

Propagation of Probabilities, Means and Variances in Mixed Graphical Association Models

by

Steffen L. Lauritzen

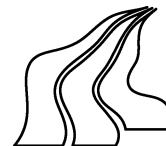
December 1992

To be referenced as:

Lauritzen, S. L. (1992). Propagation of probabilities, means and variances in mixed graphical association models. *Journal of the American Statistical Association*, 87:1098-1108.

INSTITUTE FOR ELECTRONIC SYSTEMS
DEPARTMENT OF MATHEMATICS AND COMPUTER
SCIENCE

Fredrik Bajers Vej 7E — DK 9220 Aalborg Ø — Denmark
Tel.: +45 98 15 85 22 — TELEX 69 790 aub dk



Propagation of Probabilities, Means and Variances in Mixed Graphical Association Models

Steffen L. Lauritzen*

December 1992

Abstract

A scheme is presented for modelling and local computation of exact probabilities, means and variances for mixed qualitative and quantitative variables. The models assume that the conditional distribution of the quantitative variables, given the qualitative, is multivariate Gaussian.

The computational architecture is set up by forming a tree of belief universes, and the calculations are then performed by local message passing between universes. The asymmetry between the quantitative and qualitative variables sets some additional limitations for the specification and propagation structure. Approximate methods when these are not appropriately fulfilled are sketched.

Lauritzen and Spiegelhalter (1988) showed how to exploit the local structure in the specification of a discrete probability model for fast and efficient computation, thereby paving the way for exploiting probability based models as parts of realistic systems for planning and decision support. The technique was subsequently improved by Jensen, Lauritzen and Olesen (1990). The purpose of the paper is to extend this computational scheme to networks, where some vertices represent entities that are measured on a quantitative and some on a qualitative scale.

An extension has the advantage of unifying several known techniques, but allows more flexible and faithful modelling and speeds up computation as well. To handle this more general case, the properties of CG distributions introduced in Lauritzen and Wermuth (1989) are exploited.

A fictitious but simple example is used for illustration throughout the paper, concerned with the monitoring of emission from a waste incinerator: From optical measurements of the darkness of the smoke, the concentration of CO_2 – which are both on a continuous scale – and possible knowledge about qualitative characteristics such as the type of waste burned, one wants to infer about the state of the incinerator and the current emission of heavy metals.

Key words: Bayesian methods, expert system, causal network, CG distribution, recursive model, strongly triangulated graph.

*Steffen L. Lauritzen is Professor of Mathematics and Statistics at the Department of Mathematics and Computer Science, University of Aalborg, Fredrik Bajers Vej 7, DK-9220 Aalborg Ø, Denmark. The author is indebted to Kristian G. Olesen for computational assistance and helpful technical discussions on the way. He is also grateful to Philip Dawid, David Spiegelhalter, a referee, and an associate editor for helpful comments concerning presentation of the results. This research has been supported in part by a grant from the SCIENCE programme of the EEC.

1. INTRODUCTION

Recent developments have shown that graphical models provide a flexible framework for specification and computation in probabilistic expert systems. We abstain from a detailed survey of the literature in the area, but refer to the bibliographies in Lauritzen and Spiegelhalter (1988), Pearl (1988) and Jensen, Lauritzen and Olesen (1990) as well as the volumes Oliver and Smith (1990) and Shafer and Pearl (1990).

For illustrative purposes we discuss a fictitious example throughout the paper. This is taken from a problem connected with control of the emission of heavy metals from a waste incinerator:

The emission from a waste incinerator differs because of compositional differences in incoming waste. Another important factor is the waste burning regime which can be monitored by measuring the concentration of CO_2 in the emission. The filter efficiency depends on the technical state of the electrofilter and the amount and composition of waste. The emission of heavy metal depends both on the concentration of metal in the incoming waste and the emission of dust particulates in general. The emission of dust is monitored through measuring the penetrability of light.

Here we have ignored the obvious time aspect of the monitoring problem and concentrated on a single point in time, for the sake of simplicity. The essence of the description is represented in the network of Figure 1.

Insert Figure 1 about here.

The described network could in principle be used for several purposes. Typically the emission of dust and heavy metal, the filter efficiency, as well as the actual concentration of heavy metal in the waste would normally not be directly available. The filter state might or might not be known as is also the case for the type of waste.

From the measurements and knowledge available at any time, the emission of heavy metal can be predicted, in particular the mean and standard deviation of the predictive distribution for that emission is of interest. Diagnostic probabilities for stability of the burning regime and/or the state of the filter could be required. Finally the network can be used for management purposes, in that the influence of filter efficiency, burning regime etc. on objective variables, such as the emission of heavy metals, can be computed.

The distributional theory of graphical models with both quantitative and qualitative variables are fundamental to the computational methods. A brief account of the basic elements of this is contained in Section 2.

Another formal element of the computational structure is the decompositional theory of marked graphs which are graphs with two types of vertices, here

corresponding to the discrete and continuous variables. The necessary concepts are explained in Section 3.

The remaining sections describe in some detail the model specification and the elements of the computational architecture. The example above is used to illustrate the various phases of the process.

We conclude with some remarks on computational complexity and approximative modifications when conditions for the exact results are not satisfied.

2. CG DISTRIBUTIONS AND POTENTIALS

The models behind the computations described in the present paper are based on the assumption that the conditional distribution of the continuous variables given the discrete is multivariate Gaussian. We shall briefly review some standard notation but refer the reader to Lauritzen and Wermuth (1989) or Whittaker (1990) for further details and derivations of formulae.

The set of variables V is partitioned as $V = \Delta \cup \Gamma$ into variables of *discrete* (Δ) and *continuous* (Γ) type. A typical element of the joint state space is denoted as in one of the possibilities below

$$x = (x_\alpha)_{\alpha \in V} = (i, y) = ((i_\delta)_{\delta \in \Delta}, (y_\gamma)_{\gamma \in \Gamma}),$$

where i_δ are qualitative and y_γ are real valued. A particular combination $i = (i_\delta)_{\delta \in \Delta}$ is referred to as a *cell* and the set of cells is denoted by \mathcal{I} . The joint distribution of the variables is supposed to have a density f with

$$f(x) = f(i, y) = \chi(i) \exp \left\{ g(i) + h(i)^\top y - y^\top K(i) y / 2 \right\},$$

where $\chi(i) \in \{0, 1\}$ indicates whether f is positive at i and A^\top is the transpose of the matrix A . We then say that X follows a *CG distribution* which is equivalent to the statement

$$\mathcal{L}(X_\Gamma \mid X_\Delta = i) = \mathcal{N}_{|\Gamma|}(\xi(i), \Sigma(i)) \quad \text{whenever} \quad p(i) = P\{X_\Delta = i\} > 0$$

where $X_A = (X_\alpha)_{\alpha \in A}$ and so on, and

$$\xi(i) = K(i)^{-1} h(i), \quad \Sigma(i) = K(i)^{-1}, \quad (1)$$

the latter being positive definite. The triple (g, h, K) – only defined for $\chi(i) > 0$ – constitutes the *canonical characteristics* of the distribution and $\{p, \xi, \Sigma\}$ are the *moment characteristics*.

Note that there is a slight difference between the notation used here and in Lauritzen and Wermuth (1989) in that we allow $p(i)$ to be equal to 0 for some entries i . Also, strictly speaking, χ belongs to the characteristics of the distribution, but we assume this to be implicitly represented through the domain where the other components are well-defined.

In the case where we have only one kind of variable the undefined components are denoted by zeros, i.e. $(g, 0, 0)$ or $(0, h, K)$.

A basic part of computational task consists of updating the joint distribution in the light of evidence, corresponding to a conditioning process. A simple way of doing this is by computing with unnormalised density functions. It is also an important part of the computational process to recognize and exploit a product structure in the joint density, with the factors not necessarily being densities themselves.

For this reason we extend the notion of a CG distribution to that of a *CG potential* which is any function ϕ of the form

$$\phi(x) = \phi(i, y) = \chi(i) \exp\{g(i) + h(i)^\top y - y^\top K(i)y/2\},$$

where $K(i)$ is only assumed to be a symmetric matrix. Thus ϕ is not necessarily a density. We still use the triple (g, h, K) as canonical characteristics for the potential ϕ .

A basic difference is that the moment characteristics for a CG potential are only well defined when K is positive definite for all i with $\chi(i) > 0$. Then Σ and ξ are given as in (1), whereas

$$p(i) \propto \{\det \Sigma(i)\}^{\frac{1}{2}} \exp\{g(i) + h(i)^\top \Sigma(i)h(i)/2\},$$

where \propto means ‘proportional to’. Conversely, if the moment characteristics $\{p, \xi, \Sigma\}$ are given, we can calculate the canonical characteristics as $K(i) = \Sigma(i)^{-1}$, $h(i) = K(i)\xi(i)$, and

$$g(i) = \log p(i) + \left\{ \log \det K(i) - |\Gamma| \log(2\pi) - \xi(i)^\top K(i)\xi(i) \right\} / 2.$$

3. MARKED GRAPHS AND JUNCTION TREES

In the present section we give a brief exposition of the graphtheoretic notions used in the paper. Many of the graphtheoretic terms have suggestive names that are really self-evident and the reader might want to skip this section at the first reading of the paper. When a more accurate understanding of the graphtheoretic details is needed, the reader can return.

3.1 Notation and Terminology

First we need to establish the terminology, in particular to ensure accurate understanding of the details in future developments. A section of this type must necessarily be somewhat terse, so we ask the reader to be patient.

In the paper a *network* or *graph* is formally a pair $\mathcal{G} = (V, E)$, where V is a finite set of *vertices* and the set of *edges* E is a subset of the set $V \times V$ of ordered pairs of distinct vertices. Thus our graphs have no multiple edges and no loops. Edges $(\alpha, \beta) \in E$ with both (α, β) and (β, α) in E are called *undirected*, whereas an edge (α, β) without its *opposite* (β, α) being contained in E is called *directed*.

In particular we need to work with graphs where the vertices are *marked* in the sense that they are partitioned into two groups, Δ and Γ . We use the term *marked graph* for a graph of this type.

The vertices in the set Δ are to represent qualitative variables and those in Γ quantitative variables. Therefore we say that the vertices in Δ are *discrete* and those in Γ are *continuous*.

A marked graph is conveniently represented by a picture, where we use a *dot* for a *discrete* vertex and a *circle* for a *continuous*. Further a *line* joining α to β represents an undirected edge, whereas an *arrow* from α , pointing towards β is used for a directed edge (α, β) .

If the graph has only undirected edges (lines) it is an *undirected* graph and if all edges are directed (arrows), the graph is said to be *directed*.

A subset C is *complete* if all vertices in C are joined by an arrow or a line. A complete subset that is maximal in the sense that no other vertex in the graph is connected to all its elements is called a *clique*. The cliques of the graph can be considered the fundamental blocks of the structure that the graph describes.

If there is an arrow from α pointing towards β , α is said to be a *parent* of β and β a *child* of α . The set of parents of β is denoted as $\text{pa}(\beta)$. If there is a line or an arrow between α and β , α and β are said to be *neighbours*. A selection of graphtheoretic concepts are illustrated in Figure 2.

Insert Figure 2 about here.

A *path* of length n from α to β is a sequence $\alpha = \alpha_0, \dots, \alpha_n = \beta$ of distinct vertices such that $(\alpha_{i-1}, \alpha_i) \in E$ for all $i = 1, \dots, n$. An undirected graph is a *tree* if there is a unique path between any two vertices.

An *n-cycle* is a path of length n with the modification that $\alpha = \beta$, i.e. it begins and ends in the same point. A graph is *acyclic* if it has no cycles. In particular *directed acyclic graphs* will be of interest in that it is the basic structure to be used in the model specification.

For a directed graph \mathcal{G} we define its *moral graph* \mathcal{G}^m as the undirected graph with the same vertex set as \mathcal{G} but with α and β adjacent in \mathcal{G}^m if and only if either α is a parent of β or conversely or if α and β have a common child γ .

The moral graph plays the same role in the present paper as in Lauritzen and Spiegelhalter (1988) in that it identifies groups of variables that enter simultaneously into the factors of the expression for the joint density. Figure 3 displays the moral graph of the network corresponding to the basic example studied.

Insert Figure 3 about here.

3.2 Decomposition of Marked Graphs

The basic trick enabling the computational task to be performed locally is the decomposition of a suitably modified network into partly independent components formed by the cliques of that graph.

The inherent asymmetry between discrete and continuous variables in the CG distributions implies that one also needs to take proper account of the behaviour of the markings of the graph. We refer the interested reader to Leimer (1989) for a detailed graphtheoretic study of the problems as well as all proofs. Here we introduce the notion of a decomposition by stating formally that

Definition 1 A triple (A, B, C) of disjoint subsets of the vertex set V of an undirected, marked graph \mathcal{G} is said to form a (strong) *decomposition* of \mathcal{G} if $V = A \cup B \cup C$ and the three conditions below all hold

- (i) C separates A from B
- (ii) C is a complete subset of V
- (iii) $C \subseteq \Delta \vee B \subseteq \Gamma$

When this is the case we say that (A, B, C) *decomposes* \mathcal{G} into the *components* $\mathcal{G}_{A \cup C}$ and $\mathcal{G}_{B \cup C}$.

If only (i) and (ii) hold, we say that (A, B, C) form a *weak decomposition*. Thus weak decompositions ignore the markings of the graph.

In the pure cases (iii) holds automatically and all weak decompositions are also decompositions. Note that what we have chosen to call a decomposition (without a qualifier) is what Leimer (1989) calls a strong decomposition. Figure 4 illustrates the notions of (strong) and weak decompositions.

Insert Figure 4 about here.

A decomposable graph is one that can be successively decomposed into its cliques. Again we choose to state this formally through a recursive definition as

Definition 2 An undirected, marked graph is said to be *decomposable* if it is complete, or if there exists a decomposition (A, B, C) with A and B both non-empty, into decomposable subgraphs $\mathcal{G}_{A \cup C}$ and $\mathcal{G}_{B \cup C}$.

Note that the definition makes sense because both subgraphs $\mathcal{G}_{A \cup C}$ and $\mathcal{G}_{B \cup C}$ must have fewer vertices than the original graph \mathcal{G} .

Decomposable, unmarked, graphs are characterised by being triangulated, i.e. not having cycles of length greater than three. Decomposable marked graphs are characterised by further not having any path of a particular type. More precisely

Proposition 1 *An undirected, marked graph is decomposable if and only if it is triangulated and does not contain any path $(\delta_1 = \alpha_0, \dots, \alpha_n = \delta_2)$ between two discrete vertices passing through only continuous vertices, with the discrete vertices not being neighbours.*

We have illustrated the typical forbidden path in Figure 4.

3.3 Junction Trees with Strong Roots

The construction of the computational structure in Lauritzen and Spiegelhalter (1988) as well as in Jensen, Lauritzen and Olesen (1990) begins with a directed acyclic graph, forms its moral graph, adds links to make it triangulated and forms a junction tree of the cliques of the triangulated graph. This procedure is to be generalised and modified.

The first part of the manipulations is unchanged in that we begin with a directed acyclic graph. Then we form the moral graph by adding undirected edges between parents that are not already linked and drop directions to obtain an undirected graph. Finally we add further links such that we obtain a decomposable, marked graph. In our example we have made this modification in Figure 3. Note that the link between B and F is necessary to remove the forbidden path (B, E, F) and make the graph decomposable, whereas it is (weakly) triangulated even without this.

In this particular example the discrete variables end up forming a complete subset, but this is not always the case.

The next step is the construction of the junction tree. A *junction tree* is an organization of a collection of subsets of the set of variables V into a tree that satisfies the condition that $A \cap B$ is a subset of all sets on the path in the tree between A and B .

When a collection of subsets is organised in a junction tree one can show that it must be a set of complete subsets of a triangulated graph containing the cliques. For any two sets C and D that are neighbours in the junction tree, their intersection $S = C \cap D$ is called their *separator* because it separates $C \setminus D$ from $D \setminus C$ in the graph. When we make a picture of the junction tree, the separators are drawn as rectangles, see Figure 5.

The junction tree is the basic computational structure, but the asymmetry between continuous and discrete variables make a further condition necessary for the propagation scheme to work properly. Again we formally define

Definition 3 A subset R on a junction tree is a *strong root* if any pair A, B of neighbours on the tree with A closer to R than B satisfies

$$(B \setminus A) \subseteq \Gamma \vee (B \cap A) \subseteq \Delta. \quad (2)$$

The condition (2) is equivalent to the triple $(A \setminus B, B \setminus A, A \cap B)$ forming a strong decomposition of $\mathcal{G}_{A \cup B}$. In words it expresses that when a separator between two neighbouring cliques is not purely discrete, the clique furthest away from the root has only continuous vertices beyond the separator.

The statement iii)' of Theorem 2' of Leimer (1989) ensures that *the cliques of a decomposable marked graph can be organised in a junction tree with at least one strong root*. We assume henceforth that this has been done.

Figure 5 displays a junction tree for our example. The clique $\{W, E, B, F\}$ could be used as a strong root. For example, $\{W, M_{\text{in}}, D\}$ has only the continuous variable M_{in} beyond the separator $\{W, D\}$.

4. MODEL SPECIFICATION

As in the discrete case the qualitative part of the model is initially specified by a directed acyclic graph such as the one in Figure 1.

The graph specifies the basic dependencies among the variables by assuming that the joint distribution of these has the directed Markov property with respect to the graph. In other words, we assume that the density is equal to the product of the conditional densities of the variables attached to each vertex in the graph, given the states at their parent vertices, see Kiiveri, Speed and Carlin (1984) as well as Lauritzen, Dawid, Larsen and Leimer (1990).

To exploit the properties of CG potentials we need to further assume that the graph satisfies the constraint that *no continuous vertices have discrete children*. If this assumption is not fulfilled, we have to use approximate methods in the specification phase, see Section 7. In our example the assumption is clearly satisfied.

Next we specify, for each discrete variable A , the conditional distribution at A given the states at its parent vertices (which are then all discrete). In the case where A is continuous we assume the conditional distribution of the response variable Y associated with A to be of the type

$$\mathcal{L}(Y \mid \text{pa}(A)) = \mathcal{N}(\alpha(i) + \beta(i)^\top z, \gamma(i)).$$

Here $\text{pa}(A)$ is a short notation for the combination of discrete and continuous states (i, z) of the variables that are parents of A . In this formula $\gamma(i) > 0$, $\alpha(i)$ is a real number, and $\beta(i)$ is a vector of the same dimension as the continuous part z of the parent variables.

Note that we assume that the mean depends linearly on continuous parent variables and that the variance does not depend on the continuous part of the parent variables. The linear function as well as the variance is allowed to depend on the discrete part of the parent variables.

The conditional density then corresponds to a CG potential ϕ_A with defined on the combination (i, z, y) of parent variables (i, z) and response variable y with canonical characteristics (g_A, h_A, K_A) , where

$$g_A(i) = -\frac{\alpha(i)^2}{2\gamma(i)} - [\log\{2\pi\gamma(i)\}] / 2 \quad (3)$$

$$h_A(i) = \frac{\alpha(i)}{\gamma(i)} \begin{pmatrix} 1 \\ -\beta(i) \end{pmatrix} \quad (4)$$

$$K_A(i) = \frac{1}{\gamma(i)} \begin{pmatrix} 1 & -\beta(i)^\top \\ -\beta(i) & \beta(i)\beta(i)^\top \end{pmatrix}. \quad (5)$$

This follows from direct calculation using the expression for the normal density.

We simply write

$$\phi(i, z, y) = \{2\pi\gamma(i)\}^{-1/2} \exp \left[-\{y - \alpha(i) - \beta(i)^\top z\}^2 / \{2\gamma(i)\} \right],$$

resolve the parentheses, take logarithms and identify terms. Note that $K_A(i)$ has rank one and is therefore typically not positive definite.

In our basic example we specify the conditional distributions as follows

Burning regime. This variable is discrete and denoted by B . We let

$$P(B = \textit{stable}) = .85 = 1 - P(B = \textit{unstable}).$$

Filter state. This is discrete and denoted by F . We let

$$P(F = \textit{intact}) = .95 = 1 - P(F = \textit{defect}).$$

Type of waste. This is discrete and denoted by W . We let

$$P(W = \textit{industrial}) = 2/7 = 1 - P(W = \textit{household}).$$

Filter efficiency. This is represented on a logarithmic scale and denoted by E . We assume the relation $\text{dust}_{\text{out}} = \text{dust}_{\text{in}} \times \rho$ and get in the logarithmic scale that $\log \text{dust}_{\text{out}} = \log \text{dust}_{\text{in}} + \log \rho$. We let $E = \log \rho$ and admit that filter inefficiency might be a better word for the variable E . We then specify

$$\begin{aligned} \mathcal{L}(E | \textit{intact}, \textit{household}) &= \mathcal{N}(-3.2, .00002) \\ \mathcal{L}(E | \textit{defect}, \textit{household}) &= \mathcal{N}(-.5, .0001) \\ \mathcal{L}(E | \textit{intact}, \textit{industrial}) &= \mathcal{N}(-3.9, .00002) \\ \mathcal{L}(E | \textit{defect}, \textit{industrial}) &= \mathcal{N}(-.4, .0001) \end{aligned}$$

This corresponds to filter efficiencies $1 - \rho$ on about 96%, 39%, 98% and 33% respectively. For example, when the filter is defect and household waste is burnt, the filter removes a fraction of $1 - \exp(-.5) = .39$ of the dust.

Emission of dust. This is represented on a logarithmic scale as a variable D . We let

$$\begin{aligned} \mathcal{L}(D | \textit{stable}, \textit{industrial}, e) &= \mathcal{N}(6.5 + e, .03) \\ \mathcal{L}(D | \textit{stable}, \textit{household}, e) &= \mathcal{N}(6.0 + e, .04) \\ \mathcal{L}(D | \textit{unstable}, \textit{industrial}, e) &= \mathcal{N}(7.5 + e, .1) \\ \mathcal{L}(D | \textit{unstable}, \textit{household}, e) &= \mathcal{N}(7.0 + e, .1). \end{aligned}$$

Thus, on a day when household waste is burned under a stable regime and the filter works perfectly, the typical concentration will be $\exp(6.0 - 3.2) = 16.4\text{mg}/\text{Nm}^3$. Similarly, if the filter is defect on a day with industrial waste and the burning regime is unstable, we will typically see an output concentration of dust on $\exp(7.5 - 0.4) = 1212\text{mg}/\text{Nm}^3$.

Concentration of CO₂. This is represented on a logarithmic scale as a variable C . We let

$$\mathcal{L}(C | \textit{stable}) = \mathcal{N}(-2, .1), \quad \mathcal{L}(C | \textit{unstable}) = \mathcal{N}(-1, .3).$$

Thus the concentration of CO₂ is typically around 14% under a stable regime and 37% when the burning process is unstable.

Penetrability of light. Represented on a logarithmic scale as a variable L . We let

$$\mathcal{L}(L | D = d) = \mathcal{N}(3 - d/2, .25).$$

This corresponds to the penetrability being roughly inversely proportional to the square root of the concentration of dust.

Metal in waste. The concentration of heavy metal in the waste is represented as a continuous variable M_i on a logarithmic scale. We let

$$\mathcal{L}(M_{\text{in}} | \textit{industrial}) = \mathcal{N}(0.5, 0.01), \quad \mathcal{L}(M_{\text{in}} | \textit{household}) = \mathcal{N}(-0.5, 0.005).$$

The precise interpretation is unit dependent, but the main point is that industrial waste tends to contain heavy metal in concentrations that are about three times as high as in household waste. Also the variability of the metal concentrations is higher in industrial waste.

Emission of metal. A continuous variable M_{out} on a logarithmic scale. We let

$$\mathcal{L}(M_{\text{out}} | d, m_{\text{in}}) = \mathcal{N}(d + m_{\text{in}}, 0.002).$$

Thus we simply assume that the concentration of emitted metal is about the same in the dust emitted as in the original waste.

The numbers have been constructed with a view to information from Hansen and Dalager (1985) but are otherwise purely fictitious.

5. BASIC OPERATIONS ON CG POTENTIALS

The basis for the computational scheme consists partly of a set of fundamental operations on the CG potentials, partly on a message passing scheme in the junction tree. The latter is to be described in Section 6. Here we describe the elements of the local computations.

Recall from Section 2 that a CG potential is any function ϕ of the form

$$\phi(x) = \phi(i, y) = \chi(i) \exp\{g(i) + h(i)^\top y - y^\top K(i)y/2\},$$

where $K(i)$ is a symmetric matrix. The triple (g, h, K) are the canonical characteristics for the potential ϕ .

5.1 Extension

If (g, h, K) are the characteristics of a CG potential ϕ defined on the variables (i, y) , we need sometimes to operate on this as if it was defined on a larger set (i, j, y, z) of variables. This is formally done by extending it to $\bar{\phi}$, letting $\bar{\phi}(i, j, y, z) = \phi(i, y)$. Clearly the corresponding characteristics are

$$\bar{g}(i, j) = g(i), \quad \bar{h}(i, j) = \begin{pmatrix} h(i) \\ 0 \end{pmatrix}, \quad \bar{K}(i, j) = \begin{pmatrix} K(i) & 0 \\ 0 & 0 \end{pmatrix}.$$

Hence the extension essentially amounts to adjoining zeros to the characteristics such as to give them the desired dimensions.

5.2 Multiplication and Division

Multiplication of two functions is defined the natural way, after the functions have been extended as above to be defined on the same space of variables. Expressed in terms of the canonical characteristics, multiplication becomes simple addition:

$$(g_1, h_1, K_1) * (g_2, h_2, K_2) = (g_1 + g_2, h_1 + h_2, K_1 + K_2).$$

Division is likewise defined in the obvious way, but special care has to be taken when dividing by zero. Thus for $x = (i, y)$ we let

$$(\phi/\psi)(x) = \begin{cases} 0 & \text{if } \phi(x) = 0 \\ (\phi(x)/\psi(x)) & \text{if } \psi(x) \neq 0 \\ \text{undefined} & \text{otherwise} \end{cases}$$

5.3 Marginalization

An essential difference between the pure discrete case and the situation in the present paper is due to fact that adding two CG potentials in general will result in a function of a different structure. Hence there will be some complications.

We distinguish several cases. First we discuss *marginals over continuous variables*. In this case we simply integrate. Let

$$y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}, \quad h = \begin{pmatrix} h_1 \\ h_2 \end{pmatrix}, \quad K = \begin{pmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{pmatrix}$$

with y_1 having dimension p and y_2 dimension q . We then have

Lemma 1 *The integral $\int \phi(i, y_1, y_2) dy_1$ is finite if and only if K_{11} is positive definite. It is then equal to a CG potential $\tilde{\phi}$ with canonical characteristics given as*

$$\begin{aligned} \tilde{g}(i) &= g(i) + \left\{ p \log(2\pi) - \log \det K_{11}(i) + h_1(i)^\top K_{11}(i)^{-1} h_1(i) \right\} / 2 \\ \tilde{h}(i) &= h_2(i) - K_{21}(i) K_{11}(i)^{-1} h_1(i) \\ \tilde{K}(i) &= K_{22}(i) - K_{21}(i) K_{11}(i)^{-1} K_{12}(i). \end{aligned}$$

Proof Let

$$\mu(i) = -K_{11}(i)^{-1}K_{12}(i)y_2 + K_{11}(i)^{-1}h_1(i).$$

Then we find by direct calculation that

$$\begin{aligned}\phi(i, y) &= \exp \left\{ -(y_1 - \mu(i))^{\top} K_{11}(i)(y_1 - \mu(i))/2 \right\} \\ &\times \exp \left\{ y_2^{\top} (h_2(i) - K_{21}(i)K_{11}(i)^{-1}h_1(i)) \right\} \\ &\times \exp \left\{ -y_2^{\top} (K_{22}(i) - K_{21}(i)K_{11}(i)^{-1}K_{12}(i))y_2/2 \right\} \\ &\times \exp \left\{ g(i) + h_1(i)^{\top} K_{11}(i)^{-1}h_1(i)/2 \right\}.\end{aligned}$$

Now y_1 only appears in the first factor. This can be integrated by letting $z = y_1 - \mu(i)$ and recalling that if $z \in \mathcal{R}^p$ and K is positive definite, then

$$\int e^{-z^{\top} K z/2} dz = (2\pi)^{\frac{p}{2}} (\det K)^{-\frac{1}{2}}.$$

The result follows. \square

When calculating *marginals over discrete variables* we distinguish two cases. First, if h and K do not depend on j , i.e $h(i, j) \equiv h(i)$ and $K(i, j) \equiv K(i)$, we define the marginal $\tilde{\phi}$ of ϕ over j the direct way

$$\begin{aligned}\tilde{\phi}(i, y) &= \sum_j \phi(i, j, y) = \sum_j \chi(i, j) \exp \{ g(i, j) + h(i)^{\top} y - y^{\top} K(i) y/2 \} \\ &= \exp \{ h(i)^{\top} y - y^{\top} K(i) y/2 \} \sum_j \chi(i, j) \exp g(i, j)\end{aligned}$$

which leads to the following canonical characteristics for the marginal

$$\tilde{g}(i) = \log \sum_{j: \chi(i, j)=1} \exp g(i, j), \quad \tilde{h}(i) = h(i), \quad \tilde{K}(i) = K(i).$$

If either of h or K depends on j , the marginalization process is more subtle since simple addition of CG potentials will not result in a CG potential. The procedure we shall then use is only well defined for $K(i, j)$ positive definite and is best described in terms of the moment characteristics $\{p, \xi, \Sigma\}$. The marginal $\tilde{\phi}$ is defined as the potential with moment characteristics $\{\tilde{p}, \tilde{\xi}, \tilde{\Sigma}\}$ where

$$\tilde{p}(i) = \sum_j p(i, j), \quad \tilde{\xi}(i) = \sum_j \xi(i, j) p(i, j) / \tilde{p}(i),$$

and

$$\tilde{\Sigma}(i) = \sum_j \Sigma(i, j) p(i, j) / \tilde{p}(i) + \sum_j (\xi(i, j) - \tilde{\xi}(i))^{\top} (\xi(i, j) - \tilde{\xi}(i)) p(i, j) / \tilde{p}(i).$$

The ‘marginalized’ density will then have the correct moments, i.e.

$$P(I = i) = \tilde{p}(i), \quad \mathbf{E}(Y | I = i) = \tilde{\xi}(i), \quad \mathbf{V}(Y | I = i) = \tilde{\Sigma}(i),$$

where expectations are taken with respect to the CG distribution determined by ϕ . This is a direct consequence of the familiar relations

$$\mathbf{E}(Y | I = i) = \mathbf{E} \{ \mathbf{E}(Y | (I, J)) | I = i \} \quad (6)$$

$$\mathbf{V}(Y | I = i) = \mathbf{E} \{ \mathbf{V}(Y | (I, J)) | I = i \} + \mathbf{V} \{ \mathbf{E}(Y | (I, J)) | I = i \}. \quad (7)$$

When *marginalizing over both continuous and discrete variables* we first marginalize over the continuous variables and then over the discrete. If in the second of these stages we have (h, K) independent of j , we say that we have a *strong marginalization*. In the other case we must use the marginalization process just described and we speak of a *weak marginalization*. In both cases we use the symbol $\sum_{W \setminus V} \phi_W$ for the marginalized potential, where V denotes the set of variables marginalized to, and $W \setminus V$ the set of variables marginalized over.

We leave to the reader to verify that the weak marginalization satisfies the standard composition rule such that when $U \subset V \subset W$ then

$$\sum_{V \setminus U} \left(\sum_{W \setminus V} \phi_W \right) = \sum_{W \setminus U} \phi_W. \quad (8)$$

However, only the strong marginalisations behave well when products are involved. In general we have

$$\sum_{W \setminus V} (\phi_W h_V) \neq h_V \left(\sum_{W \setminus V} \phi_W \right). \quad (9)$$

A consequence is that the axioms of Shenoy and Shafer (1990) are not fulfilled. Hence we have to establish correctness of our propagation scheme directly without exploiting their general computational theory.

In the special case of *strong* marginalisations equality holds in (9). This follows by elementary calculations since strong marginalizations are just ordinary integrations.

6. OPERATING IN THE JUNCTION TREE

When the model has been specified, the handling of incoming evidence and calculation of specific probabilities is done in the junction tree representation using the elementary operations described in the previous section.

Essentially the junction tree representation of the cliques in the strongly triangulated, moralised graph captures the computationally interesting aspects of the product structure in the joint density of the variables involved. Then the computations can be performed locally within the cliques and between the cliques that are neighbours in the junction tree.

Hence we assume that a junction tree with strong root has been established on the basis of the original graph such as discussed in Section 3.

Each subset of variables in the junction tree is referred to as a *belief universe* and the set of all variables as the *total universe*. The collection of belief universes in the junction tree is denoted by \mathcal{C} , to indicate that it is the set of cliques in a strongly decomposable graph.

Recall from Section 3 that the intersections of neighbours in the junction tree are called separators. The collection of separators is denoted by \mathcal{S} , where this may involve multiple copies of the same separator set. Both the belief universes and the separators can have belief potentials ϕ_W attached to them and these are all assumed to be CG potentials defined on the corresponding spaces of variables. The *joint system belief* ϕ_U associated with the given attachment of potentials is defined as

$$\phi_U = \frac{\prod_{V \in \mathcal{C}} \phi_V}{\prod_{S \in \mathcal{S}} \phi_S}, \quad (10)$$

and is assumed to be proportional to the joint density of all the variables. Since all potentials involved are CG potentials, the joint density will be a CG density itself.

We always assume that the tree is *supportive* meaning that for any universe V with neighbouring separator S we have $\phi_S(x) = 0 \implies \phi_V(x) = 0$. This enables us to deal correctly with cases where some states are ruled out as impossible by having potentials equal to zero.

6.1 Initializing the Junction Tree

As a first step, the junction tree with strong root has to be initialized according to the model specification, to make sure that the tree is supportive and the joint system belief given by (10) is the joint density specified by the model as in Section 4. This is done as follows.

First we assign each vertex A in the original graph to a universe V in the tree. This has to be done in such a way that $(A \cup \text{pa}(A)) \subseteq V$ but is otherwise arbitrary. This ensures the universe to be so large that the CG potential ϕ_A obtained from the conditional density of A given $\text{pa}(A)$, can be extended to V .

For each universe V we then let ϕ_V be the product of all the (extensions of) potentials ϕ_A for vertices assigned to it. On the separators we let $\phi_S \equiv 1$, i.e. the potential with canonical characteristic $(0, 0, 0)$. This is also the potential on universes with no vertices assigned to them.

In our basic example, there are several possibilities for initializing the junction tree in Figure 5. An assignment of vertices to universes could be B and C to $\{B, C\}$, F , W and E to $\{B, F, W, E\}$, D to $\{B, W, E, D\}$, L to $\{L, D\}$, M_{in} to $\{W, D, M_{\text{in}}\}$ and M_{out} to $\{D, M_{\text{in}}, M_{\text{out}}\}$.

The potentials in the belief universes are then obtained as follows: The assignment of B to the universe gives for the value *stable* a potential with characteristics $(\log .85, 0, 0) = (-.16252, 0, 0)$. The assignment of C gives for

the same value a potential where from (3)

$$g_{\{C\}}(stable) = -\frac{(-2)^2}{2 \times .1} - \{\log(2\pi \times .1)\}/2 = -20 + .23235 = -19.76765,$$

and from (4) and (5) we get $h_{\{C\}}(stable) = -2/.1 = 20$, and $k_{\{C\}}(stable) = (1)/.1 = 10$. Adding these numbers and rounding off leads for the stable case to the potentials

$$g_{\{B,C\}}(stable) = -19.930, h_{\{B,C\}}(stable) = -20, k_{\{B,C\}}(stable) = 10.$$

Analogously for the unstable case we find

$$g_{\{B,C\}}(unstable) = -3.881, h_{\{B,C\}}(unstable) = -3.333, k_{\{B,C\}}(unstable) = 3.333.$$

Only the variable L was assigned to the universe $\{L, D\}$ and this has no discrete parents. Hence only $h_{\{L,D\}}$ and $K_{\{L,D\}}$ are needed. From (4) we find

$$h_{\{L,D\}} = \frac{3}{.25} \begin{pmatrix} 1 \\ -.5 \end{pmatrix} = \begin{pmatrix} 12 \\ 6 \end{pmatrix}$$

and from (5) that

$$K_{\{L,D\}} = \frac{1}{.25} \begin{pmatrix} 1 & .5 \\ .5 & .25 \end{pmatrix} = \begin{pmatrix} 4 & 2 \\ 2 & 1 \end{pmatrix}.$$

Similar calculations have to be performed for the remaining belief universes.

The basic computational structure is now established. The belief universes are objects carrying information in the form of potentials and the separators are communication channels through which information can flow. The computational scheme which is described below is a combination of entering evidence to relevant universes and a message passing algorithm.

6.2 Entering Evidence

Incoming evidence is envisaged to be of the type that certain states are impossible for particular discrete variables or combinations of these, and that certain continuous variables are in specific states.

To be able to handle evidence in our local computational scheme, each item of evidence must be concerned with groups of variables that are members of the same universe in the junction tree. Thus an *item of evidence* is one of the following

- a function $\chi_W(i_W) \in \{0, 1\}$, where W is a set of discrete variables that is a subset of some universe V in the junction tree;
- a statement that $Y_A = y_A^*$ for a particular continuous variable A .

The first type of evidence – that we shall term *discrete* evidence – is entered simply by multiplying χ_W onto the potential ϕ_V . Then it holds that the joint system belief will be proportional to the conditional density, given that the states for which χ_W is equal to zero, are impossible, *i.e.* it represents the conditional belief, given the evidence.

If the second type of evidence – *continuous* evidence – is entered, the potentials have to be modified in all universes V containing A . We have to modify the potentials to become those, where y_A becomes fixed at the value y_A^* . If the potential ϕ has canonical characteristics (g, h, K) with

$$h(i) = \begin{pmatrix} h_1(i) \\ h_A(i) \end{pmatrix}, \quad K(i) = \begin{pmatrix} K_{11}(i) & K_{1A}(i) \\ K_{A1}(i) & K_{AA}(i) \end{pmatrix},$$

the transformed potentials ϕ^* will have canonical characteristics (g^*, h^*, K^*) given as

$$\begin{aligned} K^*(i) &= K_{11}(i) \\ h^*(i) &= h_1(i) - y_A^* K_{A1}(i) \\ g^*(i) &= g(i) + h_A(i) y_A^* - K_{AA}(i) (y_A^*)^2 / 2. \end{aligned}$$

Note that a continuous item of evidence has to be entered to all universes and separators, of which A is a member.

In the example we could know that the waste burned was of industrial type and enter this information as the function χ_W with $\chi_W(\text{industrial}) = 1$ and $\chi_W(\text{household}) = 0$. Similarly we might have measured the light penetration to be 1.1 and the concentration of CO_2 to -0.9 , both on the logarithmic scale applied when specifying the conditional distributions. The latter translates to a concentration of 41% CO_2 in the emission. Then the potentials from the initialization are modified to become, for example

$$g_{\{B\}}^*(\text{stable}) = -19.930 + 18 - 4.050 = -5.980$$

and

$$g_{\{B\}}^*(\text{unstable}) = -3.881 + 3 - 1.350 = -2.231$$

as well as

$$h_{\{D\}}^* = 6 - 1.1 \times 2 = 3.8, \quad k_{\{D\}}^* = 1.$$

6.3 Absorption

The fundamental process in the message passing algorithm is that of a universe absorbing information from its neighbours in the junction tree.

So consider a tree of belief universes with collection \mathcal{C} and separators \mathcal{S} assumed to be supportive. Let $V \in \mathcal{C}$ and let W_1, \dots, W_m be neighbours of V with separators S_1, \dots, S_m respectively. The universe V is said to *absorb* from

W_1, \dots, W_m if the following operations are performed on the belief potentials

$$\begin{aligned}\phi'_{S_i} &= \sum_{W_i \setminus V} \phi_{W_i} \\ \phi'_V &= \phi_V * (\phi'_{S_1} / \phi_{S_1}) * \dots * (\phi'_{S_m} / \phi_{S_m}).\end{aligned}$$

In words, the potentials of all the neighbours are marginalised to the separators and the ratio between the new and old separator potential is passed on as a ‘likelihood ratio’ and multiplied onto the potential at V .

We note that after an absorption, the belief potential for S_i is the marginal of W_i with respect to S_i and the tree remains supportive.

We also have that $\phi_V / (\phi_{S_1} * \dots * \phi_{S_m}) = \phi'_V / (\phi'_{S_1} * \dots * \phi'_{S_m})$, whence the joint system belief is unchanged by the absorption process.

In the particular case where $m = 1$ the universes V and W will, under certain circumstances, after absorption ‘contain the same information’ on common variables. More precisely we say that V and W are *consistent* if $\sum_{V \setminus S} \phi_V \propto \phi_S \propto \sum_{W \setminus S} \phi_W$, and a tree of belief universes is said to be *locally consistent* if all mutual neighbours in the tree are consistent. We then have

Lemma 2 *If V absorbs from W and ϕ_S is the strong marginal of ϕ_V , then V and W are consistent after absorption. In fact $\sum_{V \setminus W} \phi'_V = \phi'_S = \sum_{W \setminus V} \phi'_W$.*

Proof Since the marginalisation over $V \setminus W$ is strong it is composed of integrations and summations only. Hence we find that

$$\sum_{V \setminus W} \phi'_V = \sum_{V \setminus W} \phi_V * (\phi'_S / \phi_S) = (\phi'_S / \phi_S) * \sum_{V \setminus W} \phi_V = \phi'_S.$$

The other equality is trivial. \square

We emphasize, that the corresponding result is false when ϕ_S is only the weak marginal of ϕ_V , see (9). The necessity of using junction trees with strong roots to obtain exact propagation, is a consequence of this fact. In the situation described in the lemma we say that V has *calibrated* to W .

6.4 Collecting Evidence

Based on the notion of absorption the propagation scheme can now be constructed exactly as in the discrete case.

Each $V \in \mathcal{C}$ is given the action COLLECTEVIDENCE: When COLLECTEVIDENCE in V is called from a neighbour W , then V calls COLLECTEVIDENCE in all its other neighbours and when they have finished their COLLECTEVIDENCE, V absorbs from them (see Figure 6).

Insert Figure 6 about here.

We note that since COLLECTEVIDENCE is composed of absorptions only, after COLLECTEVIDENCE, the joint system belief is unchanged and the tree remains supportive.

The idea is now to evoke COLLECTEVIDENCE from a strong root R in the junction tree. A flow of activation of neighbours will move through the tree and a flow of absorptions towards the root will take place. When the flow terminates, the root R will have absorbed the information available from all parts of the tree.

If COLLECTEVIDENCE is evoked from a strong root R and W and W^* are neighbours with separator S such that W is closer in the tree to R than W^* , then the COLLECTEVIDENCE from R has caused W to absorb from W^* . Thus, after COLLECTEVIDENCE, the belief potential for S is the marginal of W^* with respect to S . Since the root is strong, the marginal will be strong. This can be exploited for a second flow through the tree, to be described subsequently.

6.5 Distributing Evidence

After COLLECTEVIDENCE the root R has absorbed all information available. Next it must pass this information on to the remaining universes in the tree, formalised as the operation DISTRIBUTEVIDENCE.

Each $V \in \mathcal{C}$ is given the action DISTRIBUTEVIDENCE: When DISTRIBUTEVIDENCE is called in V from a neighbour W then V absorbs from W and calls DISTRIBUTEVIDENCE in all its other neighbours.

The activation of DISTRIBUTEVIDENCE from the root R will create an outward flow of absorptions that will stop when it has reached the leaves of the tree. Again the joint system belief and supportiveness remain unchanged under DISTRIBUTEVIDENCE. But it even holds that when DISTRIBUTEVIDENCE has terminated, the resulting tree of belief universes will be locally consistent.

This follows since after COLLECTEVIDENCE all separator potentials will be strong marginals of potentials further away from the strong root. When DISTRIBUTEVIDENCE is subsequently performed, Lemma 2 ensures that all absorptions are calibrations.

We shall now argue that this locally consistent tree is the representation we are aiming for.

Ideally we would want is a tree of belief universes such that the probability distributions can be directly inferred from the local belief potentials without having to calculate the joint system belief.

This is clearly too much to demand, since the true marginal distribution at any universe would be a mixture of CG distributions and not a CG distribution itself. What we can hope to get is an equality of moments or, in other words, that each local potential is the correct (weak) marginal of the joint system belief. That this in fact is true is the main result of the paper:

Theorem 1 *Let T be a locally consistent junction tree of belief universes with a strong root R and collection \mathcal{C} . Let ϕ_U be the joint system belief for T and let $V \in \mathcal{C}$. Then*

$$\sum_{U \setminus V} \phi_U \propto \phi_V. \quad (11)$$

Proof Let n denote the number of universes in the collection \mathcal{C} . We first realise that it is enough to consider the case $n = 2$: If $n=1$ the statement obviously holds. If $n > 2$ we can find a leaf L in the tree and use the case $n = 2$ on the junction tree with strong root $R' = \cup_{V \in \mathcal{C} \setminus \{L\}} V$ and one leaf L . By induction the case gets reduced to $n = 2$.

So assume that $U = R \cup L$ where R is a strong root and let $S = R \cap L$ be the separator. The marginal to R is a strong marginal and we find by Lemma 2 that

$$\sum_{L \setminus R} \phi_U = \sum_{L \setminus R} \frac{\phi_L \phi_R}{\phi_S} = \frac{\phi_R}{\phi_S} \sum_{L \setminus R} \phi_L = \phi_R.$$

If S contains only discrete variables the marginal to L is also strong and the same calculation applies.

Else $L \setminus R$ contains only continuous variables. Then $L \setminus S \subseteq \Gamma$, i.e. only continuous vertices are in the external part of the leaf. Denote the states of the variables in S by (i, y) and those in $L \setminus S$ by z . Since ϕ_S is the weak marginal of ϕ_U , the moments $p(i)$, $\mathbf{E}(Y | I = i)$ and $\mathbf{V}(Y | I = i)$ are correct when calculated according to ϕ_S or, since these are identical, also according to ϕ_V . That the remaining moments now are correct follows since

$$\begin{aligned} \mathbf{E}(Z | I = i) &= \mathbf{E}\{\mathbf{E}(Z | Y, I = i)\} = \mathbf{E}\{A(i) + B(i)Y | I = i\} \\ &= A(i) + B(i)\mathbf{E}(Y | I = i), \end{aligned}$$

where $A(i)$ and $B(i)$ are determined from ϕ_L/ϕ_S alone. Similarly

$$\begin{aligned} \mathbf{E}(Z^\top Y | I = i) &= \mathbf{E}\{\mathbf{E}(Z^\top Y | Y, I = i)\} = \mathbf{E}\{(A(i) + B(i)Y)^\top Y | I = i\} \\ &= A(i)^\top \mathbf{E}(Y | I = i) + \mathbf{E}(Y^\top B(i)Y | I = i) \\ &= A(i)^\top \mathbf{E}(Y | I = i) + \mathbf{E}(Y | I = i)^\top B(i) \mathbf{E}(Y | I = i) + \text{tr}\{B(i)\mathbf{V}(Y | I = i)\}. \end{aligned}$$

Finally

$$\begin{aligned} \mathbf{V}(Z | I = i) &= \mathbf{E}\{\mathbf{V}(Z | Y, I = i)\} + \mathbf{V}\{\mathbf{E}(Z | Y, I = i)\} = C(i) + B(i)^\top C(i)B(i) \end{aligned}$$

where also the conditional covariance $C(i)$ is determined from ϕ_L/ϕ_S alone, whence the moments are correct. \square

In summary, after entering evidence the junction tree can be made consistent by evoking `COLLECTEVIDENCE` and then `DISTRIBUTEVIDENCE` from a strong root. The weak marginal of the belief at a vertex A can subsequently be obtained from any universe (or even separator) containing A by further weak marginalisation.

In particular this gives the correct updated probabilities of the states of any discrete variable and the correct updated mean and variance of any continuous variable.

If the full marginal density is required for a continuous variable, further computations are needed which typically involve all variables on the path between a strong root and a universe containing the variable in question. In general both the density itself and the problem of its computation can be forbiddingly complex.

We want to point out that, although the marginal density of variables cannot be obtained explicitly in practice, the tree still contains a fully correct representation of the joint system belief, given the evidence. No information is lost during the flow of evidence. Hence, the system remains ready for a correct, exact updating of beliefs when more evidence is obtained.

In the example we have displayed the initial and updated marginal probabilities, the means and the standard deviations in Table 1.

Insert Table 1 about here.

As it was to be expected from measuring a CO₂ emission of 41%, there is strong evidence for an unstable burning regime whereas the filter must be intact to explain the penetrability. This, combined with the fact that industrial waste is burned, means that the expected emission of heavy metal has been increased with a factor of $\exp 1.3 \approx 3.7$.

7. FURTHER TOPICS

We shall first briefly touch upon the issues involving feasibility of the computations. The most complex operation is the weak marginalization over a given clique. If the clique contains discrete variables $\delta \in d$ with state spaces of cardinality n_δ and q continuous variables, the computational complexity of marginalization is of the order of magnitude $q^3 \prod_{\delta \in d} n_\delta$ whereas the storage requirements are about $q^2 \prod_{\delta \in d} n_\delta$. This is because matrix inversion of a $q \times q$ matrix takes about q^3 operations and about q^2 space, and this has to be performed for every cell in the table of discrete configurations. These quantities should be compared with $2^q \prod_{\delta \in d} n_\delta$, which is the complexity, both of computation and storage, when the q variables are discretized as much as to a binary variable. Thus when q is large, dramatic savings are possible.

However, it should be remembered that extra links may have to be added to make the graph strongly triangulated instead of just triangulated. These extra links may in particular cases increase clique size so much that the savings are thereby lost.

However, another possibility is to ignore the constraint that the junction tree needs a strong root and use an ordinary triangulation for the construction of the tree. Then also COLLECTEVIDENCE would involve weak marginalization and after propagation the tree would only be approximately consistent. The quality of such an approximation is to be explored. In particular the approximative methods could give rise to pathologies such as non positive definite covariance matrices etc.

In the case where the original directed graph has continuous variables with discrete children, the initialization will have to be done approximately rather than exactly to take advantage of the CG distributions in the above computational scheme.

So let i denote a typical state for a discrete variable with discrete parent states j and continuous parent states z , where $z \in \mathcal{R}^q$. We then have to approximate $\log p(i|j, z)$ with a second degree polynomial in z for such pairs (i, j) where $p(i|j, z)$ is strictly positive. In particular the positivity is not allowed to depend on z .

An obvious suggestion is to use a CG regression model for the conditional probabilities and let

$$\log p(i|j, z) = a(i|j) + b(i|j)^\top z + z^\top C(i|j)z - \kappa(j, z)$$

where

$$\kappa(j, z) = \log \sum_i \exp \left\{ a(i|j) + b(i|j)^\top z + z^\top C(i|j)z \right\}.$$

This seems natural since CG regressions occur as conditional distributions in CG distributions, see Lauritzen and Wermuth (1989).

Since κ is not a quadratic in z , $\log p$ is to be approximated by its second order Taylor expansion around its maximal value.

Various optimization methods can be used to find this maximum and we abstain from discussing this point in detail here.

When the approximate initialization has been performed and one recalculates the conditional distribution of i given the parent states, this will be of the same type as initially specified, but the coefficients can have changed slightly. This can be used indirectly to indicate the quality of the approximation. We believe that the error of approximation is negligible compared to the general uncertainty involved in the model building itself.

REFERENCES

- Hansen, J. A. and Dalager, S. (eds.) (1985), *Emission fra affaldsforbrændingsanlæg*. (In Danish). Copenhagen: DAKOFA.
- Jensen, F. V., Lauritzen, S. L. and Olesen, K. G. (1990), Bayesian updating in causal probabilistic networks by local computations. *Computational Statistics Quarterly* **4**, 269-282.
- Kiiveri, H., Speed, T. P. and Carlin, J. B. (1984), Recursive causal models. *Journal of the Australian Mathematical Society Series A* **36**, 30-52.
- Lauritzen, S. L. (1989), Mixed graphical association models (with discussion). *Scandinavian Journal of Statistics* **16**, 273-306.
- Lauritzen, S. L., Dawid, A. P., Larsen, B. N. and Leimer, H.-G. (1990), Independence properties of directed Markov fields. *Networks* **20**, 491-505.

- Lauritzen, S. L. and Spiegelhalter, D. J. (1988), Local computations with probabilities on graphical structures and their application to expert systems (with discussion). *Journal of the Royal Statistical Society Series B* **50**, 157-224.
- Lauritzen, S. L. and Wermuth, N. (1989), Graphical models for associations between variables, some of which are qualitative and some quantitative. *The Annals of Statistics* **17**, 31-57.
- Leimer, H.-G. (1989), Triangulated graphs with marked vertices. In *Graph Theory in Memory of G. A. Dirac* (L. D. Andersen et al., eds.) *Annals of Discrete Mathematics* **41**, 311-324.
- Oliver, R. and Smith, J. Q. (eds.) (1990), *Influence Diagrams, Belief Nets and Decision Analysis*. Chichester: Wiley.
- Pearl, J. (1988), *Probabilistic Inference in Intelligent Systems*. San Mateo: Morgan Kaufmann.
- Shafer, G. R. and Pearl, J. (eds.) (1990), *Readings in Uncertain Reasoning*. San Mateo: Morgan Kaufmann.
- Shenoy, P. P. and Shafer, G. R. (1990), Axioms for probability and belief-function propagation. In: *Uncertainty in Artificial Intelligence* **4**, Shachter et al. (eds.). North-Holland, Amsterdam, pp. 169-198.
- Whittaker, J. (1990), *Graphical Models in Applied Multivariate Analysis*. Chichester: Wiley.

	F	W	B	M_{in}	M_{out}	E	D	C	L
Status	$p(i)$	$p(h)$	$p(s)$	means and standard deviations					
Initial values	0.95	0.71	0.85	-0.21	2.83	-3.25	3.04	-1.85	1.48
				0.46	0.86	0.71	0.77	0.51	0.63
Updated values	0.9995	0	0.01	0.50	4.11	-3.90	3.61	-0.90	1.10
				0.10	0.35	0.07	0.33	0	0

Table 1: The probabilities, means and standard deviations of single variables in our example before and after evidence has been entered. The information strongly suggests that the filter is intact but the burning regime is unstable. As a consequence, there is increased emission of dust and metal.

Figure 1: Graphical representation of the emission problem. The variables Filter State (F), Type of Waste (W), and Burning Regime (B) – corresponding to filled circles – are conceived as qualitative variables with states $\{intact, defect\}$, $\{industrial, household\}$, and $\{stable, unstable\}$ respectively. The remaining variables are measured on a quantitative scale. These are: Metal in Waste (M_{in}), Emission of Metal (M_{out}), Filter Efficiency (E), Emission of Dust (D), CO₂ in Emission (C), and Light Penetrability (L).

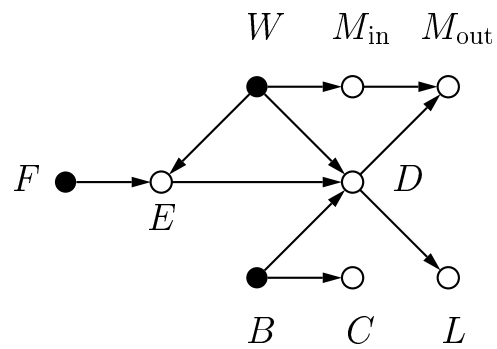
Figure 2: Illustration of graph theoretic concepts. The vertices α and β are continuous and the remaining discrete. We have $pa(\chi) = \{\gamma\}$. The vertices ϵ and ϕ are neighbours. The set $\{\alpha, \beta, \delta\}$ is a clique. A cycle is formed by $(\alpha, \delta, \phi, \epsilon)$.

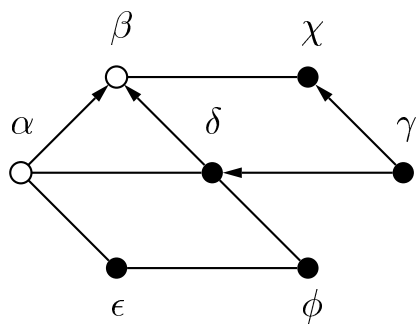
Figure 3: Modified graphs for the waste example. The graph obtained after marrying of parents and dropping directions is shown in (I). When further a link between B and F is added the strongly decomposable graph in (II) is obtained.

Figure 4: Illustration of decomposability. In (a) we see a decomposition with $C \subseteq \Delta$ and in (b) with $B \subset \Gamma$. In (c) the decomposition is only weak because none of these two conditions are fulfilled. In (d) we do not have a decomposition because the separator C is not complete. The graph (e) displays a path which is forbidden in a decomposable graph.

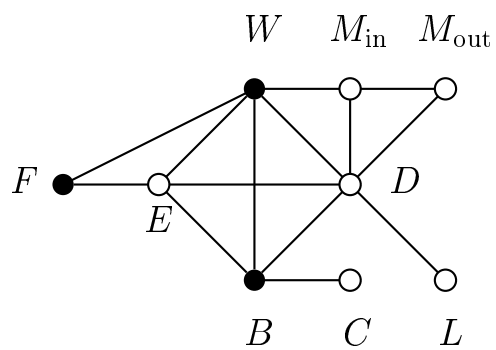
Figure 5: A junction tree of the waste example. The separators are drawn as rectangular boxes. Possible strong roots are $\{B, C\}$ and $\{W, E, B, F\}$.

Figure 6: The calls and message passing in COLLECTEVIDENCE.

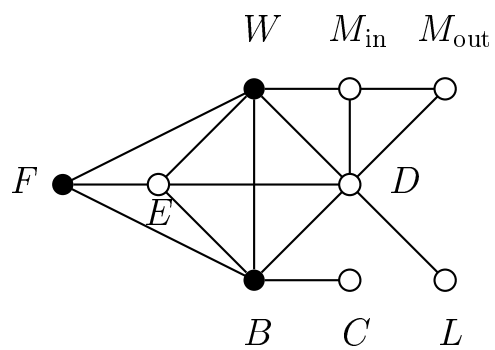


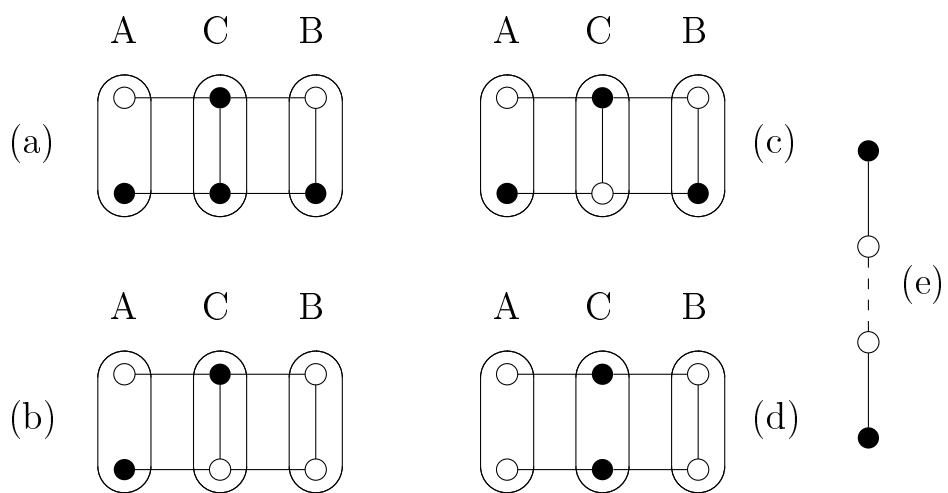


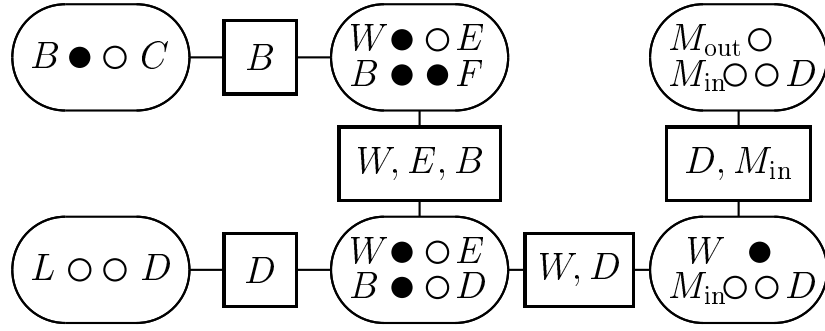
(I)

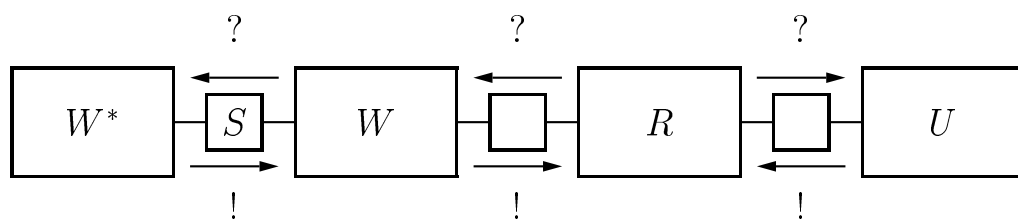


(II)









1	INTRODUCTION	2
2	CG DISTRIBUTIONS AND POTENTIALS	3
3	MARKED GRAPHS AND JUNCTION TREES	4
3.1	Notation and Terminology	4
3.2	Decomposition of Marked Graphs	5
3.3	Junction Trees with Strong Roots	6
4	MODEL SPECIFICATION	8
5	BASIC OPERATIONS ON CG POTENTIALS	10
5.1	Extension	11
5.2	Multiplication and Division	11
5.3	Marginalization	11
6	OPERATING IN THE JUNCTION TREE	13
6.1	Initializing the Junction Tree	14
6.2	Entering Evidence	15
6.3	Absorption	16
6.4	Collecting Evidence	17
6.5	Distributing Evidence	18
7	FURTHER TOPICS	20
	REFERENCES	21