



Building Bayesian Network Models in Medicine: The MENTOR Experience

SUBRAMANI MANI

Department of Electrical Engineering and Computer Science, University of Wisconsin-Milwaukee, Milwaukee, WI 53201, USA

MARCO VALTORTA

Department of Computer Science and Engineering, University of South Carolina, Columbia, SC 29208, USA
mgv@cse.sc.edu.

SUZANNE MCDERMOTT

Department of Family and Preventive Medicine, University of South Carolina, USA Columbia, SC 29208

Abstract. An experiment in Bayesian model building from a large medical dataset for Mental Retardation is discussed in this paper. We give a step by step description of the practical aspects of building a Bayesian Network from a dataset. We enumerate and briefly describe the tools required, address the problem of missing values in big datasets resulting from incomplete clinical findings and elaborate on our solution to the problem. We advance some reasons why imputation is a more desirable approach for model building than some other ad hoc methods suggested in literature. In our experiment, the initial Bayesian Network is learned from a dataset using a machine learning program called CB. The network structure and the conditional probabilities are then modified under the guidance of a domain expert. We present validation results for the unmodified and modified networks and give some suggestions for improvement of the model.

Keywords: Bayesian networks, machine learning, artificial intelligence in medicine

1. Introduction

A large quantity of non-experimental data is generated in Medicine from studies of the natural history of disease, case reports and epidemiological surveys¹. If experiments are well-designed, it is comparatively easy to analyze and interpret the data obtained. But, making sense of non-experimental data is a difficult task and involves a huge investment of time, effort and expertise. However, data collected for one purpose can often be used to answer other questions. Federally funded research projects make datasets available after the original study is completed. These datasets often are underutilized. This type of data is also referred to as archival data and is basically available to the investigators in “as is” condition [1]. Techniques based on Bayesian

networks hold great promise in the task of detecting associations which can be interpreted (with great caution!) as causal relationships using non-experimental data [2, 3].

We developed a model to answer the question—“What is the risk of Mental Retardation (MR) for a particular pregnancy or infant based on information from the prenatal, perinatal or postnatal period?” We do not have a diagnostic model in mind. We expect our model to quantify the risk of MR outcome, which in the early prenatal period can be used as a guideline for seeking invasive procedures such as amniocentesis for arriving at a definitive diagnosis and recommendation about the desirability of sustaining the pregnancy. During infancy the model may be used to screen children who are at greater risk

for MR to plan special educational or environmental interventions.

The prevalence of MR is estimated to be about 2.5 per cent of the population [4, 5]. When the category of borderline mental retardation is included in population estimates, over 16 percent of children have an IQ score less than 85, one standard deviation below the mean. MR is a developmental disability with a complex etiology, and the causative factors and mechanisms are not well understood. “Mental Retardation is characterized by significantly subaverage intellectual functioning” [6, p. 5]. The American Association on Mental Retardation (AAMR) quantifies the identification of people as those scoring below two Standard Deviations (SD) in a standardized IQ test [6, p. 5]. These tests are usually normalized to a mean of 100 with a SD of 15. Those with scores below 50 are classified as having severe mental retardation. Scores in the category of 50–69 fall in the classification of Mild Mental Retardation (MMR). AAMR suggests inclusion of limitation of adaptive skills for individual diagnosis [6, p. 6], but many studies have used cognitive tests (IQ scores) for classification [5, 7].

We shall use IQ scores and include the additional category of Borderline Mental Retardation (BMR, scores falling between one and two standard deviations). For severe MR a cause can be found in the majority of cases. In MMR, which forms 85% of MR, a cause cannot be identified in half the cases [4].

So here we have a complex web of unknown causal mechanisms, disagreement among experts, controversies (the large literature of nature versus nurture) and serious gaps in the experts’ understanding of the etiological factors. A Bayesian modeling approach may shed some light on the causal mechanisms, give us a tool for prediction of MR and open up avenues for early intervention—medical and social.

A companion publication in the developmental disabilities literature [8] discusses our model further from a medical perspective. In this paper, we discuss the techniques used in model building and validation from an applied artificial intelligence perspective.

2. Model Building Methodology

We refer the reader to [9, Section 5.3] for a precise and thorough definition of Bayesian network and to [10–15] for extended presentations of related concepts. We only give a sketch of the definition with a brief example.

A *Bayesian network* consists of a directed acyclic graph (DAG), prior marginal probability tables for the nodes in the DAG that have no parents, and conditional probability tables for the nodes in the DAG given their parents. The network and the probability tables define a joint probability distribution on all variables corresponding to the nodes, with the defining property that the conditional probability of any variable v given any set of variables that includes only the parents of v and any subset of nodes that are not descendant of v is equal to the conditional probability of v given only its parents. From this property, it follows that the joint probability of the variables in a Bayesian network decomposes in a multiplicative fashion; more precisely, if V is the set of the nodes in the DAG, the following equality (the *chain rule for Bayesian networks*) holds: $P(V) = \prod_{v \in V} P(v \mid \text{parents}(v))$. In turn, this decomposition allows for very efficient computation of marginal posterior probabilities upon observation of evidence.

As an example, the graph in Fig. 1 models a small portion of the mental retardation domain. We do not claim that this model is accurate or sensible: it is provided only for the sake of illustration. At the depth of understanding required for the example, the names of the variables should be considered self-explanatory. Recall that a Bayesian network is composed of two parts: an acyclic directed graph and the numerical specification of conditional and prior probability tables. Three features of Bayesian networks are worth mentioning.

First, the directed graph constrains the possible joint probability distributions represented by a Bayesian network. For example, in any distribution consistent with the graph of Fig. 1, *Chld_Ravn* (the IQ score of the child) is conditionally independent of *Fam_Inc* (Family

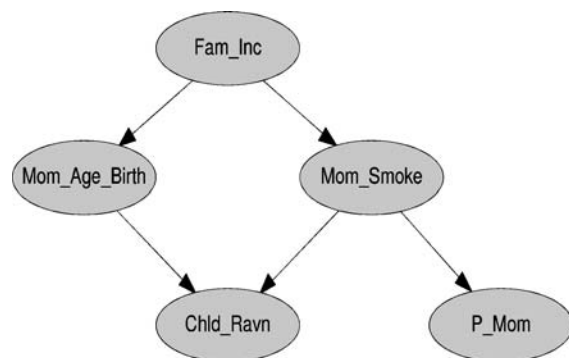


Figure 1. A microscopic model of MR.

Table 1. Values of the five micro-mentor variables.

<i>Fam_Inc</i>	$\geq 10000, < 10000$
<i>Mom_Age_Birth</i>	14–19+, 20–34, ≥ 35
<i>Mom_Smoke</i>	yes, no
<i>Child_Ravn</i>	mild, border, normal, super
<i>P_Mom</i>	mild, border, normal, super

Income) given *Mom_Age_Birth* (the age of the mother at birth) and *Mom_Smoke* (whether the mother smokes); also, *P_Mom* (the IQ score of the mother) is conditionally independent of any subset of the other variables given *Mom_Smoke*.

Secondly, the explicit representation of constraints about conditional independence allows a substantial reduction in the number of parameters to be estimated. In the example, assume that the possible values of the five variables are given in Table 1. Then, the joint probability table $P(\text{Fam_Inc}, \text{Mom_Age_Birth}, \text{Mom_Smoke}, \text{Child_Ravn}, \text{P_Mom})$ has $2 \times 3 \times 2 \times 4 \times 4 = 192$ entries. It would be very difficult to assess 191 independent parameters.² However, the independence constraints encoded in the graph permit the factorization

$$\begin{aligned}
 P(\text{Fam_Inc}, \text{Mom_Age_Birth}, \text{Mom_Smoke}, \text{Child_Ravn}, \\
 \text{P_Mom}) = & P(\text{Fam_Inc}) \times P(\text{Mom_Age_Birth} \mid \\
 & \text{Fam_Inc}) \times P(\text{Mom_Smoke} \mid \text{Fam_Inc}) \\
 & \times P(\text{Child_Ravn} \mid \text{Mom_Age_Birth}, \text{Mom_Smoke}) \\
 & \times P(\text{P_Mom} \mid \text{Mom_Smoke}),
 \end{aligned}$$

which reduces the number of parameters to be estimated to $1 + 4 + 2 + 18 + 6 = 31$. The second term in the product corresponds to the conditional probability table for *Mom_Age_Birth* given *Fam_Inc*, which is given in Table 2; note that there are only four independent parameters to be estimated, since the sum of values by column is one. Again, we emphasize that these numbers are fictitious.

Table 2. Conditional probability tables for *Mom_Age_Birth* Given *Fam_Inc*.

	<i>Fam_Inc</i>	
	≥ 10000	< 10000
14–19	0.1	0.3
20–34	0.7	0.6
≥ 35	0.2	0.1

Thirdly, the Bayesian network representation allows a substantial (usually, dramatic) reduction in the time needed to compute marginals for each variable in the domain. The explicit representation of constraints on independence relations is exploited to avoid the computation of the full joint probability table in the computation of marginals, both prior and conditioned on observations. Space prevents the description of the relevant algorithms. See, e.g., [15, Ch.5] for a discussion of the junction tree algorithm, the most widely used one.

There are two methods of building a Bayesian network for a particular application domain. The first method consists of asking the domain expert to construct the network (DAG) and assign the prior marginal probabilities for nodes without parents and the conditional probabilities for the other nodes. The second method consists in building the network from data. There are several algorithms available to accomplish this learning task—for example, BIFROST [16], K2 [17] and CB [18, 19]. The marginal and conditional probabilities can also be computed from data. The models are validated by comparing with the performance of an expert [12]. We use a combination of the two strategies—capture the skeleton network from data using the CB algorithm and refine the DAG with the help of the expert and published literature. Prior and conditional probabilities are obtained from data and fine-tuned by the expert.

3. Datasets Used in Model Construction

We obtained the Child Health and Development Study (CHDS) data set, which was developed in a study concerning pregnant mothers and their children [20]. The children were followed through their teen years and included numerous questionnaires, physical and psychological exams, and special tests. The study was conducted by the University of California at Berkeley and the Kaiser Foundation. It started in 1959 and continued into the 1980's. There are approximately 6000 children and 3000 mothers with IQ scores in the data set. The children were either 5 years old or 9 years old when their IQs were tested. The IQ test used for the children was the Raven Progressive Matrices Test. The mothers' IQs were also tested, and the test used was the Peabody Picture Vocabulary Test.

We identified about 50 variables scattered among several CHDS files that are thought to play a role in the causal mechanism of MR. Under the guidance of the domain expert this set of fifty variables was reduced to

Table 3. The variables used in MENTOR.

Variable	What the variable represents
MOM_RACE	Mother's race classified as White (European or White and American Indian or others considered to be of white stock) or non-White (Mexican, Black, Oriental, interracial mixture, South-East Asian).
MOMAGE_BR	Mother's age at time of child's birth categorized as 14–19 years, 20–34 years, or ≥ 35 years.
MOM_EDU	Mother's education categorized as ≤ 12 and did not graduate, high school, graduated high school, and $>$ high school (attended college or trade school).
DAD_EDU	Father's education categorized same as mother's.
MOM_DIS	Yes if mother had one or more of lung trouble, heart trouble, high blood pressure, kidney trouble, convulsions, diabetes, thyroid trouble, anemia, tumors, bacterial disease, measles, chicken pox, herpes simplex, eclampsia, placenta previa, any type of epilepsy, or malnutrition; no otherwise.
FAM_INC	Family income categorized as $< \$10,000$ or $\geq \$10,000$.
MOM_SMOK	Yes if mother smoked during pregnancy; no otherwise.
MOM_ALC	Mother's alcoholic drinking level classified as mild (0–6 drinks per week), moderate (7–20), or severe > 20 .
PREV_STILL	Yes if mother previously had a stillbirth; no otherwise.
PN_CARE	Yes if mother had prenatal care; no otherwise.
MOM_XRAY	Yes if mother had been X-rayed in the year prior to or during the pregnancy; no otherwise.
GESTATN	Period of gestation categorized as premature (≤ 258 days), or normal (259–294 days), or postmature (≥ 295 days)..
FET_DIST	Fetal distress classified as yes if there was prolapse of cord, mother had a history of uterine surgery, there was uterine rupture or fever at or just before delivery, or there was an abnormal fetal heart rate; no otherwise.
INDUCE_LAB	Yes if mother had induced labor; no otherwise.
C_SECTION	Yes if delivery was a caesarean section; no if it was vaginal.
CHLD_GEND	Gender of child (male or female).
BIRTH_WT	Birth weight categorized as low < 2500 g) or normal (≥ 2500 g).
RESUSCITN	Yes if child had resuscitation; no otherwise.
HEAD_CIRC	Normal if head circumference is 20 or 21; abnormal otherwise.
CHLD_ANOM	Child anomaly classified as yes if child has cerebral palsy, hypothyroidism, spina bifida, Down's syndrome, chromosomal abnormality, anencephaly, hydrocephalus, Turner's syndrome, cerebellar ataxia, speech defect, Klinefelter's syndrome, or convulsions; no otherwise.
CHLD_HPRB	Child's health problem categorized as having a physical problem, having a behavior problem, having both a physical and a behavioral problem, or having no problem.
CHLD_RAVN	Child's cognitive level, measured by the Raven test, categorized as mild MR, borderline MR, normal, or superior.
P_MOM	Mother's cognitive level, measured by the Peabody test, categorized as mild MR, borderline MR, normal, or superior.

a set of twenty-three, resulting in the datasets described in 3.1. The subject expert thought that this set of variables was sufficient to capture the domain knowledge. Only one child of the mother is included in each of the datasets. Table 3 contains a list of the twenty-three variables used in the final Bayesian network. (The files used in network construction include a twenty-fourth variable, MAR_STAT, indicating marital status of the mother, which was removed at a late stage.)

3.1. Datasets Used for Network Construction

RAVN6X24. This dataset contains 5985 cases and 24 variables. In this dataset many of the IQ scores of

mothers are missing. The percentage of missing values is 12. This dataset is the total relevant dataset **RAVN2X24**. This dataset contains 2212 cases and 24 variables. The IQ scores of mothers and children are present. There are no missing values for the IQ scores. This is a subset of the RAVN6X24 dataset, with all the rows which did not have IQ scores for the mother and child removed. The percentage of missing values is 4.

RAVN6X23. This dataset contains 5985 cases and 23 variables. As only about 3000 mothers were given IQ tests, this dataset was created without the variable P.MOM (IQ score of the mother). This is also a subset of the RAVN6X24 dataset with the variable

mother's IQ deleted. The percentage of missing values is 10.

All three datasets were used for network construction, as explained in Section 6.1.

4. Tools for Model Building

The *CB algorithm* takes as input a dataset with no missing values and outputs a Bayesian network structure. The network structure, when augmented with suitable conditional probability tables constitutes a Bayesian network, as defined in Section 2 that models the data, in the sense that the data can be taken to be a sample of the distribution encoded by the network. Moreover the network structure output by CB has usually only a few edges, because it exploits independence relations among variables well. The network is therefore appropriate for use by inference algorithms and for visual inspection.

The CB algorithm works in two phases. In the first phase, CB uses Conditional Independence tests (χ^2 tests) for ordering the nodes. In the second phase, which is based on the K2 algorithm [17], CB computes greedily an approximation to the most likely network structure given the dataset [19]. Given a dataset and network, **CondProb** computes the prior marginal and conditional probabilities using the formulas in [17]. It has been observed that constraint-based algorithm for learning Bayesian network structures cannot orient many edges. Some authors are highly skeptical of any attempt to distinguish between Bayesian network structures that encode the same conditional independence information. It is important to note that CB is a hybrid structure learning algorithm that uses both a constraint-based approach and a scoring approach, and that CB does not assign the same score to equivalent networks. An implementation of CB with a user-friendly graphical user interface is available by contacting the corresponding author.³

HUGIN provides a graphical interface for representing the nodes (domain variables) and the directed edges (usually interpretable as causal relationships between the variables). A user-friendly mechanism for naming the variables, entering the states of the variables and assigning the conditional probabilities is also provided. *HUGIN* implements the Lauritzen and Spiegelhalter method of probability propagation in DAGs [22], with some improvements. The *HUGIN* shell was developed

by Andersen, Olesen, F.V. Jensen and F. Jensen in Denmark [23].

The *IMP* program analyzes the given dataset and predicts missing values. We use statistical, case matching and randomization methods. A random guess is attempted when case matching fails. The method is expected to succeed in domains where there is good interdependency between variables. Fortunately most real world data and medical data in particular have many interdependent variables. We have not analyzed the theoretical properties of *IMP*, but we consider it to be a practical and useful method particularly for purposes of model building.

CAP-CPN is an application written in C to call Bayesian Networks using *HUGIN*.⁴ It provides modules to use the *HUGIN-API* C library in an organized way. *CAP-CPN* converts an ASCII dataset to the format required by *HUGIN* for batch validation. It also provides functions to perform simple statistical tests on the data gathered by sampling the outcome node when a batch file containing cases is processed.

5. Handling Missing Values

Real world data contains missing values. This is particularly true of medical datasets. The general practice in the analysis of missing data is deleting cases (records) with missing data. But when there are numerous variables such a policy can mean that most records will have to be disregarded from analysis or many variables will have to be sacrificed. It will help if we can come up with a scheme to predict and assign missing values. To start with, this strategy will be very useful for model building and validation from datasets. We do not discuss the merits and demerits of imputing for data analysis here.

We decided against the easy way of making a separate category for the missing values, as done in the original *MUNIN* system [24]. We believe that it is not a satisfactory procedure as in most cases it is hard to trace a causal pathway between the missing category of one variable and the missing category of another variable. Treating missing value as a separate category is also likely to create serious problems in computing the conditional probabilities from data. For example, assigning the conditional probability of a variable with 4 states which has 5 parents having 2 states each results in a table of 128 entries. Now if a missing category is included, the table space grows to 1215 entries. And for this example (which is by no means an extreme case) we have more than 1000 junk entries. Not only is

the size of the conditional probability table a problem, but we also encounter semantic difficulties computing conditional probabilities for the occurrences including missing states. Hence it is desirable to come up with a scheme to avoid missing categories. Then in the quantitative modeling stage only the valid categories of the variables in the network and the conditional probabilities will have to be entered. Another method which has been used [16, p. 94] is to assign one of the valid categories to all the missing values of a particular variable. It may be suitable for variables where the domain expert can predict with a high probability which category the missing values should have.

We developed and implemented an algorithm (IMP) for predicting and imputing missing values. Before describing it, we introduce a piece of notation. Let \mathcal{V}_j be the set of cases that have the same known⁵ values as the known values of case c_j .

The basic step of the IMP algorithm is to assign to the missing value x_i in case c_j the mode of x_i in the set of cases \mathcal{V}_j . If there are no cases with known values in set \mathcal{V}_j , then IMP drops some of the variables from the cases. Since it would be too computationally expensive to consider all subsets of variables, IMP only considers subsets in which a constant (k) number of variables are dropped and averages the value of the mode of x_i by a weight proportional to the number of cases corresponding to each subset of the variables. Finally, if there is no case in D in which x_i has a known value, even after dropping k variables, x_i is assigned a value at random in case c_j .

The accuracy of IMP can be validated using datasets across domains. Datasets without any missing values were used for validation. By random number generation a fixed percentage (say ten percent) of data values are assigned missing, thus obtaining a dataset on which IMP is run to impute missing values. The output dataset is compared with the original dataset. Our validation tests using LED, ALARM and SOYBEAN which are small to large artificial datasets used for Machine Learning research and available from the University of California at the Irvine Machine Learning repository [25] gave a mean accuracy of 80% over ten runs. The range was from 67% to 95%.

Another possibility is to impute a dataset using the algorithm. This imputed dataset has no missing values. Now we assign missing values (we can assign the same percentage of missing values originally present) generating random numbers, imputing and comparing with the dataset we created originally by imputing.

This technique, called *customized validation*, gives the predictive accuracy for the particular dataset in question with its given percentage of missing values. Even though this takes into account the size and other peculiarities of the dataset for validation purposes, it may introduce a small error for the estimate as we are using IMP twice for validation.

An alternative to IMP for inferring missing values is expectation maximization learning (the EM algorithm). In its basic form, EM finds a maximum likelihood estimate for the parameters, rather than the most probable (a posteriori) values. Some approaches to Bayesian network structure learning in the presence of missing data bypass the imputation of missing values for incomplete data and are presented in [26]. A particularly interesting technique is Friedman's Structural EM Algorithm [26, 27]. In this paper, we concentrate on our techniques. A detailed comparison of the quality and complexity of these techniques would be a good topic for future work.

Our datasets were imputed using IMP. For our datasets RAVN2X24, RAVN6X23, and RAVN6X24, we obtained an accuracy of 79%, 82% and 83% respectively. The accuracy of the imputed values were judged by the technique of customized validation.

6. Network Generation and Refinement

6.1. Network Generation

The CB algorithm was run on the three imputed datasets described in Section 3.1 for generating the networks. The datasets were randomly partitioned into two—a major part and a minor part. The bigger partition was used for constructing the network and the smaller part was set apart for validation. For RAVN2X24, we used the first 2000 cases for generating the network and for the other two, the first 5000. The network generated from the RAVN6X24 dataset is shown in Fig. 2. The networks obtained are given in Tables 4–6.

6.2. Network Refinement

We defined three rules to characterize the inadequacies of the generated networks.

6.2.1. Rule of Chronology. Events occurring later in time cannot be the parents of earlier incidents. For instance a child health problem cannot be the parent of maternal disease.

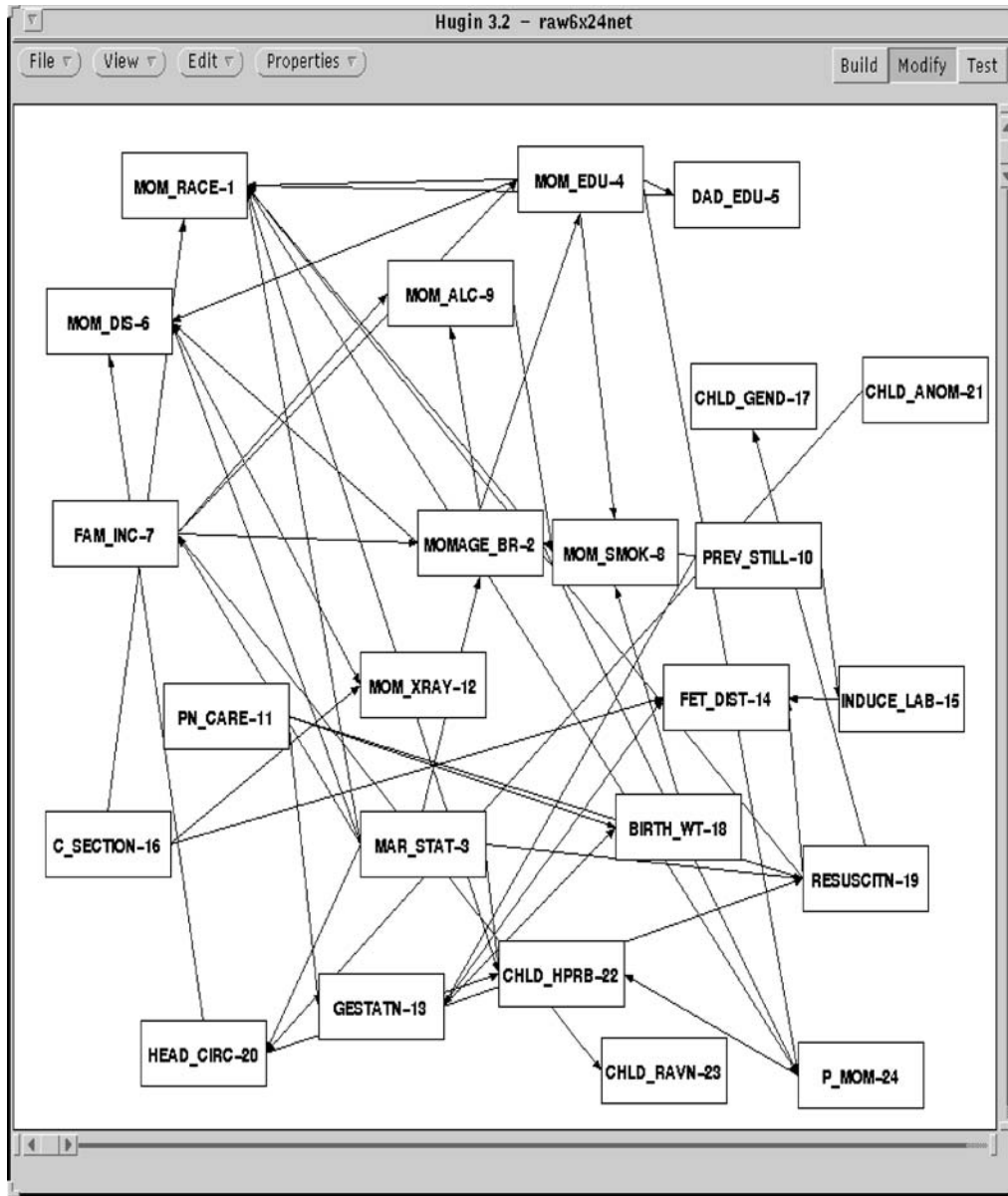


Figure 2. Network generated by CB from RAVN6X24 dataset.

6.2.2. Rule of Common Sense. The directed edges of the network should not go against common sense. For instance, Father’s education cannot be a cause of Mother’s race.

6.2.3. Domain Rule. This rule has been referred as the *Rule of Biological Plausibility* in the medical and biological science literature. This rule states that a causal explanation is tenable in terms of existing knowledge

level about the variables involved. This level is what we obtain from an intelligent review of the relevant literature. The directed edges should not violate established domain rules. For example, pre-natal care cannot be put down as a cause of maternal smoking. Mausner and Kramer strike a note of caution here: “The development of biological knowledge often introduces new factors that previous studies have not taken into account. In the existing studies, the major causal factors may have

Table 4. RAVN2X24 network.

Variable	Parents
MOM_RACE	
MOM_AGE_AT_BIRTH	MAR_STAT, MOM_EDU, FAM_INC, PREV_STILLBRTH
MAR_STAT	
MOM_EDU	MOM_RACE
DAD_EDU	MOM_EDU ² , P_MOM ¹
MOM_DIS	MOM_AGE_AT_BIRTH
FAM_INC	MOM_EDU
MOM_SMOKE	MOM_RACE, MOM_EDU, MOM_ALC, PN_CARE ³
MOM_ALC	FAM_INC
PREV_STILLBRTH	
PN_CARE	
MOM_XRAY	MOM_DIS, C_SECTION
GESTATN	MOM_RACE, FET_DIST
FET_DIST	INDUCE_LAB ¹ , C_SECTION ¹ , RESUSCITN ¹
INDUCE_LAB	
C_SECTION	
CHLD_GEND	CHLD_HPROB ¹
BIRTH_WT	PN_CARE, GESTATN
RESUSCITN	MOM_RACE ²
HEAD_CIRC	MAR_STAT, INDUCE_LAB ³ , CHLD_ANOM
CHLD_ANOM	
CHLD_HPROB	CHLD_ANOM
CHLD_RAVN	MOM_EDU, CHLD_ANOM
P_MOM	MOM_RACE, MOM_EDU

¹Violates law of chronology.²Goes against commonsense.³Violates domain rules.

been missed because their importance was not appreciated” [28, p. 187]. This point is well taken and if there is a strong case, such a directed edge should be investigated further. But for our network construction purposes, if an edge clearly violated established domain constraints, it was removed. The directed edges of the network in Fig. 2 are given in Table 6 with annotations describing examples of rules that are broken. So also new edges were incorporated to capture the knowledge of the known domain causal mechanisms. The variable MAR_STAT was removed as the expert felt that it was not playing a useful role in representing domain relations. See Tables 4–6 for examples of rules that are broken. The expert refined network is a synthesis

Table 5. RAVN6X23 network.

Variable	Parents
MOM_RACE	MAR_STAT ¹ , MOM_EDU ¹ , DAD_EDU ² , MOM_SMOKE ¹ , MOM_ALC ¹ , C_SECTION ¹
MOM_AGE_AT_BIRTH	MAR_STAT, FAM_INC, PREV_STILLBRTH
MAR_STAT	
MOM_EDU	MOM_AGE_AT_BIRTH ² , FAM_INC
DAD_EDU	
MOM_DIS	MOM_RACE, MOM_AGE_AT_BIRTH, MAR_STAT ²
FAM_INC	MAR_STAT
MOM_SMOKE	MOM_EDU, MOM_ALC, BIRTH_WT
MOM_ALC	FAM_INC, MOM_AGE_AT_BIRTH
PREV_STILLBRTH	
PN_CARE	
MOM_XRAY	MOM_RACE ² , MOM_DIS, C_SECTION
GESTATN	MOM_AGE_AT_BIRTH, PREV_STILLBRTH, PN_CARE
FET_DIST	PN_CARE, GESTATN, INDUCE_LAB ¹ , C_SECTION ¹ ,
INDUCE_LAB	PREV_STILLBRTH
C_SECTION	
CHLD_GEND	CHLD_HPROB ¹
BIRTH_WT	PN_CARE, GESTATN
RESUSCITN	MOM_RACE ² , MAR_STAT ² , FET_DIST
HEAD_CIRC	MAR_STAT, CHLD_ANOM
CHLD_ANOM	
CHLD_HPROB	MAR_STAT, HEAD_CIRC
CHLD_RAVN	FAM_INC

¹Violates law of chronology.²Goes against commonsense.³Violates domain rules.

and refinement of the three raw networks. The expert-modified network is shown in Fig. 3. The changes made to the networks are significant. This is because we have encoded the causal relations, which made the networks sparser. As a consequence, the expert-modified networks generalize better to new cases, as the results of Section 7 show.

6.3. Refinement of Conditional Probabilities

The prior marginal and conditional probabilities were computed using the program CondProb. For nodes

Table 6. RAVN6X24 network.

Variable	Parents
MOM_RACE	MAR_STAT ¹ , MOM_EDU ¹ , DAD_EDU ² , MOM_SMOKE ¹ , C_SECTION ¹ , RESUSCITN ¹
MOM_AGE_AT_BIRTH	MAR_STAT, FAM_INC, PREV_STILLBRTH
MAR_STAT	
MOM_EDU	MOM_AGE_AT_BIRTH ² , FAM_INC
DAD_EDU	
MOM_DIS	MOM_AGE_AT_BIRTH, MAR_STAT ² , MOM_EDU, HEAD_CIRC ²
FAM_INC	MAR_STAT
MOM_SMOKE	MOM_EDU, MOM_ALC, BIRTH_WT
MOM_ALC	FAM_INC, MOM_AGE_AT_BIRTH
PREV_STILLBRTH	
PN_CARE	
MOM_XRAY	MOM_DIS, C_SECTION
GESTATN	MOM_AGE_AT_BIRTH, PREV_STILLBRTH, PN_CARE
FET_DIST	GESTATN, INDUCE_LAB ¹ , C_SECTION ¹ , RESUSCITN ¹
INDUCE_LAB	PREV_STILLBRTH
C_SECTION	
CHLD_GEND	RESUSCITN ²
BIRTH_WT	PN_CARE, GESTATN
RESUSCITN	MAR_STAT ² , PN_CARE, GESTATN
HEAD_CIRC	MAR_STAT, CHLD_ANOM
CHLD_ANOM	
CHLD_HPROB	MOM_RACE, MAR_STAT, HEAD_CIRC, P_MOM
CHLD_RAVN	FAM_INC
P_MOM	MOM_RACE, MOM_AGE_AT_BIRTH, MOM_EDU

¹Violates law of chronology.

²Goes against commonsense.

³Violates domain rules.

without parents prior marginal probabilities of the various states calculated from the RAVN6X24 dataset were assigned.

For the nodes with one or more parents, the conditional probabilities calculated using the same dataset was assigned. The Conditional Probabilities of the outcome variable CHLD_RAVN (See paragraph 3) were refined by the expert. There were many pos-

sible instantiations that were not represented in the dataset RAVN6X24. A reasonable conditional probability was assigned by the expert for these. For the raw networks probabilities were assigned from the RAVN2X24 dataset using the program CondProb.

7. Validation of the Model

7.1. Validation by the Expert

As ours is a model for risk assessment and risk prediction of mental retardation, it is different from a classification or diagnostic problem. In a typical diagnostic approach we consider a set of differential diagnoses and the attempt is to assign probabilities to them and order them on that basis. In risk assessment we are interested in the change in magnitude of a particular category of interest even though it may still occupy a low position in an ordering of the variable levels. We have a prior probability of 5.6% for mild and 12.4% for borderline MR. Hence if the risk of both mild and borderline doubles, still we get a combined probability of only 36%. That leaves a probability of 64% for normal and superior. We would consider intervention for a case like this, because the cost of intervention is outweighed by the potential benefit accruing from it. Most of the actual cases from the dataset with mild or borderline MR give a >50% probability for normal outcome. This is because there are more normal outcome cases with similar instantiations of variables than outcomes that result in mental retardation. Hence we decided first on a strategy of validation by comparing with the expert. We generated nine cases with instantiation for a subset of variables. We ran these cases on the model and computed the probabilities. The expert was asked to score the results as *agree* or *disagree*. The expert was in agreement with the model's assessment in eight out of nine cases used for validation. Three of the cases are depicted in Table 7, while the conditional probabilities of the values of CHLD_RAVN for those cases are shown in Table 8.

7.2. Validation using RAVN2X24

7.2.1. Risk Means of Cases and Controls. All the cases in the dataset RAVN2X24 (unimputed) were run through the models—the expert network, and the two raw networks that have twenty-four variables (raw2x24net and raw6x24net) using CAP-CPN. The

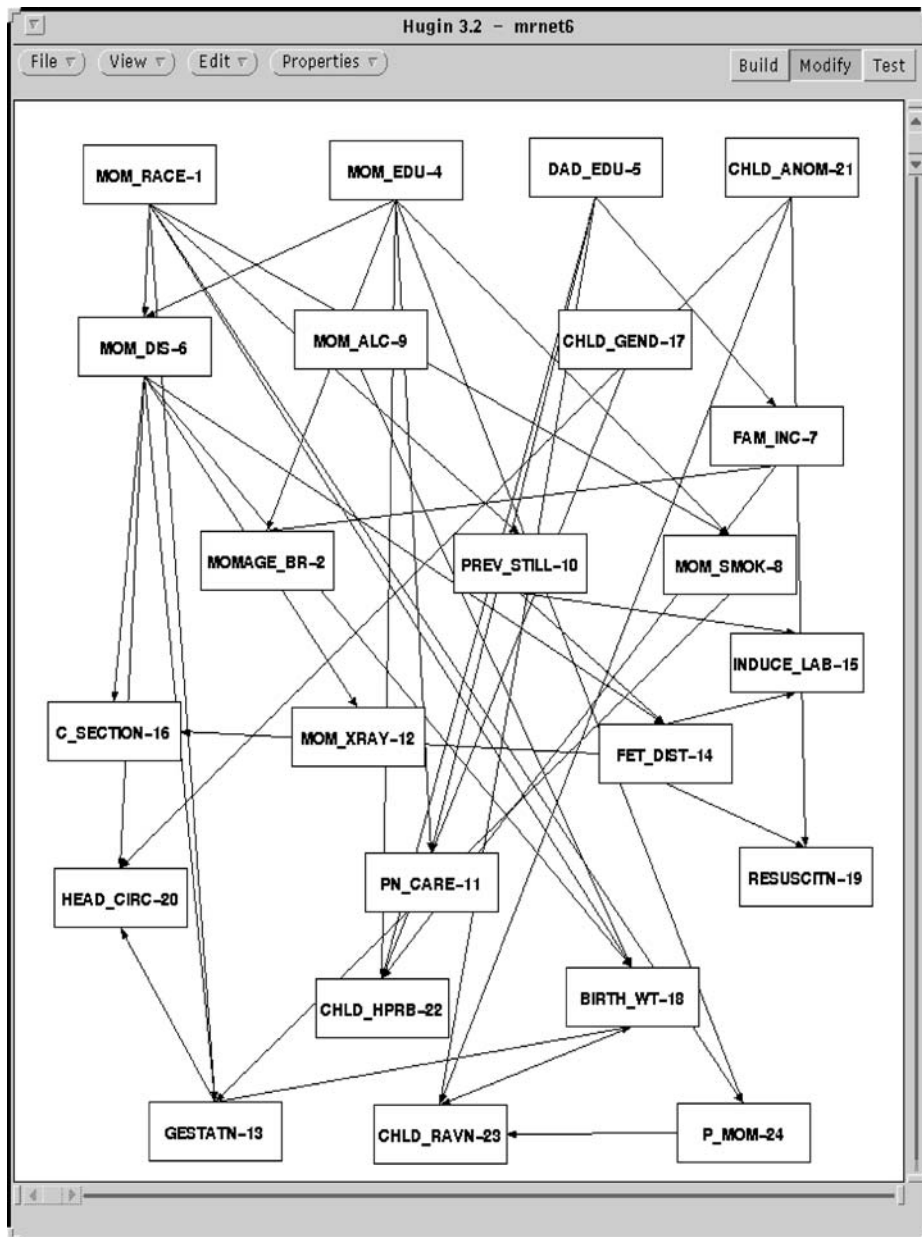


Figure 3. Network modified by the expert.

results showing the relative risk for controls (children with normal outcome) and cases (children with mild or borderline MR) are given in Table 9 for the three nets.

A *t*-test procedure was performed to assess the statistical significance of the predicted risks. The Prob $> |T|$ was 0.0001, 0.0000 and 0.1878 for the expert net, raw2x24net and raw6x24net respectively. (We do not have a good explanation for the fact that the sec-

ond number is smaller than the first; since the first two numbers are both very small, we dismiss this as an aberration.) This shows that there is significant difference in the mean risk scores between the cases and controls for the expert net and raw2x24net ($P < 0.05$). Note that there are fewer violations of the rules described in Section 6.2 for raw2x24net (9) than for raw6x24net (14). There are several possible explanations for the

Table 7. Generated values for three cases.

Variable	Case 1 Variable value	Case 2 Variable value	Case 3 Variable value
MOM_RACE	Non-White	White	White
MOMAGE_BR	14-19		≥ 35
MOM_EDU	≤12	>High school	≤12
DAD_EDU	≤12	>High school	High school
MOM_DIS			No
FAM_INC	<10,000		<10,000
MOM_SMOK			Yes
MOM_ALC			Moderate
PREV_STILL			
PN_CARE		Yes	
MOM_XRAY			Yes
GESTATN	Normal	Normal	Premature
FET_DIST		No	Yes
INDUCE_LAB			
C_SECTION			
CHLD_GEND			
BIRTH_WT	Low	Normal	Low
RESUSCITN			
HEAD_CIRC			Abnormal
CHLD_ANOM		No	
CHILD_HPRB			Both
CHLD_RAVN			
P_MOM	Normal	Superior	Borderline

Table 8. Posterior probabilities for three cases.

Value of CHLD_RAVN and Prior probability	Case 1 Posterior probability	Case 2 Posterior probability	Case 3 Posterior probability
Mild MR (.056)	.101	.010	.200
Borderline MR (.124)	.300	.040	.400
Normal (.731)	.559	.690	.380
Superior (.089)	.040	.260	.200

poor performance of the raw2x24 dataset. It might be that the data or, more precisely, the missing data, are not distributed evenly across the domain. Another (related) possibility is poor performance of IMP. The performance of the expert net was the best of the three based on the mean risks.

7.2.2. Evaluation Using a Risk Threshold. In the initialized state, the expert network gives a resting MR risk of 0.18 if the risks for mild and borderline retardation are added together, as shown in Fig. 4. (Recall from the Introduction that the prevalence of mental retardation, mild mental retardation, and borderline mental retardation is much larger than that of mental retardation alone, which is 0.025.)

If we take twice the resting state risk as our threshold for significant risk, our threshold can be set at the value of 0.36. Using this threshold we find that twenty nine per cent of cases are flagged correctly. This also results in eighteen per cent of controls being flagged as significant risk for MR. (Fig. 5 shows the increase in risk of MR for an example case with some known risk factors.) These results are contained in Table 10, which also presents the same type of results for raw2x24net, whose resting MR risk is 0.16. (The results for raw6x24 are very poor.)

We considered using n -fold cross-validation to estimate the error rates of our methods, but we did not employ it, because it can be applied only to the raw networks. In n -fold cross-validation, the cases are divided into n (typically, ten) groups of roughly equal size. All except one of the groups are used in learning, while the group that is left out is used to estimate an error rate. This process is carried out n times, each time leaving out a different group, and the overall error rate for the learning algorithm is the average of the n error rates [29]. Since the expert modified network was crafted from the three raw networks by removing and adding edges to reflect domain knowledge and the conditional

Table 9. Mean risk of MR predicted for cases and controls by the three nets.

Level	Expert net		raw2x24net		raw6x24net	
	Controls $n = 1863$	Cases $n = 349$	Controls $n = 1863$	Cases $n = 349$	Controls $n = 1863$	Cases $n = 349$
	Mean risk		Mean risk		Mean risk	
Mild MR	0.05	0.07	0.02	0.02	0.02	0.02
Borderline MR	0.11	0.14	0.14	0.16	0.14	0.15
Mild + Border	0.16	0.21	0.16	0.18	0.16	0.17

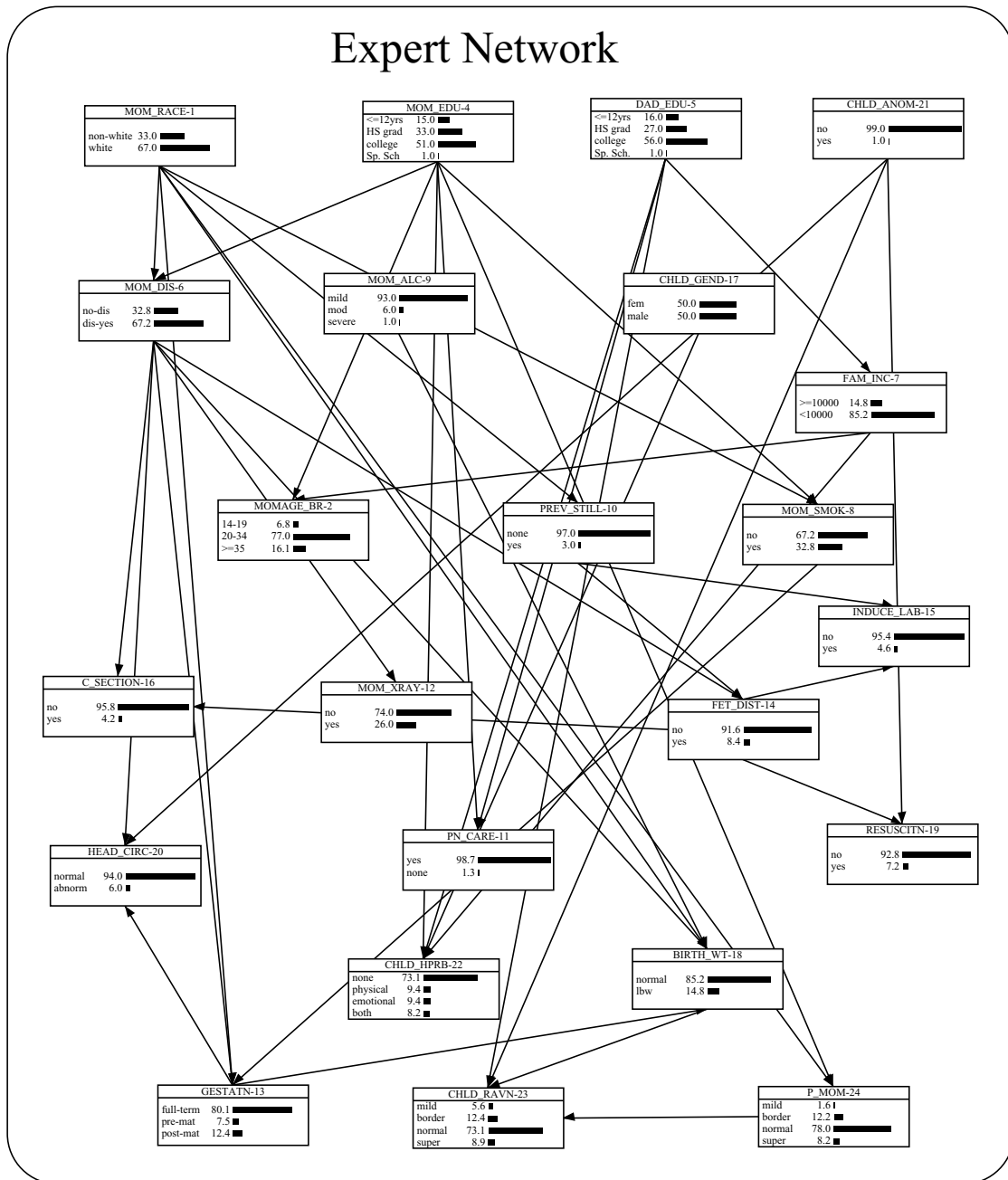


Figure 4. Initial probabilities in the expert network.

probabilities were also modified by the expert, cross-validation cannot be used to compare the performance of the raw and refined networks.

7.2.3. Validation Using a Separate Data Set. The National Collaborative Perinatal Project (NCP), of

the National Institute of Neurological and Communicative Disorders and Strokes, developed a data set containing information on pregnancies between 1959 and 1974 and 8 years of follow-up for live-born children. For each case in the data set, the values of all 22 variables except CHLD_RAVN (child's cognitive level

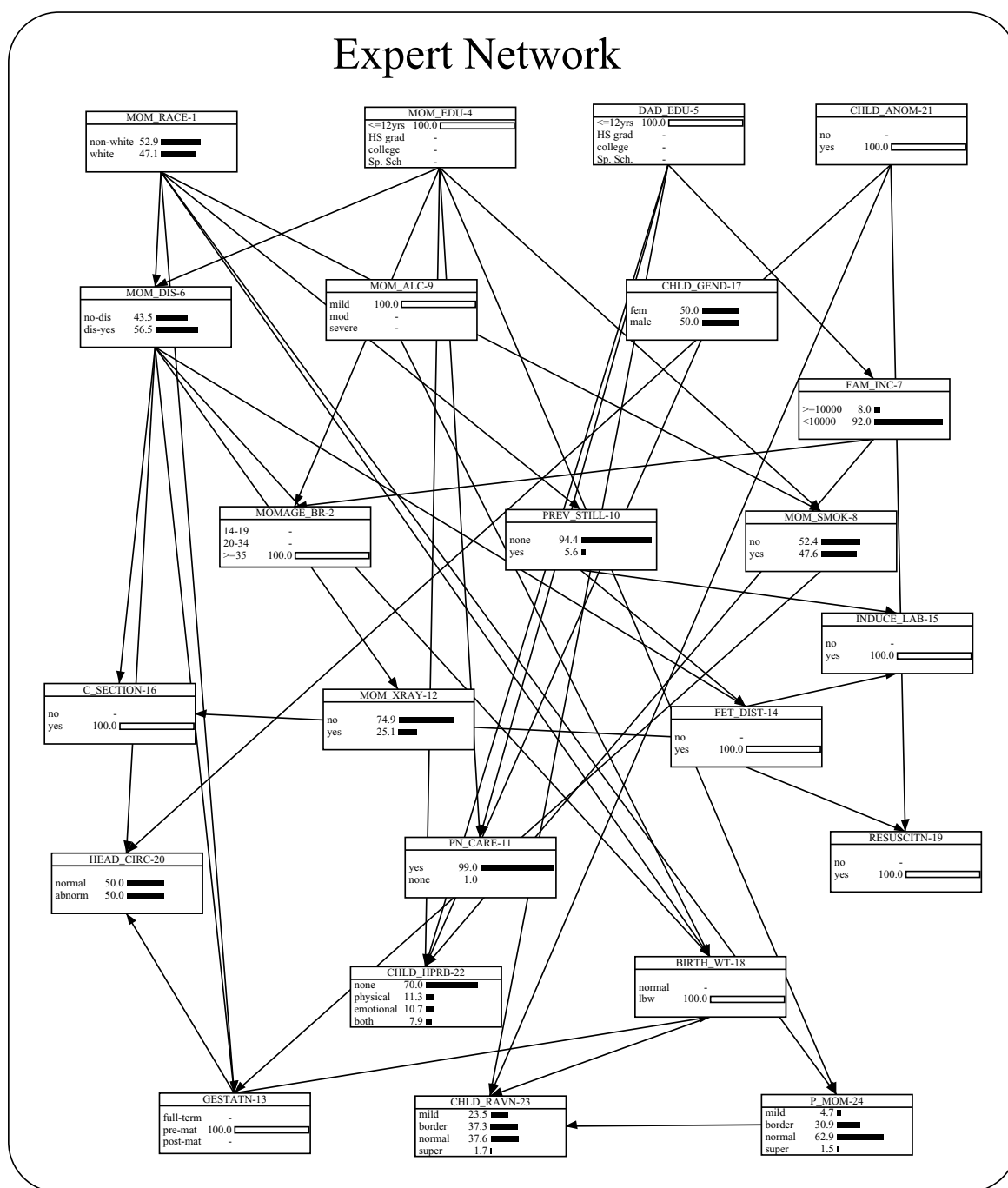


Figure 5. Probabilities in the expert network for a high risk case.

as measured by the Raven test) were entered, and the conditional probabilities of each of the four values of CHLD_RAVN were computed. Table 11 shows the average values of $P(\text{CHLD_RAVN} = \text{mildMR} \mid \mathbf{d})$

and $P(\text{CHLD_RAVN} = \text{borderlineMR} \mid \mathbf{d})$, where \mathbf{d} is the set of values of the other 22 variables, for both the controls (children in the study with normal cognitive function at age 8) and the subjects

Table 10. Cases flagged for different risk thresholds.

Risk Threshold for MR (Mild + border)	Expert net		raw2x24net	
	Controls <i>n</i> = 1863	Cases <i>n</i> = 349	Controls <i>n</i> = 1863	Cases <i>n</i> = 349
Resting value	434 (23%)	122 (35%)	901 (48%)	240 (69%)
1.5 × Resting value	370 (20%)	111 (32%)	242 (13%)	73 (21%)
2 × Resting value	342 (18%)	101 (29%)	9 (<1%)	5 (1%)

Table 11. Average probabilities, as determined by MENTOR, of having mental retardation for controls (children identified as having normal cognitive functioning at age 8) and subjects (children identified as having mild or borderline MR at age 8).

Cognitive level	Avg. probability for controls (<i>n</i> = 13019)	Avg. probability for subjects (<i>n</i> = 3598)
Mild MR	.06	.09
Borderline MR	.12	.16
Mild or Borderline MR	.18	.25

(children in the study with mild or borderline MR at age 8).

8. Discussion and Future Work

8.1. Discussion

The validation results are significant but not dramatic. We feel that this is due to the incomplete state of knowledge of the etiological factors of MR. This results in datasets where some of the relevant variables (not yet recognized as contributory or causative) have not been collected. Hence our model is constrained by the state of domain knowledge existing at this point in time.

Throughout the development of MENTOR, we emphasized the causal interpretation of the links. While this is not in any way necessary, it seemed to be a good decision for two reasons. First, there is widespread belief that ordering variables in a causal direction simplifies modeling, or, as Russell and Norvig put it, “If we stick to a causal model, we end up specifying fewer numbers, and the numbers will often be easier to come up with” [30, p. 443]. Second, it is easier to involve the expert in validating the edges if the model is causal.

Still, there is a serious problem with a causal interpretation in the case of MENTOR. It is quite likely that there are many hidden (unmeasured or even unknown)

variables playing a role in the causal pathway of MR. Despite this, we attempted to build a DAG model and assign to it a causal interpretation. This is clearly suspect. There are two possible approaches to dealing with this problem. The first is to attempt to discover hidden variables using a data analysis technique, such as TETRAD [31, Ch. 2] that purports to discover such variables. A second approach is to explicitly model correlations that have no causal interpretation by using undirected links as in chain graphs (graphical models that include both directed and undirected links). It is commonly (and somewhat simplistically) believed that the undirected links can be used to model associative, non-causally interpretable information, while the directed links are used to model the causally interpretable information [32]. For our research we did not address this issue further.

8.2. Future Work

The networks generated from the different datasets using the CB algorithm had many nodes that violated the *rule of chronology*. A facility for inputting the chronological order can be incorporated. Likewise, if some rules could be incorporated in the network generation stage to take care of domain-specific constraints, directed edges violating *domain rules* would be avoided. In other words, the ordering of nodes built by the first phase of the CB algorithm would be forced to be consistent with *chronology* and *domain* rules. Another mechanism to incorporate these rules in the network generation phase would be to set appropriate priors on the network structures, which would favor networks compatible with the rules. This is a more drastic change to the current approach, in that the scoring metric used by CB assumes prior equivalence of all network structures. (Other metrics, while forgoing prior equivalence, are also insensitive to domain peculiarities.)

In some cases, the values of variables in the original dataset have been discretized. In many cases (e.g., for head circumference), this has been done according to accepted practice in epidemiology. Still, it may be interesting to challenge accepted wisdom and attempt different ways of discretizing variable ranges, for example by using a decision tree building algorithm [33]. (We observe, incidentally, that the variable FAM_INC, which represents family income, was already normalized in the original dataset.) In order to enable other researchers to validate and extend our

results, we provide the following pointer to the datasets we used: <http://www.cse.sc.edu/~mgv/MentorData/MentorData.zip>.

(The original CHDS dataset is widely available from other sources.)

Finally, since there are many variables in the Mental Retardation domain, it may be advantageous to attempt attribute selection [34], and use only a subset of variables.

Acknowledgments

This work was supported in part by the Advanced Research and Development Activity (ARDA), an entity of the U.S. Government. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the U.S. Government. We thank Shanthi Krishnaswami for preprocessing of the data sets used in this study and for help with the t -test. We thank Doug Fisher for many detailed comments on an early draft of this paper, Finn V. Jensen and other colleagues at the Department of Computer Science at the University of Aalborg for comments on a presentation of this work, and Richard Neapolitan for comments on a more recent draft. We thank several participants in the Fifth International Workshop on AI and Statistics for useful comments. We thank Hugin Ltd. Finally, we thank the anonymous reviewers for comments that improved the paper.

Notes

1. A similar situation exists in many other fields, both in the social and in the natural sciences; consider the tremendous amount of non-experimental data sent by spacecraft for an example outside the social sciences.
2. Probabilities sum to 1, so one of the 192 parameters is dependent on the other 191.
3. This implementation, which includes **CondProb**, is called **Visual CB** and is described in [21].
4. Another name for Bayesian Networks is Causal Probabilistic Networks; hence the second part of the acronym.
5. We use “known” as the opposite of missing.

References

1. J.M. Zytow and J. Baker, “Interactive mining of regularities in databases,” in: *Knowledge Discovery in Databases*, edited by Gregory Piatetsky-Shapiro and William J. Frawley, The AAAI Press: 445 Burgess Drive, Menlo Park, CA 94025, 1991, pp. 31–54.
2. J. Pearl and T. Verma, “A theory of inferred causation,” in *Principles of Knowledge Representation and Reasoning: Proceedings of the Second Conference*, edited by J.A. Allen, R. Fikes, and E. Sandewall, Morgan-Kaufmann: San Mateo, CA, 1991, pp. 441–452.
3. J. Pearl, *Causality: Models, Reasoning, and Inference*, Cambridge University Press: Cambridge, UK, 2000.
4. M.L. Batshaw, “Mental retardation,” in *Pediatric Clinics of North America—The Child With Developmental Disabilities*, vol. 40(3), edited by M.L. Batshaw, W.B. Saunders Company: Philadelphia, PA, 1993, pp. 507–522.
5. Z.A. Stein and M.W. Susser, “Mental retardation,” in *Public Health and Preventive Medicine*, 13th edn., M.J. Last and R.B. Wallace, Appleton & Lange: San Mateo, CA, 1992.
6. American Association on Mental Retardation, *Mental Retardation* American Association on Mental Retardation, 1719 Kalamazoo Road, NW Washington DC, 1992.
7. S.W. McDermott, A.L. Coker, and R.E. McKeown, “Low birth-weight and risk of mild mental retardation by ages 5 and 9 to 11,” *Pediatric and Perinatal Epidemiology*, vol. 7, pp. 195–204, 1993.
8. S. Mani, S. McDermott, and M. Valtorta, “MENTOR: A Bayesian model for prediction of mental retardation in newborns,” *Research in Developmental Disabilities*, vol. 18, no. 5, pp. 303–318, 1997.
9. R. Neapolitan, *Probabilistic Reasoning in Expert Systems: Theory and Algorithms*. Wiley: New York, 1990.
10. J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, Morgan Kaufmann: San Mateo, California, 1988.
11. E. Charniak, “Bayesian networks without tears,” *AI Magazine*, vol. 12, no. 4, pp. 50–63, 1991.
12. D.J. Spiegelhalter, A.P. Dawid, S.L. Lauritzen, and R.G. Cowell, “Bayesian analysis in expert systems,” *Statistical Science*, vol. 8, no. 3, pp. 219–283, 1993.
13. R.G. Cowell, A.P. Dawid, S.L. Lauritzen, and D.J. Spiegelhalter, *Probabilistic Networks and Expert Systems*, Springer: New York, 1999.
14. F.V. Jensen, *An Introduction to Bayesian Networks*, Springer: New York, 1996.
15. F.V. Jensen, *Bayesian Networks and Decision Graphs*, Springer: New York, 2001.
16. S.L. Lauritzen, B. Thiesson, and D. Spiegelhalter, “Diagnostic systems created by model selection methods—A case study,” *Preliminary Papers of the Fourth International Workshop on Artificial Intelligence and Statistics*, Ft. Lauderdale, FL, January 3–6, pp. 93–105, 1993.
17. G.F. Cooper, and E. Herskovits, “A Bayesian method for the induction of probabilistic networks from data,” *Machine Learning*, vol. 9, pp. 309–347, 1992.
18. M. Singh and M. Valtorta, “A new algorithm for the construction of bayesian network structures from data,” in *Proceedings of the 9th Conference on Uncertainty in Artificial Intelligence (UAI-93)*, Washington, DC, July 1993, pp. 259–264.
19. M. Singh and M. Valtorta, “Construction of bayesian belief networks from data: A brief survey and an efficient algorithm,”

- International Journal of Approximate Reasoning*, vol. 12, pp. 111–131, 1995.
20. *Data Archive and Users Manual of the Child Health and Development Studies*, Vols. 1 & 2. School of Public Health, University of California at Berkeley, CA, March 1, 1987.
 21. B. Xia, “An algorithm to learn bayesian probabilistic network structures from data,” M.S. Thesis, Department of Computer Science and Engineering, University of South Carolina, Columbia, SC, 2002.
 22. S.L. Lauritzen and D.J. Spiegelhalter, “Local computation with probabilities on graphical structures and their applications to expert systems,” *J.R. Statist. Soc., Series B*, vol. 50, pp. 157–224, 1988.
 23. S.K. Andersen, K.G. Olesen, F.V. Jensen, and F. Jensen. “HUGIN—A shell for building bayesian belief universes for expert systems,” in *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence*, Detroit, 1989, pp. 1080–1085.
 24. S. Andreassen, M. Woldbye, B. Falck, and S.K. Andersen, “MUNIN—A causal probabilistic network for interpretation of electromyographic findings,” in *Proceedings of the Tenth International Joint Conference on Artificial Intelligence*, Milan, 1987, pp. 366–372.
 25. P.M. Murphy and D.W. Aha, UCI Repository of Machine Learning Databases. [<http://www.ics.uci.edu/mllearn/MLRepository.html>]. University of California, Dept. of Information and Computer Science: Irvine, CA.
 26. R. Neapolitan, *Learning Bayesian Networks*, Pearson Prentice Hall: Upper Saddle River, NJ, 2004.
 27. N. Friedman, “The Bayesian structural EM algorithm,” in *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence (UAI-98)*, Madison, WI, July 1998, pp. 129–138.
 28. J.S. Mausner and S. Kramer, “The concept of causality and steps in the establishment of causal relationships,” in *Epidemiology—An introductory Text*, Ch 7. Company: Philadelphia, PA, W.B Saunders, 1985.
 29. S.M. Weiss and C.A. Kulikowski, *Computer Systems that Learn*. San Mateo, Morgan-Kaufmann: CA, 1991.
 30. S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*. Englewood Cliffs, Prentice Hall: New Jersey, 1995.
 31. P. Spirtes, C. Glymour, and R. Scheines, *Causation, Prediction, and Search*, 2nd edn., MIT Press: Cambridge, MA, 2000.
 32. S.L. Lauritzen and R. Thomas, “Chain graph models and their causal interpretation,” Research Report R-01-2003, Department of Mathematical Sciences, Aalborg University. (The revised version of August 2001 is available at <http://www.math.auc.dk/steffen/papers/rss.ps>. The paper is to appear in *Journal of the Royal Statistical Society, Series B*.)
 33. S. Monti and G.F. Cooper, “A latent variable model for multivariate discretization,” in *Proceedings of the Seventh International Workshop on Artificial Intelligence and Statistics*, January 1999 (<http://www-2.cs.cmu.edu/~smonti/publications/ais99.html>).
 34. G. Provan and M. Singh, “Learning bayesian networks using feature selection,” in *Learning from Data, Lecture Notes in Statistics* edited by D. Fisher, and H. Lenz, Springer Verlag: New York, vol. 112, pp. 291–300.