



ELSEVIER

International Journal of Approximate Reasoning 28 (2001) 23–50

INTERNATIONAL JOURNAL OF
APPROXIMATE
REASONING

www.elsevier.com/locate/ijar

Performance evaluation of compromise conditional Gaussian networks for data clustering

J.M. Peña, J.A. Lozano, P. Larrañaga

*Department of Computer Science and Artificial Intelligence, Intelligent Systems Group,
University of Basque Country, P.O. Box 649, E-20080 Donostia – San Sebastián, Spain*

Received 1 May 2000; accepted 1 April 2001

Abstract

This paper is devoted to the proposal of two classes of compromise conditional Gaussian networks for data clustering as well as to their experimental evaluation and comparison on synthetic and real-world databases. According to the reported results, the models show an ideal trade-off between efficiency and effectiveness, i.e., a balance between the cost of the unsupervised model learning process and the quality of the learnt models. Moreover, the proposed models are very appealing due to their closeness to human intuition and computational advantages for the unsupervised model induction process, while preserving a rich enough modeling power. © 2001 Elsevier Science Inc. All rights reserved.

Keywords: Data clustering; Conditional Gaussian networks; Naive Bayes models; Tree augmented naive Bayes models; Extended naive Bayes models

1. Introduction

A basic problem that arises in a variety of fields, such as pattern recognition, machine learning, and statistics, is the so-called *data clustering problem* [1,2,7,8,16,18]. From the point of view adopted in this paper, the data clustering problem may be defined as the inference of a generalized

E-mail address: ccbpepaj@si.ehu.es (J.M. Peña).

joint probability distribution from a database. We assume that, in addition to the observed random variables or predictive attributes, there is a hidden random variable. This last unobserved random variable would reflect the cluster membership for every case in the database. Thus, the clustering problem is also referred to as an example of learning from unlabeled data due to the existence of such a hidden random variable. In this paper, we focus on learning conditional Gaussian networks for data clustering.

Due to the difficulty involved in learning densely connected conditional Gaussian networks (even more when working with large databases in terms of both number of cases and number of random variables), and the painfully slow probabilistic inference when working with them, it is necessary to develop methods for learning simple conditional Gaussian networks while preserving the quality of the learnt models. A symmetrical problem is found when working in discrete domains, and some of the proposed solutions are based on achieving such a balance between efficiency and effectiveness [7,10–12,19,24,27–34].

Keeping this idea in mind, we propose two classes of compromise conditional Gaussian networks for data clustering. The models of the first class belong to what in [10–12,19] are called tree augmented naive Bayes models (more recently referred to as mixtures of trees with shared structure [24]). These models are used in the referred works as Bayesian classifiers, while in [31,33] they are evaluated in the framework of data clustering for discrete domains. In this work, we evaluate them with respect to the data clustering problem for continuous domains. Tree augmented naive Bayes models for data clustering are defined by the following condition: Predictive attributes may have at most one other predictive attribute as a parent. The second class of models that we present and evaluate with respect to the data clustering problem for continuous domains is proposed by Pazzani [27,28] as Bayesian classifiers. In [30,32] this class of models are applied to data clustering in discrete domains. These models are similar to naive Bayes models [7], which is why we name them extended naive Bayes models, with the exception that the number of nodes in the structures can be smaller than the original number of random variables in the database as some of them can be grouped together under the same node as fully correlated.

When applied to data clustering, the compromise conditional Gaussian networks proposed are a weaker representation of some domains than general conditional Gaussian networks for data clustering. However, our experimental results on both synthetic and real-world databases show that the performance and modeling power of these models is still significant, while they offer obvious advantages from the point of view of the cost of the unsupervised model learning process.

The remainder of this paper is organized as follows. In Section 2, we introduce general conditional Gaussian networks for data clustering. Section 3 is

dedicated to present in detail our proposal of compromise conditional Gaussian networks for data clustering. Some experimental results on several synthetic and real-world databases are compiled in Section 4. Finally, in Section 5 we draw conclusions.

2. Unsupervised learning of conditional Gaussian networks

This section starts introducing the notation used throughout this paper. Then, we give a formal definition of conditional Gaussian networks for data clustering. We also present the Bayesian structural EM algorithm [9], which is used for explanatory purposes as well as in our experiments due to its good performance in unsupervised learning of conditional Gaussian networks.

2.1. Notation

We follow the usual convention of denoting unidimensional random variables by uppercase letters and their states by the same letters in lowercase. We use a letter or letters in boldface uppercase to designate a multidimensional random variable and the same boldface lowercase letter or letters to denote an assignment of state to the multidimensional random variable. The generalized joint probability distribution for \mathbf{X} is represented as $\rho(\mathbf{x})$. Additionally, $\rho(\mathbf{x}|\mathbf{y})$ denotes the generalized conditional probability distribution for \mathbf{X} given $\mathbf{Y} = \mathbf{y}$. If \mathbf{X} is a multidimensional discrete random variable, then $\rho(\mathbf{x}) = p(\mathbf{x})$ is the joint probability mass function for \mathbf{X} . Thus, $p(\mathbf{x}|\mathbf{y})$ denotes the conditional probability mass function for \mathbf{X} given $\mathbf{Y} = \mathbf{y}$. On the other hand, if \mathbf{X} is a multidimensional continuous random variable, then $\rho(\mathbf{x}) = f(\mathbf{x})$ is the joint probability density function for \mathbf{X} . Thus, $f(\mathbf{x}|\mathbf{y})$ denotes the conditional probability density function for \mathbf{X} given $\mathbf{Y} = \mathbf{y}$.

2.2. Conditional Gaussian networks for data clustering

As already mentioned, when facing a data clustering problem we assume the existence of a $(n + 1)$ -dimensional mixed random variable \mathbf{X} partitioned as $\mathbf{X} = (\mathbf{Y}, C)$ into a n -dimensional continuous random variable \mathbf{Y} and a unidimensional discrete *hidden* random variable C . It is usual to refer to the unidimensional random variables in \mathbf{Y} as *observed random variables* or *predictive attributes*. When dealing with a data clustering problem, \mathbf{X} is said to have a *conditional Gaussian distribution* [5,20–22] if the distribution for \mathbf{Y} , conditioned on each state of C , is a multivariate normal distribution. That is,

$$f(\mathbf{y}|C = c) = f_c(\mathbf{y}) \equiv \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}(c), \boldsymbol{\Sigma}(c)), \quad (1)$$

whenever $p(c) = p(C = c) > 0$. Given $C = c$, $\boldsymbol{\mu}(c)$ is the n -dimensional mean vector, and $\boldsymbol{\Sigma}(c)$, the $n \times n$ variance matrix, is positive definite.

We define a conditional Gaussian network (CGN) for \mathbf{X} for data clustering as a graphical model that encodes a conditional Gaussian distribution for \mathbf{X} . Essentially, CGNs for data clustering belong to a class of mixed probabilistic graphical models introduced for the first time by Lauritzen and Wermuth [22], and further developed in [20,21]. This class groups models in which both discrete and continuous random variables can be present and for which the conditional distribution for the continuous random variables given the discrete random variables is restricted to be multivariate Gaussian. More recently, CGNs have been successfully applied to data clustering [29,33].

Concretely, a CGN for \mathbf{X} for data clustering is defined by a directed acyclic graph s (model structure) determining the conditional (in)dependencies among the predictive attributes in \mathbf{Y} , a set of local probability density functions for these predictive attributes, and a multinomial distribution for the cluster random variable C . The model structure yields to a graphical factorization of the generalized joint probability distribution for \mathbf{X} as follows:

$$\rho(\mathbf{x}) = \rho(\mathbf{y}, c) = p(c)f(\mathbf{y}|c) = p(c)f_c(\mathbf{y}) = p(c) \prod_{i=1}^n f_c(y_i | \mathbf{pa}(s)_i), \quad (2)$$

where $\mathbf{pa}(s)_i$ denotes the configuration of the parents of Y_i , $\mathbf{Pa}(s)_i$, for all i . The local probability density functions for the predictive attributes and the multinomial distribution for the cluster random variable are those in the previous equation, and we assume that they depend on a finite set of parameters $\boldsymbol{\theta}_s \in \boldsymbol{\Theta}_s$. Therefore, Eq. (2) can be rewritten as follows:

$$\begin{aligned} \rho(\mathbf{x} | \boldsymbol{\theta}_s) &= \rho(\mathbf{y}, c | \boldsymbol{\theta}_s) \\ &= p(c | \boldsymbol{\theta}_s) f(\mathbf{y} | c, \boldsymbol{\theta}_s) \\ &= p(c | \boldsymbol{\theta}_s) f_c(\mathbf{y} | \boldsymbol{\theta}_s^c) \\ &= p(c | \boldsymbol{\theta}_s) \prod_{i=1}^n f_c(y_i | \mathbf{pa}(s)_i, \boldsymbol{\theta}_i^c), \end{aligned} \quad (3)$$

where $\boldsymbol{\theta}_s^c = (\boldsymbol{\theta}_1^c, \dots, \boldsymbol{\theta}_n^c)$ denotes the parameters for the local probability density functions for the predictive attributes when $C = c$.

If s^h denotes the hypothesis that the true generalized joint probability distribution for \mathbf{X} can be factorized according to the conditional independence assertions reflected in s , then we obtain from Eq. (3) that,

$$\begin{aligned}
 \rho(\mathbf{x} | \boldsymbol{\theta}_s, \mathbf{s}^h) &= \rho(\mathbf{y}, c | \boldsymbol{\theta}_s, \mathbf{s}^h) = p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) f(\mathbf{y} | c, \boldsymbol{\theta}_s, \mathbf{s}^h) \\
 &= p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) f_c(\mathbf{y} | \boldsymbol{\theta}_s^c, \mathbf{s}^h) \\
 &= p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) \prod_{i=1}^n f_c(y_i | \mathbf{pa}(s)_i, \boldsymbol{\theta}_i^c, \mathbf{s}^h).
 \end{aligned} \tag{4}$$

In order to encode a conditional Gaussian distribution for \mathbf{X} , the local probability density function for every predictive attribute Y_i of a CGN for \mathbf{X} for data clustering must be the *linear-regression model*. Thus, when $C = c$,

$$f_c(y_i | \mathbf{pa}(s)_i, \boldsymbol{\theta}_i^c, \mathbf{s}^h) \equiv \mathcal{N} \left(y_i; m_i^c + \sum_{Y_j \in \mathbf{Pa}(s)_i} b_{ji}^c (y_j - m_j^c), v_i^c \right), \tag{5}$$

where $\mathcal{N}(y; \mu, \sigma^2)$ is an univariate normal distribution with mean μ and standard deviation σ ($\sigma > 0$). When $C = c$, the parameters of the local probability density function for every Y_i are given by $\boldsymbol{\theta}_i^c = (m_i^c, \mathbf{b}_i^c, v_i^c)$. The interpretation of the components of the local parameters $\boldsymbol{\theta}_i^c$ is as follows for all i and all c state of C : m_i^c is the unconditional mean of Y_i when $C = c$, $\mathbf{b}_i^c = (b_{1i}^c, \dots, b_{i-1i}^c)^t$ is a column vector where every b_{ji}^c is a linear coefficient reflecting the strength of the relationship between Y_j and Y_i when $C = c$ if $Y_j \in \mathbf{Pa}(s)_i$, otherwise $b_{ji}^c = 0$, and v_i^c is the conditional variance of Y_i given $\mathbf{Pa}(s)_i$ when $C = c$. See Fig. 1 for an example of a CGN for data clustering with three continuous predictive attributes and one binary cluster random variable.

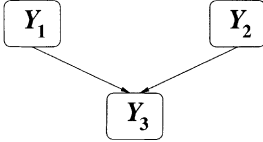
Note that the model structure is independent of the value of the cluster random variable C , thus, the model structure is the same for all the values of C . However, the parameters of the local probability density functions for the predictive attributes in \mathbf{Y} do depend on the value of C and they may be different for the different values of the random variable C .

2.3. Bayesian structural EM algorithm for unsupervised learning of conditional Gaussian networks

One of the methods for learning CGNs from unlabeled data is the well-known Bayesian structural EM (BS-EM) algorithm developed by Friedman in [9]. Due to its good performance, this algorithm has received special attention in the literature and has motivated several variants of itself [25,30–32,37]. We use the BS-EM algorithm for explanatory purposes as well as in our experiments presented in Section 4.

When applying the BS-EM algorithm in a data clustering problem, we assume that we have a database of N cases, $\mathbf{d} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, where every case is represented by an assignment to n of the $n + 1$ unidimensional random variables involved in the problem domain. In this case, it is convenient to see \mathbf{d} as a partial assignment to a set of $(n + 1)N$ unidimensional random variables that

• **Model structure**



• **Local probability density functions**

$$\begin{aligned}
 \theta_1^{c_1} &= (m_1^{c_1}, 0, v_1^{c_1}) & f_{c_1}(y_1 | \theta_1^{c_1}, \mathbf{s}^h) &\equiv \mathcal{N}(y_1; m_1^{c_1}, v_1^{c_1}) \\
 \theta_2^{c_1} &= (m_2^{c_1}, 0, v_2^{c_1}) & f_{c_1}(y_2 | \theta_2^{c_1}, \mathbf{s}^h) &\equiv \mathcal{N}(y_2; m_2^{c_1}, v_2^{c_1}) \\
 \theta_3^{c_1} &= (m_3^{c_1}, \mathbf{b}_3^{c_1}, v_3^{c_1}) & f_{c_1}(y_3 | y_1, y_2, \theta_3^{c_1}, \mathbf{s}^h) &\equiv \mathcal{N}(y_3; m_3^{c_1} + b_{13}^{c_1}(y_1 - m_1^{c_1}) + b_{23}^{c_1}(y_2 - m_2^{c_1}), v_3^{c_1}) \\
 \mathbf{b}_3^{c_1} &= (b_{13}^{c_1}, b_{23}^{c_1})^t \\
 \theta_1^{c_2} &= (m_1^{c_2}, 0, v_1^{c_2}) & f_{c_2}(y_1 | \theta_1^{c_2}, \mathbf{s}^h) &\equiv \mathcal{N}(y_1; m_1^{c_2}, v_1^{c_2}) \\
 \theta_2^{c_2} &= (m_2^{c_2}, 0, v_2^{c_2}) & f_{c_2}(y_2 | \theta_2^{c_2}, \mathbf{s}^h) &\equiv \mathcal{N}(y_2; m_2^{c_2}, v_2^{c_2}) \\
 \theta_3^{c_2} &= (m_3^{c_2}, \mathbf{b}_3^{c_2}, v_3^{c_2}) & f_{c_2}(y_3 | y_1, y_2, \theta_3^{c_2}, \mathbf{s}^h) &\equiv \mathcal{N}(y_3; m_3^{c_2} + b_{13}^{c_2}(y_1 - m_1^{c_2}) + b_{23}^{c_2}(y_2 - m_2^{c_2}), v_3^{c_2}) \\
 \mathbf{b}_3^{c_2} &= (b_{13}^{c_2}, b_{23}^{c_2})^t
 \end{aligned}$$

• **Multinomial distribution**

$$\begin{aligned}
 p(c_1 | \boldsymbol{\theta}_s, \mathbf{s}^h) \\
 p(c_2 | \boldsymbol{\theta}_s, \mathbf{s}^h) &= 1 - p(c_1 | \boldsymbol{\theta}_s, \mathbf{s}^h)
 \end{aligned}$$

Fig. 1. Model structure, multinomial distribution, and local probability density functions for a CGN for data clustering with three continuous predictive attributes and one binary cluster random variable.

we represent as $\mathbf{D} = \{X_1, \dots, X_N\}$. Let $\mathbf{D}^Y = \{Y_1, \dots, Y_N\}$ denote the set of observed random variables or predictive attributes in \mathbf{D} , that is, the nN unidimensional random variables that have assigned values in \mathbf{d} . Similarly, let $\mathbf{D}^C = \{C_1, \dots, C_N\}$ denote the set of unobserved or cluster random variables in \mathbf{D} , that is, the N unidimensional random variables that reflect the unknown cluster membership of each case of \mathbf{d} . In the forthcoming, \mathbf{d}^Y and \mathbf{d}^C represent an assignment to \mathbf{D}^Y and \mathbf{D}^C , respectively.

For learning CGNs from unlabeled data, the BS-EM algorithm performs a search over the space of CGNs for data clustering based on the well-known EM algorithm [6,23] and the direct optimization of a Bayesian score. As shown in Fig. 2, the BS-EM algorithm follows the basic intuition of the EM algorithm. To take advantage of the best estimate of the generalized joint probability distribution for the data found so far (current CGN for data clustering) in order to estimate missing values and then, to perform a structural search by using a procedure for complete data. Specifically, the BS-EM algorithm alternates between two steps. A step that finds the maximum a posteriori (MAP)

```

loop  $l = 0, 1, \dots$ 
  1. Run the EM algorithm to compute the MAP parameters  $\tilde{\theta}_{s_l}$  for  $s_l$  given  $d^Y$ 
  2. Perform a search over model structures, evaluating every one by
      $Score(s : s_l) = E[\log L(d | s^h) | d^Y, \tilde{\theta}_{s_l}, s_l^h] = \sum_{d^C} \log L(d^Y, d^C | s^h) L(d^C | d^Y, \tilde{\theta}_{s_l}, s_l^h)$ 
  3. Let  $s_{l+1}$  be the model with the highest score among these encountered during the search
  4. if  $Score(s_l : s_l) = Score(s_{l+1} : s_l)$ 
     then return  $s_l$ 

```

Fig. 2. Schematic of the BS-EM algorithm.

parameters for the current CGN structure for data clustering, usually by means of the EM algorithm, and a step that searches over CGN structures for data clustering guided by a Bayesian score. At each iteration, the BS-EM algorithm attempts to maximize the expected Bayesian score instead of the true Bayesian score.

To completely specify the BS-EM algorithm, we have to decide on the structural search procedure (step 2 in Fig. 2). The usual approach is to perform a greedy hill-climbing search over CGN structures for data clustering considering all the possible additions, removals, and reversals of a single arc at each point in the search. This structural search procedure is desirable as it exploits the decomposition properties of CGNs for data clustering and the factorization properties of the Bayesian score. However, any structural search procedure that exploits these properties can be used.

However, the direct application of the BS-EM algorithm, as it appears depicted in Fig. 2, may result in an inefficient and unrealistic solution in order to perform unsupervised CGN selection due to the fact that the computation of $Score(s : s_l)$ implies a huge computational expense as it takes account of every possible completion of the original unlabeled database. This drawback or bottleneck may be very harmful when large databases in number of cases, number of predictive attributes, or both are considered. This is often our case as we are interested in solving data clustering problems of considerable size (medium and large size databases). To overcome this problem, we usually work with a relaxed version of the BS-EM algorithm that only considers what we call the MAP *completion* of the original unlabeled database d to compute $Score(s : s_l)$, instead of considering every possible completion. The MAP completion of d in the l th iteration of the BS-EM algorithm, here denoted by d'_l , is obtained by completing every case of d with the label of the cluster where the maximum of the posterior probability distribution for the cluster random variable C given the case is achieved according to the best model found so far. Thus, this relaxed version of the BS-EM algorithm is comprised of the iteration between a parametric optimization for the current CGN for data clustering, and a structural search after (hidden) cluster random variable completion by means of probabilistic inference with the best estimate of the generalized joint

probability distribution so far (current CGN for data clustering). In this case, the score that guides the selection of the CGN structure for data clustering at each iteration of the BS-EM algorithm reduces to

$$\text{Score}(s : s_l) = \log L(\mathbf{d}'_l | \mathbf{s}^h), \quad (6)$$

where $\log L(\mathbf{d}'_l | \mathbf{s}^h)$ denotes the the log *marginal likelihood* of the MAP completion of the original unlabeled data \mathbf{d} given the structure of the current CGN for data clustering. This quantity can be calculated as follows.

According to [13], under the assumptions that (i) the database restricted to the values of the cluster random variable C , \mathbf{d}^C , is a multinomial sample, (ii) the database \mathbf{d} is complete, and (iii) the parameters of the univariate multinomial distribution for C are independent and follow a Dirichlet distribution, we have that,

$$\begin{aligned} L(\mathbf{d} | \mathbf{s}^h) &= \prod_{l=1}^N \rho(\mathbf{x}_l | \mathbf{x}_1, \dots, \mathbf{x}_{l-1}, \mathbf{s}^h) \\ &= \prod_{l=1}^N p(c_l | \mathbf{x}_1, \dots, \mathbf{x}_{l-1}, \mathbf{s}^h) f(\mathbf{y}_l | c_l, \mathbf{x}_1, \dots, \mathbf{x}_{l-1}, \mathbf{s}^h) \\ &= L(\mathbf{d}^C | \mathbf{s}^h) \prod_{l=1}^N f(\mathbf{y}_l | c_l, \mathbf{y}_1, \dots, \mathbf{y}_{l-1}, \mathbf{s}^h) \\ &= L(\mathbf{d}^C | \mathbf{s}^h) \prod_{c \in \text{Val}(C)} L(\mathbf{d}^{Y,c} | c, \mathbf{s}^h), \end{aligned} \quad (7)$$

where $\mathbf{d}^{Y,c}$ is the database \mathbf{d} restricted to the values for the continuous random variable Y and to cases where $C = c$, and $\text{Val}(C)$ is the set of values that the cluster random variable C can have. The term $L(\mathbf{d}^C | \mathbf{s}^h)$ corresponds to the marginal likelihood of a trivial Bayesian network having only a single node, C . It can be calculated in closed form under reasonable assumptions according to [4]. Moreover, each term of the form $L(\mathbf{d}^{Y,c} | c, \mathbf{s}^h)$ represents the marginal likelihood of a domain containing only continuous random variables under the assumption that the data is sampled from a multivariate normal distribution. Then, these terms can be evaluated in factorable closed form under some reasonable assumptions according to [13,14,17].

3. Compromise conditional Gaussian networks for data clustering

After introducing unrestricted CGNs for data clustering in the previous section, we propose two classes of compromise CGNs for data clustering. These compromise models represent an appealing balance between the cost of the unsupervised model learning process (efficiency) and the quality of the learnt models (effectiveness). In addition to this trade-off between efficiency and effectiveness, some other motives to focus on compromise CGNs for data

clustering include (i) their closeness to human intuition, which results in a better readability and understandability than more general and complex CGNs for data clustering, (ii) the well-known difficulty involved in learning densely connected CGNs for data clustering, and (iii) the painfully slow probabilistic inference and simulation when working with densely connected CGNs for data clustering. Moreover, these drawbacks of densely connected CGNs for data clustering are aggravated when large databases in terms of number of cases, number of predictive attributes, or both are considered for unsupervised model learning or simulation. Thus, the development of methods for learning simple CGNs for data clustering while preserving the expressive power is completely justified.

3.1. *Tree augmented naive Bayes models for data clustering*

The first class of compromise CGNs for data clustering that we consider groups the so-called tree augmented naive Bayes (TANB) models for data clustering. These models are introduced for the first time in [11] as Bayesian classifiers, and are further developed in [10,12,19]. Later, TANB models are used in discrete data clustering problems [31,33]. TANB models for data clustering are defined by the following condition: Predictive attributes may have at most one other predictive attribute as a parent. See Fig. 3 (top) for an explanatory example of the structure of a TANB model for data clustering.

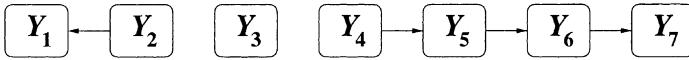
TANB models for data clustering typically imply computational advantages for the unsupervised model learning process when compared with general CGNs for data clustering as $\mathbf{Pa}(s)_i$ is either empty or a set containing only one node for all i . That is, the use of TANB models for data clustering results in a reduction of the model search space which usually decreases the learning cost, whereas the learnt models are still expressive enough. Another advantage of TANB models for data clustering is based on the fact that there is at most one path connecting every pair of predictive attributes. This is very appealing to human intuition due to the obvious conditional (in)dependence statements between subsets of predictive attributes which may be not easy to be read in densely connected CGNs for data clustering.

Notice should be taken that TANB models for data clustering are also known as *mixtures of trees with shared structure* in [24]. The interested reader should turn to this work to find out more about some advantages of these models that are out of the scope of the present work (probabilistic inference, simulation, marginalization, etc.).

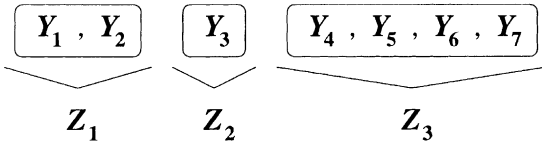
3.2. *Extended naive Bayes models for data clustering*

The well-known class of naive Bayes (NB) models [7,34] conforms one of the simplest examples of compromise CGNs that is usually considered in both

• **TANB model structure for data clustering**



• **ENB model structure for data clustering**



• **Equivalent CGN structure for data clustering in standard representation**

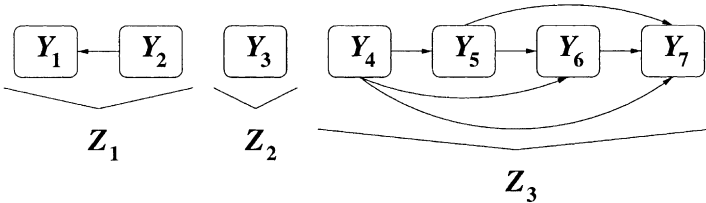


Fig. 3. Structures of a TANB model (top), and an ENB model (middle) and one of its equivalent CGNs in standard representation (bottom), when applied to data clustering.

data classification and data clustering problems. Its name comes from the naive assumption that all the predictive attributes are conditionally independent given the class (data classification) or cluster (data clustering) random variable C . Consequently, the generalized joint probability distribution for X encoded by a NB model for data clustering graphically factorizes as follows:

$$\begin{aligned}
 \rho(\mathbf{x} | \boldsymbol{\theta}_s, \mathbf{s}^h) &= \rho(\mathbf{y}, c | \boldsymbol{\theta}_s, \mathbf{s}^h) \\
 &= p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) f(\mathbf{y} | c, \boldsymbol{\theta}_s, \mathbf{s}^h) \\
 &= p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) f_c(\mathbf{y} | \boldsymbol{\theta}_s^c, \mathbf{s}^h) \\
 &= p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) \prod_{i=1}^n f_c(y_i | \theta_i^c, \mathbf{s}^h).
 \end{aligned}
 \tag{8}$$

Although NB models for data clustering present a behavior rather dependent on the domain to which they are applied, because the assumption that all the predictive attributes are conditionally independent given C is often violated in practice, they perform surprisingly well in many domains. Hence, NB

models for data clustering represent an example of models that provide us with the ideal balance between efficiency and effectiveness that we pursue.

Keeping this idea in mind, extended naive Bayes (ENB) models are introduced by Pazzani [27,28] as Bayesian classifiers and, later, used by Peña et al. [30,32] for data clustering in discrete domains. ENB models for data clustering are very similar to NB models for data clustering as all the random variables represented by the nodes of the structure of an ENB model for data clustering are conditionally independent given the random variable C . The only difference with NB models for data clustering is that the number of nodes in the structure of an ENB model for data clustering can be shorter than the original number of predictive attributes in the database. The reason is that some of these predictive attributes can be grouped together under the same node as fully correlated predictive attributes given C (we call such nodes *supernodes*). In the forthcoming, we refer to the set of nodes and supernodes of an ENB model for data clustering simply as nodes. See Fig. 3 (middle) for an example of an ENB model structure for data clustering.

The structure of an ENB model for data clustering lends itself to a graphical factorization of the generalized joint probability distribution for \mathbf{X} as follows:

$$\begin{aligned}
 \rho(\mathbf{x} | \boldsymbol{\theta}_s, \mathbf{s}^h) &= \rho(\mathbf{y}, c | \boldsymbol{\theta}_s, \mathbf{s}^h) \\
 &= p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) f(\mathbf{y} | c, \boldsymbol{\theta}_s, \mathbf{s}^h) \\
 &= p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) f_c(\mathbf{y} | \boldsymbol{\theta}_s^c, \mathbf{s}^h) \\
 &= p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) \prod_{i=1}^r f_c(\mathbf{z}_i | \boldsymbol{\theta}_i^c, \mathbf{s}^h), \tag{9}
 \end{aligned}$$

where $\{\mathbf{Z}_1, \dots, \mathbf{Z}_r\}$ is a partition of \mathbf{Y} , and r is the number of nodes in \mathbf{s} (including the special nodes referred to as supernodes). Every \mathbf{z}_i is the set of values in \mathbf{y} for the original predictive attributes that are grouped together under the node \mathbf{Z}_i . In this case, $\boldsymbol{\theta}_i^c$ represents the set of parameters for the local probability density functions for the original predictive attributes grouped together under the node \mathbf{Z}_i when $C = c$ for all i . Moreover, \mathbf{s} is a directed acyclic graph that encodes the conditional (in)dependencies among the random variables represented by $\{\mathbf{Z}_1, \dots, \mathbf{Z}_r\}$, and \mathbf{s}^h is defined as usual. Note the similarity between Eqs. (8) and (9). It should be also noticed that in our case we have that for each \mathbf{Z}_i ,

$$f_c(\mathbf{z}_i | \boldsymbol{\theta}_i^c, \mathbf{s}^h) \equiv \mathcal{N}(\mathbf{z}_i; \boldsymbol{\mu}(c), \boldsymbol{\Sigma}(c)), \tag{10}$$

whenever $p(c) = p(C = c) > 0$. Given $C = c$, $\boldsymbol{\mu}(c)$ is the multidimensional mean vector, and $\boldsymbol{\Sigma}(c)$, the variance matrix, is positive definite.

ENB models for data clustering can be considered as having an intermediate place between NB models for data clustering and models with all the predictive

attributes fully correlated for data clustering. Thus, they keep the main features of both extremes. Simplicity from the former models and a better performance from the latter models. Note that both extremes are ENB models for data clustering as well.

To appreciate that ENB models for data clustering constitute a subset of CGNs for data clustering, the reader should note that every ENB model for data clustering can be translated into at least one equivalent CGN for data clustering in standard representation (see Section 2). This can be done by simply taking into account the notion of supernode. Every predictive attribute Y_i is fully correlated with the rest of the predictive attributes grouped together under the same supernode as Y_i given the cluster random variable C , but fully uncorrelated with the predictive attributes grouped together under the remaining supernodes given C . Thus, every supernode represents a local structure for the predictive attributes grouped together under the supernode that encodes no conditional independence assertions among them given C . Fig. 3 (bottom) shows one of the possible translations of the ENB model structure for data clustering depicted in the same figure (middle). Due to the fact that every supernode of a ENB model for data clustering represents a local complete structure for the predictive attributes grouped together under the supernode, the parameters of the ENB model for data clustering can be easily translated into the parameters of its equivalent CGN for data clustering in standard representation according to the work of Shachter and Kenley [35]. See also [13,14] for more details about this translation.

According to our particular objectives, the most interesting advantages derived from the use of ENB models for data clustering are: the saving of computational effort when performing unsupervised model learning due to the obvious reduction of the search space of CGNs for data clustering, and their readability as they keep the simplicity of NB models for data clustering. It is out of the scope of our present work the study of other possible advantages that would appear in the probabilistic inference, simulation, marginalization, etc., when working with ENB models for data clustering.

It should be noticed that not all the sets of conditional (in)dependence assertions encoded by TANB models for data clustering can be encoded by ENB models for data clustering and vice versa. For example, only when every supernode of a given ENB model for data clustering groups at most two original predictive attributes, all the conditional (in)dependencies encoded by the ENB model for data clustering can be encoded by an equivalent TANB model for data clustering. On the other hand, only when the length of every path between predictive attributes in a given TANB model for data clustering is at most one, all the conditional (in)dependence assertions encoded by this TANB model for data clustering can be also represented by an equivalent ENB model for data clustering. Thus, the different expressivity of TANB models and ENB models

for data clustering might imply that the former models appear to be more suitable for some domains while the latter models show a better behavior for some other domains.

4. Experimental evaluation

This section is devoted to the experimental evaluation and comparison of the two proposed classes of compromise CGNs for data clustering using both synthetic and real-world databases. In our experiments, we use the relaxed version of the BS-EM algorithm that just considers the MAP completion of the original unlabeled database in order to compute $\text{Score}(s : s_i)$ at each iteration.

We also provide the reader with insight into the performance of the BS-EM algorithm for learning the classes of compromise CGNs for data clustering that we study. This is done through the performance of several non-parametric statistical tests in order to show statistically significant (in)dependence between the performance criteria considered and several factors.

To completely specify the BS-EM algorithm we have to decide on the structural search procedure (step 2 in Fig. 2). In our experimental evaluation, the BS-EM algorithm performs a greedy hill-climbing search over the proposed compromise model structures. When learning TANB models for data clustering, the structural search considers all the possible additions and removals of one arc at each point in the search. On the other hand, when learning ENB models for data clustering, the structural search considers all possible joints of two nodes and splits of one supernode at each point in the search. The initial model for the BS-EM algorithm is a NB model for data clustering with randomly chosen parameters.

4.1. Performance criteria and general considerations

Table 1 summarizes the criteria that we use to evaluate and compare the TANB models and ENB models for data clustering learnt by the BS-EM algorithm. The log marginal likelihood of the MAP completion of the database given the structure of the initial and the learnt TANB model or ENB model for data clustering is used in our comparison. For every synthetic database, where the original compromise CGN is available, we give special importance to the number of times that the structure of the original model is completely recovered. We also pay attention to the predictive ability of the learnt models measured as the logarithmic scoring rule of good [15]:

$$L_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \log f(y_i | \theta_s, s^h), \quad (11)$$

Table 1
Performance criteria

Expression	Comment
$L_{\text{initial}} \pm \text{S.D.}$	Mean \pm S.D. of the initial score (log marginal likelihood of the MAP completion of the database given the structure of the initial TANB model or ENB model for data clustering)
$L_{\text{final}} \pm \text{S.D.}$	Mean \pm S.D. of the final score (log marginal likelihood of the MAP completion of the database given the structure of the learnt TANB model or ENB model for data clustering)
Recovered	Number of times over all the independent runs for a given synthetic database that the structure of the original TANB model or ENB model is completely recovered
$L_{\text{test}} \pm \text{S.D.}$	Mean \pm S.D. of the multiple predictive accuracy of the learnt TANB model or ENB model for data clustering (logarithmic scoring rule of good)
Time \pm S.D.	Mean \pm S.D. of the runtime (in seconds)

where N_{test} is the number of cases in the testing database \mathbf{d}_{test} , and \mathbf{y}_i denotes the i th case of \mathbf{d}_{test} restricted to the values for the predictive attributes. Moreover, $f(\mathbf{y}_i | \boldsymbol{\theta}_s, \mathbf{s}^h)$ can be calculated as $\sum_{c \in \text{Val}(c)} p(c | \boldsymbol{\theta}_s, \mathbf{s}^h) f_c(\mathbf{y}_i | \boldsymbol{\theta}_i^c, \mathbf{s}^h)$. The higher the value for L_{test} , the higher the multiple predictive accuracy of the elicited compromise CGN for data clustering. In addition to this, we consider the runtime of the unsupervised model learning process as valuable information.

All the average and standard deviation values reported for the performance criteria are over ten independent runs. The values for the performance criterion recovered are also over ten independent runs. All the experiments are run on a Pentium 233 MHz computer. In all the experiments, we assume that the real number of clusters is known beforehand, thus, we do not perform a search to identify the true number of clusters in the databases considered.

The convergence criterion for the EM algorithm that the BS-EM algorithm runs at each iteration is satisfied when either the relative difference between successive values for the log marginal likelihood of the MAP completion of the database given a TANB model or ENB model structure for data clustering is less than 10^{-6} or 150 iterations are reached.

4.2. Results on synthetic databases

In this section, we describe our experimental results on synthetic data. The purpose of evaluating TANB models and ENB models for data clustering on synthetic databases is to show their effectiveness (measured in terms of difference between L_{final} and L_{initial} , ability to completely recover the original structure, recovered, and predictive ability of the learnt models, L_{test}), as well as their efficiency (measured by time).

We considered three factors or dimensions in the construction of each of the synthetic databases:

- Complexity of the structure of the model to be simulated, measured by the number of directed edges in the case of a TANB model and by the number of supernodes in the case of an ENB model.
- Values of the linear coefficients reflecting the strength of the relationships between parents and children in the model to be simulated (see Eq. (5)). In the case of an ENB model, where there is not such relationships, we simulated an equivalent CGN in standard representation (remember Section 3.2).
- Degree of overlapping between the joint probability density functions of the clusters encoded by the model to be simulated. This is a function of the complexity of the model to be simulated and the values of the unconditional means (fixed in our case), linear coefficients, and conditional variances of the predictive attributes involved (see Eq. (5)). Then, we focus on the values for the conditional variances because the complexity of the model and the linear coefficients of the predictive attributes correspond to the other two factors that we study.

There were involved 20 continuous predictive attributes and one three-valued cluster random variable in the original compromise CGNs. We randomly chose four TANB model structures of different complexity. The number of directed edges of each of these four TANB model structures was 0, 5, 10, and 15. Also, we randomly chose four ENB model structures of different complexity. The number of nodes of each of these four ENB model structures was 20, 15, 10, and 5. Therefore, in the last three models some of the 20 original predictive attributes had to be randomly grouped together under supernodes. In fact, these randomly chosen ENB models implied 0, 9, 14, and 31 arcs, respectively, when translated into equivalent CGNs in standard representation as indicated in Section 3.2. We considered three possible values for the linear coefficients reflecting the strength of the relationships between parents and children, these were 1, 0.5, and 0.25. Moreover, the unconditional mean of every predictive attribute was fixed to zero for the first value of the cluster random variable, to four for the second, and to eight for the third. The values considered for the conditional variances of every predictive attribute were 1, 2, and 3. Thus, when the linear coefficients of every predictive attribute were set to one and the conditional variances to three, the degree of overlapping between the joint probability density functions of the clusters achieved the maximum for each of the model complexities considered.

Taking account of all the possible combinations of the values of the three referred factors (four degrees of complexity, three values for the linear coefficients, and three values for the conditional variances), we constructed 36 TANB models and 36 ENB models to be simulated. These 72 models are referred to as the original or underlying models in the forthcoming discussion. The local probability distribution for the cluster random variable was uniform for these 72 compromise CGNs. We sampled 8000 cases from each of these models for the synthetic learning databases and 2000 cases for the testing

databases in order to compute L_{test} . Obviously, we discarded all the entries corresponding to the cluster random variable in the 72 learning databases and 72 testing databases.

We find it convenient to use a regular grammar to clearly identify each of the 72 learning databases. Then, $\text{TANB}\beta_\gamma^\delta$ denotes the database created from the original TANB model with β arcs, $\beta = 0, 5, 10, 15$, when all the linear coefficients take value γ , $\gamma = 1, 0.5, 0.25$, and all the conditional variances are δ , $\delta = 1, 2, 3$. Similarly, $\text{ENB}\beta_\gamma^\delta$ denotes the database generated from the original ENB model with β nodes, $\beta = 20, 15, 10, 5$, being γ and δ as already explained.

For the non-parametric statistical tests performed in order to elicit statistically significant (in)dependence between the performance criteria involved in the evaluation and the three factors considered in the construction of the synthetic databases, we assume $\alpha = 0.05$ as the significance level. However, we report on the P -values obtained in order to allow the reader to judge whether or not the differences are statistically significant. Concretely, χ^2 and Kruskal–Wallis statistical tests are considered.

As Table 2 shows, when learning TANB models for data clustering from the unlabeled synthetic TANB databases, the original model structures are completely recovered most of the times. Moreover, when the structure of the underlying TANB model for a given synthetic database is discovered, the parameters of the original model are also recovered with high accuracy. Thus, the BS-EM algorithm exhibits an effective behavior for learning this class of compromise CGNs for data clustering. According to the χ^2 statistical tests performed between recovered and each of the three factors considered in the construction of the synthetic databases (denoted by complexity, coefficients, and variances in Table 3), the obtained P -values do not support the rejection of the null hypothesis in any case. That is, there is not statistically significant difference between the distributions of the experimental results obtained for recovered with respect to any of the three factors. Thus, not statistically significant dependence between this measure of the effectiveness of TANB models for data clustering and any of the three factors is found.

It clearly appears from Table 2 the great improvement of the log marginal likelihood achieved by the elicited TANB models for data clustering, as it is shown by the difference between L_{final} and L_{initial} . This fact reinforces our assertions about the modeling expressivity of TANB models for data clustering. According to the P -values reported in Table 3, the Kruskal–Wallis statistical tests carried out show that there are statistically significant differences between the distributions of the values of the difference $L_{\text{final}} - L_{\text{initial}}$ with respect to complexity and coefficients. Thus, there exists statistically significant dependence between this measure of the effectiveness and the complexity and the linear coefficients of the underlying model. The improvement of the log marginal likelihood increases when the complexity or the linear coefficients increase(s). On the other hand, the performed Kruskal–Wallis statistical tests do

Table 2
Performance achieved when learning TANB models for data clustering from the 36 synthetic TANB databases

Database	$L_{\text{initial}} \pm \text{S.D.}$	$L_{\text{final}} \pm \text{S.D.}$	Recovered	$L_{\text{test}} \pm \text{S.D.}$	Time \pm S.D.
TANB0 ¹	-145407 \pm 17723	-129407 \pm 5551	8	-13.13 \pm 0.67	372 \pm 583
TANB0 ¹ _{0,5}	-137937 \pm 17253	-128021 \pm 4168	9	-12.96 \pm 0.51	346 \pm 810
TANB0 ¹ _{0,25}	-143052 \pm 15123	-126632 \pm 0	10	-12.79 \pm 0.00	84 \pm 35
TANB0 ²	-174686 \pm 10899	-156398 \pm 4503	8	-16.08 \pm 0.56	378 \pm 606
TANB0 ² _{0,5}	-173937 \pm 12302	-155296 \pm 3446	9	-15.94 \pm 0.42	286 \pm 571
TANB0 ² _{0,25}	-173544 \pm 10578	-156468 \pm 4642	8	-16.08 \pm 0.56	449 \pm 767
TANB0 ³	-181671 \pm 8306	-172230 \pm 3971	8	-17.80 \pm 0.48	307 \pm 457
TANB0 ³ _{0,5}	-181296 \pm 9251	-171213 \pm 2906	9	-17.68 \pm 0.36	230 \pm 448
TANB0 ³ _{0,25}	-189797 \pm 13783	-171226 \pm 2944	9	-17.68 \pm 0.36	188 \pm 280
TANB5 ¹	-146001 \pm 13975	-126671 \pm 0	10	-12.79 \pm 0.00	123 \pm 15
TANB5 ¹ _{0,5}	-147909 \pm 15852	-127808 \pm 3420	9	-12.93 \pm 0.42	201 \pm 252
TANB5 ¹ _{0,25}	-168741 \pm 25245	-126667 \pm 0	10	-12.79 \pm 0.00	127 \pm 12
TANB5 ²	-176656 \pm 10439	-155022 \pm 2498	9	-15.90 \pm 0.31	273 \pm 429
TANB5 ² _{0,5}	-172351 \pm 13249	-155131 \pm 2832	9	-15.92 \pm 0.34	212 \pm 268
TANB5 ² _{0,25}	-171420 \pm 10298	-156265 \pm 4157	8	-16.05 \pm 0.50	394 \pm 546
TANB5 ³	-191028 \pm 8035	-172423 \pm 3260	7	-17.82 \pm 0.39	306 \pm 265
TANB5 ³ _{0,5}	-183830 \pm 9192	-171084 \pm 2390	9	-17.66 \pm 0.28	255 \pm 372
TANB5 ³ _{0,25}	-186011 \pm 10597	-172946 \pm 4062	7	-17.88 \pm 0.48	424 \pm 454
TANB10 ¹	-173421 \pm 16617	-126719 \pm 0	10	-12.79 \pm 0.00	176 \pm 21
TANB10 ¹ _{0,5}	-162582 \pm 11058	-129465 \pm 4206	7	-13.12 \pm 0.50	438 \pm 446
TANB10 ¹ _{0,25}	-140031 \pm 17182	-128996 \pm 4572	8	-13.06 \pm 0.54	377 \pm 451
TANB10 ²	-187000 \pm 6669	-155813 \pm 2404	7	-15.99 \pm 0.29	505 \pm 456

(continued on next page)

Table 2 (Continued)

Database	$L_{\text{initial}} \pm \text{S.D.}$	$L_{\text{final}} \pm \text{S.D.}$	Recovered	$L_{\text{rest}} \pm \text{S.D.}$	Time \pm S.D.
TANB10 _{0.5} ²	-175807 \pm 16556	-155695 \pm 2918	8	-15.97 \pm 0.35	387 \pm 452
TANB10 _{0.25} ²	-166660 \pm 12249	-156109 \pm 3748	8	-16.02 \pm 0.44	330 \pm 365
TANB10 ₁ ³	-195496 \pm 3861	-170338 \pm 0	10	-17.56 \pm 0.00	282 \pm 36
TANB10 _{0.5} ³	-184310 \pm 11060	-170937 \pm 1797	9	-17.63 \pm 0.22	307 \pm 328
TANB10 _{0.25} ³	-183411 \pm 7326	-171128 \pm 2371	9	-17.65 \pm 0.28	302 \pm 449
TANB15 ₁ ¹	-165994 \pm 11718	-127382 \pm 1250	8	-12.87 \pm 0.15	327 \pm 119
TANB15 _{0.5} ¹	-142337 \pm 13188	-127419 \pm 2015	9	-12.87 \pm 0.24	260 \pm 236
TANB15 _{0.25} ¹	-143321 \pm 15604	-126745 \pm 0	10	-12.79 \pm 0.00	179 \pm 17
TANB15 ₁ ²	-194388 \pm 9055	-154266 \pm 0	10	-15.80 \pm 0.00	346 \pm 49
TANB15 _{0.5} ²	-168293 \pm 8935	-154275 \pm 0	10	-15.80 \pm 0.00	242 \pm 113
TANB15 _{0.25} ²	-166754 \pm 10621	-154274 \pm 0	10	-15.80 \pm 0.00	182 \pm 9
TANB15 ₁ ³	-206382 \pm 6555	-170341 \pm 0	10	-17.58 \pm 0.00	950 \pm 146
TANB15 _{0.5} ³	-192429 \pm 5549	-171188 \pm 1623	8	-17.66 \pm 0.20	490 \pm 289
TANB15 _{0.25} ³	-188984 \pm 10148	-172390 \pm 3073	7	-17.80 \pm 0.37	473 \pm 368

All the average and standard deviation values are over ten independent runs. The values for recovered are also over ten independent runs.

Table 3

P-values for the non-parametric statistical tests performed when evaluating the differences in the performance criteria values obtained from the 36 synthetic TANB databases

	Recovered	Complexity	Coefficients	Variances
Recovered	–	0.595 χ^2	0.975 χ^2	0.498 χ^2
$L_{\text{final}} - L_{\text{initial}}$	0.318 K–W	0.000 K–W	0.000 K–W	0.052 K–W
L_{test}	0.000 K–W	0.001 K–W	0.920 K–W	0.000 K–W
Time	0.000 K–W	0.000 K–W	0.031 K–W	0.016 K–W

χ^2 indicates the performance of a χ^2 statistical test and K–W indicates the performance of a Kruskal–Wallis statistical test.

not show statistically significant dependence between the improvement of the log marginal likelihood and recovered.

Table 2 also depicts the results obtained for the predictive ability of the learnt TANB models for data clustering on the testing databases, L_{test} . Statistically significant dependence exists between L_{test} and the ability to recover the underlying model, the complexity of the original TANB model, and the conditional variances of the predictive attributes, according to the Kruskal–Wallis statistical tests performed (see Table 3). The increase of the ability to discover the underlying model for a given database implies the increase of the predictive ability of the learnt TANB model for data clustering. On the other hand, the predictive ability of the learnt model decreases with the increase of the conditional variances of the predictive attributes of the underlying TANB model. As a result, we could conclude that, from the point of view of the effectiveness, TANB models for data clustering offer a good performance level as well as high reliability to recover the underlying models from the synthetic databases considered.

From the point of view of the efficiency (measured as the runtime of the unsupervised model learning process, time), our experimental results show empirical evidence of the trade-off between effectiveness and efficiency that TANB models for data clustering provide with: the runtime of the BS-EM algorithm is reasonable, taking into account the considerable dimensionality and number of cases of the synthetic databases in addition to the, sometimes, complex underlying models. According to the *P*-values reported in Table 3 for the Kruskal–Wallis statistical tests performed, statistically significant dependence between the cost of the unsupervised model learning process, i.e., runtime of the BS-EM algorithm, and each of the three factors exists. Also, the results of the tests reveal statistically significant dependence between the efficiency and the ability to recover the underlying model of a given database. As one might expect, the cost of the unsupervised model learning process increases with the increase of the complexity, the linear coefficients, or the conditional variances of the predictive attributes of the original TANB model, while decreasing as recovered increases.

Table 4 shows the performance achieved when learning ENB models for data clustering from the unlabeled synthetic ENB databases by means of the BS-EM algorithm. Roughly speaking, the results achieved on the synthetic ENB databases are similar to those presented in Table 2 for the models elicited from the synthetic TANB databases. Thus, TANB models for data clustering as well as ENB models for data clustering provide us with an appealing balance between effectiveness and efficiency. However, in the case of learning ENB models for data clustering the results appear to be more noticeable as the complexity of the original ENB models was greater than the complexity of the original TANB models. Recall that the original ENB models had equivalent CGNs in standard representation with 0, 9, 14, and 31 arcs, while the original TANB models had only 0, 5, 10, and 15 arcs.

The only ENB databases where the results are not satisfying at all are $ENB5_1^2$ and $ENB5_1^3$. The results reflect the inability of ENB models for data clustering to capture the underlying models as the overlapping between the joint probability density functions of the three existing clusters was considerable in these databases. This can be appreciated from the complexity of the models that were used to generate the databases $ENB5_1^2$ and $ENB5_1^3$ (only five nodes and all of them supernodes), the linear coefficients (equal to one), and the conditional variances of the predictive attributes (equal to two for the former database and to three for the latter). Remember Eq. (5) to appreciate the difficulty involved in recovering the original models from the databases $ENB5_1^2$ and $ENB5_1^3$.

With regard to the non-parametric statistical tests performed for the results obtained when learning ENB models for data clustering from the ENB databases via the BS-EM algorithm, we can say that they are very similar to those that we report in Table 3. However, some differences appear from the comparison between this table and Table 5. It is noticeable that the statistically significant dependence that exists between recovered and each of the three factors considered in the generation of the synthetic databases, when learning ENB models for data clustering. Remember that such statistically significant dependence is not observed when TANB models for data clustering are induced from the TANB databases. Concretely, the rate of recovered decreases as the conditional variances of the predictive attributes of the underlying model increase. Another difference between Tables 3 and 5 is the statistically significant dependence between the improvement of the log marginal likelihood achieved and the ability to recover the underlying model of a given database. To be more exact, as the ability to recover the original ENB model increases, the improvement of the log marginal likelihood reduces. The reason is that an initial model that has a poor log marginal likelihood may make the unsupervised model learning algorithm unable to discover the underlying model. However, the elicited model, despite being different from the original model, can improve substantially the initial model in terms of log marginal likelihood.

Table 4
Performance achieved when learning ENB models for data clustering from the 36 synthetic ENB databases

Database	$L_{initial} \pm S.D.$	$L_{final} \pm S.D.$	Recovered	$L_{test} \pm S.D.$	Time \pm S.D.
ENB20 ¹	-139968 \pm 16709	-126631 \pm 0	10	-12.79 \pm 0.00	170 \pm 25
ENB20 ¹ _{0.5}	-132789 \pm 12195	-126632 \pm 0	10	-12.79 \pm 0.00	155 \pm 20
ENB20 ¹ _{0.25}	-135483 \pm 18267	-126631 \pm 0	10	-12.79 \pm 0.00	169 \pm 15
ENB20 ²	-168420 \pm 14391	-154603 \pm 1367	9	-15.84 \pm 0.10	222 \pm 142
ENB20 ² _{0.5}	-167999 \pm 15185	-154147 \pm 0	10	-15.80 \pm 0.00	187 \pm 27
ENB20 ² _{0.25}	-165552 \pm 13372	-154604 \pm 1371	9	-15.84 \pm 0.10	197 \pm 71
ENB20 ³	-184454 \pm 9316	-171068 \pm 1647	8	-17.62 \pm 0.12	392 \pm 514
ENB20 ³ _{0.5}	-181481 \pm 9833	-171067 \pm 1643	8	-17.62 \pm 0.11	253 \pm 135
ENB20 ³ _{0.25}	-185655 \pm 8843	-170650 \pm 1217	9	-17.59 \pm 0.09	265 \pm 125
ENB15 ¹	-153876 \pm 20567	-127074 \pm 1430	9	-12.86 \pm 0.12	306 \pm 199
ENB15 ¹ _{0.5}	-140303 \pm 14803	-127078 \pm 1450	9	-12.86 \pm 0.13	307 \pm 237
ENB15 ¹ _{0.25}	-151038 \pm 20620	-127593 \pm 1997	8	-12.90 \pm 0.18	533 \pm 610
ENB15 ²	-184477 \pm 15610	-154525 \pm 1227	9	-15.85 \pm 0.09	307 \pm 106
ENB15 ² _{0.5}	-175527 \pm 12872	-154530 \pm 1245	9	-15.86 \pm 0.10	391 \pm 408
ENB15 ² _{0.25}	-171009 \pm 13400	-154883 \pm 1552	8	-15.89 \pm 0.13	314 \pm 175
ENB15 ³	-191186 \pm 8337	-170949 \pm 1467	8	-17.64 \pm 0.10	657 \pm 872
ENB15 ³ _{0.5}	-184510 \pm 11663	-171349 \pm 1733	7	-17.67 \pm 0.12	656 \pm 664
ENB15 ³ _{0.25}	-182646 \pm 8536	-170993 \pm 1557	8	-17.64 \pm 0.11	573 \pm 884
ENB10 ¹	-165082 \pm 12643	-127231 \pm 1447	9	-12.84 \pm 0.13	425 \pm 514
ENB10 ¹ _{0.5}	-147749 \pm 13386	-126741 \pm 0	10	-12.79 \pm 0.00	243 \pm 18
ENB10 ¹ _{0.25}	-151382 \pm 15171	-127209 \pm 1410	9	-12.84 \pm 0.12	371 \pm 383
ENB10 ²	-192655 \pm 7961	-155099 \pm 1652	8	-15.87 \pm 0.13	347 \pm 128

(continued on next page)

Table 4 (Continued)

Database	$L_{\text{initial}} \pm \text{S.D.}$	$L_{\text{final}} \pm \text{S.D.}$	Recovered	$L_{\text{est}} \pm \text{S.D.}$	Time \pm S.D.
ENB1 $0^2_{0,5}$	-166961 \pm 10491	-154269 \pm 0	10	-15.80 \pm 0.00	243 \pm 34
ENB1 $0^2_{0,25}$	-171971 \pm 13195	-155068 \pm 1600	8	-15.87 \pm 0.12	389 \pm 261
ENB1 0^3_1	-204529 \pm 5266	-171498 \pm 1715	7	-17.64 \pm 0.12	506 \pm 261
ENB1 $0^3_{0,5}$	-187853 \pm 7283	-170373 \pm 0	10	-17.57 \pm 0.00	328 \pm 46
ENB1 $0^3_{0,25}$	-179607 \pm 8690	-170731 \pm 1077	9	-17.59 \pm 0.07	306 \pm 171
ENB5 1_1	-200981 \pm 7884	-127384 \pm 1695	9	-12.87 \pm 0.16	521 \pm 408
ENB5 $^1_{0,5}$	-165111 \pm 11959	-126861 \pm 0	10	-12.82 \pm 0.00	274 \pm 19
ENB5 $^1_{0,25}$	-153217 \pm 18286	-127589 \pm 1585	8	-12.89 \pm 0.14	876 \pm 1450
ENB5 2_1	-222545 \pm 5611	-154675 \pm 0	0	-15.83 \pm 0.00	536 \pm 106
ENB5 $^2_{0,5}$	-189427 \pm 6562	-154656 \pm 941	9	-15.85 \pm 0.07	394 \pm 84
ENB5 $^2_{0,25}$	-172131 \pm 12697	-154664 \pm 973	9	-15.85 \pm 0.07	340 \pm 207
ENB5 3_1	-236717 \pm 4453	-171492 \pm 0	0	-17.60 \pm 0.00	707 \pm 81
ENB5 $^3_{0,5}$	-200792 \pm 4858	-170455 \pm 0	10	-17.59 \pm 0.00	543 \pm 118
ENB5 $^3_{0,25}$	-191012 \pm 14548	-170738 \pm 855	9	-17.61 \pm 0.06	484 \pm 251

All the average and standard deviation values are over ten independent runs. The values for recovered are also over ten independent runs.

Table 5

P-values for the non-parametric statistical tests performed when evaluating the differences in the performance criteria values obtained from the 36 synthetic ENB databases

	Recovered	Complexity	Coefficients	Variances
Recovered	–	0.001 χ^2	0.000 χ^2	0.005 χ^2
$L_{\text{final}} - L_{\text{initial}}$	0.000 K–W	0.000 K–W	0.000 K–W	0.388 K–W
L_{test}	0.000 K–W	0.000 K–W	0.778 K–W	0.000 K–W
Time	0.000 K–W	0.000 K–W	0.023 K–W	0.000 K–W

χ^2 indicates the performance of a χ^2 statistical test and K–W indicates the performance of a Kruskal–Wallis statistical test.

A closer look at the learnt compromise CGNs for data clustering reveals the decisive influence of the quality of the initial models on the quality of the learnt models (see the sometimes large standard deviation values for the performance criteria that we report in Table 2 and 4). Hence, it appears to be specially relevant the study of the initial conditions in order to improve the results in terms of all the performance measures considered in this work.

4.3. Results on real-world databases

Another source of data for our evaluation consists of three well-known real-world databases from the UCI repository of machine learning databases [26]. These databases provide us with a more realistic framework for the evaluation. Moreover, they allow us to compare the performance achieved by TANB models and ENB models for data clustering in the same domains. The databases considered are the following ones:

- Waveform [3] which is an artificial database consisting of 21 predictive attributes. There are three classes. Instances of each class represent a random convex combination of two of three base waves with noise added. We used the data set generator from the UCI repository to create a learning database with 4000 cases and a testing database with 1000 cases.
- Waveform21 + 19 [3] which is an artificial database consisting of 40 predictive attributes. The first 21 predictive attributes of the instances of each class represent a random convex combination of two of three base waves with noise added. The last 19 predictive attributes are noise attributes which turn out to be irrelevant for describing the underlying three classes. We used the data set generator from the UCI repository to create a learning database with 4000 cases and a testing database with 1000 cases.
- Pima [36] which is a real database containing 768 cases and eight predictive attributes. There are two classes. Every instance represents the record of a female patient tested for diabetes. We used the first 600 cases for learning and the last 168 cases for testing.

The first two databases were chosen due to our interest in working with databases of considerable size (thousands of cases and tens of predictive

attributes). The third database, considerably shorter in both number of cases and number of predictive attributes, was chosen to get feedback on the scalability of the results. Obviously, we deleted all the class entries for the three learning and three testing databases.

As in the evaluation on synthetic data, we assume $\alpha = 0.05$ as the significance level for the non-parametric statistical tests performed in order to elicit statistically significant differences in the results for the performance criteria involved in the evaluation with respect to the domains considered. However, we report on the P -values obtained in order to allow the reader to judge whether or not the differences are statistically significant. Concretely, Mann–Whitney statistical tests are considered.

Table 6 summarizes the results obtained by the BS-EM algorithm when learning TANB models and ENB models for data clustering from the three real-world databases. This table reveals the good performance of both classes of compromise CGNs for data clustering when applied to real-world domains, as well as the suitability of the BS-EM algorithm for unsupervised model induction in these domains. By comparing the performance reached by the elicited TANB models and ENB models for data clustering, we can conclude that in the two largest and, thus, more interesting databases, the waveform and waveform21 + 19 databases, the learnt ENB models for data clustering outperform the obtained TANB models for data clustering in terms of effectiveness (L_{final} and L_{test}). In the third database, however, the results address a better behavior of the induced TANB models for data clustering over the obtained ENB models for data clustering. Thus, ENB models for data clustering appear to be more suitable than TANB models for data clustering for the two largest real-world databases considered, while in the third one the latter models behave in a better way than the former models.

According to the P -values that Table 7 reports for the Mann–Whitney statistical tests performed for the results obtained, there exist statistically significant differences between the distributions of L_{test} for the TANB models

Table 6

Performance achieved when learning TANB models and ENB models for data clustering from the three real-world databases

Database	Model	$L_{\text{initial}} \pm \text{S.D.}$	$L_{\text{final}} \pm \text{S.D.}$	$L_{\text{test}} \pm \text{S.D.}$	Time \pm S.D.
Waveform	TANB	-72069 ± 1765	-67531 ± 3	-14.21 ± 0.00	651 ± 265
	ENB	-72729 ± 2518	-67255 ± 7	-14.03 ± 0.01	479 ± 72
Waveform21 + 19	TANB	-127109 ± 2032	-120981 ± 0	-25.90 ± 0.00	672 ± 79
	ENB	-127259 ± 1791	-120713 ± 0	-25.72 ± 0.00	1455 ± 111
Pima	TANB	-8372 ± 65	-7706 ± 393	-11.56 ± 0.91	9 ± 3
	ENB	-8324 ± 52	-7946 ± 317	-12.10 ± 0.72	10 ± 2

All the average and standard deviation values are over ten independent runs.

Table 7

P-values for the different Mann–Whitney statistical tests performed when comparing TANB models and ENB models for data clustering according to the differences in the performance criteria values obtained from the three real-world databases

	$L_{\text{final}} - L_{\text{initial}}$	L_{test}	Time
Waveform	0.326	0.000	0.199
Waveform21 + 19	0.496	0.000	0.000
Pima	0.070	0.067	0.496

and ENB models for data clustering learnt from the waveform and waveform21 + 19 databases via the BS-EM algorithm. Also, statistically significant difference between the distributions of the cost of the unsupervised model induction process when learning TANB models and ENB models for data clustering from the waveform21 + 19 database is found. Summarizing, the increase in the predictive ability reached by the elicited ENB models for data clustering with respect to the predictive ability achieved by the induced-TANB models for data clustering for the waveform and waveform21 + 19 databases is statistically significant. Also the increase of the cost of learning ENB models for data clustering with respect to the cost of learning TANB models for data clustering from the waveform21 + 19 database is statistically significant. Thus, there are statistically significant differences between the effectiveness and efficiency of the TANB models and ENB models for data clustering elicited from these databases. Remember our comments about the different expressivity of TANB models and ENB models for data clustering that appear in Section 3.2.

As done in the case of the synthetic databases, it seems interesting to inspect the standard deviation values for the different performance criteria that we use for the real-world databases. Contrary to what happened with the synthetic domains, the standard deviations reported in Table 6 have very small values, which reflects the robustness of the BS-EM algorithm in these databases.

5. Conclusions

We have proposed two classes of compromise CGNs for data clustering: TANB models and ENB models for data clustering. Some of their advantages have been addressed: trade-off between efficiency and effectiveness, computational advantages, readability, interpretability, closeness to human intuition, etc. Moreover, an experimental evaluation and comparison of both classes of compromise CGNs for data clustering on synthetic and real-world databases have been carried out. The results obtained on the synthetic databases have shown the modeling power of TANB models and ENB models for data

clustering as, most of the times, the underlying model was completely discovered while the cost of the unsupervised model learning process was kept at a reasonable level. From the results achieved on real-world databases, it is noticeable the slight outperformance of the ENB models for data clustering learnt from the two largest databases over the obtained TANB models for data clustering.

In addition to this, we have provided the reader with insight into the performance of the BS-EM algorithm for learning the classes of compromise CGNs for data clustering that we have studied in this work. This has been achieved by performing several non-parametric statistical tests in order to show statistically significant (in)dependence between the performance criteria considered and several factors.

Some lines of future research could be the study of some other advantages of TANB models and ENB models for data clustering (probabilistic inference, simulation, marginalization, etc.), a proper study of the influence of the initial conditions, and/or the evaluation of the performance of TANB models and ENB models for data clustering with respect to some other factors (number of predictive attributes, number of clusters, number of cases, etc.).

Acknowledgements

The authors would like to thank the anonymous referees for their useful comments related to this work. This work was supported by the Spanish Ministry of Education and Culture (Ministerio de Educación y Cultura) under AP97 44673053 grant.

References

- [1] M.R. Anderberg, *Cluster Analysis for Applications*, Academic Press, New York, 1973.
- [2] J. Banfield, A. Raftery, Model-based Gaussian and non-Gaussian clustering, *Biometrics* 49 (1993) 803–821.
- [3] L. Breiman, J.H. Friedman, R.A. Olshen, C.J. Stone, *Classification and Regression Trees*, Wadsworth International Group, Belmont, CA, 1984.
- [4] G. Cooper, E. Herskovits, A Bayesian method for the induction of probabilistic networks from data, *Mach. Learn.* 9 (1992) 309–347.
- [5] R.G. Cowell, A.P. Dawid, S.L. Lauritzen, D.J. Spiegelhalter, *Probabilistic Networks and Expert Systems*, Springer, New York, 1999.
- [6] A. Dempster, N. Laird, D. Rubin, Maximum likelihood from incomplete data via the EM algorithm, *J. R. Stat. Soc. B* 39 (1977) 1–38.
- [7] R. Duda, P. Hart, *Pattern Classification and Scene Analysis*, Wiley, New York, 1973.
- [8] D. Fisher, Knowledge acquisition via incremental conceptual clustering, *Mach. Learn.* 2 (1987) 139–172.

- [9] N. Friedman, The Bayesian structural EM algorithm, in: *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann, San Francisco, CA, 1998, pp. 129–138.
- [10] N. Friedman, D. Geiger, M. Goldszmidt, Bayesian network classifiers, *Mach. Learn.* 29 (1997) 131–163.
- [11] N. Friedman, M. Goldszmidt, Building classifiers using Bayesian networks, in: *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, AAAI Press, Menlo Park, CA, 1996, pp. 1277–1284.
- [12] N. Friedman, M. Goldszmidt, T. Lee, Bayesian network classification with continuous attributes: getting the best of both discretization and parametric fitting, in: *Proceedings of the Fifteenth National Conference on Machine Learning*, 1998.
- [13] D. Geiger, D. Heckerman, Learning Gaussian networks, Technical report MSR-TR-94-10, Microsoft Research, Redmond, WA, 1994.
- [14] D. Geiger, D. Heckerman, Learning Gaussian networks, in: *Proceedings of the Tenth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann, Seattle, WA, 1995, pp. 235–243.
- [15] I. Good, Rational decisions, *J. R. Stat. Soc. B* 14 (1952) 107–114.
- [16] J.A. Hartigan, *Clustering Algorithms*, Wiley, Canada, 1975.
- [17] D. Heckerman, D. Geiger, Likelihoods and parameter priors for Bayesian networks, Technical report MSR-TR-95-54, Microsoft Research, Redmond, WA, 1995.
- [18] L. Kaufman, P. Rousseeuw, *Finding Groups in Data*, Wiley, New York, 1990.
- [19] E.J. Keogh, M.J. Pazzani, Learning augmented Bayesian classifiers: a comparison of distribution-based and classification-based approaches, in: *Proceedings of the Seventh International Workshop on Artificial Intelligence and Statistics*, 1999, pp. 225–230.
- [20] S.L. Lauritzen, Propagation of probabilities, means and variances in mixed graphical association models, *J. Am. Stat. Assoc.* 87 (1992) 1098–1108.
- [21] S.L. Lauritzen, *Graphical Models*, Clarendon Press, Oxford, 1996.
- [22] S.L. Lauritzen, N. Wermuth, Graphical models for associations between variables, some of which are qualitative and some quantitative, *Ann. Stat.* 17 (1989) 31–57.
- [23] G.J. McLachlan, T. Krishnan, *The EM Algorithm and Extensions*, Wiley, New York, 1997.
- [24] M. Meilä, *Learning with Mixtures of Trees*, Ph.D. thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, 1999.
- [25] M. Meilä, M.I. Jordan, Estimating dependency structure as a hidden variable, *Neural Inf. Process. Syst.* 10 (1997) 584–590.
- [26] C. Merz, P. Murphy, D. Aha, UCI repository of machine learning databases, Department of Information and Computer Science, University of California, Irvine, CA, 1997. <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
- [27] M.J. Pazzani, Constructive induction of cartesian product attributes, *Inf. Stat. Induction Sci.* (1996).
- [28] M.J. Pazzani, Searching for dependencies in Bayesian classifiers, *Learning from Data: Artificial Intelligence and Statistics V*, Springer, New York, 1996, pp. 239–248.
- [29] J.M. Peña, I. Izarzugaza, J.A. Lozano, E. Aldasoro, P. Larrañaga, Geographical clustering of cancer incidence by means of Bayesian networks and conditional Gaussian networks, in: *Proceedings of the Eighth International Workshop on Artificial Intelligence and Statistics*, Morgan Kaufmann, San Francisco, CA, 2001, pp. 266–271.
- [30] J.M. Peña, J.A. Lozano, P. Larrañaga, Learning Bayesian networks for clustering by means of constructive induction, *Pattern Recognition Lett.* 20 (11–13) (1999) 1219–1230.
- [31] J.M. Peña, J.A. Lozano, P. Larrañaga, An improved Bayesian structural EM algorithm for learning Bayesian networks for clustering, *Pattern Recognition Lett.* 21 (8) (2000) 779–786.

- [32] J.M. Peña, J.A. Lozano, P. Larrañaga, Learning recursive Bayesian multinets for data clustering by means of constructive induction, *Mach. Learn.* (2001) to appear.
- [33] J.M. Peña, J.A. Lozano, P. Larrañaga, I. Inza, Dimensionality reduction in unsupervised learning of conditional Gaussian networks, *IEEE Trans. Pattern Anal. Machine Intell.* 23 (6) (2001) 590–603.
- [34] M. Peot, Geometric implications of the Naive Bayes assumption, in: *Proceedings of the Twelfth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann, San Francisco, CA, 1996, pp. 414–419.
- [35] R. Shachter, C. Kenley, Gaussian influence diagrams, *Manage. Sci.* 35 (1989) 527–550.
- [36] J.W. Smith, J.E. Everhart, W.C. Dickson, W.C. Knowler, R.S. Johannes, Using the ADAP learning algorithm to forecast the onset of diabetes mellitus, in: *Proceedings of the Symposium on Computer Applications and Medical Care*, IEEE Computer Society Press, Silver Spring, MD, 1988, pp. 261–265.
- [37] B. Thiesson, C. Meek, D.M. Chickering, D. Heckerman, Learning mixtures of DAG models, in: *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann, San Francisco, CA, 1998, pp. 504–513.