# Model Selection and Accounting for Model Uncertainty in Graphical Models using Occam's Window

David Madigan and Adrian E. Raftery
University of Washington *

July 14, 1992
Revised: June 22, 1993; October 10, 1993

## Abstract

We consider the problem of model selection and accounting for model uncertainty in high-dimensional contingency tables, motivated by expert system applications. The approach most used currently is a stepwise strategy guided by tests based on approximate asymptotic $P$-values leading to the selection of a single model; inference is then conditional on the selected model. The sampling properties of such a strategy are complex, and the failure to take account of model uncertainty leads to underestimation of uncertainty about quantities of interest. In principle, a panacea is provided by the standard Bayesian formalism which averages the posterior distributions of the quantity of interest under each of the models, weighted by their posterior model probabilities. Furthermore, this approach is optimal in the sense of maximising predictive ability. However, this has not been used in practice because computing the posterior model probabilities is hard and the number of models is very large (often greater than $10^{11}$).

We argue that the standard Bayesian formalism is unsatisfactory and we propose an alternative Bayesian approach that, we contend, takes full account of the true model uncertainty by averaging over a much smaller set of models. An efficient search algorithm is developed for finding these models. We consider two classes of graphical models that arise in expert systems: the recursive causal models and the decomposable

1

log-linear models. For each of these, we develop efficient ways of computing exact Bayes factors and hence posterior model probabilities. For the decomposable log-linear models, this is based on properties of chordal graphs and hyper-Markov prior distributions and the resultant calculations can be carried out locally. The end product is an overall strategy for model selection and accounting for model uncertainty that searches efficiently through the very large classes of models involved.

Three examples are given. The first two concern data sets which have been analysed by several authors in the context of model selection. The third addresses a urological diagnostic problem. In each example, our model averaging approach provides better out-of-sample predictive performance than any single model which might reasonably have been selected.

KEYWORDS: Chordal graph; Contingency table; Decomposable log-linear model; Expert system; Hyper-Markov distribution; Recursive causal model.

# Contents

# 1  Introduction

Fruitful approaches to inference in high-dimensional contingency tables all involve choosing a broad class of models to be considered and then comparing them on the basis of how well they predict the data. Typically, the model classes are huge and inference in the presence of the many competing models is not easy.

Here we consider two classes of graphical models: the recursive causal models of Kiiveri *et al.* (1984) and the decomposable log-linear models introduced by Goodman (1970) and Haberman (1974). This work is motivated by applications in expert systems which use a belief network to represent knowledge and perform inference (Lauritzen and Spiegelhalter, 1988). These are the two model classes that arise in such applications. Potentially the most important advantage of constructing expert systems in this fashion is the system's ability to modify itself as data becomes available. In a series of recent papers, Spiegelhalter and Lauritzen (1990a,1990b), Dawid and Lauritzen (1993) and Spiegelhalter and Cowell (1991) have addressed the issue of updating the quantitative layer of such models. Building on this work, we address the issue of updating the qualitative layer—how can the graphical structure itself be updated as data becomes available?

Currently, the most used approach to model selection in contingency tables is a stepwise one, adapted from stepwise regression by Goodman (1971); see also Bishop, Fienberg and Holland (1975, Section 4.5 and Chapter 9). This consists of sequentially adding and deleting terms on the basis of approximate asymptotic likelihood ratio tests, leading to the selection of a single model. Inference about the quantities of interest is then made conditionally on the selected model.

There are several difficulties with this approach. The sampling properties of the overall strategy are complex because it involves multiple tests and, at least implicitly, the comparison of non-nested models (Fenech and Westfall, 1988). The use of $P$-values themselves is controversial, even when there are only two models to be compared, because of the so-called "conflict between $P$-values and evidence" discussed by Berger and Sellke (1987) and Berger and Delampady (1987). One aspect of this is that tests based on $P$-values tend to reject even apparently satisfactory models when the sample size is large; a dramatic example of this was discussed by Raftery (1986b). On the other hand, when the sample size is small and the table sparse, the asymptotic approximations on which the $P$-values are based tend to break down.

Perhaps most fundamentally, conditioning on a single selected model ignores model un-

3

certainty and so leads to underestimation of the uncertainty about the quantities of interest. This underestimation can be large, as was shown by Regal and Hook (1991) in the contingency table context and by Miller (1984) in the regression context. One bad consequence is that it can lead to decisions that are too risky (Hodges, 1987).

In principle, the standard Bayesian formalism provides a panacea for all these difficulties. If $\Delta$ is the quantity of interest, such as a parameter, a future observation, or the utility of a course of action, then its posterior distribution given data $D$ is

$$\mathrm{pr}(\Delta \mid D) = \sum_{k=1}^{K} \mathrm{pr}(\Delta \mid M_k, D)\mathrm{pr}(M_k \mid D). \tag{1}$$

This is an average of the posterior distributions under each of the models, weighted by their posterior model probabilities. In equation (1), $M_1, \ldots, M_K$ are the models considered and

$$\mathrm{pr}(M_k \mid D) = \frac{\mathrm{pr}(D \mid M_k)\mathrm{pr}(M_k)}{\sum_{\ell=1}^{K} \mathrm{pr}(D \mid M_\ell)\mathrm{pr}(M_\ell)}, \tag{2}$$

where

$$\mathrm{pr}(D \mid M_k) = \int \mathrm{pr}(D \mid \theta_k, M_k)\mathrm{pr}(\theta_k \mid M_k)d\theta_k \tag{3}$$

is the marginal likelihood of model $M_k$, $\theta_k$ is the (vector) parameter of $M_k$, $\mathrm{pr}(\theta_k \mid M_k)$ is the prior distribution of $\theta_k$, $\mathrm{pr}(D \mid \theta_k, M_k)$ is the likelihood, and $\mathrm{pr}(M_k)$ is the prior probability of $M_k$.

Furthermore, averaging over *all* the models in this fashion provides better predictive ability, as measured by a logarithmic scoring rule, than using any single model $M_j$:

$$- E\left[\log\left\{\sum_{k=1}^{K} \mathrm{pr}(\Delta \mid M_k, D)\mathrm{pr}(M_k \mid D)\right\}\right] \leq -E\left[\log\{\mathrm{pr}(\Delta \mid M_j, D)\}\right] \quad (j = 1, \ldots, K), \tag{4}$$

where $\Delta$ is the observable to be predicted and the expectation is with respect to $\sum_{k=1}^{K} \mathrm{pr}(\Delta \mid M_k, D)\mathrm{pr}(M_k \mid D)$. This follows from the non-negativity of the Kullback-Leibler information divergence. The logarithmic scoring rule was suggested by Good (1952) and assigns to each event $A$ which occurs a score of $-\log\{\mathrm{pr}(A)\}$. See Dawid (1986) for further discussion and Kass and Raftery (1993) for a review of the general approach.

Cooper and Herskovits (1992) present this approach in the context of recursive causal models. However, the approach in general has not been adopted in practice. This appears to be because (a) the posterior model probabilities $\mathrm{pr}(M_k \mid D)$ are hard to compute since they involve the very high-dimensional integrals in equation (3), and (b) the number of models in the sum in equation (1) can be huge. For example, with just 10 variables (small by expert

4

system standards) there are approximately $4 \times 10^{18}$ recursive causal models and $2 \times 10^{11}$ decomposable models.

One might hope that most of the posterior probability would be accounted for by a small number of models so that the sum in equation (1) would be well approximated by a small number of terms. Unfortunately, this is not typically the case because, although a small number of models do have much higher posterior probabilities than all the others, the very many models with small posterior probabilities contribute substantially to the sum. For example, Moulton (1991) reported a regression example with $2^{12} = 4096$ models where about 800 models were needed to account for 90% of the posterior probability.

We argue that the standard Bayesian formalism of equation (1) is flawed. Adopting standard methods of scientific investigation, we contend that accounting for the true model uncertainty involves averaging over a much smaller set of models. We present simple and efficient ways of computing the exact posterior model probabilities for the two model classes considered. Our approach is to take advantage of the graphical structure to calculate the required probabilities very quickly, while representing prior opinion in an easily elicitable form. We also describe an efficient algorithm for searching the very large model space.

Putting all this together gives us a simple and computationally efficient way of selecting the best models and accounting for model uncertainty in recursive causal models and decomposable log-linear models. To demonstrate the generality of our approach, our discussion will be in the context of conventional statistical model selection rather than expert systems. In Section 2 we describe the principles underlying our approach to model selection. In Section 3 we apply those principles to the recursive causal models, while in Section 4 we consider the decomposable models.

## 2 Model Selection Strategy

### 2.1 General Principles and Occam's Razor

We argue that equation (1) does not accurately represent model uncertainty. Science is an iterative process in which competing models of reality are compared on the basis of how well they predict what is observed; models that predict much less well than their competitors are discarded. Most of the models in equation (1) have been discredited in the sense that they predict the data far less well than the best models and so they should be discarded. Hence they should not be included in equation (1).

In our approach, if a model predicts the data far less well than the best model in the

5

class it will be discarded, so that initially we exclude from equation (1) those models not belonging to the set

$$\mathcal{A}' = \left\{ M_k : \frac{\max_l \{\text{pr}(M_l \mid D)\}}{\text{pr}(M_k \mid D)} \leq c \right\}, \tag{5}$$

for some constant $c$. The value of $c$ used will depend on the context. In our examples we used $c = 20$, by analogy with the popular .05 cutoff for $P$-values; Jeffreys (1961, Appendix B) would suggest some number between 10 and 100, while Evett (1991) suggests a value of 1000 for forensic evidence in criminal cases. (Note that we use $\text{pr}(M_k \mid D)$ rather than $\text{pr}(D \mid M_k)$ as the measure of how well the model predicts the data. In this way the likelihood is weighted by the prior model probability $p(M_k)$, assumed to reflect past data. This results in a composite predictive probability for both past and present data.)

Next we appeal to one of the most widely accepted norms of scientific investigation, namely Occam's razor. Let $E$ represent the evidence and $\text{pr}(H|E)$ the probability of a specified hypothesis $H$ given the evidence $E$. Occam's razor states that if:

$$\text{pr}(H_1|E) = \text{pr}(H_2|E) = ... = \text{pr}(H_k|E)$$

for hypotheses $H_1, ..., H_k$, then the simplest among $H_1, ..., H_k$ is to be preferred (Kotz and Johnson, 1985). Thus we also exclude from equation (1) models belonging to the set

$$\mathcal{B} = \left\{ M_k : \exists M_l \in \mathcal{A}', M_l \subset M_k, \frac{\text{pr}(M_l \mid D)}{\text{pr}(M_k \mid D)} > 1 \right\} \tag{6}$$

and equation (1) is replaced by

$$\text{pr}(\Delta \mid D, \mathcal{A}) = \sum_{M_k \in \mathcal{A}} \text{pr}(\Delta \mid M_k, D)\text{pr}(M_k \mid D, \mathcal{A}) \tag{7}$$

where

$$\mathcal{A} = \mathcal{A}' \backslash \mathcal{B}.$$

This considerably reduces the number of models in the sum in equation (1) and hence simplifies the model uncertainty problem a great deal. Note that our argument is *not* an approximation adopted for computational convenience, but rather a solution based on accepted scientific methodology. Also, note also that our approach in equation (7) will not necessarily give an answer close to that given by equation (1) because, due to the very large number of models in the class, the models discarded may have a large total posterior probability $\sum_{M_k \notin \mathcal{A}} \text{pr}(M_k \mid D)$, even though each individual model discarded has a very small posterior probability. Similarly, excluding the models not in the set $\mathcal{A}$ may result in violations of the
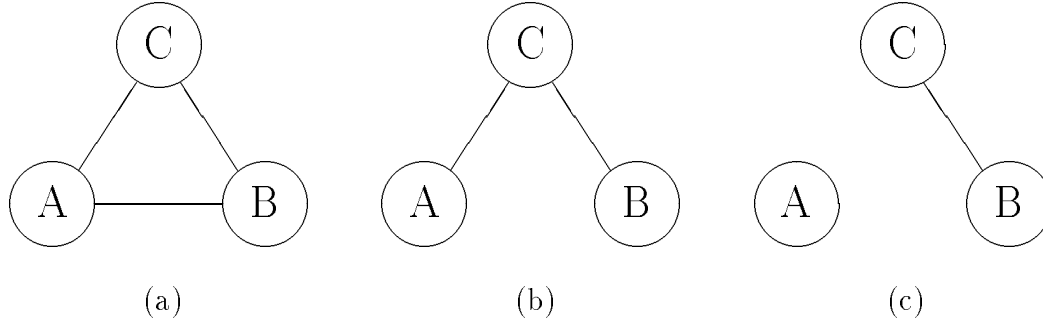
Figure 1: Model Selection Strategy - A Simple Example

inequality (4). However, experience to date suggests that this will rarely happen and that averaging over the models in $\mathcal{A}$ provides better predictive performance than conditioning on a single model. See Section 5 for further discussion.

The problem thus reduces to finding the set $\mathcal{A}$, and we now outline a computational strategy for doing this.

## 2.2   Model Selection Strategy

Our approach is heuristic in nature, and is a variant of the greedy-search algorithm. The essentials of the approach are the same for the recursive causal models and the decomposable models and could readily be applied to more general graphical models. Posterior model probabilities are used as a metric to guide the search. The strategy proceeds out into model space away from the opening set of models, comparing models via ratios of posterior model probabilities in a series of nested comparisons. In what follows, $M_0$ will denote the smaller of the two models being compared and $M_1$ will denote the larger. In fact, $M_0$ and $M_1$ will differ by just one link throughout.

Our basic rule is that if $M_0$ is rejected, then so are all its submodels. Here we define $M$ to be a submodel of $M_0$ if all the links in $M$ are also in $M_0$. To see this, consider the (undirected) example in Figure 1. Suppose that we start with the saturated model $[ABC]$ of Figure 1(a), and that when we compare it with the model of conditional independence $[AC][BC]$ of Figure 1(b), we reject the smaller model decisively. Then we are precisely rejecting the conditional independence of $A$ and $B$ given $C$. This conditional independence also holds in all the submodels of $[AC][BC]$ and so we reject all of those as well, including the model $[A][BC]$ of Figure 1(c). Thus, if we reject a model, we reject all its submodels.

7

In the algorithm described in Section 2.3, even this rule is relaxed in the sense that $[A][BC]$ may be subsequently considered as a submodel of a different model.

Our basic rule is the first of two "coherence" rules proposed by Gabriel (1969) for sequential testing procedures based on monotone test statistics. His second rule was that if a model is not rejected, then no model that includes it is considered rejected, but this second rule does not apply here because posterior model probabilities are not monotone (i.e., unlike deviance, for example, the posterior probability for a particular model can be smaller than the posterior probability of its submodels). The model selection strategy of Edwards and Havránek (1985) is based on both these rules, while that of Havránek (1984) is based on the first rule alone.

## 2.3  Occam's Window

A crucial aspect of the strategy concerns the interpretation of the ratio of posterior model probabilities when comparing two models. Again we appeal to Occam's razor which we implement as follows:

- If the log posterior odds is positive, i.e., the data provides evidence for the smaller model, then we reject $M_1$ and consider $M_0$. We could generalize this by requiring the log posterior odds to be greater than some positive constant $O_R$ before rejecting $M_1$.

- If the log posterior odds is small and negative, providing evidence against the smaller model which is not very strong, then we consider both models.

- If the log posterior odds is large and negative, i.e. smaller than $O_L = -\log(c)$ where $c$ is defined by equation (5), we reject $M_0$ and consider $M_1$.

Thus there are three possible actions following each comparison—see Figure 2.

Now that the various elements of the strategy are in place, we outline the search technique. The search can proceed in two directions: "Up" from each starting model by adding links, or "Down" from each starting model by dropping links. When starting from a non-saturated, non-empty model, we first execute the "Down" algorithm. Then we execute the "Up" algorithm, using the models from the "Down" algorithm as a starting point. Experience to date suggests that the ordering of these operations has little impact on the final set of models. Let $\mathcal{A}$ and $\mathcal{C}$ be subsets of model space $\mathcal{M}$, where $\mathcal{A}$ denotes the set of "acceptable" models and $\mathcal{C}$ denotes the models under consideration. For both algorithms, we begin with $\mathcal{A} = \emptyset$ and $\mathcal{C} =$ set of starting models.
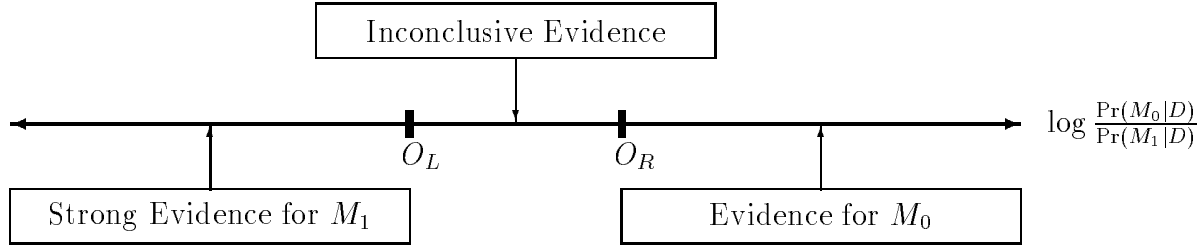
8

Figure 2: Occam's Window: Interpreting the log posterior odds, $\log \frac{\Pr(M_0|D)}{\Pr(M_1|D)}$, where $M_0$ is a submodel of $M_1$

### BGMS–Down Algorithm

1. Select a model $M$ from $\mathcal{C}$

2. $\mathcal{C} \leftarrow \mathcal{C} \setminus \{M\}$ and $\mathcal{A} \leftarrow \mathcal{A} \cup \{M\}$

3. Select a submodel $M_0$ of $M$ by removing a link from $M$

4. Compute $B = \log \frac{\mathrm{pr}(M_0|D)}{\mathrm{pr}(M|D)}$

5. If $B > O_R$ then $\mathcal{A} \leftarrow \mathcal{A} \setminus \{M\}$ and if $M_0 \notin \mathcal{C}, \mathcal{C} \leftarrow \mathcal{C} \cup \{M_0\}$

6. If $O_L \leq B \leq O_R$ then if $M_0 \notin \mathcal{C}, \mathcal{C} \leftarrow \mathcal{C} \cup \{M_0\}$

7. If there are more submodels of $M$, go to 3

8. If $\mathcal{C} \neq \emptyset$, go to 1

### BGMS–Up Algorithm

1. Select a model $M$ from $\mathcal{C}$

2. $\mathcal{C} \leftarrow \mathcal{C} \setminus \{M\}$ and $\mathcal{A} \leftarrow \mathcal{A} \cup \{M\}$

3. Select a supermodel $M_1$ of $M$ by adding a link to $M$

4. Compute $B = \log \frac{\mathrm{pr}(M|D)}{\mathrm{pr}(M_1|D)}$

5. If $B < O_L$ then $\mathcal{A} \leftarrow \mathcal{A} \setminus \{M\}$ and if $M_1 \notin \mathcal{C}, \mathcal{C} \leftarrow \mathcal{C} \cup \{M_1\}$

6. If $O_L \leq B \leq O_R$ then if $M_1 \notin \mathcal{C}, \mathcal{C} \leftarrow \mathcal{C} \cup \{M_1\}$

7. If there are more supermodels of $M$, go to 3

8. If $\mathcal{C} \neq \emptyset$, go to 1

Upon termination, $\mathcal{A}$ contains the set of potentially acceptable models. Finally, we remove all the models which satisfy equation (6), where 1 is replaced by $\exp(O_R)$, and those models $M_k$ for which

$$\frac{\max_l\{\text{pr}(M_l \mid D)\}}{\text{pr}(M_k \mid D)} > c. \tag{8}$$

The set $\mathcal{A}$ now contains the acceptable models.

# 3 The Directed Case—Recursive Causal Model Selection

## 3.1 Implementation

Implementation for the recursive causal models proceeds in a straightforward fashion. Consider a recursive causal model for a set of random variables $X_v, v \in V$. The model is represented by a directed graph where each variable in $V$ is represented by a node in the graph. For each variable $v \in V$ we define $\text{pa}(v)$ to be the set of parent nodes of $v$, i.e. nodes $w$ for which there exists a directed link from $w$ to $v$. The assumptions of the model imply that the joint distribution of $X_v, v \in V$, which we denote $\text{pr}(V)$, is given by

$$\text{pr}(V) = \prod_{v \in V} \text{pr}(v|\text{pa}(v)).$$

Here, and in what follows, we are using node labels like $v$, to represent the random variables which correspond to the nodes. In early implementations, $\text{pr}(v|\text{pa}(v))$ was assumed to be fully specified for all $v$ by the expert/data analyst. Spiegelhalter and Lauritzen (1990a) introduced a parameterization for $\text{pr}(v|\text{pa}(v))$ whereby the relationship between a node $v$ and its parents $\text{pa}(v)$ is fully specified by a vector parameter $\theta_v \in \Theta_v$. This leads to a conditional distribution for $V$:

$$\text{pr}(V|\theta) = \prod_{v \in V} \text{pr}(v|\text{pa}(v), \theta_v).$$

where $\theta$ is a general parameter with components $\theta_v$.

10

Spiegelhalter and Lauritzen (1990a) make two key assumptions which greatly simplify subsequent analysis. The first assumption is that of *global independence* whereby the parameters $\theta_v$ are assumed mutually independent *a priori*. This assumption alone allows us to calculate the likelihood for a single case:

$$\mathrm{pr}(V) = \int \mathrm{pr}(V, \theta)d\theta = \int \prod_v \mathrm{pr}(v|\mathrm{pa}(v), \theta_v)\mathrm{pr}(\theta_v)d\theta_v = \prod_v \mathrm{pr}(v|\mathrm{pa}(v))$$

where

$$\mathrm{pr}(v|\mathrm{pa}(v)) = \int \mathrm{pr}(v|\mathrm{pa}(v), \theta_v)\mathrm{pr}(\theta_v)d\theta_v.$$

The second assumption is that of *local independence* whereby the parameter $\theta_v$ breaks into components corresponding to the levels of the factors in $\mathrm{pa}(v)$. These components are assumed to be mutually independent *a priori*.

Now consider a conditional probability distribution $\mathrm{pr}(v|\mathrm{pa}(v)^+, \theta_v^+) = \theta_v^+$ for a specific set of levels, $\mathrm{pa}(v)^+$, of $\mathrm{pa}(v)$. We assume that $\theta_v^+$ has a Dirichlet distribution $\mathcal{D}[\lambda_1^+, ..., \lambda_k^+]$ where $k$ is the number of levels of $v$. Then we can show that

$$\mathrm{pr}(v = j|\mathrm{pa}(v)^+) = \lambda_j^+ / \sum_i \lambda_i^+, j = 1, 2, \ldots, k.$$

If we observe $v$ to be at level $x_{v_j}$ and the parent state to be $\mathrm{pa}(v)^+$, we have

$$\theta_v^+|v \sim \mathcal{D}[\lambda_1^+, ..., \lambda_j^+ + 1, ..., \lambda_k^+].$$

This provides a method for *sequentially* calculating the required ratios of posterior model probabilities and is simpler than the non-sequential approach. Furthermore, the sequential approach allows for efficient incorporation of new evidence. The elicitation of the required Dirichlet priors is feasible provided the cardinality of $\mathrm{pa}(v)$ is not too large. Computer-based methods for eliciting Dirichlet prior distributions have been described by Chaloner and Duncan (1987). If $\mathrm{pa}(v)$ is not observed the updating becomes more complex—see Spiegelhalter and Lauritzen (1990a,1990b) for details. A Jeffreys prior density was used in the examples, i.e. $\lambda_i^+ = 0.5$, for $i = 1, 2, \ldots, k$. A uniform prior typically selects identical models.

A considerable computational saving is obtained by noting that the sequential updating of the distribution of $\theta_v$ depends on the levels of $v$ and $\mathrm{pa}(v)$ *only*. Therefore the likelihood for all qualitative layers (graphs) having the same set $\mathrm{pa}(v)$ of parent nodes of $v$ will have identical contributions from $v$. For example, consider the two recursive causal models of Figure 3. When calculating the likelihood for the model of Figure 3(a), we store the likelihood of each
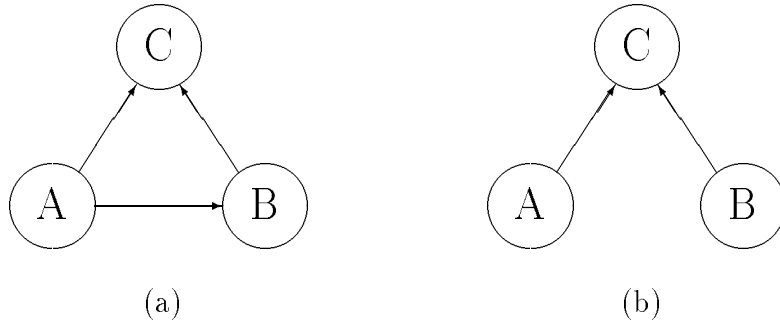
11

Figure 3: Simplifying the Likelihood Computations

node/parent combination separately. Now when subsequently calculating the likelihood for the model of Figure 3(b), only the likelihood for node $B$ requires recalculation as the sets of parent nodes of $A$ and $C$ have not changed.

To implement the model selection strategy described in Section 2 for the recursive causal models an *ordering* of the nodes must be pre-specified by the expert/data analyst. If $v_i$ precedes $v_j$ in the ordering, then a directed link from $v_j$ to $v_i$ is prohibited. In certain applications it may be possible to search over all possible orderings but this will typically not be the case. Pearl's IC-algorithm (Pearl and Verma, 1991) induces directed "causal" structures from data. An ordering of the nodes is not required, but for each pair of nodes $v_i$ and $v_j$, the algorithm does involve searching amongst all subsets of $V \setminus \{v_i, v_j\}$ for cutsets between $v_i$ and $v_j$ (sets which when conditioned on, render $v_i$ and $v_j$ independent.) Cooper and Herskovits (1992) provide a review of other approaches.

## 3.2   Examples

### 3.2.1   Coronary Heart Disease Risk Factors

Firstly we consider a data set which has been previously analysed by Edwards and Havránek (1985). The data concerns 1,841 men cross-classified according to six coronary heart disease risk factors. The risk factors are as follows: $A$, smoking; $B$, strenuous mental work; $C$, strenuous physical work; $D$, systolic blood pressure; $E$, ratio of $\beta$ and $\alpha$ proteins; $F$, family anamnesis of coronary heart disease.

Their likelihood ratio-based model selection strategy selected two graphical log-linear models: $[AC][ADE][BC][BE][F]$ which is not decomposable and therefore is not equivalent to any recursive causal model, and $[ACE][ADE][BC][F]$ which is decomposable. A striking
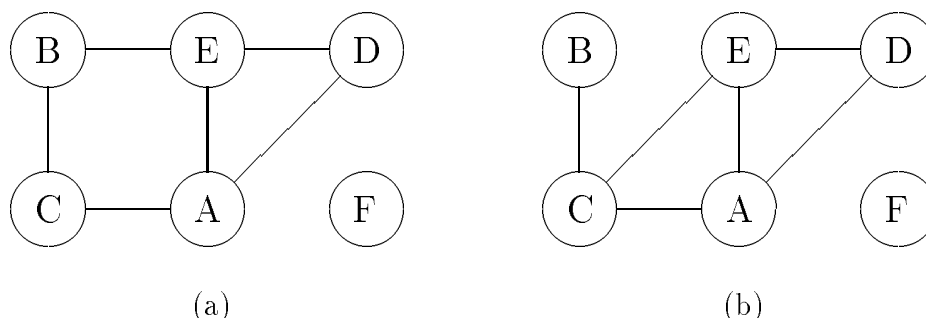
12

Figure 4: Models Selected by Edwards and Havránek

Table 1: Coronary Heart Disease: Posterior Model Probabilities for Recursive Causal Models

| Figure | Posterior probability % |
|--------|--------------------------|
| 5(a) | 52 |
| 5(b) | 40 |
| 5(c) | 5 |
| 5(d) | 4 |

feature of both models is the independence of $F$, family anamnesis. The models are shown in Figure 4.

To implement the Bayesian graphical model selection procedure, we started from the saturated model and used the "Down" algorithm only (starting from the empty model and using the "Up" algorithm produced the same set of models). All qualitative structures were assumed equally likely *a priori*. A natural partial ordering of the variables suggests itself: $F, (B, C), A, (E, D)$. The variables $B, F$ or $C$ could not be "influenced" by the other factors and must be exogenous, although the ordering of $B$ and $C$ is unclear. Similarly, $D$ or $E$ could hardly influence $A$, although the ordering of $E$ and $D$ is unclear. The four corresponding complete orderings produced strong evidence for the precedence of $E$ over $D$, and weak evidence for the precedence of $C$ over $B$. Several further orderings were tried, but this "natural" ordering resulted in the models with highest posterior probabilities. The selected models are shown in Figure 5 and their posterior probabilities in Table 1.

The two most likely models are shown in Figures 5(a) and 5(b). They are rather similar in that both contain the $CB$, $CA$, $AE$, $ED$ and $AD$ links. The main difference between
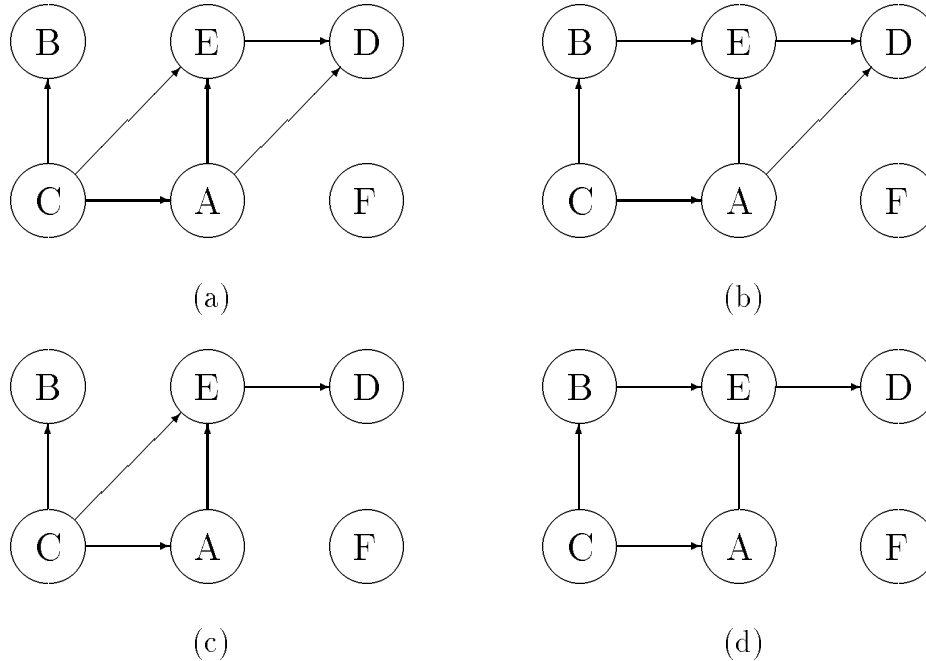
Figure 5: Coronary Heart Disease: Recursive Causal Models Selected

them lies in the way they describe the effect of strenuous mental work $(B)$ and strenuous physical work $(C)$ on the ratio of $\beta$ and $\alpha$ proteins $(E)$. Model 5(a) says that $C$ affects $E$ both directly and indirectly via $A$, whereas model 5(b) says that the effect of $C$ on $E$ is solely indirect, being mediated by $B$ and $A$. There is also some uncertainty about the presence of a link from smoking $(A)$ to systolic blood pressure $(D)$. The evidence favors the marginal independence of family anamnesis of coronary heart disease $(F)$.

The four models selected are similar to the models selected by Edwards and Havránek (1985) and shown in Figure 4. We note that the $AD$ link (smoking and systolic blood pressure) is present in both of the models of Figure 4 and also in models (a) and (b) of Figure 5, but it is absent from models (c) and (d) of Figure 5. In fact, the exact test for zero partial association of $A$ and $D$ reported by Edwards and Havránek (1985) had a significance level of 0.04 which was the largest of any of the links whose absence was rejected at the 5% level.

14

### 3.2.2  Women and Mathematics

Our second example concerns a survey which was reported in Fowlkes *et al.* (1988) concerning the attitudes of New Jersey high-school students towards mathematics. The data has been further analysed by Upton (1991). A total of 1190 students in eight schools took part in the survey. Data on six dichotomous variables was collected:

A. Lecture Attendance; attended or did not attend;

B. Sex; female or male;

C. School Type; suburban or urban;

D. "I'll need mathematics in my future work"; agree or disagree;

E. Subject Preference; maths/science or liberal arts;

F. Future Plans; college or job;

Upton (1991) reports that a model selection procedure based on the AIC criterion (Akaike, 1973) selects $[ABCE][CDF][BCD][DEF]$ while a procedure based on the BIC criterion (Raftery, 1986a) selects the much simpler $[A][BE][CE][CF][BD][DE][DF]$. Clearly an important difference between these two models is the treatment of $A$.

The Bayesian graphical model selection procedure started from the empty model and used the "Up" algorithm. It is clear that $B$ (Sex) cannot be influenced by other variables and must be exogenous. Initially it was also assumed that $C$ (School Type) was exogenous. An exhaustive search over all consequent orderings produced the single model shown in Figure 6.

The selected model is similar to the model selected by Upton's BIC procedure. The model selected by AIC clearly over-fits the data (Upton, 1991). It is of interest to note the direction of the link from $D$ to $F$. Both Upton (1991) and Fowlkes *et al.* (1988) treat $D$ as a response variable and Upton's path diagram shows a directed link from $F$ to $D$. However, the data provides strong evidence that the direction of the influence is from $D$ to $F$, i.e. that students' attitudes towards mathematics influence their future plans, rather than the other way around. The ability of the selected model to predict is unaffected by the direction of the $ED$ link.

Further analysis removed the restriction that $C$ be exogenous. The data now provides some support for the presence of a link from $E$ to $C$ although its interpretation is somewhat unclear.
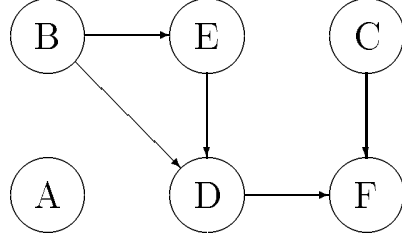
Figure 6: Women and Mathematics: Recursive Causal Model Selected

# 4 The Undirected Case—Decomposable Model Selection

## 4.1 Implementation

To implement the strategy for the decomposable models, we rely heavily on a recent fundamental paper by Dawid and Lauritzen (1993), hereafter DL. We consider three issues which are specific to model selection for the decomposable models.

First, how should we add and remove links whilst efficiently ensuring that all the models created are decomposable? Here we use a result which follows from Lemma 3 of Frydenberg and Lauritzen (1989): *Let $\mathcal{G} = (V, E)$ be a chordal graph with vertices $V$ and edges $E$ and let $\mathcal{G}' = (V, E')$ be a chordal subgraph of $\mathcal{G}$ with exactly one edge, $e$, less. Then $e$ is contained in exactly one clique of $\mathcal{G}$.* Therefore, the model selection strategy must remove only links which are members of a single clique. When adding links, the strategy must not create any chordless four-cycles.

Second, given any two decomposable models $M$ and $M^*$, is it possible to generate $M^*$ from $M$, adding or removing only one edge at a time but staying within the class of decomposable models? Lemma 5 of Frydenberg and Lauritzen (1989) shows this to be the case.

Finally, how do we calculate the required posterior model probabilities? Following DL, we consider a decomposable model $M$ for a set of random variables $X_v, v \in V$, whose joint distribution is specified by a vector parameter, $\theta$. $\theta$, in turn, is determined by the clique marginal probability tables $\theta_C = (\theta_C)_{C \in \mathcal{C}}$ where $\mathcal{C}$ denotes the set of cliques of $M$:

$$\theta(i) = \frac{\prod_{C \in \mathcal{C}} \theta_C(i_C)}{\prod_{S \in \mathcal{S}} \theta_S(i_S)}, i \in \mathcal{I},$$

16

where $\mathcal{S}$ denotes the system of separators in an arbitrary perfect ordering of $\mathcal{C}$, and $\mathcal{I}$ denotes the set of possible configurations of $X$.

For each clique $C \in \mathcal{C}$, let

$$\lambda_C = (\lambda_C(i_C))_{i_C \in \mathcal{I}_C}$$

be a given table of arbitrary positive numbers and let $\mathcal{D}(\lambda_C)$ denote the Dirichlet distribution for $\theta_C$ with density

$$\pi(\theta_C|\lambda_C) \propto \prod_{i_C \in \mathcal{I}_C} \theta_C(i_C)^{\lambda_C(i_C)-1},$$

where $\sum_{i_C} \theta_C(i_C) = 1$ and $\theta(i_C) > 0$.

Now let us suppose that the collection of specifications $\mathcal{D}(\lambda_C), C \in \mathcal{C}$ are constructed in such a way that for any two cliques $C$ and $D$ in $\mathcal{C}$ we have:

$$\lambda_C(i_{C \cap D}) = \lambda_D(i_{C \cap D}).$$

Then DL show that there exists a unique *strong hyper-Markov* distribution for $\theta$ over $M$ that has density $\mathcal{D}(\lambda_C)$ for all $C \in \mathcal{C}$. DL call this the hyper-Dirichlet distribution for $\theta$. A distribution for $\theta$ is strong hyper-Markov if and only if $\theta_{A|B}, \theta_{B|A}$ and $\theta_{A \cap B}$ are mutually independent whenever $A \cap B$ is complete and separates $A$ from $B$. It follows that by letting $\lambda_0 = \sum_{i \in \mathcal{I}} \lambda_i$, the likelihood for a single case is given by:

$$\mathrm{pr}(v) = \frac{\prod_{C \in \mathcal{C}} \lambda_C}{\lambda_0(\prod_{S \in \mathcal{S}} \lambda_S)}.$$

From Proposition 1 we have that updating can be carried out one clique at a time:

PROPOSITION 1 *If the prior distribution $\mathcal{L}(\theta)$ is strong hyper-Markov, the posterior distribution of $\theta$ is the unique hyper-Markov distribution $\mathcal{L}^*$ specified by the clique-marginal distributions $\{\mathcal{L}_C^* : C \in \mathcal{C}\}$, where $\mathcal{L}_C^*$ is the posterior distribution of $\theta_C$ based on its prior distribution $\mathcal{L}_C$ and the clique-specific data $X_C = x_C$.*

*Proof.* This is Corollary 9 of DL.

The posterior distribution for $\theta_C$ given data $n_C$ from the marginal table corresponding to clique $C$ is $\mathcal{D}(\lambda_C + n_C)$.

Consider the Bayes factor

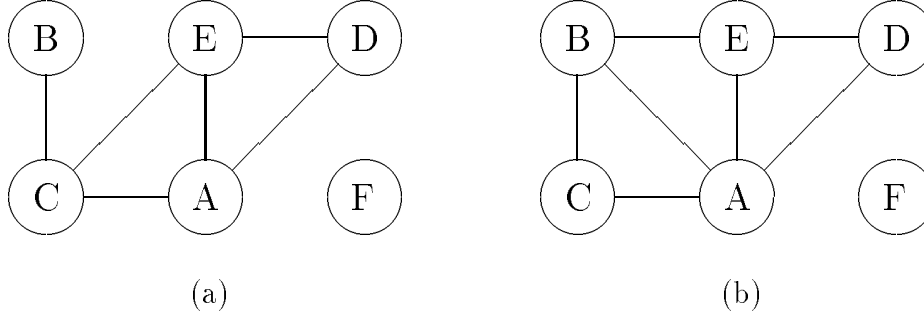$$B_{01} = \frac{\mathrm{pr}(D|M_0)}{\mathrm{pr}(D|M_1)}$$

17

Figure 7: Coronary Heart Disease: Decomposable Models Selected

where $M_0$ and $M_1$ are decomposable and $M_0$ is obtained from $M_1$ by deleting one edge $e$ linking $u$ with $v$. Since both models are decomposable we have that $e$ is contained in a single clique, $C$ say, of $M_1$. Let $C_u = C \setminus \{v\}, C_v = C \setminus \{u\}, C_0 = C \setminus \{u, v\}$. Then DL show that the Bayes factor is given by:

$$B_{01} = \frac{p_{C_u}(D_{C_u})p_{C_v}(D_{C_v})}{p_{C_0}(D_{C_0})p_C(D_C)}.$$

Thus, the required decomposable model comparisons can be carried out very rapidly with calculations local to single cliques.

## 4.2 Examples

### 4.2.1 Coronary Heart Disease Risk Factors

Firstly we consider again the coronary heart disease risk factor data of Edwards and Havránek (1985). We note that the model of Figure 4(a) which was selected by the Edwards and Havránek procedure is not decomposable and hence will not be selected by our procedure.

The selection procedure started from the saturated model and used the "Down" algorithm. All qualitative structures were assumed equally likely *a priori*. A standard Jeffreys prior was adopted for $\theta_C, C \in \mathcal{C}$. Just two models were selected and they are shown in Figure 7. Starting from the empty model and using the "Up" algorithm resulted in the same two models. The corresponding posterior probabilities are shown in Table 2.

The model of Figure 7(a) was also selected by the directed model selection procedure and by Edwards and Havránek (1985). The model of Figure 7(b) is essentially a decomposable version of the directed model of Figure 5(b) and Edwards and Havránek's model of Figure 4(a).

18

Table 2: Coronary Heart Disease: Posterior Model Probabilities for Decomposable Models

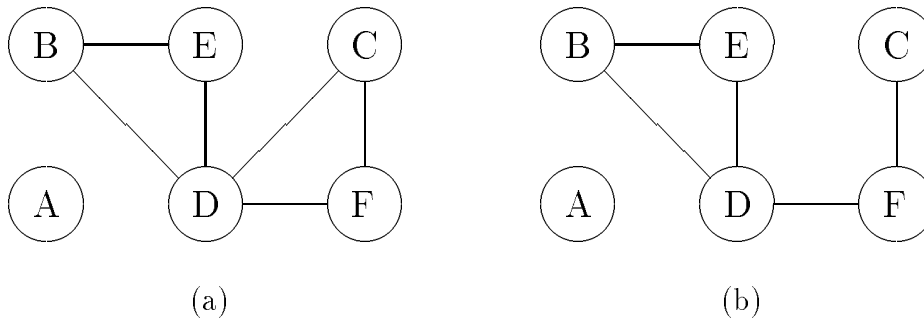| Figure | Model | Posterior probability % |
|--------|-------|------------------------|
| 7(a) | [BC][ACE][ADE][F] | 92 |
| 7(b) | [ABC][ABE][ADE][F] | 8 |



Figure 8: Women and Mathematics: Decomposable Models Selected

Overall, the model selection exercise indicates that there is very strong evidence for the $BC$, $AC$, $AE$ and $DE$ links, with evidence for the $AD$ link that is strong but somewhat less so. There is also some evidence for the $CE$ and $BE$ links, but it seems that one of these alone is enough to describe the data, and it is not fully clear which one is better. Of course the interpretation of these links is different in the two model classes. Again, as in the directed case, there is evidence for the marginal independence of $F$.

### 4.2.2 Women and Mathematics

We consider again the survey data previously analysed by Fowlkes *et al.* (1988) and Upton (1991). We note that the models selected in Upton (1991) are not graphical and hence will not be selected by our procedure. The procedure adopted was identical to that adopted for the example of Section 4.2.1. The two models selected are shown in Figure 8 and the corresponding posterior probabilities are shown in Table 3.

As in the directed case, the selected models are close to the models selected by the BIC model selection procedure carried out by Upton (1991). However there is uncertainty about the $CD$ link (School Type and "I'll need mathematics in my future work") which is not apparent in Upton's analysis. The odds in favor of the inclusion of the $CD$ link are 3 to 1,

Table 3: Women and Mathematics: Posterior Model Probabilities for Decomposable Models

| Figure | Model | Posterior probability % |
|--------|-------|-------------------------|
| 8(a) | $[A][BDE][CDF]$ | 75 |
| 8(b) | $[A][BDE][DF][CF]$ | 25 |

which Jeffreys (1961) would call evidence "not worth more than a bare mention". The data strongly supports the marginal independence of $A$.

### 4.2.3 Scrotal Swellings

Our final example concerns the diagnosis of scrotal swellings. Data on 299 patients was gathered at the Meath Hospital, Dublin, Ireland under the supervision of Mr. Michael R. Butler. We consider a cross-classification of the patients according to one disease class, Hernia ($H$), and 7 binary indicants as follows: $A$, possible to get above the swelling; $B$, swelling transilluminates; $C$, swelling separate from testes; $D$, positive valsalva/stand test; $E$, tender; $F$, pain; $G$, evidence of other urinary tract infections. The data is shown in Table 4. There are 28 possible links to be considered by the selection procedure in this example. In the absence of prior expert opinion, computation times can be prohibitive. Clearly, if the starting point for the selection procedure were close to the models for which the data provides evidence, this problem could be overcome.

With this objective we adopted the following heuristic procedure. First, Bayes factors for each of the 28 links are calculated by comparing the saturated model with the 28 sub-models generated by removing single links. The model consisting of the links for which the data provides evidence in this manner is then used as a starting point for the selection procedure. If this model is not decomposable, some of the links may be removed or additional ones may be added. A similar approach was suggested by Goodman (1973). The starting model is shown in Figure 9.

Now the "Up" algorithm is executed, followed by the "Down" algorithm (or *vice versa*). Note that if the starting links are badly chosen, the complete procedure has the opportunity to remove them, although, in this example, the final model contains all the links from the starting model. Two models were selected by this procedure and they are shown in Figure 10. The corresponding posterior probabilities are shown in Table 5.

The result of primary interest here is the importance of $A$ (possible to get above swelling)

20

Table 4: Scrotal Swelling data

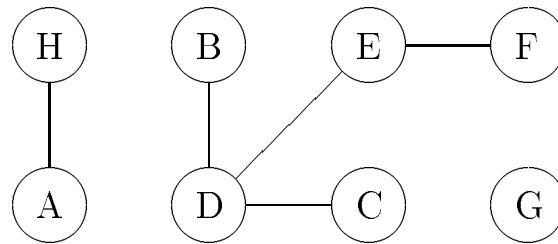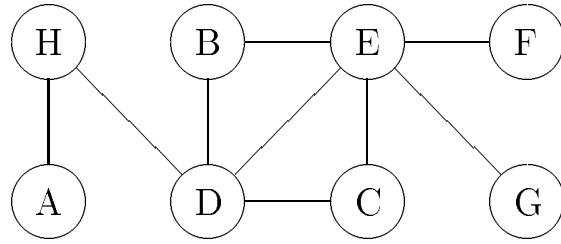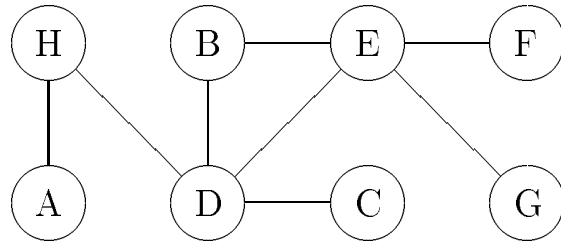| | Indicants | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Hernia | $A$ | $B$ | $C$ | $D$ | $E$ | $F$ | $G$ | Count |
| N | N | N | N | N | Y | Y | N | 1 |
| N | Y | N | N | N | N | N | N | 16 |
| N | Y | N | N | N | N | Y | N | 3 |
| N | Y | N | N | N | Y | Y | N | 51 |
| N | Y | N | N | N | Y | Y | Y | 17 |
| N | Y | N | Y | N | N | N | N | 30 |
| N | Y | N | Y | N | N | N | Y | 1 |
| N | Y | N | Y | N | N | Y | N | 3 |
| N | Y | N | Y | N | Y | N | N | 1 |
| N | Y | N | Y | N | Y | Y | N | 20 |
| N | Y | N | Y | N | Y | Y | Y | 4 |
| N | Y | N | Y | Y | N | N | N | 36 |
| N | Y | N | Y | Y | N | Y | N | 3 |
| N | Y | Y | N | N | N | N | N | 38 |
| N | Y | Y | N | N | N | N | Y | 1 |
| N | Y | Y | N | N | N | Y | N | 3 |
| N | Y | Y | N | N | Y | Y | N | 3 |
| N | Y | Y | Y | N | N | N | N | 21 |
| N | Y | Y | Y | N | Y | Y | N | 2 |
| Y | N | N | Y | Y | N | N | N | 39 |
| Y | N | N | Y | Y | N | Y | N | 5 |
| Y | Y | N | Y | Y | N | N | N | 1 |



Figure 9: Starting Model for Scrotal Swelling Example

(a)



(b)

Figure 10: Scrotal Swellings: Decomposable Models Selected

Table 5: Scrotal Swellings: Posterior Model Probabilities for Decomposable Models

| Figure | Model | Posterior probability % |
|--------|-------|--------------------------|
| 8(a) | $[AH][DH][BDE][CDE][EF][EG]$ | 75 |
| 8(b) | $[AH][DH][BDE][CD][EF][EG]$ | 25 |

and $D$ (valsalva/stand test) with respect to Hernia diagnosis. Both indicants can be established through simple procedures at physical examination. The only real model uncertainty which is exhibited concerns the relationship between $C$ (swelling separate from testes) and $E$ (tender). The odds in favor of the inclusion of the $CE$ link are 3 to 1 (evidence not worth more than a bare mention). Analysis of further cross-classifications extracted from this database also yield similarly sparse models.

# 5   Performance

Following Dawid (1984), we contend that one of the primary purposes of statistical analysis is to make forecasts for the future. Therefore, one way we can judge the efficacy of a model selection strategy, is to measure how well the resulting models predict future observations. In the case of Occam's Window, our specific objective is to compare the quality of the predictions based on model averaging against those based on any single model that an analyst might reasonably have selected.

We examined the predictive performance for each of the examples considered previously as follows: we randomly split the complete data sets into two subsets. One subset, $D^S$, containing 25% of the data, was used to select models, while $D^T = D \setminus D^S$, was used as a set of test cases. We measured performance by the logarithmic scoring rule of Good (1952). Specifically, we measured the predictive ability of an individual model, $M$, with:

$$- \sum_{d \in D^T} \log \mathrm{pr}(d \mid M, D^S).$$

We measured the predictive performance of model averaging with:

$$- \sum_{d \in D^T} \log \{ \sum_{M \in \mathcal{A}} \mathrm{pr}(d \mid M, D^S) \mathrm{pr}(M \mid D^S) \},$$

where $\mathcal{A}$ is the set of selected models.

We present results in Tables 6, 7 and 8 for each of the undirected examples of Section 4. In each case, we give the models selected and the performance measure (up to a normalising constant) for each individual model and for model averaging. For the Coronary Heart Disease example, we also include the score for the model selected by Whittaker (1990) on the basis of the full data set. The models selected by Upton (1991) and Fowlkes *et al.* (1988) are not included because they are not decomposable.

In each case, the method that averages over the models selected provides predictive performance which is superior to the performance resulting from basing the inference on any

Table 6: Coronary Heart Disease: Predictive Performance

| Model | Posterior probability % | Logarithmic Score |
|---|---|---|
| [AE][BC][BE][DE][F] | 26 | 4984.4 |
| [AC][BC][BE][DE][F] | 16 | 4990.2 |
| [AC][AE][BC][DE][F] | 13 | 4992.2 |
| [A][BC][BE][DE][F] | 9 | 4990.1 |
| [AE][BC][BE][D][F] | 8 | 4981.7 |
| [AE][BC][DE][F] | 7 | 4983.7 |
| [AC][BC][BE][D][F] | 5 | 4981.6 |
| [AC][BC][DE][F] | 4 | 4989.5 |
| [AC][AE][BC][D][F] | 4 | 4987.4 |
| [A][BC][BE][D][F] | 3 | 4989.4 |
| [A][BC][DE][F] | 2 | 4981.0 |
| [AE][BC][D][F] | 2 | 4980.9 |
| [AC][BC][D][E][F] | 1 | 4986.7 |
| [ABCE][ADE][BF] | Whittaker | 4984.8 |
| Model Averaging | | 4953.6 |

Table 7: Women and Mathematics: Predictive Performance

| Model | Posterior probability % | Logarithmic Score |
|---|---|---|
| [A][B][CDF][DE] | 75 | 3318.9 |
| [A][B][CF][DE][DF] | 21 | 3317.3 |
| [A][B][CF][DE] | 4 | 3320.4 |
| Model Averaging | | 3313.9 |

24

www.manaraa.com

Table 8: Scrotal Swellings: Predictive Performance

| Model | Posterior probability % | Logarithmic Score |
|---|---|---|
| $[AH][AD][BDE][CD][EF][FG]$ | 3 | 605.3 |
| $[AH][DH][BDE][CD][EF][FG]$ | 3 | 599.6 |
| $[AH][DH][BDE][CDE][EF][FG]$ | 5 | 600.6 |
| $[AH][AD][BDE][CDE][EF][FG]$ | 5 | 606.3 |
| $[AH][AD][BDE][CD][EF][EG]$ | 15 | 603.4 |
| $[AH][DH][BDE][CD][EF][EG]$ | 15 | 597.7 |
| $[AH][DH][BDE][CDE][EF][EG]$ | 27 | 598.7 |
| $[AH][AD][BDE][CDE][EF][EG]$ | 27 | 604.4 |
| Model Averaging | | 594.2 |

single model which might reasonably have been selected. In the coronary heart disease data, for example, our model averaging method outperforms the "best" model (i.e. that with the highest posterior probability) by 31 points of log predictive probability, or 62 points on the scale of twice the log probability on which deviances are measured. Repeating the random split or varying the subset proportions produces very similar performance results.

We also carried out a ROC (receiver operating characteristic) analysis for each of the examples. In Figure 11 we show two ROC curves for the variable $E$, ratio of $\beta$ and $\alpha$ proteins, in the coronary heart disease example. The solid ROC curve shows how well the single model with the highest posterior probability predicts variable $E$ while the dashed curve shows the performance achieved by averaging over the selected models. Again we used 25% of the data to select models and the remainder for testing.

These ROC curves show the false-positive and true-positive proportions for different probability thresholds for variable $E$. The area above the curve in the unit square provides a measure of predictive ability. Model averaging provides substantially better predictive performance in this instance. The area above the curve is 0.08 for model averaging, and 0.22 for the best single model. Thus model averaging reduces the average False Positive rate for a given True Positive rate by about two-thirds, where the False Positive rates are averaged over all True Positive rates. Model averaging is *not* guaranteed to provide superior predictive performance for each variable although the situation in Figure 11 is typical.
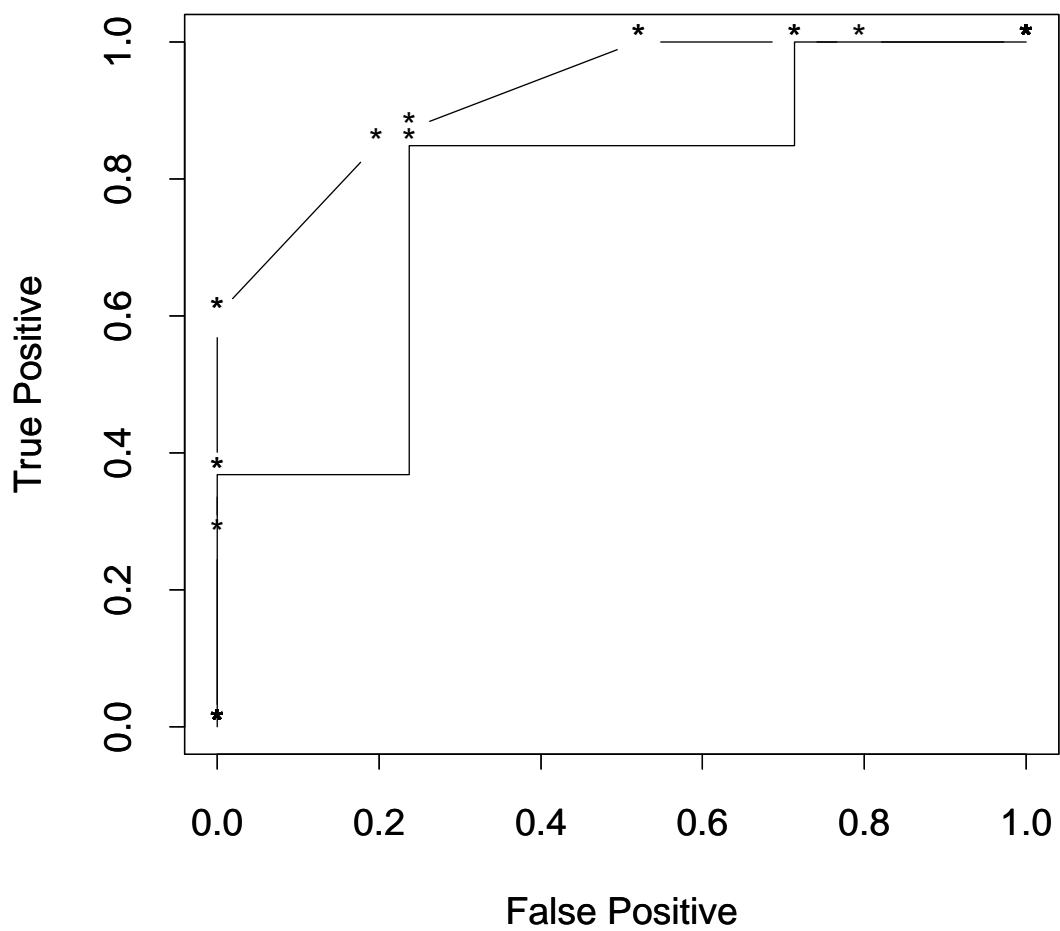
Figure 11: Coronary Heart Disease: ROC Curves for Node $E$

# 6  Discussion

## 6.1  General Comments and Other Approaches

We have outlined an overall strategy for model selection and accounting for model uncertainty in two important classes of models for high-dimensional contingency tables. This involves a redefinition of the Bayesian model uncertainty formalism, an efficient way of computing exact Bayes factors that exploits the graphical structure, and an algorithm for quickly searching through the very large model classes involved. The resulting procedure is quite efficient: for the example of Section 4.2.1, approximately 3,000 model comparisons per minute can be carried out on a Sun IPC.

There is a considerable literature on model selection for multidimensional contingency tables; this is generally concerned with the selection of a single "best" model. Most of it is based on the asymptotic properties of goodness-of-fit statistics (Goodman (1973), Wermuth (1976), Havránek (1984), Whittaker (1984), Edwards and Havránek (1985) or Fowlkes *et al.* (1988)). There are also approaches based on information criteria and discrepancy measures (Gokhale and Kullback, 1978; Sakamoto, 1984; Linhart and Zucchini, 1986). A recent review is provided by Upton (1991) who advocates the use of the BIC statistic. The calculation of Bayes factors for contingency table models has been considered by Spiegelhalter and Smith (1982), Raftery (1986a, 1988, 1993), Spiegelhalter and Lauritzen (1990a) and Spiegelhalter and Cowell (1991).

Pearl and Verma (1991) and Glymour *et al.* (1987) have proposed strategies for recovering causal structure from data. While these authors' objectives differ from ours, their procedures for selecting directed graphical structures have much in common with our recursive causal model selection strategy.

Cooper and Herskovits (1992) and Anderson *et al.* (1991) have examined model selection in the context of probabilistic expert systems. In both cases, the examples are based solely on data analysis and the incorporation of prior expert opinion is not considered. Cooper and Herskovits (1992) describe a general theory involving averaging over all models and suggest possible approximations. Their K2 strategy which seeks out the "best" recursive causal model for the qualitative layer, where "best" is taken to mean the *single* model with maximum probability. The algorithm starts with a model with no links and at each stage adds the directed link which most increases the model probability. The user must pre-specify an ordering of the nodes. Anderson *et al.* (1991) carry out their search in the undirected graphical model framework using a method introduced by Kreiner (1987). The difficulties

with large sparse tables mentioned above are avoided by using exact tests when comparing models.

## 6.2   Model Priors

In the examples considered above, the prior model probabilities $\mathrm{pr}(M)$ were assumed equal (Cooper and Herskovits, 1992, also assume that models are equally likely *a priori*). In general this can be unrealistic and may also be expensive and we will want to penalise the search strategy as it moves further away from the model(s) provided by the expert(s)/data analyst(s). Ideally one would elicit prior probabilities for all possible qualitative structures from the expert but this will be feasible only in trivial cases.

For models with fewer than 15 to 20 nodes, prior model probabilities may be approximated by eliciting prior probabilities for the presence of every possible link and assuming that the links are mutually independent, as follows. Let $\mathcal{E} = \mathcal{E}_P \cup \mathcal{E}_A$ denote the set of all possible links for the nodes of model $M$, where $\mathcal{E}_P$ denotes the set of links which are present in model $M$ and $\mathcal{E}_A$ denotes the absent links. For every link $e \in \mathcal{E}$ we elicit $\mathrm{pr}(e)$, the prior probability that link $e$ is included in $M$. The prior model probability is then approximated by

$$\mathrm{pr}(M) \propto \prod_{e \in \mathcal{E}_P} \mathrm{pr}(e) \prod_{e \in \mathcal{E}_A} (1 - \mathrm{pr}(e)).$$

Prior link probabilities from multiple experts are treated as independent sources of information and are simply multiplied together to give pooled prior model probabilities. Clearly, the contribution from each expert/data analyst could be weighted.

For applications involving a larger number of nodes or where the elicitation of link probabilities is not possible, we could assume that the "evidence" in favour of each link included by the expert(s)/data analyst(s) in the elicited qualitative structure(s) is "substantial" or "strong" but not "very strong" or "decisive" (Jeffreys, 1961). For example, we could assume that the evidence in favour of an included link lies at the center of Occam's window corresponding to a prior link probability for all $e \in \mathcal{E}_P$ of

$$\mathrm{pr}(e) = \frac{1}{1 + \exp(\frac{O_L + O_R}{2})}.$$

Similarly, the prior link probabilities for $e \in \mathcal{E}_A$ are given by

$$\mathrm{pr}(e) = \frac{\exp(\frac{O_L + O_R}{2})}{1 + \exp(\frac{O_L + O_R}{2})}.$$

In the directed case it may be possible to construct a prior distribution on the space of orderings—see Critchlow (1985) for further discussion.

## 6.3   Remaining Issues

While we believe that the methods we propose provide a workable approach to qualitative updating in expert systems, some issues remain. Spiegelhalter and Lauritzen (1990a) and other authors have expressed concerns about automatically updating the qualitative structure without reference to the domain expert. Such concerns need to be addressed in the context of real expert systems. Extension of the methods to include the more general graphical models of Wermuth and Lauritzen (1990) and Edwards (1990) will also be important. Missing data will frequently be a problem and we are currently exploring a number of techniques for the incorporation of missing data in the model selection strategy.

In the examples we have used vague priors for the model parameters that do not incorporate specific prior information. However, in expert system applications there will often be substantial prior information, and taking account of it would be expected to improve performance. How to elicit the required Dirichlet prior distributions is therefore a major issue. Direct elicitation is typically intractable; this has been a barrier to the use of the DL approach.

Madigan and Raftery (1991) outlined a simple approach to the elicitation of the required priors. They regarded the parameters of the Dirichlet prior distribution as "equivalent prior samples", which are elicited subject to constraints that ensure consistency. The priors are elicited sequentially in a way that avoids the need to store the full "equivalent prior" table.

# References

Anderson, L.R., Krebs, J.H. and Anderson J.D. (1991) STENO : An expert system for medical diagnosis based on graphical models and model search. *Journal of Applied Statistics*, **18**, 139–153.

Berger, J.O. and Delampady, M. (1987) Testing precise hypotheses (with Discussion). *Statistical Science*, **2**, 317–352.

Berger, J.O. and Sellke, T. (1987) Testing a point null hypothesis: The irreconcilability of P values and evidence (with Discussion). *Journal of the American Statistical Association* **82**, 112–139.

Bishop, Y.M.M., Fienberg, S.E. and Holland, P.W. (1975) *Discrete Multivariate Analysis.* Cambridge, Mass.: MIT Press.

Chaloner, K.M. and Duncan G.T. (1987) Some properties of the Dirichlet-Multinomial distribution and its use in prior elicitation. *Communications in Statistics–Theory and Methods*, **16**,511–523.

Cooper, G.F. and Herskovits, E. (1992) A Bayesian method for the induction of probabilistic networks from data. *Machine Learning*, **9**,309–347.

Critchlow, D.E. (1985) *Metric Methods for Analyzing Partially Ranked Data*, Springer-Verlag.

Dawid, A.P. (1984) Statistical theory–The prequential approach *Journal of the Royal Statistical Society (Series A)*, **147**,278–292.

Dawid, A.P. (1986) Probability Forecasting. In *Encyclopedia of Statistical Sciences, Volume 7*, (ed. Kotz, S. and Johnson, N.L.). Wiley: New York. 210–218.

Dawid, A.P. and Lauritzen, S.L. (1993) Hyper-Markov laws in the statistical analysis of decomposable graphical models. *Annals of Statistics*, to appear.

Edwards, D. (1990) Hierarchical mixed interaction models. *Journal of the Royal Statistical Society (Series B)*, **52**,3–20.

Edwards, D. and Havránek, T. (1985) A fast procedure for model search in multidimensional contingency tables. *Biometrika*, **72** ,339–351.

Evett, I.W. (1991) Implementing Bayesian methods in forensic science. Paper presented to the Fourth Valencia International Meeting on Bayesian Statistics.

Fenech, A. and Westfall, P. (1988) The power function of conditional log-linear model tests. *Journal of the American Statistical Association*, **83**,198–203.

Fowlkes, E.B., Freeny, A.E. and Landwehr, J.M. (1988) Evaluating logistic models for large contingency tables. *Journal of the American Statistical Association*, **83**,611–622.

Frydenberg, M. and Lauritzen, S.L. (1989) Decomposition of maximum likelihood in mixed graphical interaction models. *Biometrika*, **76**,539–555.

Gabriel, K.R. (1969) Simultaneous test procedures—Some theory of multiple comparisons. *Annals of Mathematical Statistics*, **40**,224–250.

Glymour, C., Scheines, R., Spirtes, P. and Kelly, K. (1987) *Discovering Causal Structure.* New York: Academic Press.

Gokhale, D.V. and Kullback, S. (1978) *The Information in Contingency Tables.* New York: Marcel Dekker.

Good, I.J. (1952) Rational Decisions. *Journal of the Royal Statistical Society (Series B)*, **14**,107–114.

Goodman, L.A. (1970) The multivariate analysis of qualitative data: Interaction among multiple classifications. *Journal of the American Statistical Association*, **65**,226–256.

Goodman, L.A. (1971) The analysis of multidimensional contingency tables: Stepwise procedures and direct estimation methods for building models for multiple classifications. *Technometrics* **13**, 33–61.

Goodman, L.A. (1973) Guided and unguided methods for the selection of models for a set of T multidimensional contingency tables. *Journal of the American Statistical Association*, **68**,165–175.

Haberman, S.J. (1974) *The Analysis of Frequency Data.* IMS Monographs, University of Chicago Press.

Havránek, T. (1984) A procedure for model search in multidimensional contingency tables. *Biometrics*, **40**,95–100.

Hodges, J.S. (1987) Uncertainty, policy analysis and statistics (with Discussion). *Statistical Science* **2**, 259–291.

Jeffreys, H. (1961) *Theory of Probability.* (3rd ed.), Oxford University Press.

Kass, R.E. and Raftery, A.E. (1993). Bayes factors and model uncertainty. Technical Report no. 254, Department of Statistics, University of Washington. Submitted for publication to *Journal of the American Statistical Association.*

Kotz, S. and Johnson, N.L. (eds.) (1985). *Encyclopedia of Statistical Sciences*, Volume 6, 578–579.

Kiiveri, H., Speed, T.P. and Carlin, J.B. (1984) Recursive causal models. *Journal of the Australian Mathematical Society (Series A)*, **36**,30–52.

Kreiner, S. (1987) Analysis of multidimensional contingency tables by exact conditional tests: techniques and strategies. *Scandinavian Journal of Statistics*, **14**,97–112.

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.

Linhart, H. and Zucchini, W. (1986) *Model Selection.* New York: Wiley.

Madigan, D. and Raftery, A.E. (1991) Model selection and accounting for model uncertainty in graphical models using Occam's window. Technical Report no. 213, Department of Statistics, University of Washington.

Miller, A.J. (1984) Selection of subsets of regression variables (with Discussion). *Journal of the Royal Statistical Society (Series A)*, **147**, 389–425.

Moulton, B.R. (1991) A Bayesian approach to regression selection and estimation with applica-

31

tion to a price index for radio services. *Journal of Econometrics*, **49**, 169–193.

Pearl, J. and Verma,T.S. (1991) A theory of inferred causation. In *Knowledge Representation and Reasoning: Proceedings of the Second International Conference*, (eds. J.A. Allen, Fikes, R. and Sandewall, E.). Morgan Kaufmann: San Mateo. 441–452.

Raftery, A.E. (1986a) A note on Bayes factors for log-linear contingency table models with vague prior information. *Journal of the Royal Statistical Society (Series B)*, **48**,249–250.

Raftery, A.E. (1986b). Choosing models for cross-classifications. *American Sociological Review* **51**, 145–146.

Raftery, A.E. (1988). Approximate Bayes factors for generalized linear models. Technical Report 121, Department of Statistics, University of Washington.

Raftery, A.E. (1993). Approximate Bayes factors and accounting for model uncertainty in generalized linear models. Technical Report no. 255, Department of Statistics, University of Washington. Submitted for publication to *Applied Statistics*.

Regal and Hook (1991) The effects of model selection on confidence intervals for the size of a closed population. *Statistics in Medicine* **10**, 717–721.

Sakamoto (1984) Efficient use of Akaike's information criterion on high dimensional contingency table analysis. *Metron* **40**, 257–275.

Spiegelhalter, D.J. and Cowell, R.G. (1991) Learning in probabilistic expert systems. Paper presented to the Fourth Valencia International Meeting on Bayesian Statistics.

Spiegelhalter, D.J. and Lauritzen, S.L. (1990a) Sequential updating of conditional probabilities on directed graphical structures. *Networks*, **20**,579–605.

Spiegelhalter, D.J. and Lauritzen, S.L. (1990b) Techniques for Bayesian analysis in expert systems. *Annals of Mathematics and Artificial Intelligence*, **2**,353–366.

Spiegelhalter, D.J. and Smith, A.F.M. (1982) Bayes factors for linear and log-linear models with vague prior information. *Journal of the Royal Statistical Society (Series B)*, **44**,377–387.

Upton, G.J.G. (1991) The exploratory analysis of survey data using log-linear models. *The Statistician*, **40**,169–182.

Wermuth, N. (1976) Model search among multiplicative models. *Biometrics*, **32**,253–263.

Wermuth, N. and Lauritzen, S.L. (1990) On substantive research hypotheses, conditional independence graphs and graphical chain models. *Journal of the Royal Statistical Society (Series B)*, **52**,21–50.

Whittaker, J. (1984) Fitting all possible decomposable and graphical models to multiway contingency tables. In *Compstat 84*, (ed. T.Havránek *et al.*). Physica-Verlag: Vienna. 401–406.

Whittaker, J. (1990) *Graphical Models in Applied Multivariate Statistics*. Chichester: Wiley.