probability vector, i.e.

$$p(\mathbf{x}|\boldsymbol{\pi}, C) = \prod_{i=1}^k p_i(\mathbf{x}_i|\boldsymbol{\pi}_i, C)$$

where $\boldsymbol{\pi} = \{\boldsymbol{\pi}_1, \boldsymbol{\pi}_2, \dots, \boldsymbol{\pi}_k\}$, and $\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ is the complete sample data such that each $\mathbf{x}_i = (x_{i1}, \dots, x_{ik_i})'$ is the vector of the number of units in the sample (for example, the students in Example 1) that start in stage u_i and move to the stage at the end of edge e_{ij} for $j \in \{1, \dots, k_i\}$.

If it is further assumed that $\mathbf{x}_i \perp\!\!\!\perp \mathbf{x}_j | \boldsymbol{\pi}, \forall i \neq j$ then

$$p_i(\mathbf{x}_i|\boldsymbol{\pi}_i, C) = \prod_{j=1}^{k_i} \pi_{ij}^{x_{ij}} \quad (3)$$

Thus, just as for the analogous situation with BNs, the likelihood of a random sample also separates over the components of $\boldsymbol{\pi}$. With BNs, a common modelling assumption is of local and global independence of the probability parameters [12]; the corresponding assumption here is that the parameters $\boldsymbol{\pi}_1, \boldsymbol{\pi}_2, \dots, \boldsymbol{\pi}_k$ of $\boldsymbol{\pi}$ are all mutually independent a priori. It will then follow, with the separable likelihood, that they will also be independent a posteriori.

If the probabilities $\boldsymbol{\pi}_i$ are assigned a Dirichlet distribution, $\text{Dir}(\boldsymbol{\alpha}_i)$, a priori, where $\boldsymbol{\alpha}_i = (\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{ik_i})'$, so that for values of π_{ij} such that $\sum_{j=1}^{k_i} \pi_{ij} = 1$ and $\pi_{ij} > 0$ for $1 \leq j \leq k_i$, the

density of $\boldsymbol{\pi}_i$, $q_i(\boldsymbol{\pi}_i|C)$, can be written

$$q_i(\boldsymbol{\pi}_i|C) = \frac{\Gamma(\alpha_{i1} + \dots + \alpha_{ik_i})}{\Gamma(\alpha_{i1}) \dots \Gamma(\alpha_{ik_i})} \prod_{j=1}^{k_i} \pi_{ij}^{\alpha_{ij}-1}$$

where $\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$ is the Gamma function. It then follows that $\boldsymbol{\pi}_i|\mathbf{x}$ ($= \boldsymbol{\pi}_i|\mathbf{x}_i$) also has a Dirichlet distribution, $\text{Dir}(\boldsymbol{\alpha}_i^*)$, a posteriori, where $\boldsymbol{\alpha}_i^* = (\alpha_{i1}^*, \dots, \alpha_{ik_i}^*)'$, $\alpha_{ij}^* = \alpha_{ij} + x_{ij}$ for $1 \leq j \leq k_i, 1 \leq i \leq k$. The marginal likelihood of this model can be written down explicitly as the function of the prior and posterior Dirichlet parameters:

$$p(\mathbf{x}|C) = \prod_{i=1}^k \left[\frac{\Gamma(\sum_j \alpha_{ij})}{\Gamma(\sum_j \alpha_{ij}^*)} \prod_{j=1}^{k_i} \frac{\Gamma(\alpha_{ij}^*)}{\Gamma(\alpha_{ij})} \right].$$

The computationally more useful logarithm of the marginal likelihood is therefore a linear combination of functions of α_{ij} and α_{ij}^* . Explicitly,

$$\log p(\mathbf{x}|C) = \sum_{i=1}^k [s(\boldsymbol{\alpha}_i) - s(\boldsymbol{\alpha}_i^*)] + \sum_{i=1}^k [t(\boldsymbol{\alpha}_i^*) - t(\boldsymbol{\alpha}_i)] \quad (4)$$

where for any vector $\mathbf{c} = (c_1, c_2, \dots, c_n)'$,

$$s(\mathbf{c}) = \log \Gamma\left(\sum_{v=1}^n c_v\right) \text{ and } t(\mathbf{c}) = \sum_{v=1}^n \log \Gamma(c_v) \quad (5)$$

So the posterior probability of a CEG C after observing \mathbf{x} , $q(C|\mathbf{x})$, can be calculated using Bayes' Theorem, given a prior probability $q(C)$:

$$\log q(C|\mathbf{x}) = \log p(\mathbf{x}|C) + \log q(C) + K \quad (6)$$

for some value K which does not depend on C . This is the SCORE that will be used when searching over the candidate set of CEGs for the model that best describes the data.

4. A Local Search Algorithm for Chain Event Graphs

4.1. Preliminaries

With the log marginal posterior probability of a CEG model, $\log q(C|\mathbf{x})$, as its score, searching for the highest-scoring CEG in the set of all candidate models is equivalent to trying to find the Maximum A Posteriori (MAP) model [14]. The intuitive approach for searching \mathcal{C} , the candidate set of CEGs — calculating $q(C|\mathbf{x})$ (or $\log q(C|\mathbf{x})$) for every $C \in \mathcal{C}$ and choosing $C^* := \max_C q(C|\mathbf{x}) = \max_C \log q(C|\mathbf{x})$ — is infeasible for any but the most trivial problems. We describe in this section an algorithm for efficiently searching the model space by reformulating the model search problem as a clustering problem.

As mentioned in Section 2.3, every CEG that can be formed from a given event tree can be identified exactly with a partition of the event tree's nodes into stages. The coarsest partition C_∞ has all nodes with k outgoing edges in the same stage, u_k ; the finest partition C_0 has each situation

in its own stage, except for the trivial cases of those nodes with only one outgoing edge. Defined this way, the search for the highest-scoring CEG is equivalent to searching for the highest-scoring clustering of stages.

Various Bayesian clustering algorithms exist [15], including many involving MCMC [16]. We show here how to implement an Bayesian agglomerative hierarchical clustering (AHC) exact algorithm related to that of Heard et al [17]. The AHC algorithm here is a local search algorithm that begins with the finest partition of the nodes of the underlying ET model (called C_0 above and henceforth) and seeks at each step to find the two nodes that will yield the highest-scoring CEG if combined.

Some optional steps can be taken to simplify the search, which we will implement here. The first of these involves the calculation of the scores of the proposed models in the algorithm. By assuming that the probability distributions of stages that are formed from the same nodes of the underlying ET are equal in all CEGs, i.e. $p(\mathbf{x}_i | \boldsymbol{\pi}_i, C_1) = p(\mathbf{x}_i | \boldsymbol{\pi}_i, C_2)$ when $u_i \in J(C_1), J(C_2)$, it becomes more efficient to calculate the differences of model scores, i.e. the logarithms of the relevant Bayes factors, than to calculate the two individual model scores separately. This is because, if for two CEGs their stage sets $J(C_1)$ and $J(C_2)$ differ only in that stages $u_{1a}, u_{1b} \in C_1$ are combined into $u_{2c} \in C_2$, with all other stages unchanged, then the calculation of the logarithm of their posterior Bayes factor depends only on the stages involved; using the notation of Equation (5),

$$\log \frac{q(C_1|\mathbf{x})}{q(C_2|\mathbf{x})} = \log q(C_1|\mathbf{x}) - \log q(C_2|\mathbf{x}) \quad (7)$$

$$= \log q(C_1) - \log q(C_2) + \log q(\mathbf{x}|C_1) - \log q(\mathbf{x}|C_2) \quad (8)$$

$$= \log q(C_1) - \log q(C_2) + \sum_i [s(\boldsymbol{\alpha}_{1i}) - s(\boldsymbol{\alpha}_{1i}^*)] + \sum_i [t(\boldsymbol{\alpha}_{1i}^*) - t(\boldsymbol{\alpha}_{1i})] \\ - \sum_i [s(\boldsymbol{\alpha}_{2i}) - s(\boldsymbol{\alpha}_{2i}^*)] - \sum_i [t(\boldsymbol{\alpha}_{2i}^*) - t(\boldsymbol{\alpha}_{2i})] \quad (9)$$

$$= \log q(C_1) - \log q(C_2) + s(\boldsymbol{\alpha}_{1a}) - s(\boldsymbol{\alpha}_{1a}^*) + t(\boldsymbol{\alpha}_{1a}^*) - t(\boldsymbol{\alpha}_{1a}) \\ + s(\boldsymbol{\alpha}_{1b}) - s(\boldsymbol{\alpha}_{1b}^*) + t(\boldsymbol{\alpha}_{1b}^*) - t(\boldsymbol{\alpha}_{1b}) \\ - s(\boldsymbol{\alpha}_{2c}) + s(\boldsymbol{\alpha}_{2c}^*) - t(\boldsymbol{\alpha}_{2c}^*) + t(\boldsymbol{\alpha}_{2c}) \quad (10)$$

Using the trivial result that for any three CEGs

$$\log q(C_3|\mathbf{x}) - \log q(C_2|\mathbf{x}) = [\log q(C_3|\mathbf{x}) - \log q(C_1|\mathbf{x})] - [\log q(C_2|\mathbf{x}) - \log q(C_1|\mathbf{x})],$$

it can be seen that in the course of the AHC algorithm, comparing two proposal CEGs from the current CEG can be done equivalently by comparing their log Bayes factors with the current CEG, which as shown above requires fewer calculations.

The calculation of the score for each CEG C , as shown by Equation (6), shows that it is formed of two components: the prior probability of the CEG being the true model and the marginal likelihood of the data. These must therefore be set before the algorithm can be run, and it is here that the other simplifications are made.

4.2. The prior over the CEG space

For any practical problem \mathcal{C} , the set of all possible CEGs for a given ET, is likely to be a very large set, making setting a value for $q(C), \forall C \in \mathcal{C}$ a non-trivial task. An obvious way to set a non-informative or exploratory prior is to choose the uniform prior, so that $q(C) = \frac{1}{|\mathcal{C}|}$. This has

the advantages of being simple to set and of eliminating the $\log q(C_1) - \log q(C_2)$ term in Equation (10).

A more sophisticated approach is to consider which potential clusters are more or less likely a priori, according to structural or causal beliefs, and to exploit the modular nature of CEGs by stating that the prior log Bayes factor of a CEG relative to C_0 is the sum of the prior log Bayes factors of the individual clusters relative to their components completely unclustered, and that these priors are modular across CEGs. This approach makes it simple to elicit priors over \mathbf{C} from a lay expert, by requiring the elicitation only of the prior probability of each possible stage.

A particular computational benefit of this approach is when the prior Bayes factor of any CEG C with C_0 is believed to be zero, because one or more of its clusters is considered to be impossible. This is equivalent in the algorithm to not including the CEG in its search at all, as though it was never in \mathbf{C} in the first place, with the obvious simplification of the search following.

4.3. The prior over the parameter space

Just as when attempting to set $q(C)$, the size of most CEGs in practise leads to intractability of setting $p(\mathbf{x}|C)$ for each CEG C individually. However, the task is again made possible by exploiting the structure of a CEG with judicious modelling assumptions.

Assuming independence between the likelihoods of the stages for every CEG, so that $p(\mathbf{x}|\boldsymbol{\pi}, C)$ is as determined by Equation (3), and the fact that $p(\mathbf{x}|C) = \int p(\mathbf{x}|\boldsymbol{\pi}, C)p(\boldsymbol{\pi}|C)d\boldsymbol{\pi}$, it is clear that to set the marginal likelihood for each CEG is equivalent to setting the prior over the CEG's parameters, i.e. setting $p(\boldsymbol{\pi}|C)$ for each C . With the two further structural assumptions that the stage priors are independent for all CEGs (so that $p(\boldsymbol{\pi}|C) = \prod_{i=1}^k p(\boldsymbol{\pi}_i|C)$) and that equivalent stages in different CEGs have the same prior distributions on their probability vectors, (i.e. $p(\boldsymbol{\pi}_i|C_1) = p(\boldsymbol{\pi}_i|C_2)$), it can be seen that the problem of setting $p(\mathbf{x}|\boldsymbol{\pi}, C)$ is reduced to setting the parameter priors of each non-trivial floret in C_0 ($p(\boldsymbol{\pi}_i|C_0), i = 1, \dots, k$) and the parameter priors of stages that are clusters of stages of C_0 .

The usual prior put on the probability parameters of finite discrete BNs is the product Dirichlet distribution. In [18] the surprising result was shown that a product Dirichlet prior is inevitable if local and global independence are assumed to hold over all Markov equivalent graphs on at least two variables. In this paper we show that a similar characterisation can be made for CEGs given the assumptions in the previous paragraph. We will first show that the floret parameters in C_0 must have Dirichlet priors, and second that all CEGs formed by clustering the florets in C_0 have Dirichlet priors on the stage parameters. One characterisation of C_0 is given by Theorem 18.

Theorem 18. *If it is assumed a priori that the rates at which units take the root-to-leaf paths in C_0 are independent (“path independence”) and that the probability of which edge units take after arriving at a situation v is independent of the rate at which units arrive at v (“floret independence”), then the non-trivial florets of C_0 have independent Dirichlet priors on their probability vectors.*

Proof. The proof is in the Appendix. □

Thus $p(\boldsymbol{\pi}_i|C_0)$ is entirely determined by rates $\gamma(\lambda)$ on the root-to-leaf paths $\lambda \in \Lambda(v_0, C_0)$ of C_0 . This is similar to the “equivalent sample sizes” method of assessing prior uncertainty of Dirichlet hyperparameters in BNs as discussed in Section 2 of [10]. Lemma 28 shows inter alia that the parameters of the Dirichlet distribution of $p(\boldsymbol{\pi}_i|C_0)$ are determined for each edge by the sums of the rates of the root-to-leaf paths passing through that edge.

Another way to show that all non-trivial situations in C_0 have Dirichlet priors on their parameter spaces is to use the characterisation of the Dirichlet distribution first proven by [18], repeated here as Theorem 19.

Theorem 19. *Let $\{\theta_{ij}\}, 1 \leq i \leq k, 1 \leq j \leq n, \sum_{ij} \theta_{ij} = 1$, where k and n are integers greater than 1, be positive random variables having a strictly positive pdf $f_U(\{\theta_{ij}\})$. Define $\theta_{i.} = \sum_{j=1}^n \theta_{ij}$, $\theta_{.i} = \{\theta_{.i}\}_{i=1}^{k-1}$, $\theta_{j|i} = \theta_{ij}/\sum_j \theta_{ij}$, and $\theta_{J|i} = \{\theta_{j|i}\}_{j=1}^{n-1}$. Then if $\{\theta_{i.}, \theta_{J|1}, \dots, \theta_{J|k}\}$ are mutually independent, $f_U(\{\theta_{ij}\})$ is Dirichlet.*

Proof. Theorem 2 of Geiger and Heckerman [18]. \square

Corollary 20. *If C_0 has a composite number m of root-to-leaf paths and all Markov equivalent CEGs have independent floret distributions then the vector of probabilities on the root-to-leaf paths of C_0 must have a Dirichlet prior. This means in particular that, from the properties of the Dirichlet distribution, the floret of each situation with at least two outgoing edges has a Dirichlet prior on its edges.*

Proof. Construct an event tree C'_0 with m root-to-leaf paths, where the floret of the root node v'_0 has k edges and each of the florets extending from the children of v'_0 have n edges terminating in leaf nodes, where $m = kn, k \geq 2, n \geq 2$. This will always be possible with a composite m . C'_0 describes the same atomic events as C_0 with a different decomposition.

Let the random variable associated with the root floret of C'_0 be X , and let the random variable associated with each of the other florets be $Y|X = i, i = 1, \dots, k$. Let $\theta_{ij} = P(X = i, Y = j)$. Then by the definition of event trees, $P(\theta_{ij} > 0) > 0, 1 \leq i \leq k, 1 \leq j \leq n$ and $\sum \theta_{ij} = 1$. By the notation of Theorem 19, $\theta_{i.} = P(X = i)$ and $\theta_{j|i} = P(Y = j|X = i)$.

By hypothesis the floret distributions of C'_0 are independent. Therefore the condition of Theorem 19 holds and hence $f_U(\theta_{ij})$ is Dirichlet. From the equivalence of the atomic events, the probability distribution over the root-to-leaf path probabilities of C_0 is also Dirichlet, and so by Lemma 29, all non-trivial florets of C_0 therefore have Dirichlet priors on their probability vectors. \square

To show that the stage parameters of all the other CEGs in \mathcal{C} have independent Dirichlet priors, an inductive approach will be taken. Because of the assumption of consistency – that two identically composed stages in different CEGs have identical priors on their parameter space – for any given CEG C whose stages all have independent Dirichlet priors on their parameters spaces, it is known that another CEG C^* formed by clustering two stages u_{1c}, u_{2c} from C into one stage u_{c^*} will have independent Dirichlet priors on all its stages apart from u_{c^*} . It is thus only required to show that π_{c^*} has a Dirichlet prior. We this result for a class of CEGs called REGULAR CEGS.

Definition 21. *A stage u is REGULAR if and only if every path $\lambda \in \Lambda(v_0, C)$ contains either one situation in u or none of the situations in u .*

Definition 22. *A CEG is REGULAR if and only if every situation $u \in \mathbf{u}(C)$ is regular.*

Theorem 23. *Let C be a regular CEG, and let C^* be the CEG that is formed from C by setting two of its stages, u_{1c} and u_{2c} , as being in the same stage u_{c^*} , where u_{c^*} is a regular stage, with all other attributes of the CEG unchanged from C .*

If all stages in C have Dirichlet priors, then assuming that equivalent stages in different CEGs have equivalent priors, all stages in C^ have Dirichlet priors.*

Proof. Without loss of generality, let all situations in u_{1c} and u_{2c} have s children each, and let the total number of situations in u_{1c} and u_{2c} be r . Thus there are r situations in u_{c^*} , each with s children. By the assumption of prior consistency across stages, all stages in C^* have Dirichlet priors on their parameter spaces, so it is only required to prove that u_{c^*} has a Dirichlet prior.

Consider the CEG C' formed as follows: Let the root node of C' , v_0 , have 2 children, v_1 and v' . Let v' be a terminal node, and let v_1 have r children, $\{v_1(1), \dots, v_1(r)\}$, which are equivalent to the situations in u_{c^*} , including the property that they are in the same stage $u_{c'}$. Lastly, let the children of $\{v_1(1), \dots, v_1(r)\}$, $\{v_1(1, 1), \dots, v_1(1, s), \dots, v_1(r, 1), \dots, v_1(r, s)\}$, be leaf nodes in C' .

By construction, the prior for $u_{c'}$ is the same as that for u_{c^*} .

Now construct another CEG $C^{*'}$ from C' by reversing the order of the stages v_1 and $u_{c'}$. The new CEG has root node v_0 with the same distribution as $v_0 \in C'$. v_0 now has two children v' – the same as before – and v_2 , which has s children $\{v_2(1), \dots, v_2(s)\}$ in the same stage. Each node $v_2(i)$, $i = 1, \dots, s$ has r children $v_2(i, 1), \dots, v_2(i, r)$, all of which are leaf nodes.

The two CEGs $C^{*'}$ and C' are Markov equivalent, as it is clear that $P(v_1(i, j)) = P(v_2(j, i))$, $i = 1, \dots, r, j = 1, \dots, s$. The probabilities on the floret of v_2 are thus equal to the probabilities of the situations in the stage of $u_{c'}$, and hence u_{c^*} . Because v_2 is a stage with only one situation, Theorem 18 implies that it has a Dirichlet prior. Therefore u_{c^*} has a Dirichlet prior. \square

An alternative justification for assigning a Dirichlet prior to any stage that is formed by clustering situations with Dirichlet priors on their state spaces can be obtained, which does not depend on assuming Markov equivalency between CEGs derived from different event trees, by assuming a property analogous to that of “parameter modularity” for BNs [19]. This property states that the distribution over structures common to two CEGs should be identical.

Definition 24. Let u be a stage in a CEG C composed of the situations v_1, \dots, v_n from C_0 , each of which has m children v_{i1}, \dots, v_{im} , $i = 1, \dots, n$ such that v_{ij} are the same colour for all i for each j . Then u has the property of MARGIN EQUIVALENCY if

$$\pi_{u_j} = P(v_{1j} \text{ or } v_{2j} \text{ or } \dots \text{ or } v_{nj} | v_1 \text{ or } v_2 \text{ or } \dots \text{ or } v_n) \quad (11)$$

$$= \frac{\sum_{i=1}^n P(v_{ij})}{\sum_{i=1}^n P(v_i)} \quad (12)$$

is the same for both C and C_0 for $j = 1, \dots, m$.

Definition 25. C has margin equivalency if all of its stages have margin equivalency.

Theorem 26. Let u_c be a stage as defined in Definition 24 with $m \geq 2$. Then assuming independent priors between the situations for the associated finest-partition CEG C_0 of C , $\pi_{v_i} \sim \text{Dir}(\alpha_i)$ where $\alpha_i = (\alpha_{i1}, \dots, \alpha_{im})$ for each v_i , $i = 1, \dots, n$. Furthermore, for both C and C_0 , $\pi_u \sim \text{Dir}(\alpha_u)$, where $\alpha_u = (\sum_i \alpha_{i1}, \dots, \sum_i \alpha_{im})$.

Proof. From Theorem [5] or Corollary [7], every non-trivial floret in C_0 has a Dirichlet prior on its edges, which includes in this case the situations v_1, \dots, v_n .

Let $\gamma_{ij} = \gamma \pi_{ij}$ for $i = 1, \dots, n$, $j = 1, \dots, m$ for some $\gamma \in \mathbb{R}^+$. Then it is a well-known fact that $\gamma_{ij} \sim \text{Gamma}(\alpha_{ij}, \beta)$ for all $1 \leq i \leq n, 1 \leq j \leq m$ for some $\beta > 0$, and that $\perp_j \gamma_{ij}$. As $\perp_i \pi_{v_i}$, $\perp_j \gamma_{ij}$. Then by Lemma 28, letting $I[j]$ be the set of edges $\{e_{ij} = e(v_i, v_{ij}), i = 1, \dots, n\}$ for $j = 1, \dots, m$,

$$\pi_u \sim \text{Dir}\left(\sum_i \alpha_{i1}, \dots, \sum_i \alpha_{im}\right)$$

By margin equivalency, π_u must be set the same way for C . \square

Note that the posterior of π_u for a stage u that is composed of the C_0 situations v_1, \dots, v_n is thus $\pi_u | \mathbf{x} \sim \text{Dir}(\boldsymbol{\alpha}_u^*)$ where $\boldsymbol{\alpha}_u^* = \boldsymbol{\alpha}_u + \mathbf{x}_u = \sum_{i=1}^n \boldsymbol{\alpha}_{v_i} + \sum_{i=1}^n \mathbf{x}_{v_i}$. Equation (10), therefore, becomes

$$\begin{aligned} \log \frac{q(C_1 | \mathbf{x})}{q(C_2 | \mathbf{x})} &= \log q(C_1) - \log q(C_2) + s(\boldsymbol{\alpha}_{1a}) - s(\boldsymbol{\alpha}_{1a}^*) + t(\boldsymbol{\alpha}_{1a}^*) - t(\boldsymbol{\alpha}_{1a}) \\ &\quad + s(\boldsymbol{\alpha}_{1b}) - s(\boldsymbol{\alpha}_{1b}^*) + t(\boldsymbol{\alpha}_{1b}^*) - t(\boldsymbol{\alpha}_{1b}) - s(\boldsymbol{\alpha}_{1a} + \boldsymbol{\alpha}_{1b}) \\ &\quad + s(\boldsymbol{\alpha}_{1a}^* + \boldsymbol{\alpha}_{1b}^*) - t(\boldsymbol{\alpha}_{1a}^* + \boldsymbol{\alpha}_{1b}^*) + t(\boldsymbol{\alpha}_{1a} + \boldsymbol{\alpha}_{1b}) \quad (13) \end{aligned}$$

Setting priors on the paths rather than the florets also ensures that the probabilities of the atomic events remain under different tree representations of the event space.

The path priors would in the first instance be set based on expert knowledge of the system at hand, possibly using the “equivalent sample size” heuristic to aid elicitation. In problems where there is no strong prior information, as with the analogous Dirichlet model selection issues for Bayesian networks [20, 21], the performance of the selection procedure is rather sensitive to the prior value we put on each of the components of $\boldsymbol{\alpha}$.

Within the context of the types of problem discussed here it seems natural in the absence of information to the contrary to set all the components of this vector equal to each other a priori. This implies that for the model with no stages, C_0 , we a priori believe that all the atoms — i.e. all possible root to leaf paths — are equally probable, implying that were a model with no structure true then we have no prior information to expect one path to be more likely than another.

Even if we choose to set these all equal, the equivalent sample size parameter $\alpha \triangleq \mathbf{1}^T \boldsymbol{\alpha}$ — the sum of the rate parameters — has an important role in determining the performance of the selection procedure. One default is to let $\boldsymbol{\alpha}$ be a vector of 1s. This ensures both a uniform prior over all possible combinations of path probabilities and equal expected path probabilities. Another possibility is the Jeffreys prior [22, 23] which sets every α to be equal to $\frac{1}{2}$. This has the advantage of being invariant under reparameterization.

4.4. The algorithm

The algorithm thus proceeds as follows:

1. Starting with the initial ET model, form the CEG C_0 with the finest possible partition, where all leaf nodes are placed in the terminal stage u_∞ and all nodes with only one emanating edge are placed in the same stage. Calculate $\log q(C_0 | \mathbf{x})$ using (6).
2. For each pair of situations $v_i, v_j \in C_0$ with the same number of edges, calculate $\log \frac{q(C_1^* | \mathbf{x})}{q(C_0 | \mathbf{x})}$ where C_1^* is the CEG formed by having v_i, v_j in the same stage and keeping all others in their own stage; do not calculate if $q(C_1^*) = 0$.
3. Let $C_1 = \max_{C_1^*} (\log \frac{q(C_1^* | \mathbf{x})}{q(C_0 | \mathbf{x})})$.
4. Now calculate C_2^* for each pair of stages in C_1 except where $q(C_2^*) = 0$, and record $C_2 = \max(q(C_2^* | \mathbf{x}))$.
5. Continue for C_3, C_4 and so on until the coarsest partition C_∞ has been reached.
6. Find $C = \max(C_0, C_1, \dots, C_\infty)$, and select this as the MAP model.

We note that the algorithm can also be run backwards, starting from C_∞ and splitting one cluster in two at each step. This has the advantage of making the identification of positions in the MAP model easier.

5. Examples

5.1. Simulated data

To first demonstrate the efficacy of the algorithm described above we implement the algorithm using simulated data for Example 1, where the CEG generating the data was as known and described in Section 1. Figure 4 shows the number of students in the sample who reached each situation in the tree.

In this complete dataset the progress of 1000 students has been tracked through the event tree. Half are assigned to take module A first and the other half B . By finding the MAP CEG model in the light of this data we may find out whether the three hypotheses posed in the introduction are valid. We repeat them here for convenience:

1. The chances of doing well in the second component are the same whether the student passed first time or after a resit.
2. The components A and B are equally hard.
3. The distribution of marks for the second component is unaffected by whether students passed or got a distinction for the first component.

We set a uniform prior on the CEG priors and a constant rate of 1 on the root-to-leaf paths of C_0 for illustration purposes. The algorithm is then implemented as follows.

There are only two florets with two edges; with Beta(1,3) priors on each and a Beta(2,6) prior on the combined stage, resulting from the path priors, the log Bayes factor is -1.85. Carrying out similar calculations for all the pairs of nodes with three edges, it is first decided to merge the nodes $P_{1,A}$ and $P_{1,B}$, which has a log Bayes factor of -3.76 against leaving them apart. Applying the algorithm to the updated set of nodes and iterating, the CEG in Figure 5 is found to be the MAP one. All edges in the CEG have different colours.

Under this model, it can be seen that all three hypotheses above are satisfied.

5.2. Student test data

In our second example we apply the learning algorithm to a real dataset in order to test the algorithm's efficacy in a real-life situation and to identify remaining issues with its usage. The dataset we used was an appropriately disguised set of marks taken over a 10-year period from four core modules of the MORSE degree course taught at the University of Warwick. A part of the event tree used as the underlying model for the first two modules is shown in Figure 6, along with a few illustrative data points. This is a simplification of a much larger study that we are currently investigating but large enough to illustrate the richness of inference possible with our model search.

For simplicity, the prior distributions on the candidate models and on the root-to-leaf paths for C_0 were both chosen to be uniform distributions, in the latter case by again assuming $\alpha = 1$ for each path.

The MAP CEG model was not C_0 , so that there were some non-trivial stages. In total, 170 situations were clustered into 32 stages. Some of the more interesting stages of this model are described in Table 1.

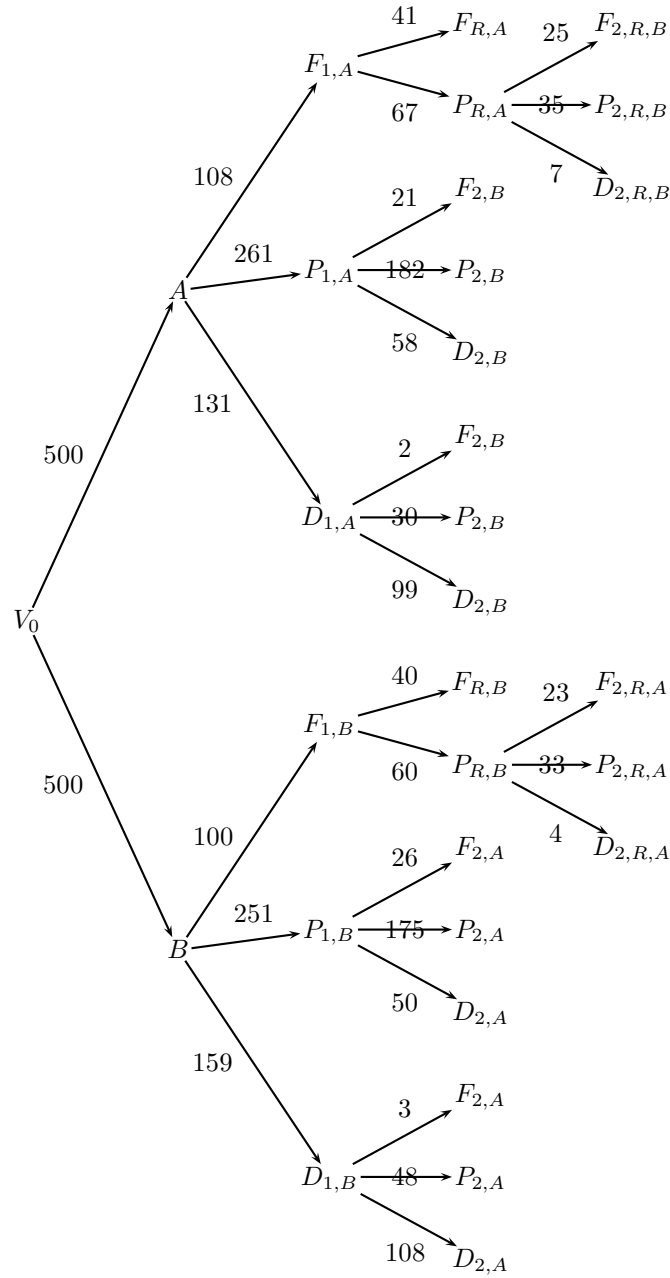


Figure 4: The event tree from Example 1 with the numbers representing the number of students in a simulated sample who reached each situation.

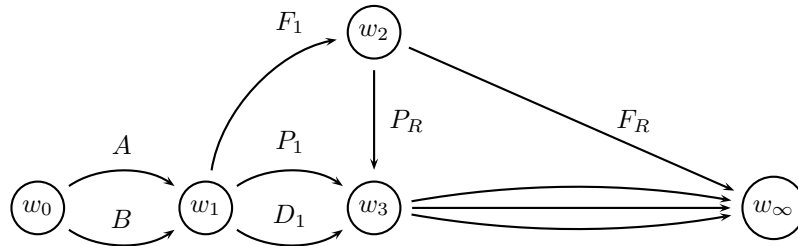


Figure 5: The MAP CEG for that event tree in Figure 4

Stage	Probability vector	Students	Situations	Locations	Comments
7	(0.47, 0.44, 0.08)	685	2	1; 1,1,1	High achievers
11	(0.22, 0.43, 0.35)	412	6	3; 1,2; 3,1; 1,1,3	Middling students
13	(0.33, 0.33, 0.33)	16	18	4; 4,2; 4,3	No students appeared in 17 of these situations
17	(0.07, 0.27, 0.66)	86	4	1,3; 3,2; 3,2,4	Struggling students
27	(0.19, 0.56, 0.25)	464	7	1,1,4; 1,2,2; 1,3,2; 1,4,2	More likely to get grade 2 than stage 11
28	(0.11, 0.51, 0.38)	436	6	1,2,3; 3,1,3; 1,2,4	More likely to get grade 3 than stage 27

Table 1: Selected stages of MAP CEG model formed from data described in Section 5.2. The columns respectively detail the stage number, posterior expectation of the probability vector of that stage (rounded to two decimal places), number of students passing through that stage in the dataset, number of situations from the original ET in that stage, examples of situations in that stage (shown as sequence of grades, where “4” means that grade is missing), and any comments or observations related to that stage.

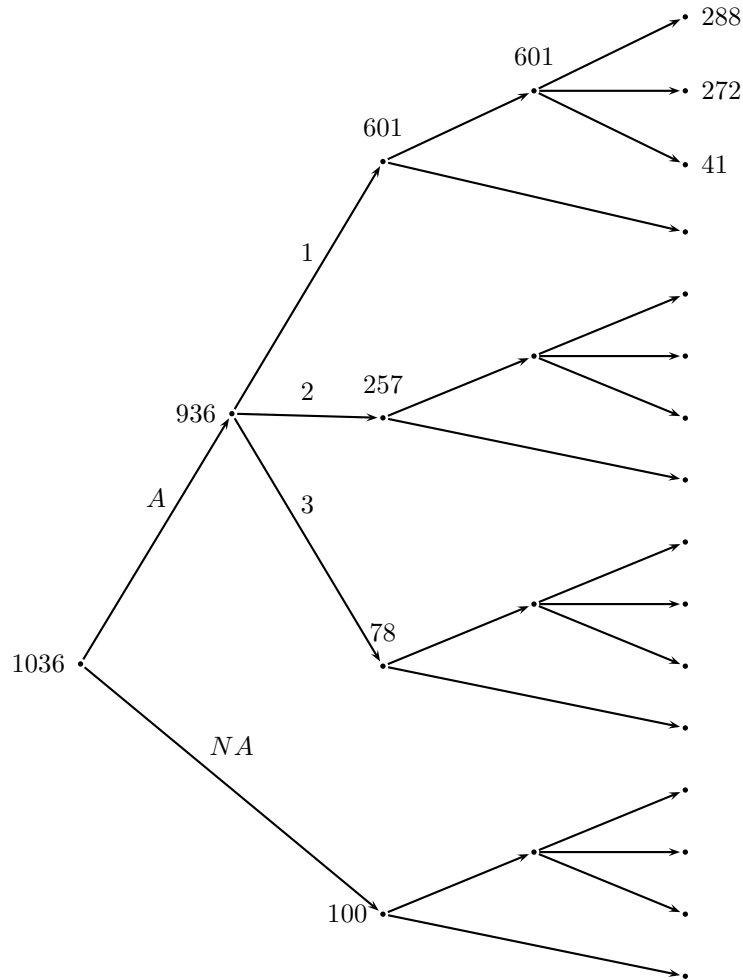


Figure 6: Sub-tree of the event tree of possible grades for the MORSE degree course at the University of Warwick. Each floret of two edges describes whether a student's marks are available for a particular module (denoted by the edge labelled *A* for the first module) or whether they are missing (*NA*). If they are available, then they are counted as grade 1 if are 70% or higher, grade 2 if they are between 50% and 69% inclusive, and grade 3 if they are below 50%. Some illustrative count data are shown on corresponding nodes.

From inspecting the membership of stages it was possible to identify various situations which were discovered to share distributions. From example, students who reach one of the two situations in stage 7 have an expected probability of 0.47 in getting a high mark, an expected probability of 0.44 of getting a middling grade, and only an expected probability of 0.08 of achieving the lowest grade. From being in a stage of their own, it can be deduced that students in these situations have qualitatively different prospects from students in any other situations. In contrast, students who reach one of the four situations in stage 17 have an expected probability of 0.66 of getting the lowest grade.

6. Discussion

In this paper we have shown that chain event graphs are not just an efficient way of storing the information contained in an event tree, but also a natural way to represent the information that is most easily elicited from a domain expert: the order in which events happen, the distributions of variables conditional on the process up to the point they are reached, and prior beliefs about the relative homogeneity of different situations. This strength is exploited when the MAP CEG is discovered, as this can be used in a qualitative fashion to detect homogeneity between seemingly disparate situations.

There are a number extensions to the theory in this paper that are currently being pursued. These fall mostly into the two categories: creating even richer model classes than those considered here; and developing even more efficient algorithms for selecting the MAP model in these model classes.

The first category includes dynamic chain event graphs. This framework can supply a number of different model classes. The simplest case involves selecting a CEG structure that is constant across time, but with a time series on its parameters. A bigger class would allow the MAP CEG structure to change over time. These larger model classes would clearly be useful in the educational setting considered in this paper, as they would allow for background changes in the students' abilities, for example.

Another important model class is that which arises from uncertainty about the underlying event tree. A similar model search algorithm to the one described in this paper is possible in this case after setting a prior distribution on the candidate event trees.

One difficulty with model selection over CEGs is simply the expressiveness and hence relative largeness of the model space, which means that to be feasible for large problems we need to add more contextual information to limit the size of the space. This is particularly the case if the underlying tree is allowed to embody different orders for when situations happen. One method we are investigating is to use the usual BN MAP search as a coarse initialisation step and then, taking a tree consistent with its corresponding CEG, refine the search using methods described in this paper. In other contexts, to allow all combinations of florets into stages would be implausible. When this is the case, the search algorithm can accommodate this information easily and therefore be carried out faster. We will report on these methods in due course.

Finally, we note that another way to search any of these model classes is to reformulate the search as a weighted MAX-SAT problem, for which very fast algorithms have been developed. This approach was used to great effect for finding a MAP BN by Cussens [24], and we plan to report on this approach in a later paper.

Appendix

Theorem 18 is based on three well-known results concerning properties of the Dirichlet distribution, which we review below.

Lemma 27. *Let $\gamma_j \sim \text{Gamma}(\alpha_j, \beta)$, $j = 1, \dots, n$ where $\alpha_j > 0$ for $j \in \{1, \dots, n\}$, $\beta > 0$ and $\gamma = \bigsqcup_{i \in \{1, \dots, n\}} \gamma_i$. Furthermore, let $\theta_j = \frac{\gamma_j}{\gamma}$ for $j \in \{1, \dots, n\}$, where $\gamma = \sum_{i=1}^n \gamma_i$.*

Then $\theta = (\theta_i)_{i \in \{1, \dots, n\}} \sim \text{Dir}(\alpha_1, \dots, \alpha_n)$.

Proof. Kotz et al [25]. □

Lemma 28. *Let $I[j] \subseteq \{1, \dots, n\}$, $\gamma(I[j]) = \sum_{i \in I[j]} \gamma_i$ and $\theta(I[j]) = \sum_{i \in I[j]} \theta_i$. Then for any partition $I = \{I[1], \dots, I[k]\}$ of $\{1, \dots, n\}$,*

$$\theta(I) = (\theta(I[1]), \theta(I[2]), \dots, \theta(I[k])) \sim \text{Dir}(\alpha(I[1]), \dots, \alpha(I[k]))$$

where $\alpha(I[j]) = \sum_{i \in I[j]} \alpha_i$.

Proof. For any $I[j] \subseteq \{1, \dots, n\}$, $\bigsqcup_{i \in I[j]} \gamma_i$, $\gamma(I[j]) \sim \text{Gamma}(\alpha(I[j]), \beta)$ (a well-known result; see, for example, Weatherburn [26]), and for any partition $I = \{I[1], \dots, I[k]\}$ of $\{1, \dots, n\}$, $\bigsqcup_{i \in \{1, \dots, k\}} \gamma(I[j])$.

Therefore, as

$$\theta(I[j]) = \sum_{i \in I[j]} \theta_i = \sum_{i \in I[j]} \frac{\gamma_i}{\gamma} = \frac{\gamma(I[j])}{\gamma}, \quad j = 1, \dots, k$$

and $\gamma = \sum_{i=1}^k \gamma(I[i])$, the result follows from Lemma 27. □

Lemma 29. *For any $I[j] \subseteq \{1, \dots, n\}$ where $|I[j]| \geq 2$,*

$$\theta_{I[j]} = \left(\frac{\theta_i}{\theta(I[j])} \right)_{i \in I[j]} \sim \text{Dir}((\alpha_i)_{i \in I[j]})$$

Proof. Wilks [27]. □

Theorem 30. *Let the rates of units along the root-to-leaf paths $\lambda_i \in \mathbb{X}$, $i \in \{1, \dots, |\mathbb{X}|\}$ of an event tree T have independent Gamma distributions with the same scale parameter, i.e. $\gamma_i = \gamma(\lambda_i) \sim \text{Gamma}(\alpha_i, \beta)$, $i \in \{1, \dots, |\mathbb{X}|\}$ and $\bigsqcup_{i \in \{1, \dots, |\mathbb{X}|\}} \gamma_i$. Then the distribution on each floret in the tree will be Dirichlet.*

Proof. Consider a floret \mathcal{F} with root node v and edge set $\{e_1, \dots, e_l\}$. The rate for each edge e_i , $\gamma(e_i)$, is equal to $\sum_{\lambda_j \in \lambda(e_i)} \gamma(\lambda_j)$, where $\lambda(e_i)$ is the set of root-to-leaf paths that contain e_i , so that $\gamma(e_i) \sim \text{Gamma}(\alpha(e_i), \beta)$ when $\bigsqcup_{i \in \{1, \dots, l\}} \gamma(e_i)$.

Let $I = \{I[\mathcal{F}], I[\overline{\mathcal{F}}]\}$ partition \mathbb{X} , where $I[\mathcal{F}] = \{\lambda_{e_1}, \dots, \lambda_{e_l}\}$ and $I[\overline{\mathcal{F}}] = I \setminus I[\mathcal{F}]$. Then by Lemma 29, the probability vector on \mathcal{F} is Dirichlet, where

$$\theta_{I[\mathcal{F}]} \sim \text{Dir}((\alpha_{e_i})_{i \in \{1, \dots, l\}})$$

□

- [1] R. G. Cowell, A. P. Dawid, S. L. Lauritzen, D. J. Spiegelhalter, Probabilistic Networks and Expert Systems, Springer, 1999.
- [2] S. L. Lauritzen, Graphical Models (Oxford Statistical Science Series), Oxford University Press, USA, 1996.
- [3] J. E. Smith, S. Holtzman, J. E. Matheson, Structuring conditional relationships in influence diagrams, *Operations Research* 41 (2) (1993) 280–297
- [4] C. Bouilrier, N. Friedman, M. Goldszmidt, D. Koller, Context-Specific independence in bayesian networks, in: E. Horvitz, F. V. Jensen (Eds.), *Proceedings of the Twelfth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann, Reed College, Portland, Oregon, USA, 1996, pp. 115–123.
- [5] D. Poole, N. L. Zhang, Exploiting contextual independence in probabilistic inference, *J. Artificial Intelligence Res.* 18 (2003) 263–313.
- [6] J. Q. Smith, P. E. Anderson, Conditional independence and chain event graphs, *Artificial Intelligence* 172 (1) (2008) 42–68.
- [7] G. Shafer, *The Art of Causal Conjecture*, Artificial Intelligence, The MIT Press, 1996.
- [8] D. G. T. Denison, C. C. Holmes, B. K. Mallick, A. F. M. Smith, *Bayesian Methods for Nonlinear Classification and Regression*, Wiley Series in Probability and Statistics, Wiley, 2002.
- [9] R. Castelo, *The discrete acyclic digraph markov model in data mining*, Ph.D. thesis, Faculteit Wiskunde en Informatica, Universiteit Utrecht (Apr. 2002).
- [10] D. Heckerman, A tutorial on learning with bayesian networks, in: M. I. Jordan (Ed.), *Learning in Graphical Models*, MIT Press, 1999, pp. 301–354.
- [11] D. B. West, *Introduction to Graph Theory*, Pearson Education Asia Limited and China Machine Press, China, 2001.
- [12] D. J. Spiegelhalter, S. L. Lauritzen, Sequential updating of conditional probabilities on directed graphical structures, *Networks* 20 (5) (1990) 579–605.
- [13] P. Thwaites, J. Q. Smith, E. Riccomagno, *Artificial Intelligence* 174 (12-13) (2010) 889–909.
- [14] J. Bernardo, A. F. M. Smith, *Bayesian Theory*, Wiley, Chichester, England, 1994.
- [15] J. W. Lau, P. J. Green, Bayesian Model-Based clustering procedures, *Journal of Computational and Graphical Statistics* 16 (3) (2007) 526–558.
- [16] S. Richardson, P. J. Green, On bayesian analysis of mixtures with an unknown number of components, *Journal of the Royal Statistical Society. Series B (Methodological)* 59 (4) (1997) 731–792.
- [17] N. A. Heard, C. C. Holmes, D. A. Stephens, A quantitative study of gene regulation involved in the immune response of anopheline mosquitoes: An application of bayesian hierarchical clustering of curves, *Journal of the American Statistical Association* 101 (473) (2006) 18–29.

- [18] D. Geiger, D. Heckerman, A characterization of the dirichlet distribution through global and local parameter independence, *The Annals of Statistics* 25 (3) (1997) 1344–1369.
- [19] D. Heckerman, M. P. Wellman, Bayesian networks, *Communications of the ACM* 38 (3) (1995) 27–30.
- [20] H. Steck, T. Jaakkola, On the dirichlet prior and bayesian regularization, in: S. Becker, S. Thrun, K. Obermayer (Eds.), *Advances in Neural Information Processing Systems 15*, MIT Press, Vancouver, British Columbia, Canada, 2003, pp. 697–704.
- [21] T. Silander, P. Kontkanen, P. Myllymäki, On sensitivity of the MAP bayesian network structure to the equivalent sample size parameter, in: R. Parr, L. van der Gaag (Eds.), *Proceedings of the The 23rd Conference on Uncertainty in Artificial Intelligence*, AUAI Press, 2007, pp. 360–367.
- [22] H. Jeffreys, An invariant form for the prior probability in estimation problems, *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences* 186 (1007) (1946) 453–461
- [23] J. M. Bernardo, Reference posterior distributions for bayesian inference, *Journal of the Royal Statistical Society. Series B (Methodological)* 41 (2) (1979) 113–147
- [24] J. Cussens, Bayesian network learning by compiling to weighted MAX-SAT, in: D. A. McAllester, P. Myllymäki (Eds.), *Proceedings of the 24th Conference in Uncertainty in Artificial Intelligence*, AUAI Press, Helsinki, Finland, 2008, pp. 105–112.
- [25] S. Kotz, N. Balakrishnan, N. L. Johnson, *Continuous Multivariate Distributions*, 2nd Edition, Wiley series in probability and statistics. Applied probability and statistics, Wiley, New York, 2000.
- [26] C. E. Weatherburn, *A first course in mathematical statistics*, 2nd Edition, CUP, 1949.
- [27] S. S. Wilks, *Mathematical Statistics*, Wiley, New York, 1962.