

University of Castilla-La Mancha



A publication of the
Department of Computing Systems

Towards consistency in general dependency networks

by

José A. Gámez and Juan L. Mateo and José M. Puerta

Technical Report #DIAB-08-04-1 April 2008

COMPUTING SYSTEMS DEPARTMENT
COMPUTER SCIENCE SCHOOL
UNIVERSITY OF CASTILLA-LA MANCHA
Campus Universitario s/n
Albacete - 02071 - Spain
Phone +34.967.599200, Fax +34.967.599224

Towards consistency in general dependency networks

José A. Gámez Juan L. Mateo José M. Puerta

April 29, 2008

Abstract

Dependency networks are a probabilistic graphical model that claim several advantages from other models like Bayesian networks and Markov networks, for instance. One of these advantages in general dependency networks, which are the object of study in this work, is the ease of learning from data. Nonetheless this easiness is also the cause of its main drawback: inconsistency. A dependency network cannot encode the probability distribution underlaid in the data but an approximation. This approximation can be enough good for some applications but not in other cases.

In this work we make a study of this inconsistency and propose a method to reduce it. From the conclusions we have taken from this analysis we have developed an algorithm that has to be run after the standard learning algorithm yields its solution. Our method is an heuristic approach so we cannot assure that the resulting model is fully consistent, however we have carried out some experiments which make us to think that it produces high quality models and therefore is advisable its use.

1 INTRODUCTION

Probabilistic graphical models (PGM) (Lauritzen, 1996; Jensen and Nielsen, 2007) have been deeply under research and have been used in many applications in the last two decades because they combine a good way to represent knowledge readable for humans with a strong theoretical foundation based on probability theory. There are several kinds of PGMs like decision graphs or Markov networks (MN), but probably the most famous and used are Bayesian networks (BN). A BN can be built from data, from expert knowledge or both, and can be used for different purposes, for instance decision making support, risk analysis, relationship visualization and interpretation, etc. The idea of BNs is based on modeling a domain in which the objects or entities are represented by random variables because they have attached uncertainty about their value. This uncertainty is modeled by means of a joint probability distribution for all the variables. The objects in the domain typically exhibit some relations or

dependencies and we use a directed acyclic graph to encode those dependencies with links in the graph where the nodes represent the objects.

Thus, with a BN we can represent qualitatively and quantitatively the relationships between the entities in the domain by means of the graphical representation and the probability measure respectively. In a BN the graph always has to satisfy the *Markov condition* given the joint probability distribution, i.e. every variable is conditionally independent of all its non-descendant given its parents, then the joint probability distribution can be recovered from the set of conditional probability distributions, one for every variable given its parents, reducing notably the complexity of the model.

In a BN we know that if we set the values for the parents for a given variable there is no influence from the variables which are above the parents in the topological order, nonetheless still this variable can be influenced by other variables. The set of variables which define entirely one given and makes it independent of all the others in the domain is called *Markov blanket set* (MB). In a BN the MB consists of the parents, the children and the parent of the children in the graph.

Dependency networks (DN) are a probabilistic graphical model proposed by (Heckerman et al., 2000) as an alternative to BN. The main difference between them is that the graph in DN does not have to be acyclic. The parametric component is the same, i.e. every variable has a conditional probability distribution given its parents. Another difference is that in DNs the parents for each variable is the set of variables which make it independent from the other variables, i.e. its MB in the Bayesian network encoding the same domain. This is the reason why the graph of a dependency network can be cyclic, we can have bi-directional links, apart of other cycles, because in a BN for every variable (X_i), every variable in its MB (X_j) has X_i in its own MB, that is $\forall X_j \in MB(X_i) \Rightarrow X_i \in MB(X_j)$.

In (Heckerman et al., 2000) are presented some tasks in which DNs can be worthwhile like probabilistic inference, collaborative filtering and visualization of relationships. Nonetheless, from the automatic learning point of view DNs have a drawback because its not easy to learn a set of conditional probability distributions which satisfied the definition given above. That is the reason the authors relaxed the definition of DNs and they defined *general dependency networks*, however now we cannot expect that with the set of conditional probability distributions we were able to recover the joint probability distribution but an approximation. Heckerman et al. (2000) argue that this approximation can be better as the amount of data used in the learning process increases, however it still is an approximation.

In this work we want to analyse how this approximation can deteriorate the performance of a DN model and we propose a way to improve the whole model with a minimum computational cost.

In section 2 we present a more formal and detailed definition of DNs. In section 3 we make an analysis of the inconsistencies that can appear in general DNs. In section 4 we explain our proposal to reduce those inconsistencies. In section 5 we describe some experiments we have carried out in order to validate

our proposal and show the results. Finally in section 6 we conclude with the final remarks.

2 DEPENDENCY NETWORKS

Dependency networks were proposed in (Heckerman et al., 2000) in response to a common complaint over BNs according to the authors. This complaint is based on the visualization properties of BNs, because they have been used as a tool which let humans visualize relationships learned from data. Using the graph of a BN we can easily interpret the knowledge encoded in that BN looking at the links for every variable.

Specially if a BN encodes causal relationships anybody can understand quickly all the interaction between the objects in the model. Nonetheless, is more difficult to learn causal relationships from data and not all BNs have this kind of representation. Moreover, an untrained individual would understand at a first glance that only the parents can give information about a given variable. Of course, is easy to learn that not only the parents, but the children and parents of the children are needed to define completely a given variable and avoid interferences of other variables.

Thus, that is the main reason why DN were created. In (Heckerman et al., 2000) authors use the example in Figure 1, in subfigure (a) is depicted a BN with three variables. If we interpret the relationships as causal ones then we can say that Age and Gender are *causes* of Income, or in other way Income is affected by our knowledge about Age and Gender. Even though this is true there are more dependencies beyond this first interpretation. Our knowledge about Income also affects the value of Age and Gender despite the link is oriented in the opposite direction, and if we do not know anything about Income then Age and Gender are independent, but if this fact changes then these two variables become dependent.

So, we can say that in a way or another all variables are related and in order to make a model more visually attractive every variable in that case should be connected with the other two. That is the aim of DNs in which every variable has as parents the MB set in the BN for the same domain, and then in this case the graph for the DN would be the one depicted in Figure 1(b).

So far we have explained the basic idea behind DNs, next we present a formal definition. We assume from now on that all variables are discrete, although the definition can be extended for continuous variables.

2.1 CONSISTENT DEPENDENCY NETWORKS

Given a set of variables $\mathbf{X} = \{X_1, \dots, X_n\}$ with a positive joint probability distribution $P(\mathbf{X})$, a *consistent dependency network* for this domain consists of a pair $(\mathcal{G}, \mathcal{P})$ where \mathcal{G} is a directed graph (not necessarily acyclic), in which every node represents a variable, and \mathcal{P} is a set of conditional probability distributions. In \mathcal{G} the set of parents for each variable X_i , denoted by \mathbf{Pa}_i , is formed by all

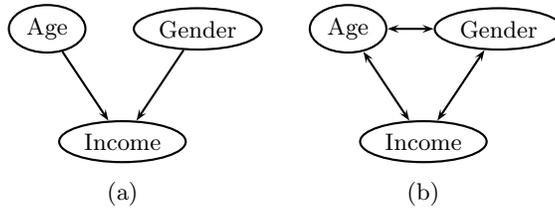


Figure 1: Example of a Bayesian and a dependency network for the same domain taken from (Heckerman et al., 2000).

those variables such that verify

$$P(X_i|\mathbf{Pa}_i) = P(X_i|X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n). \quad (1)$$

So if $P(\mathbf{X})$ is faithful to a graph, what is a common assumption in machine learning algorithms, then the parents for a given variable in a DN are its MB set. Other way to say so is that a DN has the same adjacencies than a MN.

A DN is consistent in the sense that all the conditional probability distributions in \mathcal{P} can be obtained from the joint probability distribution $P(\mathbf{x})$, i.e. we can obtain $P(\mathbf{X})$ from \mathcal{P} in a similar way than in an BN or MN.

$$P(\mathbf{X}) = \prod_{i=1}^n P(X_i|\mathbf{Pa}_i)$$

In (Heckerman et al., 2000) is shown the equivalence between DNs and MNs. The only difference is that in MNs the quantitative component is provided by potential functions whereas in DNs is provided by conditional probability distributions. Given this equivalence one approach to learn DNs from data can be to learn a MN in order to obtain the structure and then compute the set of conditional probability distributions from the MN via probabilistic inference. Other possibility suggested also in that paper is to learn another probabilistic model (a BN for instance) and translate it into a DN. Nonetheless the problem with these approaches is that the conversion can be computational expensive and inefficient in many cases. That is the reason why the authors presented another definition for DNs more relaxed in order to ease the automatic learning from data. This new definition is covered in next section.

2.2 GENERAL DEPENDENCY NETWORKS

A consistent DN is not attractive from an machine learning point of view because of the difficulties related with obtaining the set of conditional probability distributions, specially with the restriction that this set has to be consistent with the joint probability distribution for the variables in the domain. So *general* DNs, also described in (Heckerman et al., 2000), are based on the idea of removing

those restriction about consistency. Thus every single conditional probability distribution $P(X_i|\mathbf{Pa}_i)$ can be estimated independently from the others by any probabilistic classification method. Are proposed techniques such as a decision tree (Buntine, 1991), a generalized linear model (McCullagh and Nelder, 1989), a neural network (Bishop, 1995), a probabilistic support vector machine (Platt, 1999), or an embedded classification model (Heckerman and Meek, 1997). Once we have all the conditional probability distributions we can build the structure of the dependency network from the (in)dependencies that are appeared during the learning process.

This way of learning a DN can be more efficient than learning from a MN in many cases, and other advantage is that its parallelization is straightforward, what can report a great benefit if we are dealing with a domain with a large number of variables. Nonetheless this heuristic approach has a disadvantage, due mainly to the independent search over the variables and poor estimations in small datasets, the learned conditional probability distributions may not be consistent with the joint probability distribution, this can be called *parametrical inconsistency*. But also *structural inconsistencies* can appear because after learning the conditional probability distribution we can see that, for instance, X_i can be parent of X_j but not the opposite, i.e. the conditional probability distribution for X_i would not contain X_j but the conditional probability distribution for X_j would contain X_i . In (Heckerman et al., 2000), authors argue that this inconsistencies can be reduced as the amount of data used for leaning increases.

A formal definition for this new model is as follows. Given a set of variables $\mathbf{X} = \{X_1, \dots, X_n\}$, consider the set of conditional probability distributions $\mathcal{P} = \{P_1(X_1|\mathbf{X}\setminus X_1), P_2(X_2|\mathbf{X}\setminus X_2), \dots, P_n(X_n|\mathbf{X}\setminus X_n)\}$. It is not required that these distributions be consistent with $P(\mathbf{X})$, i.e. it is not required that this set be obtained via inference from the joint probability distribution. Under these conditions a *dependency network* for \mathbf{X} and \mathcal{P} is the pair $(\mathcal{G}, \mathcal{P}')$, where \mathcal{G} is a directed graph usually cyclic and \mathcal{P}' is a set of conditional probability distributions such that

$$P_i(X_i|\mathbf{Pa}_i) = P_i(X_i|\mathbf{X}\setminus X_i) \quad (2)$$

for every $P_i \in \mathcal{P}$.

2.3 INFERENCE

In any case, with a consistent or general dependency network, given the likely existence of cycles in the graph we cannot use exact inference algorithms used in BNs and some of the approximate. In the case of consistent DNs it can be converted to a MN and use standard techniques for probabilistic inference over MNs. Nonetheless a more general option is suggested in (Heckerman et al., 2000) for both models, Gibbs sampling (Geman and Geman, 1984). Basically this method works by repeatedly cycling through each variable in a fixed order during all the process, and sampling each X_i according to $P(X_i|\mathbf{Pa}_i)$. This

procedure is called *ordered Gibbs sampler* but in the case of a general DN, given that the conditional probability distributions may not be consistent with the joint probability distribution, is called *ordered pseudo-Gibbs sampler*. Besides it is developed a more efficient method which can avoid some sampling and it is called *modified ordered (pseudo-)Gibbs sampler*. This method, in order to get $P(\mathbf{Y}|\mathbf{Z})$ and the value of \mathbf{Z} is \mathbf{z} for a DN in a domain with a set of variables \mathbf{X} , is shown in Figure 2.

```

1  $\mathbf{U} = \mathbf{Y}$  (* the unprocessed variables *)
2  $\mathbf{P} = \mathbf{Z}$  (* the processed and conditioning variables *)
3  $\mathbf{p} = \mathbf{z}$  (* the values of  $\mathbf{P}$  *)
4 While  $\mathbf{U} \neq \emptyset$ 
5   Pick  $X_i \in \mathbf{U}$  s.t.  $X_i$  has no more parents in  $\mathbf{U}$  than any variable
   in  $\mathbf{U}$ 
6   If all parents of  $X_i$  are in  $\mathbf{P}$ 
7      $P(X_i|\mathbf{p}) = p(X_i|\mathbf{Pa}_i)$ 
8   Else
9     Use modified ordered Gibbs sampling to get  $P(X_i|\mathbf{p})$ 
10   $\mathbf{U} = \mathbf{U} - X_i$ 
11   $\mathbf{P} = \mathbf{P} + X_i$ 
12   $\mathbf{p} = \mathbf{p} + x_i$ 
13 Return the product of the conditionals  $P(X_i|\mathbf{p})$ 

```

Figure 2: Modified ordered Gibbs sampler

The key point is in line 6, because if are known the values for all the parents for a given variable we can avoid the sampling for that variable and just take its value from its conditional probability distribution. This algorithm is justified by Equations 1 and 2.

In (Heckerman et al., 2000) is discussed whether the use of a sampling process can increase the parametrical inconsistencies in a general DN. The conclusion is that it should not be a determinant fact but more research should be done in order to characterize better this behavior and assure a good performance of DNs.

However, given the modified ordered Gibbs sampler, there are some situations in which we can avoid completely the sampling. For instance if we use a DN as a classifier and we assume that we always know the values for all predictive variables (Gámez et al., 2006). Other case is when is needed to obtain the probability for a full configuration in a DN, i.e. $P(\mathbf{X})$ when we have fixed the value for every single variable. We can see this computation in other way

$$P(\mathbf{X}) = \prod_{i=1}^n P(X_i|\mathbf{X}\setminus X_i).$$

If we take the right part of the equation and use the modified ordered Gibbs sampler we have that $\mathbf{Y} = \{X_i\}$ and $\mathbf{Z} = \{\mathbf{X}\setminus X_i\}$, in line 5 we have only one choice and in line 6 the condition is true so the sampling is avoided. Therefore we can compute statistics such as the likelihood of a DN for a dataset, and we assume that the dataset does not contain missing data, without any sampling.

3 ANALYSIS OF PARAMETRICAL INCONSISTENCY

In this section we want to analyse some issues regarding with inconsistency in DNs. From now on we consider only general DNs.

Example 1 Consider the case in which we have two variables, X and Y , and they are dependent, then the DN for this domain, \mathcal{DN} , should have a graph with two links $X \rightarrow Y$ and $X \leftarrow Y$.

Hence $\mathcal{P}' = \{P(X|Y), P(Y|X)\}$ for \mathcal{DN} . However is clear that $P(X, Y) \neq P(X|Y) \cdot P(Y|X)$. In fact

$$\hat{P}(X, Y) = P(X|Y) \cdot P(Y|X) = \frac{P(X, Y) \cdot P(X, Y)}{P(X) \cdot P(Y)}.$$

So we can define the estimated joint probability distribution defined by \mathcal{DN} in this way

$$\hat{P}(X, Y) = f(X, Y) \cdot P(X, Y), \quad (3)$$

where

$$f(X, Y) = \frac{P(X, Y)}{P(X) \cdot P(Y)}$$

Therefore, even in a situation so simple like this one, we cannot expect to have a DN without inconsistencies, when it is learned from data of course. Moreover, looking at Equation 3 we can say that the inconsistency is smaller as the dependency between X and Y is weaker and $f(X, Y)$ tends to 1.

In (Heckerman et al., 2000) they propose learn DNs by means of probabilistic decision trees. This model are very good to encode contextual dependencies. Encoding the conditional probabilities distribution in this way can help to reduce the inconsistencies because a decision tree tries to represent a more general probability distribution by pruning some branches which are similar. Then the dependence between the variables can be smoothed and thus $f(X, Y)$ is closer to 1. However if this happens we have a poorer estimation for the joint probability distribution even though is less inconsistent.

In order to illustrate that we can use the next example.

Example 2 Considering the model in Example 1 and both variables are discrete with three states and their joint probability distribution is defined by this table:

	Y=0	Y=1	Y=2
X=0	0.12	0.04	0.04
X=1	0.06	0.18	0.06
X=2	0.10	0.10	0.30

When we compute $P(X|Y)$ and $P(Y|X)$ from $P(X, Y)$ in the way of a probability table or a full expanded probabilistic decision trees (see Figure 3 (a) and (b)) we obtain an estimation $\hat{P}_1(X, Y)$ which is shown in Figure 3(c):

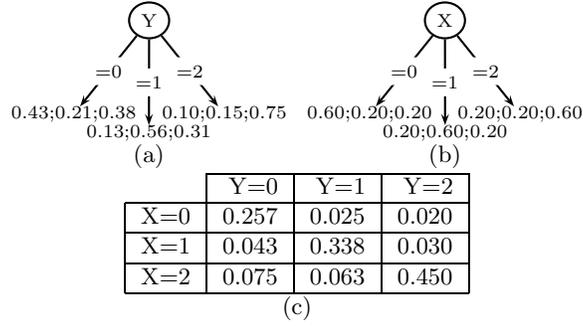


Figure 3: Joint probability distribution $\hat{P}_1(X, Y)$ (c) obtained using full decision trees for conditional probability distributions $P(X|Y)$ (a) and $P(Y|X)$ (b).

We can see large differences between the true figures and the estimation. Besides we can check that $\hat{P}_1(X, Y)$ is not a probability distribution because $\sum_{x,y} \hat{P}_1(x, y) = 1.301 \neq 1$. If, for instance, the learning procedure decides to change the representation of $P(X|Y)$ for a probabilistic decision tree in which branches for values 0 and 1 are merged (Figure 4(a)), because they are the most similar, then we obtain a new estimation $\hat{P}_2(X, Y)$ which is shown in Figure 4(b). $\hat{P}_2(X, Y)$ still differs from $P(X, Y)$ but is closer to it than $\hat{P}_1(X, Y)$ in average, and also is closer to be a probability distribution because $\sum_{x,y} \hat{P}_2(x, y) = 1.170$.

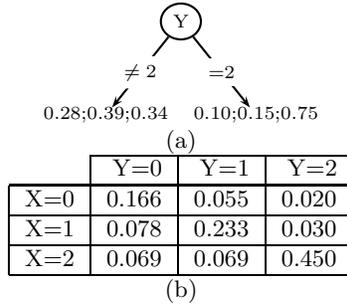


Figure 4: Joint probability distribution $\hat{P}_2(X, Y)$ (b) obtained using a simpler decision tree for $P(X|Y)$ (a).

Therefore in spite of the use of probabilistic decision trees we still have an approximation which could not be good enough for some applications. In next section we present a simple heuristic method that can reduce inconsistencies in DNs improving its accuracy.

4 HOW TO IMPROVE CONSISTENCY

As it has been seen in the previous Section even in a simple case like Example 2 we cannot expect to get a consistent DN if the conditional probability distributions are learned independently. If two variables are dependent this equation $P(X, Y) = P(X|Y) \cdot P(Y|X)$ will never be true, nonetheless this expression $P(X, Y) = P(X) \cdot P(Y|X) = P(X|Y) \cdot P(Y)$ is always true and does not matter if both variable are dependent or not. Bearing this in mind, in Example 2 we can ensure consistency if at the end of the learning process we realize that X is a predictive variable for Y and vice versa and then instead of maintaining both conditional probability distributions we replace $P(X|Y)$ by $P(X)$ or $P(Y|X)$ by $P(Y)$. This is the basic idea of our proposal, but there is not so easy when there are more variables involved. In that case we do not expect to obtain the best set of probabilities whose composition yield the right joint probability distribution, but a good approximation and, more important, more consistent.

More precisely the proposal consists in estimating a set of conditional probability distributions of a BN that encode the same (in)dependencies that the learned DN. We have to point out that our proposal only changes the set of probability distributions but not the graph, so the model learned still has the same advantages about visualization. However, given that the relationships represented in a DN can be encoded by several BNs with different factorizations of the joint probability distribution, and that the conversion from DN to BN can not be attractive from the computational point of view, this proposal is based on a heuristic approach whose complexity order is linear in the number of dependencies found. The method proposed is shown in Figure 5.

```
1 For each variable  $X_i$ 
2   For each predictive variable  $Y_j$  of  $X_i$ 
3     If  $X_i$  is also a predictive variable of  $Y_j$ 
4       If the conditioning set of  $X_i$  is grater that  $Y_j$ 's
5          $Y_j$  is removed as predictive variable of  $X_i$ 
6       Else
7          $X_i$  is removed as predictive variable of  $Y_j$ 
```

Figure 5: Proposed method to obtain a more consistent set of conditional probability distributions.

We can call this new step in the learning process as parametric reduction. An important point in this procedure is in line 4. With this condition we want avoid large conditioning sets, what can reduce overfitting in the parameters estimation. The order in which the links can be traversed can be any although not all of then will yield the same solution. The reason for that is the heuristic nature of this algorithm and that a more sophisticated search would not be interesting for practical reasons. One of the benefits of DNs is ease of learning so we do not want to change that by introducing a complicated post-learning algorithm.

After performing this step is needed to re-compute every probability distribution which has been modified. In the case that these probability distributions are in form of probability trees, if the removed variables are in the leaves the only thing to do is to aggregate its values to the up node in the tree, otherwise the entire tree should be re-built. However, if in the learning process we have cached the statistics the new tree can be built without computational cost. In the case of probability tables we can differ the learning of these tables after that step.

5 EXPERIMENTAL RESULTS

This section is devoted to evaluate our proposal with some experiments. Our testing framework is base on the one used in (Heckerman et al., 2000) for testing probabilistic inference with real data. We use the same score function for a test dataset with N instances $\{d_1, \dots, d_N\}$ and n variables:

$$score(d_1, \dots, d_N | model) = - \frac{\sum_{i=1}^N \ln P(d_i | model)}{nN}. \quad (4)$$

However, instead of using real dataset we prefer using data sampled from known BNs. The reason is that we want focus only in parametrical learning and inference so if the real dependencies are known we can give this information to the different algorithms in order to avoid that the results were affected by the structural learning. Next we present a detailed description of our experimentation.

5.1 DESCRIPTION OF THE EXPERIMENTS

We have selected seven BNs from different sources: **alarm** (Beinlich et al., 1989), **asia** (Lauritzen and Spiegelhalter, 1988), **car-starts** and **headache** (Elvira Consortium, 2002), **insurance** (Binder et al., 1992), **credit** (DSL) and **water** (Jensen et al., 1989), which is a dynamic network and we have use only the two first slices. Some details of these networks can be seen in Table 1. From each of these networks we have sampled two datasets with 5000 instances each one, one for training and one for testing.

We have defined eight models to make a comparison between them. First one is the reference model and is a BN in which the structure is fixed with the real links (BN-f). Second model is an empty network (Empty). Next we have three dependency networks models, one with probability tables in which links have been fixed from the real MB for each variable in the network (PT-f), other with probabilistic decision trees learned from data (PDT), and other with probabilistic decision trees but in which the search space for each PDT have been restricted to the real MB (PDT-f). In both cases we use the sugested value for $\kappa = 0.1$. For any of these three models we have another version in which we have used our method for reducing the conditional probability distributions.

Table 1: Set of Bayesian networks used in our experiments.

network	Num. vars	States range	Aver. states	MB range	Aver. MB
alarm	37	2-4	2.84	1-12	3.89
asia	8	2-2	2.00	1-5	2.50
car-starts	18	2-3	2.06	1-9	3.44
credit	12	2-4	2.83	2-6	3.67
headache	12	1-4	2.92	1-4	2.67
insurance	27	2-5	3.30	1-16	6.22
water	16	3-4	3.63	1-12	6.00

These new models are labeled with an asterisk (PT-f*, PDT* and PDT-f*). In all cases parameters are learned from data by using Laplace smoothing.

Every model has been learned with each training dataset. For all of them it has been computed their score (Equation 4) with the test dataset. As the model BN-f is the reference one we have also obtained the absolute difference of score between each model and BN-f. This value is more informative because we are looking for models closer to the true probability distribution what is represented by BN-f. Besides, we have computed also the summation of all possible configurations, i.e. total joint probability, which should be equal to 1, but only for those models with a tractable number of configurations (**asia**, **car-starts**, **credit**, **headache**).

5.2 RESULTS

In Table 2 we report the score value for every model and dataset. At the bottom line we show the average value for each model. Lower values should indicate a better model, so all pure DN models should be taken as the best ones. However that does not make sense because they are even better than our reference model (BN-f) which represents the true joint probability distribution. The reason is that, as we have seen in Section 3, inconsistent DNs tend to have greater probability values in average so their score is lower. That is the reason why we prefer pay more attention to the difference with respect to the reference model.

Thus, these new results are shown in Table 3. There we can see that always the model closer to BN-f is the one in which we have applied our proposal. Specially the model based on probability tables is always the best one but in two datasets. Also is important to notice that our proposal improves the original model in every dataset for PT-f model. However, in PDT model our proposal deteriorates the accuracy in **alarm** and **headache** dataset, although in average its application improves the global accuracy.

Another interesting point is that PDT models without our proposal are much better than PT-f. That corroborate the idea that for DNs the use of more general

Table 2: Score for each model and dataset.

	BN-f	Empty	PT-f	PT-f*	PDT	PDT*	PDT-f	PDT-f*
alarm	0.282	0.397	0.173	0.298	0.253	0.342	0.242	0.338
asia	0.287	0.342	0.224	0.289	0.225	0.287	0.225	0.289
car-starts	0.127	0.175	0.070	0.127	0.070	0.136	0.070	0.127
credit	0.879	0.959	0.765	0.886	0.807	0.888	0.807	0.900
headache	0.435	0.609	0.214	0.435	0.419	0.585	0.419	0.593
insurance	0.490	0.651	0.399	0.519	0.420	0.556	0.420	0.550
water	0.401	0.410	0.417	0.410	0.388	0.409	0.388	0.408
	0.414	0.506	0.323	0.423	0.369	0.458	0.367	0.458

Table 3: Absolute score difference between BN-f and the other models.

	Empty	PT-f	PT-f*	PDT	PDT*	PDT-f	PDT-f*
alarm	0.115	0.110	0.015	0.029	0.060	0.040	0.056
asia	0.055	0.062	0.002	0.062	0.000	0.062	0.002
car-starts	0.048	0.057	0.000	0.057	0.009	0.057	0.000
credit	0.080	0.114	0.007	0.071	0.009	0.071	0.021
headache	0.174	0.222	0.000	0.017	0.150	0.017	0.158
insurance	0.161	0.092	0.029	0.070	0.066	0.071	0.059
water	0.009	0.016	0.010	0.013	0.008	0.013	0.007
	0.092	0.096	0.009	0.046	0.043	0.047	0.043

encoding for the conditional probability distributions is advisable despite that this encoding is also an approximation in many cases.

Previous results give us an idea about the quality of those model. We can suppose that the increment in accuracy must be related with the reduction in the inconsistency. Additionally we have checked if the models encode a real probability distribution, i.e. whether the total joint probability for a given model is equal to one. This computation has been only done for the models learned with the smaller networks because this computation is computationally unfeasible for the others. The result is shown in Table 4. According with the table is clear that the pure DN models are quite far from being a probability distribution, but our proposal achieve that condition for all of them.

6 CONCLUSIONS

In this paper we have presented a novel method which aims to improve DNs. The main advantage of (general) DNs is that they can be learned from data easily, in fact easier than BNs because of the lack of restrictions about cyclicity and easier than MNs because conditional probability distributions can be learn

Table 4: Total joint probability for tested models.

	BN-f	Empty	PT-f	PT-f*	PDT	PDT*	PDT-f	PDT-f*
asia	1.00	1.00	3.60	1.00	3.42	1.00	3.42	1.00
car-starts	1.00	1.00	20.04	1.08	11.40	1.00	11.40	1.00
credit	1.00	1.00	6.41	1.00	4.26	1.00	4.26	1.00
headache	1.00	1.00	29.68	1.00	5.70	1.00	5.70	1.00

independently. Nonetheless this is also the main problem, the independent learning can lead to inconsistencies. These inconsistencies can be both structural and parametrical, however the later are more important. Whereas structural inconsistencies can be interesting for a better interpretation of the model (strong and weak dependencies), parametrical ones deteriorate model performance.

Thus our proposal is based on improving DNs accuracy by reducing conditional probability distributions, because, as it has been seen in Section 3, the use of full distributions is the cause of that inconsistencies. Is worthy to point out that our proposal does not change the qualitative component of a model, i.e. its links. This new method, that can be seen as a post-learning stage, works by trying to recover a set of conditional probability distributions similar to a BN which represents the same relationships between variables. In order not to lose the computational advantage of the DNs learning we have chosen a heuristic approach. This approach has a linear complexity order in the number of links. Its heuristic nature can be object of complaint, nonetheless in our experimentation we have made clear its benefit.

We plan to extend this work in two lines as future work. First we plan to make a deeper analysis of our proposal checking the performance with different sample sizes and different ordering in the reduction step and see whether it affects the results. Besides we want to test probabilistic queries for different set of variables with evidence in which we have to use Gibbs sampling. The second line of work is applying this method to scenarios where DNs have been use in order to improve their results, such as classifiers or Estimation of Distribution Algorithms.

Acknowledgements

This work has been partially supported by Spanish Ministerio de Educación y Ciencia (project TIN2004-06204-C03-03); Junta de Comunidades de Castilla-La Mancha (project PBI-08-048) and FEDER funds.

References

I.A. Beinlich, H.J. Suermondt, R.M. Chavez, and G.F. Cooper. The ALARM monitoring system: A case study with two probabilistic inference techniques

- for belief networks. In *Proc. of the 2nd European Conf. on Artificial Intelligence in Medicine*, pages 247–256, 1989.
- J. Binder, D. Koller, S. Russell, and K. Kanazawa. Adaptive probabilistic networks with hidden variables. *Machine Learning*, 29(2):213–244, 1992.
- C. M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
- W. Buntine. Theory refinement on bayesian networks. In *Uncertainty in Artificial Intelligence (UAI91)*, pages 52–60, 1991.
- Decision Systems Laboratory DSL. Genie. <http://genie.sis.pitt.edu/>.
- Elvira Consortium. Elvira: An Environment for Creating and Using Probabilistic Graphical Models. In *Proceedings of the First European Workshop on Probabilistic Graphical Models*, pages 222–230, 2002. <http://leo.ugr.es/elvira>.
- J. A. Gámez, J. L. Mateo, and J. M. Puerta. Dependency networks based classifiers: learning models by using independence test. In *Third European Workshop on Probabilistic Graphical Models (PGM06)*, pages 115–122, 2006.
- S. Geman and D. Geman. Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:147–156, 1984.
- D. Heckerman and C. Meek. Models and Selection Criteria for Regression and Classification. Technical Report MSR-TR-97-08, Microsoft Research (MSR), May 1997.
- D. Heckerman, D. M. Chickering, and C. Meek. Dependency networks for inference, collaborative filtering and data visualization. *Journal of Machine Learning Research*, 1:49–75, 2000.
- F. V. Jensen and T. D. Nielsen. *Bayesian networks and decision graphs*. Springer, 2007.
- F. V. Jensen, U. Kjærulff, K. G. Olesen, and J. Pedersen. Et forprojekt til et ekspertsystem for drift af spildevandsrensning (an expert system for control of waste water treatment — a pilot project). Technical report, Judex Datasystemer A/S, Aalborg, Denmark, 1989.
- S. L. Lauritzen. *Graphical Models*. Oxford University Press, 1996.
- S. L. Lauritzen and D. J. Spiegelhalter. Local computations with probabilities on graphical structures and their application to expert systems. *Journal of the Royal Statistical Society, Series B*, 50:157–194, 1988.
- P. McCullagh and J. Nelder. *Generalized Linear Models*. Chapman and Hall, 1989.

J. Platt. *Advances in Kernel Methods – Support Vector Learning*, chapter Fast training of support vector machines using sequential minimal optimization, pages 185–208. MIT Press, 1999.