

# Discrete Bayesian Network Classifiers: A Survey

CONCHA BIELZA and PEDRO LARRAÑAGA, Universidad Politécnica de Madrid

We have had to wait over 30 years since the naive Bayes model was first introduced in 1960 for the so-called Bayesian network classifiers to resurge. Based on Bayesian networks, these classifiers have many strengths, like model interpretability, accommodation to complex data and classification problem settings, existence of efficient algorithms for learning and classification tasks, and successful applicability in real-world problems. In this article, we survey the whole set of discrete Bayesian network classifiers devised to date, organized in increasing order of structure complexity: naive Bayes, selective naive Bayes, seminaive Bayes, one-dependence Bayesian classifiers,  $k$ -dependence Bayesian classifiers, Bayesian network-augmented naive Bayes, Markov blanket-based Bayesian classifier, unrestricted Bayesian classifiers, and Bayesian multinets. Issues of feature subset selection and generative and discriminative structure and parameter learning are also covered.

Categories and Subject Descriptors: I.5.1 [Pattern Recognition]: Models

General Terms: Algorithms, Design, Performance

Additional Key Words and Phrases: Supervised classification, Bayesian network, naive Bayes, Markov blanket, Bayesian multinets, feature subset selection, generative and discriminative classifiers

## ACM Reference Format:

Concha Bielza and Pedro Larrañaga. 2014. Discrete Bayesian network classifiers: A survey. *ACM Comput. Surv.* 47, 1, Article 60 (April 2014), 43 pages.

DOI: <http://dx.doi.org/10.1145/2576868>

## 1. INTRODUCTION

Bayesian network classifiers are special types of Bayesian networks designed for classification problems. Supervised classification aims at assigning labels or categories to instances described by a set of predictor variables or features. The classification model that assigns labels to instances is automatically induced from a dataset containing labeled instances or sometimes by hand with the aid of an expert. We will focus on learning models from data, favored by the large amount of data collected and accessible nowadays.

Bayesian network classifiers have many advantages over other classification techniques, as follows: (1) They offer an explicit, graphical, and interpretable representation of uncertain knowledge. Their semantics is based on the sound concept of conditional independence since they are an example of a probabilistic graphical model. (2) As they output a probabilistic model, decision theory is naturally applicable for dealing with cost-sensitive problems, thereby providing a confidence measure on the chosen predicted label. (3) Thanks to the model expressiveness of Bayesian network classifiers,

Research partially supported by the Spanish Ministry of Economy and Competitiveness, projects TIN2010-20900-C04-04 and Cajal Blue Brain.

Author's addresses: Concha Bielza and Pedro Larrañaga, Departamento de Inteligencia Artificial, Universidad Politécnica de Madrid.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2014 ACM 0360-0300/2014/04-ART60 \$15.00

DOI: <http://dx.doi.org/10.1145/2576868>

Q1

36 they can easily accommodate feature selection methods and handle missing data in  
37 both learning and inference phases. Also, they fit more complex classification problems  
38 in any type of domain (discrete, continuous, and mixed data), with undetermined la-  
39 bels, partial labels, many class variables to be simultaneously predicted, new flows of  
40 streaming data, and so forth. (4) There is an active research field developing a plethora  
41 of learning from data algorithms, covering different frequentist and Bayesian, expert,  
42 and/or data-based viewpoints. Besides, the induced models can be organized hierar-  
43 chically according to their structure complexity. (5) Bayesian network classifiers can  
44 be built with computationally efficient algorithms whose learning time complexity is  
45 linear on the number of instances and linear, quadratic, or cubic (depending on model  
46 complexity) on the number of variables, and whose classification time is linear on the  
47 number of variables. (6) These algorithms are easily implemented, although most of the  
48 available software only contains the simplest options (naive Bayes and tree-augmented  
49 naive Bayes), focusing instead on learning general-purpose Bayesian networks. (7) Nu-  
50 merous successful real-world applications have been reported in the literature, with  
51 competitive performance results against state-of-the-art classifiers.

52 This article offers a comprehensive survey of the state of the art of the Bayesian  
53 network classifier in discrete domains. Unlike other reviews mentioned later, this arti-  
54 cle covers many model specificities: (1) for naive Bayes, its weighted version, inclusion  
55 of hidden variables, metaclassifiers, special situations like homologous sets, multiple  
56 instances, cost-sensitive problems, instance ranking, imprecise probabilities, text cat-  
57 egorization, and discriminative learning of parameters; (2) for selective naive Bayes,  
58 univariate and multivariate filter approaches and wrapper and embedded methods;  
59 (3) the not-so-well-known seminaive Bayes classifier; (4) for one-dependence Bayesian  
60 classifiers, wrapper approaches, metaclassifiers based on tree-augmented naive Bayes,  
61 and discriminative learning; (5) for general Bayesian network classifiers, classifiers  
62 based on identifying the class variable Markov blanket, metaclassifiers, and discrim-  
63 inative and generative learning of general Bayesian networks used for classification  
64 problems; and (6) Bayesian multinets for encoding probabilistic relationships of asym-  
65 metric independence. Besides, we provide a clear unified notation for all models and  
66 graphical representations of their corresponding networks.

67 A recent overview of Bayesian network classifiers is Flores et al. [2012]. However,  
68 the authors only cover the basic details of naive Bayes, tree-augmented naive Bayes,  $k$ -  
69 dependence Bayesian classifiers, averaged one-dependence estimators, Bayesian multi-  
70 nets, dependency networks, and probabilistic decision graphs. Other shorter reviews  
71 of Bayesian network classifiers are Goldszmidt [2010], discussing only naive Bayes  
72 and tree-augmented naive Bayes, and Al-Aidaroos et al. [2010], focusing on variants of  
73 naive Bayes classifiers. This article is a comprehensive, methodical, and detailed sur-  
74 vey of Bayesian network classifiers ever conducted, elaborating on a variety of facets  
75 and a diversity of models.

76 The article is organized as follows. Section 2 reviews the fundamentals of Bayesian  
77 network classifiers in discrete domains. Then, different models of increasing struc-  
78 ture complexity are presented consecutively. Section 3 describes naive Bayes. Section 4  
79 addresses selective naive Bayes. Section 5 introduces seminaive Bayes. Section 6 fo-  
80 cuses on one-dependence Bayesian classifiers, like tree-augmented naive Bayes and the  
81 super-parent one-dependence estimator. Section 7 discusses  $k$ -dependence Bayesian  
82 classifiers. Section 8 sets out general Bayesian network classifiers, covering Bayesian  
83 network-augmented naive Bayes, classifiers based on identifying the Markov blanket  
84 of the class variable, unrestricted Bayesian classifiers, and discriminative learning.  
85 Section 9 discusses the broadest models, Bayesian multinets. Section 10 shows an il-  
86 lustrative example highlighting the differences between the most important classifiers.  
87 Finally, Section 11 rounds the article off with a discussion and future work.

## 2. FUNDAMENTALS

Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a vector of discrete predictor random variables or features, with  $x_i \in \Omega_{X_i} = \{1, 2, \dots, r_i\}$ , and let  $C$  be a label or class variable, with  $c \in \Omega_C = \{1, 2, \dots, r_c\}$ . Given a simple random sample  $\mathcal{D} = \{(\mathbf{x}^{(1)}, c^{(1)}), \dots, (\mathbf{x}^{(N)}, c^{(N)})\}$ , of size  $N$ , with  $\mathbf{x}^{(j)} = (x_{j1}, \dots, x_{jn})$ , drawn from the joint probability distribution  $p(\mathbf{X}, C)$ , the supervised classification problem consists of inducing a classification model from  $\mathcal{D}$  able to assign labels to new instances given by the value of their predictor variables. Common performance measures include classification accuracy, sensitivity, specificity, the F-measure, and area under the ROC curve. All these measures must be estimated using honest evaluation methods, like hold-out, k-fold cross-validation, bootstrapping, and so forth [Japkowicz and Mohak 2011].

A *Bayes classifier* assigns the most probable a posteriori (MAP) class to a given instance  $\mathbf{x} = (x_1, \dots, x_n)$ , that is,

$$\arg \max_c p(c|\mathbf{x}) = \arg \max_c p(\mathbf{x}, c), \quad (1)$$

which, under a 0/1 loss function, is optimal in terms of minimizing the conditional risk [Duda et al. 2001].

For a general *loss function*,  $\lambda(c', c)$ , where  $c'$  is the class value output by a model and  $c$  is the true class value, the Bayesian classifier can be learned by using the Bayes decision rule that minimizes the expected loss or conditional risk  $R(c'|\mathbf{x}) = \sum_{c \in \Omega_C} \lambda(c', c)p(c|\mathbf{x})$ , for any instance  $\mathbf{x}$  [Duda et al. 2001].

*Bayesian network classifiers* [Friedman et al. 1997] approximate  $p(\mathbf{x}, c)$  with a factorization according to a *Bayesian network* [Pearl 1988]. The structure of a Bayesian network on the random variables  $X_1, \dots, X_n, C$  is a directed acyclic graph (DAG) whose vertices correspond to the random variables and whose arcs encode the probabilistic (in)dependencies among triplets of variables; that is, each factor is a categorical distribution  $p(x_i|\mathbf{pa}(x_i))$  or  $p(c|\mathbf{pa}(c))$ , where  $\mathbf{pa}(x_i)$  is a value of the set of variables  $\mathbf{Pa}(X_i)$ , which are parents of variable  $X_i$  in the graphical structure. The same applies for  $\mathbf{pa}(c)$ . Thus,

$$p(\mathbf{x}, c) = p(c|\mathbf{pa}(c)) \prod_{i=1}^n p(x_i|\mathbf{pa}(x_i)). \quad (2)$$

When the sets  $\mathbf{Pa}(X_i)$  are sparse, this factorization prevents having to estimate an exponential number of parameters, which would otherwise be required.

For the special case of  $\mathbf{Pa}(C) = \emptyset$ , the problem is to maximize on  $c$ :

$$p(\mathbf{x}, c) = p(c)p(\mathbf{x}|c).$$

Therefore, the different Bayesian network classifiers explained later correspond with different factorizations of  $p(\mathbf{x}|c)$ . The simplest model is the naive Bayes, where  $C$  is the parent of all predictor variables and there are no dependence relationships among them (Sections 3 and 4). We can progressively increase the level of dependence in these relationships (one-dependence,  $k$ -dependence, etc.) giving rise to a family of augmented naive Bayes models, explained in Sections 5 through 8.1; see Figure 1.

Equation (2) states a more general case; see also Figure 1.  $p(\mathbf{x}, c)$  is factorized in different ways,  $C$  can have parents, and we have to search the Markov blanket of  $C$  to solve Equation (1) (Section 8.2). The *Markov blanket* (see Pearl [1988, p. 97]) of  $C$  is the set of variables  $MB_C$  that make  $C$  conditionally independent of the other variables in the network, given  $MB_C$ , that is,

$$p(c|\mathbf{x}) = p(c|\mathbf{x}_{MB_C}), \quad (3)$$

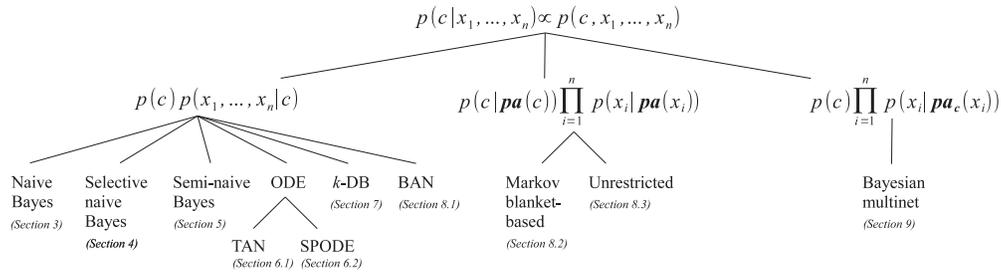


Fig. 1. Categorization of discrete Bayesian network classifiers according to the factorization of  $p(\mathbf{x}, c)$ .

129 where  $\mathbf{x}_{MB_C}$  denotes the projection of  $\mathbf{x}$  onto the variables in  $MB_C$ . Therefore, the  
 130 Markov blanket of  $C$  is the only knowledge needed to predict its behavior. A probability  
 131 distribution  $p$  is *faithful* to a DAG representing a Bayesian network if, for all triplets of  
 132 variables, they are conditionally independent with respect to  $p$  iff they are  $d$ -separated  
 133 in the DAG. For such  $p$ ,  $MB_C$  is unique and is composed of  $C$ 's parents, children, and  
 134 the children's other parents (spouses) [Pearl 1988].

135 There are two strategies for learning both the Markov blanket and the structures  
 136 for augmented naive Bayes: testing conditional independences (constraint-based tech-  
 137 niques [Spirtes et al. 1993]) and searching in the space of models guided by a score to be  
 138 optimized (score + search techniques [Cooper and Herskovits 1992]). They can also be  
 139 combined in hybrid techniques. Alternatively, we can use these strategies to learn an  
 140 *unrestricted* Bayesian network, which does not consider  $C$  as a distinguished variable,  
 141 from which only the Markov blanket of  $C$  must be extracted for classification purposes  
 142 (Section 8.3). Finally, specific conditional independence relationships can be modeled  
 143 for different  $c$  values, giving rise to different Bayesian classifiers, which are then joined  
 144 in the more complex Bayesian multinet (Section 9). The parents of  $X_i$ ,  $\mathbf{Pa}_c(X_i)$ , may be  
 145 different depending on  $c$ ; see Figure 1.

146 Apart from learning the network structure, the probabilities  $p(x_i|\mathbf{pa}(x_i))$  are esti-  
 147 mated from  $\mathcal{D}$  by standard methods like maximum likelihood or Bayesian estimation.  
 148 In Bayesian estimation, assuming a Dirichlet prior distribution over  $(p(X_i = 1|\mathbf{Pa}(X_i) =$   
 149  $j), \dots, p(X_i = r_i|\mathbf{Pa}(X_i) = j))$  with all hyperparameters equal to  $\alpha$ , then the posterior  
 150 distribution is Dirichlet with hyperparameters equal to  $N_{ijk} + \alpha$ ,  $k = 1, \dots, r_i$ , where  $N_{ijk}$   
 151 is the frequency in  $\mathcal{D}$  of cases with  $X_i = k$  and  $\mathbf{Pa}(X_i) = j$ . Hence,  $p(X_i = k|\mathbf{Pa}(X_i) = j)$   
 152 is estimated by

$$\frac{N_{ijk} + \alpha}{N_{.j} + r_i \alpha}, \tag{4}$$

153 where  $N_{.j}$  is the frequency in  $\mathcal{D}$  of cases with  $\mathbf{Pa}(X_i) = j$ . This is called the *Lindstone*  
 154 *rule*. A special case of the Lindstone rule called *Laplace estimation*, with  $\alpha = 1$  in  
 155 Equation (4), is used in Good [1965]. Also, the *Schurmann-Grassberger rule*, where  
 156  $\alpha = \frac{1}{r_i}$ , is employed in Hilden and Bjerregaard [1976] and Titterington et al. [1981].

157 Obviously, the maximum likelihood estimate is given by  $\frac{N_{ijk}}{N_{.j}}$ .

158 So far we have proceeded with only one selected Bayesian network classifier, as if that  
 159 model had generated the data, thus ignoring uncertainty in model selection. *Bayesian*  
 160 *model averaging* provides a way of accounting for model uncertainty. It uses the Bayes  
 161 rule to combine the posterior distributions under each of the models considered with  
 162 structure  $S_m$  in a space  $S$ , each weighted by its posterior model probabilities:

$$p(\mathbf{x}, c|\mathcal{D}) = \sum_{S_m \in S} p(\mathbf{x}, c|S_m, \mathcal{D})p(S_m|\mathcal{D}). \tag{5}$$

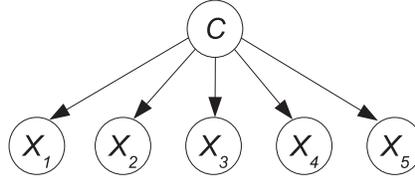


Fig. 2. A naive Bayes structure from which  $p(c|\mathbf{x}) \propto p(c)p(x_1|c)p(x_2|c)p(x_3|c)p(x_4|c)p(x_5|c)$ .

The posterior probability of model  $S_m$  is given by

163

$$p(S_m|\mathcal{D}) = \frac{p(\mathcal{D}|S_m)p(S_m)}{\sum_{S_l \in \mathcal{S}} p(\mathcal{D}|S_l)p(S_l)} \quad (6)$$

and the (marginal) likelihood of model  $S_m$  is

164

$$p(\mathcal{D}|S_m) = \int p(\mathcal{D}|\theta_m, S_m)p(\theta_m|S_m)d\theta_m, \quad (7)$$

where the vector of parameters of model  $S_m$  is  $\theta_m = (\theta_C, \theta_{X_1}, \dots, \theta_{X_n})$ , and for the case of  $\mathbf{Pa}(C) = \emptyset$ ,  $\theta_C = ((p(c))_{c=1}^c)$  and  $\theta_{X_i} = (((\theta_{ijk})_{k=1}^r)_{j=1}^{q_i})$ .  $\theta_{ijk}$  denote  $p(X_i = k | \mathbf{Pa}(X_i) = j)$  and  $q_i$  represents the total number of different configurations of  $\mathbf{Pa}(X_i)$ .

165

166

167

Since our models are Bayesian network classifiers and, according to Equation (2),  $p(\mathbf{x}, c|S_m, \mathcal{D}) = p(c) \prod_{i=1}^n \theta_{ijk}$ , Equation (5) is then simplified as

168

169

$$p(\mathbf{x}, c|\mathcal{D}) \propto \sum_{S_m \in \mathcal{S}} p(c) \left( \prod_{i=1}^n \theta_{ijk} \right) p(\mathcal{D}|S_m)p(S_m).$$

### 3. NAIVE BAYES

170

*Naive Bayes* [Maron and Kuhns 1960; Minsky 1961] is the simplest Bayesian network classifier (Figure 2), since the predictive variables are assumed to be conditionally independent given the class, transforming Equation (1) into

171

172

173

$$p(c|\mathbf{x}) \propto p(c) \prod_{i=1}^n p(x_i|c). \quad (8)$$

This assumption is useful when  $n$  is high and/or  $N$  is small, making  $p(\mathbf{x}|c)$  difficult to estimate. Even if the assumption does not hold, the model classification performance may still be good in practice (although the probabilities are not well calibrated) because the decision boundaries may be insensitive to the specificities of the class-conditional probabilities  $p(x_i|c)$  [Domingos and Pazzani 1997]; that is, variance is reduced because few parameters are required and the biased probability estimates may not matter since the aim is classification rather than accurate posterior class probability estimation [Hand and Yu 2001].

174

175

176

177

178

179

180

181

Other approaches transform the data to avoid the effects of violating the conditional independence assumption, thereby improving the probability estimates made by naive Bayes. The class dispersion problem covers distributions  $p(\mathbf{x}|c)$ , where clusters of cases that belong to the same class are dispersed across the input space. One possible solution is to transform the class distribution by applying a clustering algorithm to each subset of cases with the same label, producing a refinement (extension) on the number of labels. This is proposed in Vilalta and Rish [2003], where a naive Bayes is then learned over this new dataset, and finally the predicted (extended) labels are mapped to the original space of labels.

182

183

184

185

186

187

188

189

190

191 From a theoretical point of view, if all variables (predictors and class) are binary,  
 192 the decision boundary has been shown to be a hyperplane [Minsky 1961]. For ordinal  
 193 nonbinary predictor variables, the decision boundary is a sum of  $n$  polynomials, one  
 194 for each variable  $X_i$ , with a degree equal to  $r_i - 1$  [Duda et al. 2001]. Naive Bayes  
 195 has proved to be optimal (i.e., achieving lower zero-one loss than any other classifier)  
 196 for learning conjunctions and disjunctions of literals [Domingos and Pazzani 1997]. A  
 197 bound for the degradation of the probability of correct classification when naive Bayes  
 198 is used as an approximation of the Bayes classifier is given in Ekdahl and Koski [2006].

199 The inclusion of irrelevant (redundant) variables for the class does not (does) worsen  
 200 the performance of a naive Bayes classifier [Langley and Sage 1994]. Hence, it is  
 201 important to remove irrelevant and redundant variables, as the so-called *selective*  
 202 *naive Bayes* should ideally do (see Section 4).

203 From a practical point of view, there have been some attempts to visualize the effects  
 204 of individual predictor values on the classification decision. Most are based on an  
 205 equivalent expression for a naive Bayes model in terms of the log odds that for a binary  
 206 class ( $c$  vs.  $\bar{c}$ ) results in

$$\text{logit } p(c|\mathbf{x}) = \log \frac{p(c|\mathbf{x})}{p(\bar{c}|\mathbf{x})} = \log \frac{p(c)}{p(\bar{c})} + \sum_{i=1}^n \log \frac{p(x_i|c)}{p(x_i|\bar{c})}.$$

207 While Orange software [Možina et al. 2004] uses nomograms to represent the additive  
 208 influence of each predictor value, ExplainD [Poulin et al. 2006] uses bar-based charts  
 209 with different levels of explanation capabilities.

### 210 3.1. Parameter Estimation

211 The Bayesian probability estimate called *m-estimate* is successfully used in the naive  
 212 Bayes classifier [Cestnik 1990]. It has a tunable parameter  $m$  whereby it can adapt to  
 213 domain properties, such as the level of noise in the dataset.

214 A Bayesian bootstrap method of probability estimation is presented in Norén and  
 215 Orre [2005]. This results in sampling from the dataset of just the  $N' \leq N$  different  
 216 cases of  $\mathcal{D}$  with a Dirichlet distribution with hyperparameters related to the frequency  
 217 of these  $N'$  distinct values in  $\mathcal{D}$ . The variables in a Dirichlet random vector can never be  
 218 positively correlated and must have the same normalized variance. These constraints  
 219 deteriorate the performance of the naive Bayes classifier and motivate the introduction  
 220 of other prior distributions, like the generalized Dirichlet and the Liouville distribu-  
 221 tions [Wong 2009].

222 An estimation inspired by an iterative Hebbian rule is proposed in Gama [1999]. In  
 223 each iteration and for each of the  $N$  cases, if the case is well (incorrectly) classified by  
 224 the current naive Bayes model, then  $p(x_i|c)$  for its corresponding values  $x_i$  and its true  
 225 class  $c$  should be increased (decreased), adjusting the other conditional probabilities.

### 226 3.2. Weighted Naive Bayes

227 Adjusting the naive Bayesian probabilities during classification may significantly im-  
 228 prove predictive accuracy. A general formula is

$$p(c|\mathbf{x}) \propto w_c p(c) \prod_{i=1}^n [p(x_i|c)]^{w_i} \quad (9)$$

229 for some weights  $w_c, w_i, i = 1, \dots, n$ . In Hilden and Bjerregaard [1976],  $w_c = 1$  and  
 230  $w_i = w \in (0, 1), \forall i$ , attaching more importance to the prior probability of the class  
 231 variable.  $w$  is fixed by looking for a good performance after some trials. Also, in Hall  
 232 [2007],  $w_c = 1$  and  $w_i$  is set to  $1/\sqrt{d_i}$ , where  $d_i$  is the minimum depth at which variable

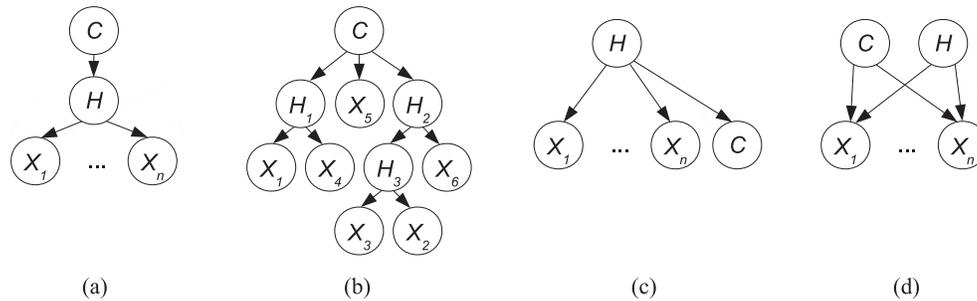


Fig. 3. (a) Naive Bayes with a hidden variable  $H$  [Kwoh and Gillies 1996]; (b) hierarchical naive Bayes [Zhang et al. 2004; Langseth and Nielsen 2006]; (c) finite mixture model, with a hidden variable as a parent of predictor variables and the class [Kontkanen et al. 1996]; (d) finite-mixture-augmented naive Bayes [Monti and Cooper 1999].

$X_i$  is tested in the unpruned decision tree constructed from the data. Fixing the root node to depth 1,  $d_i$  weighs  $X_i$  according to the degree to which it depends on the values of other variables. Finally, in Webb and Pazzani [1998], the linear adjustment  $w_c$  is found by employing a hill-climbing search maximizing the resubstitution accuracy and  $w_i = 1, \forall i$ . 233 234 235 236 237

### 3.3. Missing Data 238

When the training set is incomplete (i.e., some variable values are unknown), both classifier efficiency and accuracy can be lost. 239 240

Simple solutions for handling missing data are either to ignore the cases including unknown values or to consider unknowns to be a separate value of the respective variables [Kohavi et al. 1997]. These solutions introduce biases in the estimates. Another common solution is imputation, where likely values (mode or class-conditional mode) stand in for the missing data. Other suggestions [Friedman et al. 1997] are to use the *expectation-maximization (EM) algorithm* [Dempster et al. 1977] or gradient descent method. However, these methods rely on the assumption that data are *missing at random* (i.e., the probability that an entry will be missing is a function of the observed values in the dataset). This cannot be verified in a particular dataset, and if violated, the methods lead to decreased accuracy. 241 242 243 244 245 246 247 248 249 250

This is why the *robust Bayesian estimator* is introduced in Ramoni and Sebastiani [2001b] to learn conditional probability distributions from incomplete datasets without any assumption about the missing data mechanism. The estimation is given by an interval including all the estimates induced from all possible completions of the original dataset. A new algorithm to compute posterior probability intervals from interval-valued probabilities is then proposed in Ramoni and Sebastiani [2001a]. In the classification phase, all these intervals are ranked according to a score to decide the class with the highest-ranked interval. 251 252 253 254 255 256 257 258

### 3.4. Including Hidden Variables 259

The violation of the conditional independence assumption in naive Bayes can be interpreted as an indication of the presence of hidden or latent variables. Introducing one hidden variable in the naive Bayes model as a child of the class variable and parent of all predictor variables is the simplest solution to this problem; see Figure 3(a). This is the approach reported in Kwoh and Gillies [1996], where the conditional probabilities attached to the hidden node are determined using a gradient descent method. The objective function to be minimized is the squared error between the real class values 260 261 262 263 264 265 266

267 and the class posterior probabilities. The approach taken in Zhang et al. [2004] is more  
 268 general, since many hidden variables are arranged in a tree-shaped Bayesian network  
 269 called *hierarchical naive Bayes*. The root is the class variable, the leaves are the pre-  
 270 dictor variables, and the internal nodes are the hidden variables. An example is given  
 271 in Figure 3(b). This structure is learned using a hill-climbing algorithm that compares  
 272 candidate models with the Bayesian information criterion (BIC), whereas its param-  
 273 eters are estimated using the EM algorithm [Dempster et al. 1977]. A classification  
 274 accuracy-focused improvement is shown in Langseth and Nielsen [2006]. This strategy  
 275 is faster since latent variables are proposed by testing for conditional independencies.

276 There are other options for relaxing the conditional independence assumption. First,  
 277 the *finite mixture model* introduced in Kontkanen et al. [1996] leaves the class vari-  
 278 able as a child node, whereas the common parent for both the discrete or continuous  
 279 predictors and the class variable is a hidden variable; see Figure 3(c). This unmea-  
 280 sured discrete variable is learned using the EM algorithm and models the interaction  
 281 between the predictor variables and between the predictor variables and the class  
 282 variable. Thus, the class and the predictor variables are conditionally independent  
 283 given the hidden variable. Second, the *finite-mixture-augmented naive Bayes* [Monti  
 284 and Cooper 1999] is a combination of this model and naive Bayes. The standard naive  
 285 Bayes is augmented with another naive Bayes with a hidden variable acting as the  
 286 parent of the predictor variables; see Figure 3(d). The hidden variable models the de-  
 287 pendences among the predictor variables that are not captured by the class variable.  
 288 Therefore, it is expected to have fewer states in its domain (i.e., the mixture will have  
 289 fewer components) than the finite mixture model.

### 290 3.5. Metaclassifiers

291 We may use many rather than just one naive Bayes. Thus, the *recursive Bayesian*  
 292 *classifier* [Langley 1993] observes each predicted label (given by the naive Bayes)  
 293 separately. Whenever a label is misclassified, a new naive Bayes is induced from those  
 294 cases having that predicted label. Otherwise, the process stops. The *successive naive*  
 295 *Bayes classifier* [Kononenko 1993] repeats for a fixed number of iterations the learning  
 296 of a naive Bayes from the whole data with redefined labels: a special label  $c_0$  is assigned  
 297 to cases correctly classified by the current naive Bayes, whereas their original labels  
 298 are retained in the other instances. When classifying a new instance, the naive Bayes  
 299 learned last should be applied first. If  $c_0$  is predicted, the next latest naive Bayes  
 300 must be applied; otherwise, the predicted label will be the answer. Also, any *ensemble*  
 301 *method* can be used taking naive Bayes as the base classifier. A specific property of the  
 302 AdaBoost algorithm based on naive Bayes models is that the final boosted model is  
 303 shown to be another naive Bayes [Ridgeway et al. 1998]. Finally, two naive Bayes can  
 304 be used as the base classifier in a *random oracle classifier* [Rodríguez and Kuncheva  
 305 2007]. This is formed by two naive Bayes models and a random oracle that chooses  
 306 one of them in the classification phase. The oracle first divides the predictive variable  
 307 space into two disjoint subspaces based on some random decisions. A naive Bayes is  
 308 then learned from those instances belonging to each subspace. A possible reason for  
 309 the success of (ensembles based on) random oracle classifiers is that the classification  
 310 may be easier in each subspace than in the original space.

311 Multiclass problems are often transformed into a set of binary problems via class bi-  
 312 naryzation techniques. Prominent examples are pairwise classification and one-against-  
 313 all binarization. Training all these binary classifiers, each of which is less complex and  
 314 has simpler decision boundaries, increases the robustness of the final classifier with  
 315 probably less computational burden. The classifier resulting from an *ensemble of pair-*  
 316 *wise naive Bayes* ( $c_i$  vs.  $c_j$ ) that combines the predictions of the individual classifiers

using voting and weighted voting techniques is equivalent to a common naive Bayes. This does not hold for one-against-all binarization [Sulzmann et al. 2007].

Alternatively, naive Bayes can be hybridized with other classification models. The *NBtree* is introduced in Kohavi [1996], combining naive Bayes and decision trees. *NBtree* partitions the training data using a tree structure and builds a local naive Bayes in each leaf with nontested variables. The particular case of a tree with only one branching variable is reported in Cano et al. [2005], where several methods for choosing this variable are proposed. Optionally, for each new case to be classified, a (local) naive Bayes can be induced only from its  $k$  closest cases in the dataset. This hybrid between naive Bayes and the  $k$ -nearest neighbor model is called *locally weighted naive Bayes* [Frank et al. 2003], since the instances in the neighborhood are weighted, attaching less weight to instances that are further from the test instance. Finally, the *lazy Bayesian rule* learning algorithm [Zheng and Webb 2000] induces a rule for each example, whose antecedent is a variable-value conjunction while the consequent is a local naive Bayes with features that are not in the antecedent.

### 3.6. Special Situations

**(a) Homologous sets.** We sometimes have to classify a set of cases that belong to the same unknown class (i.e., a homologous set), for example, a set of leaves taken from the same unknown plant whose species we intend to identify. The *homologous naive Bayes* [Huang and Hsu 2002] takes this knowledge into account, where Equation (8) is now given by

$$p(c|\mathbf{x}_1, \dots, \mathbf{x}_H, \mathcal{H}) \propto p(c) \prod_{h=1}^H \prod_{i=1}^n p(x_{hi}|c),$$

since we wish to classify the homologous set  $\{\mathbf{x}_1, \dots, \mathbf{x}_H\}$ , and  $\mathcal{H}$  denotes that all cases in this set have the same unknown class label. This way, we ensure that different labels are not assigned to all these cases.

**(b) Multiple instances.** In this setting, the learner receives a set of bags that are labeled positive or negative. Each bag contains many instances. A bag is labeled positive (negative) if at least one (all) of its instances is (are) positive (negative). We are looking for a standard classification of individual instances from a collection of labeled bags, for example, learning a simple description of a person from a series of images that are positively labeled if they contain the person and negatively labeled otherwise.

The *multiple-instance naive Bayes* [Murray et al. 2005] starts by assigning negative labels to all the instances in a negative bag. In a positive bag, all the instances are assigned a negative label except one, which receives a positive label. Then a naive Bayes is applied to this dataset. For every positive bag that was misclassified (i.e., all its instances were classified as negative), the instance with the maximum a posteriori probability of being positive is relabeled as positive. A second naive Bayes is applied to this new dataset. This succession of naive Bayes models is halted when a stopping condition is met.

**(c) Cost sensitivity.** For general loss functions, a *cost-sensitive naive Bayes* selects, for each instance  $\mathbf{x}$ , the class value minimizing the expected loss [Ibáñez et al. 2014] of predictions.

We can consider the associated costs of obtaining the missing values in a new case to be classified (e.g., an X-ray test). In this respect, a *test-cost-sensitive naive Bayes classifier* is proposed in Chai et al. [2004], whose aim is to minimize the expected loss by finding how the unknown test variables should be chosen (sequentially or batch-wise). A different situation arises when we have a fixed budget and we are concerned with costs during the learning phase. Here we wish to decide sequentially which tests

364 to run on which instance subject to the budget (i.e., *budgeted learning* [Lizotte et al.  
365 2003]). Naive Bayes's conditional independence assumption simplifies the sequential  
366 process for test selection.

367 **(d) Instance ranking.** In many applications, an accurate ranking of instances is  
368 more desirable than their mere classification, for example, a ranking of candidates in  
369 terms of several aspects in order to award scholarships. Since naive Bayes produces  
370 poor probability estimates [Domingos and Pazzani 1997], an interesting question is to  
371 examine this model's ranking behavior in terms of a well-known ranking quality mea-  
372 sure, the area under the ROC curve or AUC. When all variables are binary, theoretical  
373 results on its optimality for ranking *m-of-n* concepts are given in Zhang and Su [2008],  
374 unlike for classification, where naive Bayes cannot learn all *m-of-n* concepts [Domingos  
375 and Pazzani 1997]. The ideas are extended in Zhang and Sheng [2004] to a weighted  
376 naive Bayes given by Equation (9) with  $w_c = 1$ , where weights  $w_i$  are learned using  
377 several heuristics.

378 **(e) Imprecise and inaccurate probabilities.** Unobserved or rare events, expert  
379 estimates, missing data, or small sample sizes can possibly generate imprecise and  
380 inaccurate probabilities. Using confidence intervals rather than point estimates for  
381  $p(x_i|c)$  and  $p(c)$  is an option, as in the *interval estimation naive Bayes* [Robles et al.  
382 2003]. An evolutionary algorithm can search all the possible (precise) models obtained  
383 by taking values in those confidence intervals for the most accurate model. A more  
384 general way to deal with imprecision in probabilities is by giving a credal set (i.e., the  
385 convex hull of a nonempty and finite family of probability distributions). The *naive  
386 credal classifier* [Zaffalon 2002] uses the class posterior probability intervals and a  
387 dominance criterion to obtain the output of the classification procedure, which, in this  
388 case, can be a set of labels instead of singletons. The effects of parameter inaccuracies  
389 are investigated in Renooij and van der Gaag [2008] with sensitivity analysis tech-  
390 niques. The effect of varying one parameter on the posterior probability of the class  
391 does not significantly influence the performance of the naive Bayes model. However,  
392 this article does not investigate the effect of varying more than one parameter at a  
393 time.

394 **(f) Text categorization.** In this field, documents are represented by a set of random  
395 variables  $C, X_1, \dots, X_n$ , where  $C$  denotes the class of document.  $X_i$  has a different  
396 meaning depending on the chosen model [Eyheramendy et al. 2002]. Thus, in the  
397 *binary independence model*, it represents the presence/absence of a particular term  
398 (word) in the document, and  $p(x_i|c)$  follows a Bernoulli distribution with parameter  $p_{ic}$ .  
399 In other models,  $X_i$  represents the number of occurrences of particular words in the  
400 document. The *multinomial model* assumes that the document length and document  
401 class are marginally independent, transforming Equation (8) into

$$p(c|\mathbf{x}) \propto p(c) \left( \sum_{i=1}^n x_i \right)! \prod_{i=1}^n \frac{p_{ic}^{x_i}}{x_i!}, \quad (10)$$

402 where, for each  $c$ ,  $p_{ic}$  denotes the probability of occurrence of the  $i$ th word and  $\sum_{i=1}^n p_{ic} =$   
403 1. The *Poisson naive Bayes model* assumes that, in Equation (8),  $p(x_i|c)$  follows a Poisson  
404 distribution, whereas in the *negative binomial naive Bayes model*, it is a negative  
405 binomial distribution.

### 406 3.7. Discriminative Learning of Parameters

407 All previous research models the joint probability distribution  $p(\mathbf{x}, c)$  according to what  
408 is called a *generative* approach. A *discriminative* approach [Jebara 2004], however,  
409 directly models the conditional distribution  $p(c|\mathbf{x})$ .

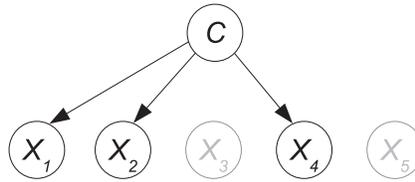


Fig. 4. A selective naive Bayes structure from which  $p(c|\mathbf{x}) \propto p(c)p(x_1|c)p(x_2|c)p(x_4|c)$ . The variables in the shaded nodes have not been selected.

When computing  $p(c|\mathbf{x})$  from the joint probability distribution given by a naive Bayes model, it has been shown [Bishop 1995] to be a linear softmax regression. The parameters of this discriminative model may be estimated by standard techniques (like the Newton-Raphson method). Another more direct way of discriminative learning of the naive Bayes parameters is given in Santafé et al. [2005]: the estimations of parameters maximizing the *conditional likelihood* are approximated using the TM algorithm [Edwards and Lauritzen 2001].

#### 4. SELECTIVE NAIVE BAYES

As mentioned in the previous section, the classification performance of naive Bayes will improve if only relevant, and especially nonredundant, variables are selected to be in the model. Generally, parsimonious models reduce the cost of data acquisition and model learning time, are easier to explain and understand, and increase model applicability, robustness, and performance. Then, a *selective naive Bayes* (Figure 4) is stated as a *feature subset selection* problem, with  $\mathbf{X}_F$  denoting the projection of  $\mathbf{X}$  onto the selected feature subset  $F \subseteq \{1, 2, \dots, n\}$ , where Equation (8) is now

$$p(c|\mathbf{x}) \propto p(c|\mathbf{x}_F) = p(c) \prod_{i \in F} p(x_i|c).$$

The exhaustive search in the space of all possible selective naive Bayes requires the computation of  $2^n$  structures. Although the induction and classification time for a naive Bayes model is short, the enumerative search for the optimal model can be prohibitive. This justifies the use of heuristic approaches for this search.

When a *filter* approach is applied for feature selection, each proposed feature subset is assessed using a scoring measure based on intrinsic characteristics of the data computed from simple statistics on the empirical distribution, totally ignoring the effects on classifier performance. A *wrapper* approach assesses each subset using the classifier performance (accuracy, AUC,  $F_1$  measure, etc.). Finally, an *embedded* approach selects features using the information obtained from training a classifier and is thereby embedded (learning and feature selection tasks cannot be separated) in and specific to a model [Saeys et al. 2007].

##### 4.1. Filter Approaches

When the feature subset is a singleton, we have *univariate filter* methods. This leads to a ranking of features from which the selected feature set is chosen once a threshold on the scoring measure is fixed. The most used scoring measure is the mutual information of each feature and the class variable  $I(X_i, C)$  [Pazzani and Billsus 1997]. Other scoring measures for a feature, like odds ratio, weight of evidence, and symmetrical uncertainty coefficient, can be used, some of which are empirically compared in Mladenic and Grobelnik [1999].

The scoring measures in *multivariate filter* methods are defined on a feature subset. The scoring measure introduced in Hall [1999], called *correlation-based feature*

447 *selection* (CFS), promotes the inclusion of variables that are relevant for classification  
 448 and, at the same time, avoids including redundant variables. Any kind of heuristic  
 449 (forward selection, backward elimination, best first, etc.) can be used to search for this  
 450 optimal subset. Another possibility is to simply select those features that the C4.5  
 451 algorithm would use in its classification tree, as in Ratanamahatana and Gunopulos  
 452 [2003]. A Bayesian criterion for feature selection proposed in Kontkanen et al. [1998]  
 453 is based on approximating the *supervised marginal likelihood* of the class value vector  
 454 given the rest of the data. This is closely related to the conditional log-likelihood (see  
 455 Section 8.4), turning the learning of the selective naive Bayes into a discriminative  
 456 approach.

#### 457 4.2. Wrapper Approaches

458 A wrapper approach outputs the feature subset with a higher computational cost than  
 459 the filter approach. The key issue is how to search the space of feature subsets of  
 460 cardinality  $2^n$ . The strategies used range from simple heuristics, like greedy forward  
 461 [Langley and Sage 1994] and floating search [Pernkopf and O’Leary 2003], to more  
 462 sophisticated population-based heuristics, like genetic algorithms [Liu et al. 2001] and  
 463 estimation of distribution algorithms [Inza et al. 2000].

464 For a large  $n$ , a wrapper approach may be impracticable even with the simplest  
 465 heuristics. This is why many researchers apply a wrapper strategy over a reduced  
 466 filtered subset, thereby adopting a filter-wrapper option [Inza et al. 2004].

#### 467 4.3. Embedded Approaches

468 Regularization techniques are a kind of embedded approach that typically sets out to  
 469 minimize the negative log-likelihood function of the data given the model plus a penalty  
 470 term on the size of the model parameters. An  $L_1$  penalty is useful for feature selection  
 471 because the size of some parameters is driven to zero. An  $L_1/L_2$ -regularized naive Bayes  
 472 for continuous and discrete predictor variables is introduced in Vidaurre et al. [2012].  
 473 In addition, a *stagewise version of the selective naive Bayes*, which can be considered a  
 474 regularized version of a naive Bayes, is also presented. Whereas the  $L_1/L_2$ -regularized  
 475 naive Bayes model only discards irrelevant predictors, the stagewise version of the  
 476 selective naive Bayes can discard both irrelevant and redundant predictors.

#### 477 4.4. Metaclassifiers

478 As with naive Bayes (Section 3.5), selective naive Bayes models can be combined in a  
 479 metaclassifier. The *random naive Bayes* [Prinzie and Van den Poel 2007] is a bagged  
 480 classifier combining many naive Bayes, each of which has been estimated from a boot-  
 481 strap sample with  $m < n$  randomly selected features. The *naive Bayesian classifier*  
 482 *committee* [Zheng 1998] sequentially generates selective naive Bayes models to be  
 483 members of the committee. The probability that a feature is used for the next model  
 484 increases if the current model performs better than the naive Bayes (with all features).  
 485 For each class, the probabilities provided by all committee members are summed up,  
 486 taking as the predicted class the one with the largest summed probability.

487 Bayesian model averaging (see Equation (5)) is an ensemble learning technique.  
 488 Applied to all selective naive Bayes models, this gives rise to a unique naive Bayes  
 489 model, as shown in Dash and Cooper [2002]. Here Dirichlet priors are assumed for  
 490  $p(\theta_m|S_m)$  in Equation (7) and uniform priors for  $p(S_m)$  in Equation (6).

### 491 5. SEMINAIVE BAYES

492 *Seminaive Bayes* models (Figure 5) aim to relax the conditional independence assump-  
 493 tion of naive Bayes by introducing new features obtained as the Cartesian product of  
 494 two or more original predictor variables. By doing this, the model is able to represent

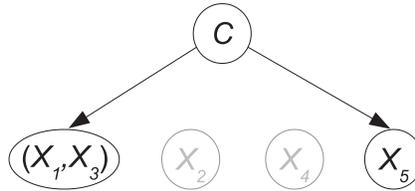


Fig. 5. A seminaive Bayes structure from which  $p(c|\mathbf{x}) \propto p(c)p(x_1, x_3|c)p(x_5|c)$ .

dependencies between original predictor variables. However, these new predictor variables are still conditionally independent given the class variable. Thus, if  $S_j \subseteq \{1, 2, \dots, n\}$  denotes the indices in the  $j$ th feature (original or Cartesian product),  $j = 1, \dots, K$ , Equation (8) is now

$$p(c|\mathbf{x}) \propto p(c) \prod_{j=1}^K p(\mathbf{x}_{S_j}|c),$$

where  $S_j \cap S_l = \emptyset$ , for  $j \neq l$ .

The seminaive Bayes model of Pazzani [1996] starts from an empty structure and considers the best option between (a) adding a variable not used by the current classifier as conditionally independent of the features (original or Cartesian products) used in the classifier, and (b) joining a variable not used by the current classifier with each feature (original or Cartesian products) present in the classifier. This is a greedy search algorithm, called *forward sequential selection and joining*, guided wrapper-wise (the objective function is the classification accuracy), that stops when there is no accuracy improvement. An alternative backward version starting from a naive Bayes, called *backward sequential elimination and joining*, is also proposed by the same author. Evolutionary computation has been used to guide the search for the best semi-naive Bayes model in Robles et al. [2003] wrapper-wise with estimation of distribution algorithms. Using a wrapper approach avoids including redundant variables in the model, since these degrade accuracy, as mentioned in Section 3.

A filter adaptation of the forward sequential selection and joining algorithm is presented in Blanco et al. [2005]. Options (a) and (b) listed previously are evaluated with a  $\chi^2$  test of independence based on the mutual information  $I(C, X_i)$  of the class and each variable not in the current model (for (a)) and on the mutual information of the class and a joint variable formed by a variable not in the current model and a feature present in the model (for (b)). We always select the variable with the smallest  $p$ -value until no more new variables can be added to the model (because they do not reject the null hypothesis of independence). Other filter approaches use alternative scoring metrics like Bayesian Dirichlet equivalence (BDe) [Heckerman et al. 1995], and leave one out and log-likelihood ratio test, as in Abellán et al. [2007]. Every time variables form a new joint variable, this approach [Abellán et al. 2007] tries to merge values of this new variable to reduce its cardinality and computation time. For imprecise probabilities, a filter *seminaive credal classifier* is given in Abellán et al. [2006].

A seminaive Bayes model (or naive Bayes or interval estimation naive Bayes) is the model built in Robles et al. [2004] at the second level of a metaclassifier following a stacked generalization scheme, taking as input data the different labels provided by different classifiers at the first level.

## 6. ONE-DEPENDENCE BAYESIAN CLASSIFIERS

*One-dependence estimators* (ODEs) are similar to naive Bayes except that each predictor variable is allowed to depend on at most one other predictor in addition to the class.

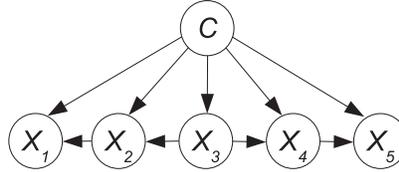


Fig. 6. A TAN structure, whose root node is  $X_3$ , from which  $p(c|\mathbf{x}) \propto p(c)p(x_1|c, x_2)p(x_2|c, x_3)p(x_3|c)p(x_4|c, x_3)p(x_5|c, x_4)$ .

533 They can improve naive Bayes accuracy when its conditional independence assumption  
534 is violated.

### 535 6.1. Tree-Augmented Naive Bayes

536 Unlike in seminaive Bayes, which introduces new features to relax the condi-  
537 tional independence assumption of naive Bayes, the *tree-augmented network* (TAN)  
538 [Friedman et al. 1997] maintains the original predictor variables and models relation-  
539 ships of at most order 1 among the variables. Specifically, a tree-shaped graph models  
540 the predictor subgraph (Figure 6).

541 Learning a TAN structure first involves constructing an undirected tree. Kruskal's  
542 algorithm [Kruskal 1956] is used to calculate the maximum weighted spanning tree  
543 (MWST), containing  $n - 1$  edges, where the weight of an edge  $X_i - X_j$  is  $I(X_i, X_j|C)$ ,  
544 which is the conditional mutual information of  $X_i$  and  $X_j$  given  $C$ . The undirected tree  
545 is then converted into a directed tree by selecting at random a variable as the root node  
546 and replacing the edges by arcs. This is the tree shaping the predictor subgraph. Finally,  
547 a naive Bayes structure is superimposed to form the TAN structure. The posterior  
548 distribution in Equation (1) is then

$$p(c|\mathbf{x}) \propto p(c)p(x_r|c) \prod_{i=1, i \neq r}^n p(x_i|c, x_{j(i)}), \quad (11)$$

549 where  $X_r$  denotes the root node and  $\{X_{j(i)}\} = \mathbf{Pa}(X_i) \setminus C$ , for any  $i \neq r$ .

550 These ideas are adapted from Chow and Liu [1968], where several trees, one for each  
551 value  $c$  of the class, were constructed rather than a single tree for the entire domain.  
552 This works like TAN, but uses only the cases from  $\mathcal{D}$  satisfying  $C = c$  to construct each  
553 tree. This collection of trees is a special case of a Bayesian multinet, a terminology  
554 introduced by Geiger and Heckerman [1996] for the first time (see Section 9).

555 From a theoretical point of view, the procedures in Chow and Liu [1968] (Figure 7(a))  
556 and Friedman et al. [1997] (Figure 7(b)) construct, respectively, the tree-based Bayesian  
557 multinet and the TAN structure that both maximize the likelihood.

558 Rather than obtaining a spanning tree, the method described in Ruz and Pham  
559 [2009] suggests that Kruskal's algorithm be stopped whenever a Bayesian criterion  
560 controlling the likelihood of the data and the complexity of the TAN structure holds.  
561 The predictor subgraph will then include  $e \leq n - 1$  arcs. This procedure has been proven  
562 to find an augmented naive Bayes classifier that minimizes the Kullback-Leibler (KL)  
563 divergence between the real joint probability distribution and the approximation given  
564 by the model, across all network structures with  $e$  arcs.

565 Two special situations are when data are incomplete and probabilities are imprecise.  
566 The *structural EM algorithm* [Friedman 1997] in the space of trees is used in François  
567 and Leray [2006] for the first case. The *tree-based credal classifier* algorithm that is  
568 able to induce credal Bayesian networks with a TAN structure is proposed in Zaffalon  
569 and Fagioli [2003] for the second case.

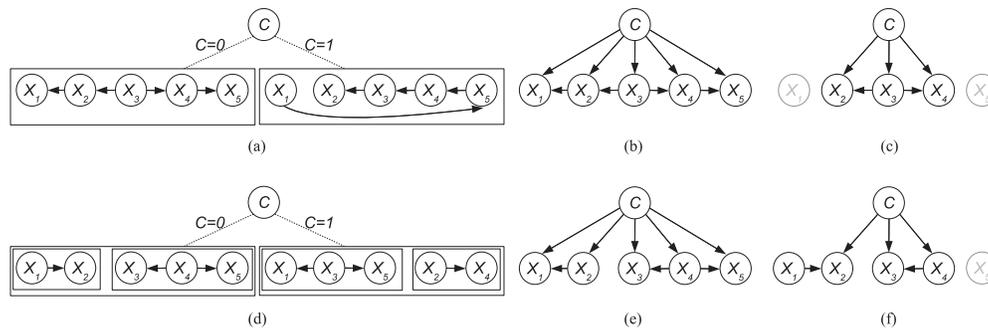


Fig. 7. (a) Bayesian multinet as a collection of trees [Chow and Liu 1968]:  $p(C = 0|\mathbf{x}) \propto p(C = 0)p(x_1|C = 0, x_2)p(x_2|C = 0, x_3)p(x_3|C = 0)p(x_4|C = 0, x_3)p(x_5|C = 0, x_4)$  and  $p(C = 1|\mathbf{x}) \propto p(C = 1)p(x_1|C = 1)p(x_2|C = 1, x_3)p(x_3|C = 1, x_4)p(x_4|C = 1, x_5)p(x_5|C = 1, x_1)$ ; (b) TAN [Friedman et al. 1997]:  $p(c|\mathbf{x}) \propto p(c)p(x_1|c, x_2)p(x_2|c, x_3)p(x_3|c)p(x_4|c, x_3)p(x_5|c, x_4)$ ; (c) selective TAN [Blanco et al. 2005]:  $p(c|\mathbf{x}) \propto p(c)p(x_2|c, x_3)p(x_3|c)p(x_4|c, x_3)$ ; (d) Bayesian multinet as a collection of forests [Pham et al. 2002]:  $p(C = 0|\mathbf{x}) \propto p(C = 0)p(x_1|C = 0)p(x_2|C = 0, x_1)p(x_3|C = 0, x_4)p(x_4|C = 0)p(x_5|C = 0, x_4)$  and  $p(C = 1|\mathbf{x}) \propto p(C = 1)p(x_1|C = 1, x_3)p(x_2|C = 1)p(x_3|C = 1)p(x_4|C = 1, x_2)p(x_5|C = 1, x_3)$ ; (e) FAN [Lucas 2004]:  $p(c|\mathbf{x}) \propto p(c)p(x_1|c, x_2)p(x_2|c)p(x_3|c, x_4)p(x_4|c)p(x_5|c, x_4)$ ; (f) selective FAN [Ziebart et al. 2007]:  $p(c|\mathbf{x}) \propto p(c)p(x_2|c, x_1)p(x_3|c, x_4)p(x_4|c)$ .

If the weights of the undirected tree based on conditional mutual information are first filtered with a  $\chi^2$  test of independence, the resulting structure is the *selective TAN* [Blanco et al. 2005] (Figure 7(c)). The predictor subgraph could be a forest rather than a tree since it may result in many root nodes.

Other authors propose following a wrapper instead of a filter approach. The next three references, again, lead to forest predictor structures (i.e., a disjoint union of trees). Thus, initializing the network to a naive Bayes, we can consider adding possible arcs from  $X_i$  to  $X_j$ , for  $X_j$  without any predictor variable as parent, and selecting the arc giving the highest accuracy improvement. This hill-climbing search algorithm is described in Keogh and Pazzani [2002]. The authors also propose another less expensive search. Finding the best arc to add is broken down into two steps. First, we consider making each node a superparent in the current classifier (i.e., with arcs directed to all nodes without a predictor parent). The best superparent yields the highest accuracy. Second, we choose one of all the superparent’s children (i.e., the favorite child that most improves accuracy) for the final structure. Also starting from a naive Bayes, a sequential floating search heuristic is used in Pernkopf and O’Leary [2003]. In Blanco et al. [2005], by initializing with an empty predictor subgraph, an algorithm greedily decides whether to add a new predictor or to create an arc between two predictors already in the model. Unlike the last two wrapper techniques, it actually performs a feature subset selection.

**Forest-augmented naive Bayes.** Rather than using a collection of trees as in Chow and Liu [1968], a collection of forests, one for each value  $c$  of the class, is built in Pham et al. [2002] (Figure 7(d)). The forests are obtained using a maximum weighted spanning forest algorithm (e.g., [Fredman and Tarjan 1987]). The *forest-augmented naive Bayes* (FAN) was first defined in Lucas [2004], with only one rather than a collection of forests in the predictor subgraph, augmented with a naive Bayes (Figure 7(e)). Therefore, the research reported in Lucas [2004] adapts Pham et al. [2002] for FAN models as Friedman et al. [1997] did with Chow and Liu [1968] for TAN. The *selective FAN* introduced in Ziebart et al. [2007] adds the novelty of allowing the predictor variables to be optionally dependent on the class variable; that is, missing arcs from  $C$  to some  $X_i$  can be found (Figure 7(f)). Moreover, the learning approach is based on maximizing

570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600

601 the likelihood of the data, which is penalized for avoiding the class variable as a  
602 parent.

603 **Metaclassifiers.** Bagging-type metaclassifiers use bootstrap samples and thus re-  
604 quire an unstable base classifier to generate diverse results from the different clas-  
605 sifiers. However, the TAN classifier is stable. A randomization is then needed in the  
606 standard TAN algorithm. Thus, the *bagging-randomTAN* in Ma and Shi [2004] takes  
607 randomTAN as base classifiers in a bagging scheme. The *randomTAN* randomly selects  
608 the edges between predictor variables whose conditional mutual information surpasses  
609 a fixed threshold. These selective TAN models vote for the final classification. Using  
610 boosting instead means sampling the original data with weights according to the clas-  
611 sification results of each data item to form a new dataset for the next classifier. This  
612 scheme is employed in the *boosted augmented naive Bayes (bAN)* [Jing et al. 2008].  
613 The base classifier is chosen by first running a trial with a naive Bayes, then greedily  
614 augmenting the current structure at iteration  $s$  with the  $s$ th edge having the highest  
615 conditional mutual information. We stop when the added edge does not improve the  
616 classification accuracy. Note that the final structure of the base classifier can be a  
617 FAN.

618 The *averaged TAN (ATAN)* [Jiang et al. 2012] takes not a random node but each  
619 predictor variable as root node and then builds the corresponding MWST conditioned  
620 to that selection. Finally, the posterior probabilities  $p(c|\mathbf{x})$  of ATAN are given by the  
621 average of the  $n$  TAN classifier posterior probabilities.

622 Bayesian model averaging (see Equation (5)) over TAN structures and parameters is  
623 carried out in Cerquides and López de Mántaras [2005b]. The authors define decompos-  
624 able (conjugate) distributions as priors for  $p(S_m)$  in Equation (6) and choose Dirichlet  
625 priors for  $p(\theta_m|S_m)$  in Equation (7). They compute the exact *Bayesian model averaging*  
626 *over TANs*. In addition, they propose an ensemble of the  $k$  most probable a posteriori  
627 TAN models.

628 **Discriminative learning.** A discriminative learning of a TAN model is proposed in  
629 Feng et al. [2007]. First, the TAN structure is learned as in Friedman et al. [1997] but  
630 replacing the conditional mutual information by the explaining away residual (EAR)  
631 criterion [Bilmes 2000], that is, using  $I(X_i, X_j|C) - I(X_i, X_j)$ . Maximizing EAR over the  
632 tree is in fact an approximation to maximizing the conditional likelihood. Second, they  
633 define an objective function based mainly on the KL divergence between the empirical  
634 distribution and the distribution given by the previous TAN structure for each value  $c$   
635 of the class to discriminatively learn the parameters.

636 A different discriminative score, the maximum margin, is proposed in Pernkopf and  
637 Wohlmayr [2013] to search for the structure of TAN with both greedy hill-climbing  
638 and simulated annealing strategies. The multiclass margin of an instance  $\mathbf{x}^{(i)}$  is  $d^{(i)} =$   
639  $\frac{p(c^{(i)}|\mathbf{x}^{(i)})}{\max_{c \neq c^{(i)}} p(c|\mathbf{x}^{(i)})}$ . Rather than searching for the structure that maximizes  $\min_{i=1, \dots, N} d^{(i)}$ ,  
640 this is relaxed with a soft margin, finally defining the *maximum margin score* of a  
641 structure as  $\sum_{i=1}^N \min\{1, \lambda \log d^{(i)}\}$ , where  $\lambda > 0$  is a scaling parameter and is set by  
642 cross-validation.

643 As in Section 3.7 with naive Bayes, the TM algorithm [Edwards and Lauritzen  
644 2001] can be adapted for the discriminatively learning parameters in a TAN classifier  
645 [Santafé et al. 2005].

## 646 6.2. SuperParent-One-Dependence Estimators

647 *SuperParent-One-Dependence Estimators (SPODEs)* are an ODE where all predictors  
648 depend on the same predictor (the superparent) in addition to the class [Keogh and  
649 Pazzani 2002] (Figure 8). Note that this is a particular case of a TAN model. The

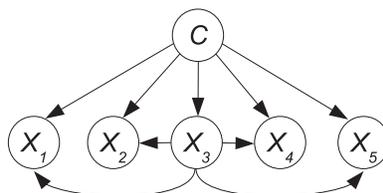


Fig. 8. A SPODE structure, with  $X_3$  as superparent, from which  $p(c|\mathbf{x}) \propto p(c)p(x_1|c, x_3)p(x_2|c, x_3)p(x_3|c)p(x_4|c, x_3)p(x_5|c, x_3)$ .

posterior distribution in Equation (1) is

650

$$p(c|\mathbf{x}) \propto p(c)p(x_{sp}|c) \prod_{i=1, i \neq sp}^n p(x_i|c, x_{sp}),$$

where  $X_{sp}$  denotes the superparent node. This equation is similar to Equation (11), particularized as  $X_r = X_{j(i)} = X_{sp}$ , for any  $i \neq sp$ .

651  
652

**Metaclassifiers.** The *averaged one-dependence estimator* (AODE) [Webb et al. 2005] averages the predictions of all qualified SPODEs, where “qualified” means that it includes, for each instance  $\mathbf{x} = (x_1, \dots, x_{sp}, \dots, x_n)$ , only the SPODEs for which the probability estimates are accurate, that is, where the training data contain more than  $m$  instances verifying  $X_{sp} = x_{sp}$ . The authors suggest fixing  $m = 30$ . The average prediction is given by

653  
654  
655  
656  
657  
658

$$p(c|\mathbf{x}) \propto p(c, \mathbf{x}) = \frac{1}{|SP_{\mathbf{x}}^m|} \sum_{X_{sp} \in SP_{\mathbf{x}}^m} p(c)p(x_{sp}|c) \prod_{i=1, i \neq sp}^n p(x_i|c, x_{sp}), \quad (12)$$

where  $SP_{\mathbf{x}}^m$  denotes for each  $\mathbf{x}$  the set of predictor variables qualified as superparents and  $|\cdot|$  is its cardinal. AODE avoids model selection, thereby decreasing the variance component of the classifier.

659  
660  
661

The AODE can be further improved by deleting  $X_j$  from the set of predictors whenever  $P(x_j|x_i) = 1$  ( $x_i$  and  $x_j$  are highly dependent predictor values) when classifying a new instance  $\mathbf{x}$ . Note that this technique introduced in Zheng and Webb [2006] is performed at classification time for each new instance, and this is why it is called *lazy elimination*. It is shown that it significantly reduces classification bias and error without undue computation.

662  
663  
664  
665  
666  
667

Another improvement is the *lazy AODE* [Jiang and Zhang 2006], which builds an AODE for each test instance. The training data is expanded by adding a number of copies (clones) of each training instance equal to its similarity to the test instance. This similarity is the number of identical predictor variables.

668  
669  
670  
671

Since AODE requires all the SPODE models to be stored in main memory, *generalized additive Bayesian network classifiers* (GABNs) defined in Li et al. [2007] propose aggregating only some SPODEs (or other simple Bayesian classifiers) within the framework of generalized additive models. SPODEs with the lowest mutual information scores  $I(X_{sp}, C)$  are not considered in the aggregation. Thus, this aggregation is given by the linear combination of  $n' \leq n$  probabilities  $p_{sp}(\mathbf{x}, c)$  obtained in the SPODE models:

672  
673  
674  
675  
676  
677

$$\sum_{sp=1}^{n'} \lambda_{sp} g_{sp}(p_{sp}(\mathbf{x}, c)),$$

where  $g_{sp}$  is the link function and  $0 \leq \lambda_{sp} \leq 1$  are parameters to be estimated such that  $\sum_{sp=1}^{n'} \lambda_{sp} = 1$ . When  $g_{sp}$  is the log function, then  $p(\mathbf{x}, c) \propto \prod_{sp=1}^{n'} p_{sp}^{\lambda_{sp}}(\mathbf{x}, c)$ . It is

678  
679

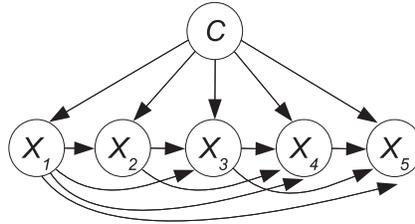


Fig. 9. An example of 3-DB structure from which  $p(c|\mathbf{x}) \propto p(c)p(x_1|c)p(x_2|c, x_1)p(x_3|c, x_1, x_2)p(x_4|c, x_1, x_2, x_3)p(x_5|c, x_1, x_3, x_4)$ .

680 easy to design a gradient-based method to optimize its associated quasi-likelihood that  
 681 outputs the combining parameters  $\lambda_{sp}$ .

682 Another way to obtain an ensemble of SPODEs in the AODE is proposed in Yang  
 683 et al. [2005] as a wrapper approach. The aim is to select SPODEs so as to maximize  
 684 classification accuracy. We need a metric (like minimum description length [MDL],  
 685 minimum message length [MML], leave-one-out classification accuracy, accuracy from  
 686 backward sequential elimination, or forward sequential addition processes) to order  
 687 the  $n$  possible SPODEs for selection, and a stopping criterion always based on the  
 688 accuracy.

689 The idea of Yang et al. [2007] is to compute the final predictions as a weighted  
 690 average in Equation (12), rather than as an average. Four different weighting schemes  
 691 are then proposed. Two of them use the posterior probability of each SPODE given  
 692 the data as its weight. The first is based on the inversion of Shannon’s law and the  
 693 second is within a Bayesian model averaging, where uniform priors over the  $n$  SPODE  
 694 structures and Dirichlet priors over the corresponding parameters are assumed. The  
 695 other two schemes use a MAP estimation to find the most probable a posteriori set of  
 696 weights for a SPODE ensemble, assuming a Dirichlet prior over the weights. These  
 697 two last schemes differ as to the posterior, generative, or discriminative models (see  
 698 Cerquides and López de Mántaras [2005a] for further details).

699 **6.3. Other One-Dependence Estimators**

700 The *weighted ODE* can be used to approximate the conditional probabilities  $p(x_i|c)$   
 701 in the naive Bayes. This was proposed by Jiang et al. [2009], resulting in

$$p(c|\mathbf{x}) \propto p(c, \mathbf{x}) \approx p(c) \prod_{i=1}^n \left( \sum_{j=1, j \neq i}^n w_{ij} p(x_i|c, x_j) \right), \quad (13)$$

702 where  $w_{ij} \propto I(X_i, X_j|C)$ . The same authors propose in Jiang et al. [2012] other weighting  
 703 schemes, based on performance measures of the different ODE models, like AUC or  
 704 classification accuracy.

705 The *hidden one-dependence estimator* classifier (HODE) [Flores et al. 2009] avoids  
 706 using any SPODE. HODE introduces, via the EM algorithm, a new variable (the hidden  
 707 variable  $H$ ), with the aim of representing the links existing in the  $n$  SPODE models.  
 708 Node  $C$  in the naive Bayes structure is replaced by the Cartesian product of  $C$  and  
 709  $H$ . Then we have to estimate the probability of  $x_i$  conditioned by  $c$  and  $h$  searching for  
 710  $\arg \max_c \sum_h p(c, h) \prod_{i=1}^n p(x_i|c, h)$ .

711 **7.  $k$ -DEPENDENCE BAYESIAN CLASSIFIERS**

712 The  *$k$ -dependence Bayesian classifier* ( $k$ -DB) [Sahami 1996] allows each predictor vari-  
 713 able to have a maximum of  $k$  parent variables apart from the class variable (Figure 9).  
 714 The inclusion order of the predictor variables  $X_i$  in the model is given by  $I(X_i, C)$ ,

starting with the highest. Once  $X_i$  enters the model, its parents are selected by choosing those  $k$  variables  $X_j$  in the model with the highest values of  $I(X_i, X_j|C)$ . The main disadvantages of the standard  $k$ -DB are the lack of feature selection (all the original predictor variables are included in the final model) and the need to determine the optimal value for  $k$ . Also, once  $k$  has been fixed, the number of parents of each predictor variable is inflexible. Obviously, naive Bayes and TAN are particular cases of  $k$ -DBs, with  $k = 0$  and  $k = 1$ , respectively.

The posterior distribution in Equation (1) is

$$p(c|\mathbf{x}) \propto p(c) \prod_{i=1}^n p(x_i|c, x_{i_1}, \dots, x_{i_k}),$$

where  $X_{i_1}, \dots, X_{i_k}$  are parents of  $X_i$  in the structure. Note that the first  $k$  variables entering the model will have fewer than  $k$  parents (the first variable entering the model has no parents, the second variable has one parent, and so on) and the remaining  $n - k$  variables have exactly  $k$  parents.

Feature subset selection is performed in Blanco et al. [2005] within a  $k$ -DB using filter and wrapper approaches. In the filter approach, an initial step selects the predictor variables that pass a  $\chi^2$  test of independence based on the mutual information  $I(C, X_i)$ . Then the standard  $k$ -DB algorithm is applied on this reduced subset, considering only those arcs that pass an analogous independence test based on conditional mutual information  $I(X_i, X_j|C)$ . In the wrapper approach, as in the wrapper TAN approach discussed in Section 6.1, the decision on whether to add a new predictor or to create an arc between two predictors already in the model is guided by accuracy, provided that the added arc does not violate the  $k$ -DB restrictions. As a consequence, all the predictors in the structures output by this wrapper approach have at most  $k$  parents, but there is no need to have  $n - k$  variables with exactly  $k$  parents. In general, graphs where each node has at most  $k$  parents are called *k-graphs*.

A  $k$ -graph as the predictor subgraph is also the result of a kind of evolutionary computation method described in Xiao et al. [2009], inspired by the so-called group method of data handling (GMDH) [Ivakhnenko 1970]. The algorithm to build *GMDH-based Bayesian classifiers* starts from a set of  $s \propto n + 1$  models with only one arc, corresponding to the pair of variables ( $C$  included) with the highest mutual information. Then a new set of  $\binom{s}{2}$  models is obtained by pairwise joining the previous structures. The best  $s$  models according to BDe or BIC are selected. This process that incrementally increases the model complexity is repeated until the new best does not improve the current best model. The number of parents is always bounded by a fixed  $k$ .

The  $k$ -graphs obtained in Carvalho et al. [2007] are obliged to be consistent with an order between the predictor variables. This order,  $\sigma$ , is based on a breadth-first search (BFS) over the TAN predictor subgraph obtained in the usual manner [Friedman et al. 1997]. This means that for any arc  $X_i \rightarrow X_j$  in the  $k$ -graph,  $X_i$  is visited before  $X_j$  in a total order completing  $\sigma$ . The learning algorithm of *BFS-consistent Bayesian network classifiers* can cope with any decomposable score, score expressible as a sum of local scores that depend only on each node and its parents.

$k$ -graphs are also induced in Pernkopf and Bilmes [2010]. They first establish an ordering of the predictor variables by using a greedy algorithm. A variable  $X$  is chosen whenever it is the most informative about  $C$  given the previous variables in the order, where informativeness is measured by the conditional mutual information,  $I(C, X|\mathbf{X}_{\text{prev}})$ . This order can alternatively use classification accuracy as a score assuming a fully connected subgraph over  $C$ ,  $X$ , and  $\mathbf{X}_{\text{prev}}$ . In any case, the best  $k$  parents for each variable among  $\mathbf{X}_{\text{prev}}$  are selected in a second step by scoring each possibility

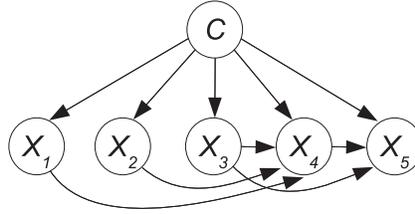


Fig. 10. A Bayesian network-augmented naive Bayes structure from which  $p(c|\mathbf{x}) \propto p(c)p(x_1|c)p(x_2|c)p(x_3|c)p(x_4|c, x_1, x_2, x_3)p(x_5|c, x_3, x_4)$ .

762 with the classification accuracy. Here a naive Bayes assumption is used for  $\mathbf{X} \setminus \{\mathbf{X}_{\text{prev}}, X\}$ ,  
 763 that is, the variables whose parents have not yet been chosen.

764 **Metaclassifiers.** A combination of  $k$ -DB models in a bagging fashion is proposed in  
 765 Louzada and Ara [2012].

## 766 8. GENERAL BAYESIAN NETWORK CLASSIFIERS

767 This section discusses more general structures. First, relaxing the structure of the  
 768 predictor subgraph but maintaining  $C$  without any parent defines a Bayesian network-  
 769 augmented naive Bayes (Section 8.1). Second, if  $C$  is allowed to have parents, its Markov  
 770 blanket is the only knowledge needed to predict its behavior (see Equation (3)), and  
 771 some classifiers have been designed to search for the Markov blanket (Section 8.2).  
 772 Finally, a very general unrestricted Bayesian network that does not consider  $C$  as a  
 773 special variable can be induced with any existing Bayesian network structure learning  
 774 algorithm. The corresponding Markov blanket of  $C$  can be used later for classification  
 775 purposes (Section 8.3). In all three cases, Equation (1) is

$$p(c|\mathbf{x}) \propto p(c|\mathbf{pa}(c)) \prod_{i=1}^n p(x_i|\mathbf{pa}(x_i)),$$

776 where  $\mathbf{Pa}(C) = \emptyset$  in Section 8.1.

### 777 8.1. Bayesian Network-Augmented Naive Bayes

778 Relaxing the fixed number of parents,  $k$ , in a  $k$ -DB, does not place any limitations  
 779 on links among predictor variables (except that they do not form a cycle); that is, a  
 780 Bayesian network structure can be the predictor subgraph (Figure 10). This model is  
 781 called *Bayesian network-augmented naive Bayes* (BAN), a term first coined by Friedman  
 782 et al. [1997]. The factorization is

$$p(c|\mathbf{x}) \propto p(c) \prod_{i=1}^n p(x_i|\mathbf{pa}(x_i)).$$

783 The first reference to a learning algorithm for this model is Ezawa and Norton [1996].  
 784 First, it ranks the  $n$  predictor variables based on  $I(X_i, C)$ , and then it selects the min-  
 785 imum number of predictor variables  $k$  verifying  $\sum_{j=1}^k I(X_j, C) \geq t_{CX} \sum_{j=1}^n I(X_j, C)$ ,  
 786 where  $0 < t_{CX} < 1$  is the threshold. Second,  $I(X_i, X_j|C)$  is computed for all pairs of  
 787 selected variables. The edges corresponding to the highest values are selected until  
 788 a percentage  $t_{XX}$  of the overall conditional mutual information  $\sum_{i < j} I(X_i, X_j|C)$   
 789 is surpassed. Edge directionality is based on the variable ranking of the first step:  
 790 higher-ranked variables point toward lower-ranked variables. Note that this algorithm  
 791 resembles the initial proposal for learning a  $k$ -DB model [Sahami 1996]; see Section 7.

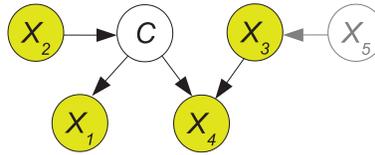


Fig. 11. A Markov blanket structure for  $C$  from which  $p(c|\mathbf{x}) \propto p(c|x_2)p(x_1|c)p(x_2)p(x_3)p(x_4|c, x_3)$ .

As explained in Section 2, a Bayesian network can be learned using conditional independence tests. This is the strategy adopted in Cheng and Greiner [1999] to obtain the predictor subgraph. This algorithm has three phases: drafting, thickening, and thinning. First, it computes  $I(X_i, X_j|C)$  as a measure of closeness and creates a draft based on this information. Second, it adds arcs (thickening) when the pairs of nodes cannot be  $d$ -separated, resulting in an independence map (I-map) of the underlying dependency model. Third, each arc of the I-map is examined using conditional independence tests and will be removed (thinning) if both nodes of the arc can be  $d$ -separated. The final result is the minimal I-map [Pearl 1988].

Also, a Bayesian network can be learned with a score + search technique. In Friedman et al. [1997], the structure is learned by minimizing the MDL score with a greedy forward search. In van Gerven and Lucas [2004], the (conditional) mutual information score and a forward greedy search is used in the *maximum mutual information* (MMI) algorithm. MMI iteratively selects the arc with the highest (conditional) mutual information from two sets of candidate arcs:  $C \rightarrow X_i$ -type arcs, chosen with  $I(X_i, C)$ , followed, as soon as  $C$  has children, by  $X_j \rightarrow X_i$ -type arcs where  $X_i$  is a child of  $C$ , chosen with  $I(X_i, X_j|\mathbf{Pa}(X_i))$ . Note that  $\mathbf{Pa}(X_i)$  can add new variables at each iteration, and the conditional mutual information should be recomputed accordingly. The parameter learning uses nonuniform Dirichlet priors to avoid spurious dependences. Another example of a score + search approach is reported in Pernkopf and O’Leary [2003], where accuracy is used as the score with a sequential floating search heuristic.

## 8.2. Bayesian Classifiers Based on Identifying the Markov Blanket of the Class Variable

**(a) Detecting conditional independences.** Finding the Markov blanket of  $C$  (Figure 11),  $MB_C$ , can be stated as a feature selection problem, where we start from the set of all the predictor variables and eliminate a variable at each step (backward greedy strategy) until we have approximated  $MB_C$ . A feature is eliminated if it gives little or no additional information about  $C$  beyond what is subsumed by the remaining features. The method in Koller and Sahami [1996] eliminates feature by feature trying to keep  $p(C|MB_C^{(t)})$ , the conditional probability of  $C$  given the current estimation of the Markov blanket at step  $t$ , as close to  $p(C|\mathbf{X})$  as possible. Closeness is defined by the expected KL divergence. The main idea is to note that eliminating a variable  $X_i^*$ , which is conditionally independent of  $C$  given  $MB_C^{(t)}$ , keeps the expected “distance” from  $p(C|MB_C^{(t)}, X_i)$  to  $p(C|MB_C^{(t)})$  close to zero. The obtained succession of  $\{MB_C^{(t)}\}_t$ , where  $MB_C^{(t)} = MB_C^{(t-1)} \setminus \{X_i^*\}$ , should converge to the true  $MB_C$ .

At each step  $t$ , the algorithm chooses which variable  $X_i^*$  to eliminate, as follows. For each  $X_i$ , we compute for any  $X_j$  not yet eliminated,  $D_{KL}(p(C|X_i = x_i, X_j = x_j), p(C|X_j = x_j))$ ,  $\forall x_i, x_j, j \neq i$ , where  $D_{KL}$  is the KL divergence. The expected  $D_{KL}$  is then computed as  $\delta(X_i|X_j) = \sum_{x_i, x_j} p(x_i, x_j) D_{KL}(p(C|X_i = x_i, X_j = x_j), p(C|X_j = x_j))$ . We select the  $K$  features  $(X_{i_1}, \dots, X_{i_K}) = \mathbf{M}_i$  for which  $\delta(X_i|X_j)$  is smallest.  $\mathbf{M}_i$  tries to capture the variables  $X_j$  for which  $X_i$  is conditionally independent of  $C$  given  $X_j$ . The process is repeated for each  $X_i$ , and then we choose the variable  $X_i^*$  to be eliminated as the one

833 with minimum

$$\sum_{\mathbf{m}_i, x_i} p(\mathbf{m}_i, x_i) D_{KL}(p(C|\mathbf{M}_i = \mathbf{m}_i, X_i = x_i), p(C|\mathbf{M}_i = \mathbf{m}_i)).$$

834 Finally, the next step  $t + 1$  is started with  $MB_C^{(t+1)} = MB_C^{(t)} \setminus \{X_i^*\}$ . The number of steps  
 835 is prespecified and is the number of variables for elimination from the approximate  
 836 Markov blanket. Note that, as mentioned in Koller and Sahami [1996], the algorithm  
 837 is suboptimal in many ways, particularly due to the very naive approximations that it  
 838 uses and the need to specify a good value for  $K$  and for the number of variables in the  
 839 Markov blanket.

840 This and the following algorithms are based on the observation that if  $X_i \notin MB_C$   
 841 then  $I_p(C, X_i|MB_C)$  holds; that is,  $C$  and  $X_i$  are conditionally independent under  $p$  given  
 842  $MB_C$ . This holds if we apply the decomposition property of the conditional independence  
 843 [Pearl 1988]

$$I_p(T, Y \cup W|Z) \Rightarrow I_p(T, Y|Z), I_p(T, W|Z) \quad (14)$$

844 to Equation (3).

845 A common assumption in all these algorithms is that  $\mathcal{D}$  is a sample from a probability  
 846 distribution  $p$  faithful to a DAG representing a Bayesian network.

847 The *grow-shrink (GS) Markov blanket algorithm* [Margaritis and Thrun 2000] starts  
 848 from an empty Markov blanket, current Markov blanket  $CMB_C$ , and adds a variable  
 849  $X_i$  as long as the Markov blanket property of  $C$  is violated, that is,  $\neg I_p(C, X_i|CMB_C)$ ,  
 850 until there are no more such variables (growing phase). Many false positives may have  
 851 entered the  $MB_C$  during the growing phase. Thus, the second phase identifies and  
 852 removes the variables that are independent of  $C$  given the other variables in the  $MB_C$   
 853 one by one (shrinking phase). In practice, it is possible to reduce the number of tests  
 854 in the shrinking phase by heuristically ordering the variables by ascending  $I(X_i, C)$   
 855 or the probability of dependence between  $X_i$  and  $C$  in the growing step. Orientation  
 856 rules are then applied to this Markov blanket to get its directed version. GS is the first  
 857 correct Markov blanket induction algorithm under the faithfulness assumption; that  
 858 is, it returns the true  $MB_C$ . GS is scalable because it outputs the Markov blanket of  
 859  $C$  without learning a Bayesian network for all variables  $\mathbf{X}$  and  $C$ . GS has to condition  
 860 on at least as many variables simultaneously as the Markov blanket size, and it is  
 861 therefore impractical, because it requires a sample that grows exponentially to this  
 862 size if the conditional independence tests are to be reliable. This means that GS is  
 863 not data efficient. A randomized version of the GS algorithm with members of the  
 864 conditioning set chosen randomly from  $CMB_C$  is also proposed as a faster and more  
 865 reliable variant.

866 The *incremental association Markov blanket (IAMB)* algorithm [Tsamardinos and  
 867 Aliferis 2003], a modified version of GS, consists of a forward phase followed by a  
 868 backward phase. Starting from an empty Markov blanket, it iteratively includes the  
 869 variable  $X_i$  that has the highest association with  $C$  conditioned on  $CMB_C$  (e.g., condi-  
 870 tional mutual information) in the first forward (admission) phase, after checking the  
 871 same condition as in GS ( $\neg I_p(C, X_i|CMB_C)$ ). We stop when this association is weak.  
 872 For each  $X_i \in CMB_C$ , we remove  $X_i$  from  $CMB_C$  if  $I_p(C, X_i|CMB_C \setminus \{X_i\})$  holds to elim-  
 873 inate the false positives in the second backward conditioning phase. IAMB scales to  
 874 high-dimensional datasets. The authors prove that the Markov blanket corresponds to  
 875 the strongly relevant features as defined by Kohavi and John [1997]. Likewise to GS,  
 876 IAMB is correct and scalable but data inefficient.

877 There have been many variants of the IAMB algorithm. The InterIAMBnPC al-  
 878 gorithm [Tsamardinos et al. 2003a] interleaves the admission phase with backward  
 879 conditioning attempting to keep the size of  $CMB_C$  as small as possible during all

the steps. It also substitutes the backward conditioning phase with the PC algorithm [Spirtes et al. 1993]. Fast-IAMB [Yaramakala and Margaritis 2005] speeds up IAMB, reducing the number of tests in the admission phase by adding not one but a number of variables at a time.

The *HITON* algorithm [Aliferis et al. 2003] consists of three steps. First, HITON-PC identifies the parents and children of  $C$ , the set  $PC$ . This is started from an empty set and includes the variable  $X_i$  that has the maximum association with  $C$  in the current  $PC$ ,  $CPC$ . Then, a variable  $X_j \in CPC$  that meets  $\neg I_p(C, X_j|S)$  for some subset  $S$  from  $CPC$  is removed from  $CPC$  and not considered again for admission. The process is repeated until no more variables are left. After outputting  $PC$ , in the second step, HITON-PC is again applied to each variable in  $PC$  to obtain  $PCPC$ , the parents and children of  $PC$ . Thus, the current  $MB_C$  is  $CMB_C = PC \cup PCPC$ . False positives, which retain just the spouses of  $C$ , are removed from  $CMB_C$ :  $X_j \in CMB_C$  is only retained if  $\nexists S \in CMB_C \setminus PC$  such that  $\neg I_p(C, X_j|S)$ . Unlike the GS and IAMB algorithms, HITON works with conditional (in)dependence statements involving any subset  $S$  in  $CMB_C$ , rather than just with  $CMB_C$ . Finally, in a third step, a greedy backward elimination approach is applied wrapper-like to the previously obtained Markov blanket. HITON is scalable and data efficient because the number of instances required to identify the Markov blanket does not depend on its size but on its topology. However, HITON is incorrect, as proved by Peña et al. [2007].

The *max-min Markov blanket* (MMMB) algorithm [Tsamardinos et al. 2003b] is similar to HITON. However, it chooses the variable  $X_i$  in  $CPC$  that exhibits the maximum association with  $C$  conditioned on the subset  $S^*$  of  $CPC$  that achieves the minimum association possible for this variable; that is,  $S^*$  is the subset  $S$  of  $CPC$  that minimizes the association of  $X_i$  and  $C$  given  $S$ . This selection method typically admits very few false positives, whereby all subsets on which we condition in the next steps have a manageable size. Also, the second step of MMMB introduces a more sophisticated criterion to identify the spouses of  $C$  than HITON. MMMB has the same properties as HITON.

The *parents- and children-based Markov boundary* (PCMB) algorithm [Peña et al. 2007] is a variant of MMMB that incorporates so-called “symmetry correction.” The parents–children relationship is symmetric in the sense that  $X_i$  belongs to the set of parents and children of  $C$ , and  $C$  should also belong to the set of parents and children of  $X_i$ . A breach of this symmetry is a sign of a false-positive member in the Markov blanket. This leads to the first algorithm that is correct, scalable, and data efficient. This symmetry correction, based on an AND operator, makes it harder for a true positive to enter the Markov blanket. This is relaxed in the MBOR algorithm [Rodrigues de Morais and Aussem 2010], which uses an OR operator and is correct and scalable but data inefficient. A faster PCMB called *breadth-first search of Markov blanket* (BFMB) [Fu and Desmarais 2007] relies on fewer data passes and conditioning on the minimum set.

The *generalized local learning framework for Markov blanket induction* algorithms is proposed in Aliferis et al. [2010]. It can be instantiated in many ways, giving rise to existing state-of-the-art (HITON and MMPC) algorithms. Both the  $PC$  set and the Markov blanket are seen as the results of searching for direct causes, direct effects, and direct causes of the direct effects of a variable  $C$ .

Table I shows a summary of the main algorithms assuming faithfulness and their properties.

Few algorithms have tried to relax the faithfulness assumption. A weaker condition is the *composition property*, which is the converse of Equation (14), which does not have the guarantee of the Markov blanket being unique. IAMB is still correct under this composition property, but because it is a deterministic algorithm, it cannot discover

Table I. Properties of the Main Algorithms for Markov Blanket Discovery under the Faithfulness Assumption

	Correct	Scalable	Data efficient
GS [Margaritis and Thrun 2000]	✓	✓	
IAMB [Tsamardinos and Aliferis 2003]	✓	✓	
HITON [Aliferis et al. 2003]		✓	✓
MMMB [Tsamardinos et al. 2003b]		✓	✓
PCMB [Peña et al. 2007]	✓	✓	✓
MBOR [Rodrigues de Morais and Aussem 2010]	✓	✓	

932 different Markov blankets. This drawback is overcome by KIAMB [Peña et al. 2007], a  
 933 stochastic version of IAMB, which is not only correct and scalable like IAMB but also  
 934 data efficient unlike IAMB. Rather than conditioning on  $CMB_C$  when searching for  
 935 the highest association in the IAMB admission phase, KIAMB conditions on a random  
 936 subset of  $CMB_C$ , whose size is proportional to  $K \in [0, 1]$ . IAMB corresponds to KIAMB  
 937 with  $K = 1$ .

938 Note that none of these algorithms takes into account arcs between either the chil-  
 939 dren of  $C$  or  $\text{Pa}(C)$  and the children of  $C$ .

940 **(b) Score + search techniques.** The *partial Bayesian network* (PBN) for the  
 941 Markov blanket around  $C$  [Madden 2002] involves three steps. In the first step, each  
 942 predictor variable is classified as either parent of  $C$ , child of  $C$ , or unconnected to  $C$ .  
 943 During the second step, the spouses of  $C$  are added from the set of parents and uncon-  
 944 nected nodes. The third step determines the dependences between the nodes that are  
 945 children of  $C$ . The three steps are guided by the K2 score [Cooper and Herskovits 1992],  
 946 thereby requiring a node ordering. The inclusion of an arc is decided with the score in  
 947 a forward greedy way. A similar idea is presented in dos Santos et al. [2011], where the  
 948 K2 algorithm [Cooper and Herskovits 1992] is applied on an ordering starting with  $C$ .  
 949 This ordering prevents  $C$  from having parents resulting in an approximated Markov  
 950 blanket of  $C$ .

951 For small sample situations, a bootstrap procedure for determining membership in  
 952 the Markov blanket is proposed in Friedman et al. [1999]. They answer the question  
 953 of how confident we can be that  $X_i$  is in  $X_j$ 's Markov blanket (in our case we would  
 954 be interested in  $X_j = C$ ). From each bootstrap sample, a Bayesian network is learned  
 955 using the BDe score with a uniform prior distribution and using a greedy hill-climbing  
 956 search. Using the procedure described in Chickering [1995], each Bayesian network  
 957 is converted into a partially directed acyclic graph (PDAG). From these PDAGs, the  
 958 final PDAG is composed of the arcs and edges whose confidence (measured by their  
 959 occurrence frequency in these networks) surpasses a given threshold. A PDAG repre-  
 960 sents an *equivalence class* of Bayesian network structures, where equivalence means  
 961 that all networks in the class imply the same set of independence statements. Thus, an  
 962 equivalence class includes equivalent networks, with the same skeleton (the undirected  
 963 version of the DAG) and the same set of immoralities or v-structures (arcs  $X \rightarrow Z$  and  
 964  $Y \rightarrow Z$  but with nonadjacent  $X$  and  $Y$ ) [Verma and Pearl 1990]. An arc in a PDAG  
 965 denotes that all members in the equivalence class contain that arc; an edge  $X_i - X_j$   
 966 in a PDAG indicates that some members contain the arc  $X_i \rightarrow X_j$  and some contain  
 967  $X_j \rightarrow X_i$ .

968 Rather than using a filter score, the search can be guided in a wrapper-wise using  
 969 classification accuracy as the score. An example is given in Sierra and Larrañaga [1998],  
 970 where the search is performed by means of a genetic algorithm. Each individual in the  
 971 population represents a Markov blanket structure for  $C$ .

972 **(c) Hybrid techniques.** A two-stage algorithm called *tabu search-enhanced Markov*  
 973 *blanket* is presented in Bai et al. [2008]. In the first stage, an initial Markov blanket is

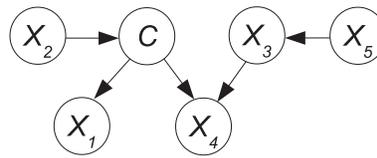


Fig. 12. An unrestricted Bayesian network classifier structure from which  $p(c|\mathbf{x}) \propto p(c|x_2)p(x_1|c)p(x_2)p(x_3)p(x_4|c, x_3)$ .

obtained based on conditional independence tests carried out according to a breadth-first search heuristic. In the second stage, tabu search enhancement, allowing four kinds of move (arc addition, arc deletion, arc switch, and arc switch with node pruning) is introduced. Each possible move is evaluated taking into account classification accuracy.

### 8.3. Unrestricted Bayesian Classifiers

This section includes the general unrestricted Bayesian classifiers where  $C$  is not considered as a special variable in the induction process (Figure 12).

The complexity of algorithms that learn Bayesian networks from data identifying high-scoring structures in which each node has at most  $k$  parents, for all  $k \geq 3$ , has been shown to be NP hard [Chickering et al. 2004]. It holds whenever the learning algorithm uses a consistent scoring criterion and is applied to a sufficiently large dataset. This justifies the use of search heuristics.

The *K2-attribute selection* (K2-AS) algorithm [Provan and Singh 1995] consists of two main steps. The node selection phase chooses the set of nodes from which the final network is built. In the network construction phase, the network is built with those nodes. Nodes are selected incrementally by adding the variable whose inclusion results in the maximum increase in accuracy (of the resulting network). Using these selected variables, the final network is built using the *CB algorithm* [Singh and Valtorta 1995]. This algorithm uses conditional independence tests to generate a “good” node ordering and then uses the K2 algorithm on that ordering to induce the Bayesian network. A variant of K2-AS is Info-AS [Singh and Provan 1996]. They differ only as to node selection being guided by a conditional information-theoretic metric (conditional information gain, conditional gain ratio, or complement of conditional distance). A simpler approach is to use a node ordering for the K2 algorithm given by the ranking of variables yielded with a score (like information gain or chi-squared score) as in Hruschka and Ebecken [2007].

Instead of searching the Bayesian classifier in the space of DAGs, we can use a reduced search space that consists of a type of PDAGs, called *class-focused restricted PDAGs* (C-RPDAGs) [Acid et al. 2005]. C-RPDAGs combine two concepts of DAG equivalence: independence equivalence and a new concept, classification equivalence. This classification equivalence means producing the same posterior probabilities for the class. Local search is performed by means of specific operators to move from one C-RPDAG to another neighboring C-RPDAG. Standard decomposable and score-equivalent (where equivalent networks have the same score) functions guide the search.

As mentioned at the beginning of this section, from the general Bayesian network obtained with all these methods, the Markov blanket of  $C$  is used for classification.

**Metaclassifiers.** Following the stacked generalization method, a general Bayesian network classifier is built in Sierra et al. [2001] from the response given by a set of classifiers. The algorithm for building this network searches for the structure that maximizes classification accuracy, guided by a genetic algorithm.

Exact Bayesian model averaging of a particular class of structures, consistent with a fixed partial ordering of the nodes and with bounded in-degree  $k$ , is considered in

Table II. Generative and Discriminative Approaches for Structure and Parameter Learning of General Bayesian Network Classifiers

		Structure learning	
		Generative	Discriminative
	Generative	Sections 8.1, 8.2, 8.3	CMDL [Grossman and Domingos 2004], CBIC [Guo and Greiner 2005], $\hat{f}$ CLL [Carvalho et al. 2011], ACL-MLE [Burge and Lane 2005], EAR [Narasimhan and Bilmes 2005], MDL-FS [Drugan and Wiering 2010], <i>Hist-dist</i> [Sierra et al. 2009]
Parameter learning	Discriminative	LR-Roos [Roos et al. 2005], LR-Feelders [Feelders and Ivanovs 2006], ELR [Greiner and Zhou 2002; Greiner et al. 2005], DFE [Su et al. 2008], ECL, ACL, and EBW [Pernkopf and Wohlmayr 2009], MCLR [Guo et al. 2005; Pernkopf et al. 2012]	CMDL-ELR [Grossman and Domingos 2004], CBIC-ELR [Guo and Greiner 2005], ACL- <i>Max</i> [Burge and Lane 2005]
	Generative-Discriminative	Normalized hybrid [Raina et al. 2004; Fujino et al. 2007], JoDiG [Xue and Titterington 2010], HBayes [Kang and Tian 2006], Bayesian blending [Bishop and Lasserre 2007]	

1016 Dash and Cooper [2004]. The authors prove that there is a single Bayesian network  
 1017 whose prediction is equivalent to the one obtained by averaging the structures of this  
 1018 particular class. Since constructing this network is computationally prohibitive, they  
 1019 provide a tractable approximation whereby approximate model-averaging probability  
 1020 calculations can be performed in linear time. Rather than starting from a fixed node  
 1021 order, which is hard to obtain and may affect classification performance, the idea of  
 1022 Hwang and Zhang [2005] is to extend Bayesian model averaging of general Bayesian  
 1023 network classifiers by averaging over several distinct node orders. The average is  
 1024 approximated using the Markov chain Monte Carlo sampling technique. This method  
 1025 performs well when the dataset is sparse and noisy.

#### 1026 8.4. Discriminative Learning of General Bayesian Network Classifiers

1027 As mentioned in Section 3.7, generative classifiers learn a model of the joint probability  
 1028 distribution  $p(\mathbf{x}, c)$  and perform classification using Bayes's rule to compute the pos-  
 1029 terior probability of the class variable. The standard approach for learning generative  
 1030 classifiers is maximum likelihood estimation, possibly augmented with a (Bayesian)  
 1031 smoothing prior. Discriminative classifiers directly model the posterior probability of  
 1032 the class variable, which is the distribution used for classification. Therefore, gener-  
 1033 ative models maximize the log-likelihood or a related function, whereas discriminative  
 1034 models maximize the conditional log-likelihood. Table II summarizes the content of  
 1035 this section.

1036 **(a) Discriminative learning of structures.** The log-likelihood of the data  $\mathcal{D}$  given  
 1037 a Bayesian network classifier  $B$ ,  $LL(\mathcal{D}|B)$ , and the conditional log-likelihood,  $CLL(\mathcal{D}|B)$ ,

are both related as follows:

1038

$$\begin{aligned}
 LL(\mathcal{D}|B) &= \sum_{i=1}^N \log p_B(c^{(i)}, x_1^{(i)}, \dots, x_n^{(i)}) \\
 &= \sum_{i=1}^N \log p_B(c^{(i)}|x_1^{(i)}, \dots, x_n^{(i)}) + \sum_{i=1}^N \log p_B(x_1^{(i)}, \dots, x_n^{(i)}) \\
 &= CLL(\mathcal{D}|B) + \sum_{i=1}^N \log p_B(x_1^{(i)}, \dots, x_n^{(i)}). \tag{15}
 \end{aligned}$$

It is the first addend that matters in classification, and a better approach would be to use  $CLL(\mathcal{D}|B)$  alone as the objective function. Unfortunately, the CLL function does not decompose into a separate term for each variable, and there is no known closed-form solution for the optimal parameter estimates.

1039  
1040  
1041  
1042

The CLL function is used in Grossman and Domingos [2004] to learn the structure of the network, where the maximum number of parents per variable is bounded, while parameters are approximated by their maximum likelihood estimates (MLEs), which is extremely fast. Also, they propose using a modified CLL, which penalizes complex structures via the number of parameters in the network, that is, a *conditional MDL* score (CMDL). A hill-climbing algorithm is used to maximize CLL and CMDL, starting from an empty network and at each step considering the addition, deletion, or reversion of an arc. Additionally, this discriminative learning of structures is extended to a discriminative learning of parameters by computing their estimates via the extended logistic regression (ELR) algorithm [Greiner and Zhou 2002], although the results were not much better.

1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053

Another way of modifying CLL is to penalize by the number of parameters in  $C$ 's Markov blanket. This results in the *conditional BIC* score (CBIC) defined in Guo and Greiner [2005] as an analog of the generative BIC criterion. This CBIC criterion can be accompanied by generative (MLE) or discriminative (ELR) parameter learning.

1054  
1055  
1056  
1057

Rather than working with CLL, other authors propose criteria similar to CLL but with better computational properties. The *factorized conditional log-likelihood* ( $\hat{f}CLL$ ) is introduced in Carvalho et al. [2011] with the properties of being decomposable and score equivalent for BAN classifiers. Note that the addends in CLL (see Equation (15)) can be expressed, for a binary  $C$  ( $c$  vs.  $\neg c$ ), as a difference of logarithms:

1058  
1059  
1060  
1061  
1062

$$\begin{aligned}
 \log p_B(c^{(i)}|x_1^{(i)}, \dots, x_n^{(i)}) &= \log p(c^{(i)}, x_1^{(i)}, \dots, x_n^{(i)}) \\
 &\quad - \log(p(c^{(i)}, x_1^{(i)}, \dots, x_n^{(i)}) + p(\neg c^{(i)}, x_1^{(i)}, \dots, x_n^{(i)})),
 \end{aligned}$$

the second one being the log of a sum of terms, whereby it is nondecomposable. Then these addends are approximated by a linear function of the log of these terms. When substituted in the  $\hat{f}CLL$  score, this can be rewritten in terms of conditional mutual information and *interaction information* [McGill 1954]. For parameter learning, the authors use MLEs.

1063  
1064  
1065  
1066  
1067

Another simpler approximation to CLL is the *approximate conditional likelihood* (ACL) [Burge and Lane 2005], where the sum mentioned earlier is replaced by a single term, that is, by  $\log p(\neg c^{(i)}, x_1^{(i)}, \dots, x_n^{(i)})$ , to avoid the nondecomposability drawback. This formulation can be applied even for complex classifiers like Bayesian multi-nets (see Section 9). This results in a decomposable (although unbounded) score. The (discriminatively learned) parameters maximizing this score (*ACL-Max*) have a closed form. Alternatively, MLEs can be used for parameter learning (*ACL-MLE*).

1068  
1069  
1070  
1071  
1072  
1073  
1074

1075 The EAR measure is the criterion maximized in Narasimhan and Bilmes [2005]  
1076 using a greedy forward algorithm with an MLE of parameters.

1077 The idea of Drugan and Wiering [2010] is to use both the Bayesian network clas-  
1078 sifier that factorizes the joint distribution  $p(c, \mathbf{x})$  and an auxiliary Bayesian network  
1079 that factorizes  $p(\mathbf{x})$ . Since the quotient between these two distributions is  $p(c|\mathbf{x})$ , the  
1080 conditional log-likelihood, CLL, of the classifier is then approximated by the differ-  
1081 ence between the unconditional log-likelihood of the classifier and the log-likelihood of  
1082 the auxiliary network; see the first three sums in Equation (15). Both structures are  
1083 learned using a generative method. A new score, called *minimum description length for*  
1084 *feature selection* (MDL-FS), is introduced to guide the search for good structures, also  
1085 allowing feature selection. MDL-FS, like MDL, penalizes the complexity of the classi-  
1086 fier and, rather than including the log-likelihood, it includes the so-called *conditional*  
1087 *auxiliary log-likelihood*, the difference between the log-likelihood of the data given the  
1088 Bayesian network classifier and that given the auxiliary Bayesian network over  $\mathbf{X}$ . In  
1089 practical applications, they propose to set a specific family of auxiliary networks before-  
1090 hand. Depending on their complexity, the MDL-FS can serve to identify and remove  
1091 redundant variables at various levels. Thus, with trees as auxiliary networks, learning  
1092 a selective TAN classifier starts with all predictor variables in both types of structures.  
1093 The corresponding MDL-FS is computed and guides the next variable to be deleted  
1094 following a backward elimination strategy. New structures are learned from the new  
1095 set of variables. MLE is used for parameter learning.

1096 A score that takes into account the posterior distribution of the class variable during  
1097 the structure learning process should in principle lead to models with higher classi-  
1098 fication capabilities. The score introduced in Sierra et al. [2009] (*Hist-dist*) uses, for  
1099 each case, the distance between the predicted posterior distribution of the class and  
1100 an approximation of the real (degenerated) posterior distribution. This is defined by  
1101 giving an  $\alpha$  value (close to 1) to the real class of the case and dividing the remain-  
1102 ing  $1 - \alpha$  evenly across the other class values. The final score to be minimized is  
1103 the mean of those distances for all cases. Different distance measures are proposed  
1104 (Euclidean, Kolmogorov-Smirnov, chi-square, etc.). The wrapper approach is based  
1105 on the greedy *Algorithm B* [Buntine 1991], which searches for the best unrestricted  
1106 Bayesian classifier.

1107 **(b) Discriminative learning of parameters.** Logistic regression can be seen as  
1108 discriminatively trained naive Bayes classifiers [Agresti 1990]. See also Ng and Jordan  
1109 [2001] for an empirical and theoretical comparison of both models, where for small  
1110 sample sizes the generative naive Bayes can outperform the discriminatively trained  
1111 naive Bayes. In general, discriminatively trained classifiers are usually more accurate  
1112 when  $N$  is high.

1113 For a fixed Bayesian network structure, finding the values  $\theta_{ijk}$  for the conditional  
1114 probability tables that maximize the CLL is NP hard for a given incomplete dataset  
1115 [Greiner et al. 2005], something more readily solved in generative models maximizing  
1116 the likelihood, which have straightforward EM methods for handling missing data.

1117 Given complete data, the complexity of maximizing the CLL for arbitrary structures  
1118 is unknown. However, the CLL does not have local maxima for structures satisfying  
1119 a certain graph-theoretic property, and the global maximum can be found by mapping  
1120 the corresponding optimization problem to an equivalent logistic regression model  
1121 [Roos et al. 2005]. This model has fewer parameters than its Bayesian network clas-  
1122 sifier counterpart and is known to have a strictly concave log-likelihood function. The  
1123 graph-theoretic property is that the structure of the Bayesian network is such that its  
1124 canonical version is perfect; that is, all nodes having a common child are connected.  
1125 The canonical version is constructed by first restricting the original structure to  $C$ 's  
1126 Markov blanket and then adding as many arcs as needed to make the parents of  $C$

fully connected. All Bayesian networks with the same canonical version are equivalent in terms of  $p(c|x_1, \dots, x_n)$ . Naive Bayes and TAN models comply with this property. The conditional distributions  $p(c|x_1, \dots, x_n)$  in the CLL expression are reparameterized using a logistic regression model where the covariates are derived from the original variables. There are two types of covariates: (a) indicator variables for each configuration  $\mathbf{pa}(c)$  and (b) indicator variables for each configuration  $(x_i, \mathbf{pa}^C(x_i))$ , where  $X_i$  denotes any children of  $C$ , and  $\mathbf{Pa}^C(X_i) = \mathbf{Pa}(X_i) \setminus \{C\}$ . The original parameters,  $\theta_{ijk}$ , are recovered via the exponential function of the logistic regression parameters. We call this approach LR-Roos, an acronym of logistic regression for perfect structures.

A different mapping for perfect graphs to an equivalent logistic regression model with fewer parameters than LR-Roos is proposed in Feelders and Ivanovs [2006]. The Broyden-Fletcher-Goldfarb-Shanno (BFGS) method (the best for simple structures) and conjugate gradient are used to optimize the CLL. We call this approach LR-Feelders.

The aforementioned ELR algorithm [Greiner et al. 2005] is the most popular approximation procedure for maximizing the CLL for a given Bayesian network structure. ELR applies to arbitrary Bayesian network structures and works effectively even with an incomplete dataset. It is often superior to classifiers produced by standard generative algorithms, especially in common situations where the given Bayesian network structure is incorrect; that is, it is not an I-map of the underlying distribution. This occurs when the learning algorithm is conservative about adding new arcs to avoid overfitting the data or because the algorithm only considers a restricted class of structures that is not guaranteed to contain the correct structure. For each conditional probability table entry, ELR is a conjugate gradient-ascent algorithm that tries to maximize CLL with respect to a softmax function of  $\theta_{ijk}$ , that is,  $\theta_{ijk} = \frac{e^{\theta_{ijk}}}{\sum_{k'} e^{\theta_{ijk'}}$ .

A different idea is to take the effect of estimating  $\theta_{ijk}$  on classification into account by adapting the appropriate frequencies from data.  $\theta_{ijk}$  is initialized as the MLE in iteration  $t = 0$ . Going through all the training data, the update at iteration  $t + 1$  consists of summing, for each instance  $\mathbf{x}$ , the difference between the true posterior probability  $p(c|\mathbf{x})$  (assumed to be 1 when  $\mathbf{x}$  has label  $c$  in the dataset) and the predicted probability generated by the current parameters  $p_t(c|\mathbf{x})$ , that is,  $\theta_{ijk}^{(t+1)} = \theta_{ijk}^{(t)} + p(c|\mathbf{x}) - p_t(c|\mathbf{x})$ . This approach was proposed in Su et al. [2008] and named *discriminative frequency estimate* (DFE). DFE can be seen as a more sophisticated approach than the one proposed in Gama [1999].

Three discriminative parameter learning algorithms are introduced in Pernkopf and Wohlmayr [2009] for naive Bayes, TAN, or 2-DB structures. First, the *exact CLL decomposition* (ECL) algorithm tries to optimize the CLL function. Second, the *approximate CLL decomposition* (ACL) algorithm aims at optimizing a lower-bound surrogate of the CLL function. Third, the *extended Baum-Welch* (EBW) algorithm is used for these three structures. All the algorithms initialize the parameters to the MLEs.

A different criterion is optimized in Guo et al. [2005]. The discriminative objective is to maximize the *minimum conditional likelihood ratio* (MCLR):

$$MCLR(\theta) = \min_{i=1, \dots, N} \min_{c \neq c^{(i)}} \frac{p(c^{(i)}|\mathbf{x}^{(i)}, \theta)}{p(c|\mathbf{x}^{(i)}, \theta)}.$$

When Bayesian networks are formulated as a form of exponential model,  $\log MCLR(\theta)$  resembles a large margin criterion of support vector machines, but subject to normalization constraints over each variable (probabilities summing 1). These restrictions are nonlinear, and this yields a difficult optimization problem. The authors solve the problem with convex relaxation for a wide range of graph topologies.

1174 A conjugate gradient algorithm is instead proposed in Pernkopf et al. [2012] and is  
1175 advantageous in terms of computational requirements.

1176 **(c) Generative-discriminative learning.** Some researchers try to take advantage  
1177 of the best of both approaches through hybrid parameter learning (partly generative  
1178 and partly discriminative) and generative modeling.

1179 Thus, in the context of text classification, the multinomial naive Bayes model of  
1180 Raina et al. [2004] divides the set of predictors into  $R$  regions. For the sake of clarity,  
1181 we will focus on  $R = 2$ , and therefore  $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$ . Equation (8) is modified as

$$p(c|\mathbf{x}) \propto p(c)p(\mathbf{x}_1|c)^{\frac{w_1}{n_1}} p(\mathbf{x}_2|c)^{\frac{w_2}{n_2}},$$

1182 where  $(w_1, w_2)$  controls the relative weighting between the regions, and  $n_1, n_2$  are their  
1183 lengths. For instance, in emails consisting of two regions, subject and body,  $n_2 \gg n_1$   
1184 since bodies are usually much longer than subjects, and the usual naive Bayes equation  
1185 will be mostly dominated by the message body (with many more factors). This model  
1186 tries instead to convey that different predictors are of different importance (words in  
1187 the subject might be more important) and counteracts the independence assumption of  
1188 naive Bayes with normalization factors  $n_1, n_2$ . The expression of  $p(c|\mathbf{x})$  is then rewritten  
1189 in a logistic regression form, where its linear combination contains parameters, gener-  
1190 atively learned functions of  $p(\mathbf{x}_i|c)$ . Parameters  $w_i$  are discriminatively learned (by  
1191 maximizing the CLL),  $i = 1, 2$ . They call this model the *normalized hybrid* algorithm,  
1192 designed for a binary class. A multiclass extension is reported in Fujino et al. [2007].

1193 The *joint discriminative-generative* (JoDiG) approach of Xue and Titterington [2010]  
1194 partitions  $\mathbf{X}$  into two subvectors:  $\mathbf{X} = (\mathbf{X}_D, \mathbf{X}_G)$ . A generative approach is applied to  $\mathbf{X}_G$   
1195 to estimate  $p(\mathbf{x}_G|c)$  and a discriminative approach is applied to  $\mathbf{X}_D$  to estimate  $p(c|\mathbf{x}_D)$ .  
1196 A data-generating process is always assumed in generative but never in discriminative  
1197 approaches. In general, when this process is well specified, the generative approach  
1198 performs better than the discriminative approach. This is the idea for finding the  
1199 partition of  $\mathbf{X}$ :  $\mathbf{X}_D$  will contain the variables that violate the assumption underlying  
1200 the data-generating process (as given by a statistical test). Finally, since  $\mathbf{X}_G$  and  $\mathbf{X}_D$   
1201 are assumed to be (block-wise) conditionally independent given  $C$ , then  $p(\mathbf{x}_D, \mathbf{x}_G, c) =$   
1202  $p(\mathbf{x}_D)p(c|\mathbf{x}_D)p(\mathbf{x}_G|c)$ , and both approaches are probabilistically combined to classify a  
1203 new instance via the MAP criterion

$$\arg \max_c p(c|\mathbf{x}_D)p(\mathbf{x}_G|c).$$

1204 The *hybrid generative/discriminative Bayesian* (HBayes) classifier [Kang and Tian  
1205 2006] uses a similar idea. The difference lies in how the partition is chosen, for which  
1206 purpose a wrapper strategy is adopted in this case: starting from  $\mathbf{X}_G = \mathbf{X}$ , the variable  
1207 producing the greatest improvement in classification performance is greedily moved  
1208 from  $\mathbf{X}_G$  to  $\mathbf{X}_D$ . Ridge logistic regression is used to estimate  $p(c|\mathbf{x}_D)$ , whereas naive  
1209 Bayes or TAN is used to estimate  $p(\mathbf{x}_G|c)$ . The Bayesian network structure is thereby  
1210 restricted (Figure 13) to reduce the computational effort.

1211 A Bayesian approach for the combination of generative and discriminative learning  
1212 of classifiers is found in Bishop and Lasserre [2007]. This is intended to find the appro-  
1213 priate tradeoff between generative and discriminative extremes. Generative and dis-  
1214 criminative models correspond to specific choices for the priors over parameters. Since  
1215 generative approaches can model unlabelled instances while discriminative approaches  
1216 do not, this *Bayesian blending* can also be applied to semisupervised classification.

## 1217 9. BAYESIAN MULTINETS

1218 Bayesian networks are unable to encode *asymmetric* independence assertions in their  
1219 topology. This refers to conditional independence relationships only held for some but

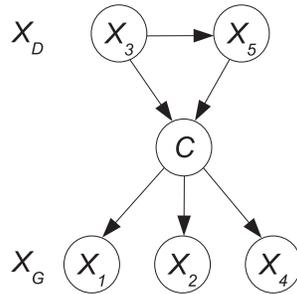


Fig. 13. A HBayes classifier structure from which  $p(c|\mathbf{x}) \propto p(c|\mathbf{x}_D)p(\mathbf{x}_G|c)$ .

not all the values of the variables involved. *Bayesian multinets* [Geiger and Heckerman 1996] offer a solution. They consist of several (local) Bayesian networks associated with a subset of a partition of the domain of a variable  $H$ , called the hypothesis or distinguished variable; that is, each local network represents a joint probability of all (but  $H$ ) variables conditioned on a subset of  $H$  values. As a result of this conditioning, asymmetric independence assertions are represented in each local network topology. Consequently, structures are expected to be simpler, with computational and memory requirement savings. Whereas the typical setting is when  $H$  is a root node, other situations are addressed in Geiger and Heckerman [1996]:  $H$  is a nonroot node, and there is more than one variable representing hypotheses.

For classification problems, the distinguished variable is naturally the class variable  $C$ . All subsets of the  $C$  domain partition are commonly singletons. Thus, conditioned on each  $c$ , the predictors can form different local networks with different structures. Therefore, the relations among variables do not have to be the same for all  $c$ . Equation (1) is, for Bayesian multinets, given by

$$p(c|\mathbf{x}) \propto p(c) \prod_{i=1}^n p(x_i|\mathbf{pa}_c(x_i)),$$

where  $\mathbf{Pa}_c(X_i)$  is the parent set of  $X_i$  in the local Bayesian network associated with  $C = c$ ; see Figure 1. Therefore, a Bayesian multinet is defined via its local Bayesian networks and the prior distribution on  $C$ .

Particular cases of multinets were explained in Section 6.1: networks reported in Chow and Liu [1968] and Pham et al. [2002] with trees and forests, respectively, as local Bayesian networks (illustrated in Figure 7(a) and (d)). Trees are also used in Kłopotek [2005], although the learning is based on a new algorithm designed for very large datasets rather than Kruskal’s algorithm. The trees in Huang et al. [2003] are learned by optimizing a function that includes a penalty term representing the divergence between the different joint distributions defined at each local network. Finally, the trees in Gurwicz and Lerner [2006] are learned from all instances, instead of learning the local structures from only those instances with  $C = c$ . The process is guided by a score that simultaneously detects class patterns and rejects patterns of the other classes. Thus, for the local network for  $C = c$ , the score of  $\mathbf{x}$  with true class value  $c$  is higher when  $p(C = c|\mathbf{x}) \geq p(C = c'|\mathbf{x}), \forall c' \neq c$  and the score of  $\mathbf{x}$  with true class value  $c' \neq c$  is higher when  $p(C = c'|\mathbf{x}) \geq p(C = c|\mathbf{x})$ . The search is based on the hill-climbing algorithm described in Keogh and Pazzani [2002] (see Section 6.1).

The local structures are general unrestricted Bayesian networks in Friedman et al. [1997] and Hussein and Santos [2004]. However, the approach taken in Hussein and Santos [2004] is different. The data are not partitioned according to  $C = c$ . The training

Table III. Mean Accuracies (%)  $\pm$  Standard Deviations of the 12 Bayesian Network Classifiers  
 “#” means the number of variables included in the model.

	All variables	#	Filter	#	Wrapper	#
Naive Bayes	71.64 $\pm$ 9.78	9	71.98 $\pm$ 11.59	5	77.20 $\pm$ 8.01	3
Tree-augmented naive Bayes	77.57 $\pm$ 8.08	9	76.50 $\pm$ 9.10	5	77.55 $\pm$ 9.35	5
Bayesian network-augmented naive Bayes	74.78 $\pm$ 8.62	9	76.83 $\pm$ 10.54	5	77.22 $\pm$ 10.14	6
Markov blanket-based Bayesian classifiers	75.16 $\pm$ 7.62	9	73.74 $\pm$ 7.67	5	76.52 $\pm$ 9.00	6

1255 data are first partitioned into clusters from which a set of rules characterizing their  
 1256 cases are derived. Then a local Bayesian network is learned from the cases satisfying  
 1257 the rules. This is why the resulting models are called *case-based Bayesian network clas-*  
 1258 *sifiers*, capturing case-dependent relationships, a generalization of hypothesis-specific  
 1259 relationships.

## 1260 10. ILLUSTRATIVE EXAMPLE

1261 This section reports the classification accuracy results of 12 different Bayesian net-  
 1262 work classifiers, according to four increasing model complexities (naive Bayes, tree-  
 1263 augmented naive Bayes, Bayesian network-augmented naive Bayes, and Markov  
 1264 blanket-based Bayesian classifiers) including all predictor variables and using two  
 1265 feature subset selection methods (a filter and a wrapper approach). The filter approach  
 1266 is univariate and based on information gain, whereas the wrapper search uses a greedy  
 1267 forward strategy in all models but the Markov blanked-based classifier, which employs  
 1268 a genetic algorithm.

1269 The classifiers were learned from the Ljubljana breast cancer dataset [Michalski  
 1270 et al. 1986] with 286 labeled instances of real patients. The classification problem was  
 1271 to predict breast cancer recurrence (yes or no) in the 5 years after surgery. Recurrence  
 1272 was observed in 85 out of the 286 patients. The nine predictor variables, measured at  
 1273 diagnosis, are:

- 1274 —age: patient age in years, discretized into three equal-width intervals
- 1275 —menopause: non-, pre-, or postmenopausal patient
- 1276 —deg-malig: degree of tumor malignancy (histological grade scored 1–3)
- 1277 —node-caps: whether or not the tumor has perforated through the lymph node capsule
- 1278 —inv-nodes: the number (range 0–26) of involved axillary lymph nodes that contain  
 1279 metastatic breast cancer visible on histological examination, discretized into three  
 1280 intervals
- 1281 —irradiation: whether or not the patient has been irradiated
- 1282 —breast: left- or right-sided breast cancer
- 1283 —breast-quad: location of the tumor according to the four breast quadrants (upper-  
 1284 outer, lower-outer, upper-inner, and lower-inner) plus the nipple as a central point
- 1285 —size: maximum excised tumor diameter (in mm), discretized into three equal-width  
 1286 intervals

1287 Table III shows the classification accuracy (%) and standard deviations of all model  
 1288 combinations. They have been estimated with 10-fold stratified cross-validation using  
 1289 WEKA [Hall et al. 2009] software.

1290 Naive Bayes and the filter-based selective naive Bayes (Figure14(a)) are the worst-  
 1291 performing algorithms ( $\approx$ 71% accuracy). However, the accuracy of selective naive Bayes  
 1292 increases considerably (up to 77%) using a wrapper-wise-guided search, with only  
 1293 three predictor variables. WEKA was parameterized to run similar algorithms to those  
 1294 proposed in the literature and reviewed within this article: Maron and Kuhns [1960]  
 1295 for naive Bayes, Pazzani and Billsus [1997] for filter-based selective naive Bayes, and  
 1296 Langley and Sage [1994] for wrapper-based selective naive Bayes.

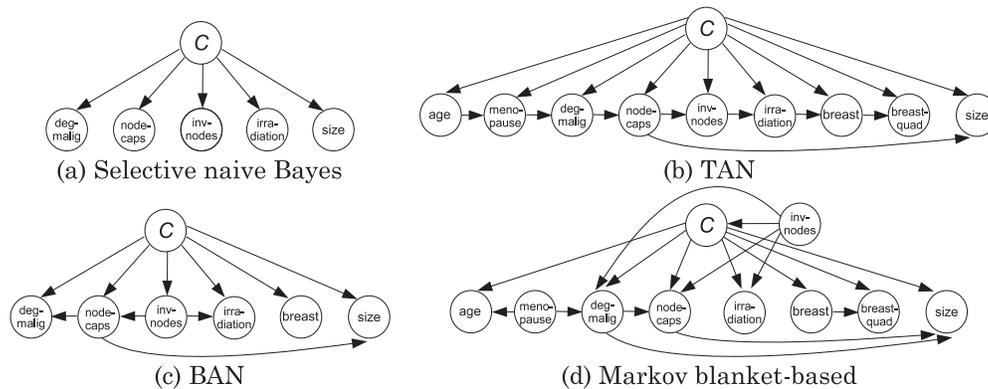


Fig. 14. Structures of (a) selective naive Bayes output using a filter approach, (b) TAN, (c) wrapper BAN, and (d) Markov blanket-based Bayesian classifier.

TAN and its selective versions (filter and wrapper) are the best-performing models on average. The TAN spanning tree (Figure 14(b)) is rooted at node age. It captures expected relationships, as specified by the arcs age→menopause, deg-malign→node-caps, and node-caps→size. Age and menopause are obviously related. There is a greater likelihood of the tumor penetrating through the lymph node capsule and invading the surrounding tissues at worse tumor grades. Tumor grade also conditions tumor size. The WEKA algorithms for these TAN models were similar to the learning algorithms described in Friedman et al. [1997] for TAN and in Blanco et al. [2005] for both selective TAN models.

BAN models (Table III, row 3) were learned by setting the maximum number of parents to 3. Selective BAN models behave similarly to their TAN counterparts. Without feature selection, BAN accuracy decreases. The best BAN, which is in fact a FAN (Figure 14(c)), is the wrapper version. This model did not select age, menopause, and breast-quad. Its structure shares two arcs with the TAN classifier (Figure 14(b)), node-caps→size and inv-nodes→irradiation. TAN also identified arcs inv-nodes→node-caps and node-caps→deg-malign, albeit reversed. The most similar algorithms to those run in WEKA are Friedman et al. [1997] for BAN, Ezawa and Norton [1996] for the filter-based BAN, and Pernkopf and O’Leary [2003] for the wrapper-based BAN.

Finally, despite the flexibility of the Markov blanket-based classifier structures, they do not exhibit very high accuracies. Without variable selection (Figure 14(d)), C has only one parent, inv-nodes. This model has many relationships in common with TAN (Figure 14(b)). However, three nodes (deg-malign, node-caps, and size) have three parents, requiring bigger conditional probability tables. Also, there is a new arc, deg-malign→size (justified by following the aforementioned reasoning), and a missing arc, C→menopause. The algorithm reported in Madden [2002] is close to the WEKA implementations of Markov blanket-based classifiers (all variables and filter), whereas we used WEKA’s genetic algorithm-guided search for the wrapper version as reported in Sierra and Larrañaga [1998].

In summary, the wrapper versions are the models that work best here. All of them include at least the inv-nodes, deg-malign, and breast variables. Filter approaches seem to improve the all-variables strategy. With only nine variables, carefully chosen by physicians to be relevant for the problem, the advantages of feature selection are limited. The best model is the wrapper-based TAN. Thus, increasing model complexity does not necessarily imply a better model. This is why it is always worthwhile to explore the whole hierarchy of Bayesian classifiers.

1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331

Table IV. Summary of Bayesian Network Classifiers and Their Most Relevant References

Name	Structure	Feature subset selection			Metaclassifiers
		Seminal paper	Filter	Wrapper	
Naive Bayes		[Maron and Kuhns 1960]	NA	NA	[Langley 1993]
Selective naive Bayes		[Langley and Sage 1994]	[Pazzani and Billsus 1997]	[Langley and Sage 1994]	[Zheng 1998]
Semi-naive Bayes		[Pazzani 1996]	[Blanco et al. 2005]	[Robles et al. 2003]	[Robles et al. 2004]
Tree-augmented naive Bayes		[Friedman et al. 1997]	[Blanco et al. 2005]	[Keogh and Pazzani 2002]	[Ma and Shi 2004]
Forest-augmented naive Bayes		[Lucas 2004]	[Ziebart et al. 2007]		
Superparent-one-dependence estimator		[Keogh and Pazzani 2002]	NA	NA	[Webb et al. 2005]
k-dependence Bayesian classifier		[Sahami 1996]	[Blanco et al. 2005]	[Blanco et al. 2005]	[Louzada and Ara 2012]
Bayesian network-augmented naive Bayes		[Friedman et al. 1997]	[Ezawa and Norton 1996]	[Pernkopf and O'Leary 2003]	
Markov blanket-based Bayesian classifiers		[Koller and Sahami 1996]	[Koller and Sahami 1996]	[Sierra and Larrañaga 1998]	
Unrestricted Bayesian classifiers		[Provan and Singh 1995]	[Singh and Provan 1996]	[Provan and Singh 1995]	[Dash and Cooper 2004]
Bayesian multinet		[Geiger and Heckerman 1996]			

## 11. DISCUSSION

This survey has shown the power of Bayesian network classifiers in terms of model expressiveness and algorithm efficiency/effectiveness for learning models from data and for use in classification. Unlike other pattern recognition classifiers, Bayesian network classifiers can be clearly organized hierarchically from the simplest naive Bayes to the most complex Bayesian multinet.

The Bayesian network classifiers are hierarchized in the rows of Table IV, whereas the columns give an example of their graphical structure, the associated seminal paper, and the first references proposing filter/wrapper approaches for feature subset selection and metaclassifiers.

We did not set out to survey the behavior of these classifiers in big real-world problems. As the no-free-lunch theorem states, this depends on the dataset. However, some relevant papers, already cited within this survey [Friedman et al. 1997; Cheng and Greiner 1999, 2001; Pernkopf 2005; Madden 2009], do include empirical comparisons of the algorithms for learning naive Bayes, TAN, BAN, unrestricted Bayesian classifiers, and Bayesian multinets. They all use datasets from the UCI repository [Bache and Lichman 2013]. Also, both discriminative and generative parameter learning on both discriminatively and generatively structured models are compared in Pernkopf and Bilmes [2005]. The general findings are that more complex structures perform better whenever the sample size is big enough to guarantee reliable probability estimates. Also, smoothing parameter estimation can significantly improve the classification rate. Discriminative parameter learning produces on average a better classifier than maximum likelihood parameter learning. In most datasets, structures learned with wrapper approaches yield the most accurate classifiers.

Since the focus of this article is on Bayesian network classifiers based on Bayesian networks, other models—models with cycles, like dependency networks, and undirected models, like Markov networks—are beyond its scope. We have not considered data-streaming situations or specific problems like multilabel or semisupervised classification or classification with probabilistic labels either. Although the survey has focused on discrete data, research on continuous and mixed data is on-going.

Research on discrete Bayesian network classifiers may in the future target more theoretical studies on determining the decision boundary for classifier types apart from the naive Bayes reviewed here. Also, the gaps in Table IV suggest that there is still room for research on metaclassifiers and feature subset selection. Metaclassifiers might also be formed by hybridizing Bayesian classifiers with different types of classifiers other than the decision trees and  $k$ -nearest neighbors mentioned in this article. Finally, we have seen how naive Bayes can tackle complex classification situations (e.g., with homologous sets, multiple instances, cost-sensitive learning, instance ranking, and imprecise probabilities). We expect to see other models dealing with these and more challenging settings soon.

## REFERENCES

- J. Abellán. 2006. Application of uncertainty measures on credal sets on the naive Bayes classifier. *International Journal of General Systems* 35 (2006), 675–686.
- J. Abellán, A. Cano, A. R. Masegosa, and S. Moral. 2007. A semi-naive Bayes classifier with grouping of cases. In *Proceedings of the 9th European Conference in Symbolic and Quantitative Approaches to Reasoning with Uncertainty (ECSQARU-2007)*. *Lecture Notes in Artificial Intelligence*, Vol. 4724. Springer, 477–488.
- S. Acid, L. M. de Campos, and J. G. Castellano. 2005. Learning Bayesian network classifiers: Searching in a space of partially directed acyclic graphs. *Machine Learning* 59, 3 (2005), 213–235.
- A. Agresti. 1990. *Categorical Data Analysis*. Wiley.
- K. M. Al-Aidaroos, A. A. Bakar, and Z. Othman. 2010. Naive Bayes variants in classification learning. In *Proceedings of the International Conference on Information Retrieval Knowledge Management (CAMP-2010)*. 276–281.

- 1384 C. F. Aliferis, A. R. Statnikov, I. Tsamardinos, S. Mani, and X. D. Koutsoukos. 2010. Local causal and Markov  
 1385 blanket induction for causal discovery and feature selection for classification. Part I: Algorithms and  
 1386 empirical evaluation. *Journal of Machine Learning Research* 11 (2010), 171–234.
- 1387 C. F. Aliferis, I. Tsamardinos, and M. S. Statnikov. 2003. HITON: A novel Markov blanket algorithm for  
 1388 optimal variable selection. In *AMIA Annual Symposium Proceedings*. 21–25.
- 1389 K. Bache and M. Lichman. 2013. UCI Machine Learning Repository. (2013). Retrieved from <http://archive.ics.uci.edu/ml>.  
 1390
- 1391 X. Bai, R. Padman, J. Ramsey, and P. Spirtes. 2008. Tabu search-enhanced graphical models for classification  
 1392 in high dimensions. *INFORMS Journal on Computing* 20, 3 (2008), 423–437.
- 1393 J. Bilmes. 2000. Dynamic Bayesian multinets. In *Proceedings of the 16th Conference in Uncertainty in  
 1394 Artificial Intelligence (UAI-2000)*. Morgan Kaufmann, 38–45.
- 1395 C. Bishop. 1995. *Neural Networks for Pattern Recognition*. Oxford University Press.
- 1396 C. M. Bishop and J. Lasserre. 2007. Generative or discriminative? Getting the best of both worlds. In *Bayesian  
 1397 Statistics*, Vol. 8. Oxford University Press, 3–23.
- 1398 R. Blanco, I. Inza, M. Merino, J. Quiroga, and P. Larrañaga. 2005. Feature selection in Bayesian classifiers  
 1399 for the prognosis of survival of cirrhotic patients treated with TIPS. *Journal of Biomedical Informatics*  
 1400 38, 5 (2005), 376–388.
- 1401 W. L. Buntine. 1991. Theory refinement on Bayesian networks. In *Proceedings of the 7th Conference on  
 1402 Uncertainty in Artificial Intelligence (UAI-1991)*. Morgan Kaufmann, 52–60.
- 1403 J. Burge and T. Lane. 2005. Learning class-discriminative dynamic Bayesian networks. In *Proceedings of the  
 1404 22nd International Conference on Machine Learning (ICML-2005)*. ACM, 97–104.
- 1405 A. Cano, J. G. Castellano, A. R. Masegosa, and S. Moral. 2005. Methods to determine the branching at-  
 1406 tribute in Bayesian multinets classifiers. In *Proceedings of the 8th European Conference in Symbolic  
 1407 and Quantitative Approaches to Reasoning with Uncertainty (ECSQARU-2005)*. *Lecture Notes in Artificial  
 1408 Intelligence*, Vol. 3571. Springer, 932–943.
- 1409 A. M. Carvalho, A. L. Oliveira, and M.-F. Sagot. 2007. Efficient learning of Bayesian network classifiers. In  
 1410 *Proceedings of the 20th Australian Joint Conference on Artificial Intelligence (AI-2007)*. *Lecture Notes in  
 1411 Computer Science*, Vol. 4830. Springer, 16–25.
- 1412 A. M. Carvalho, T. Roos, A. L. Oliveira, and P. Myllymäki. 2011. Discriminative learning of Bayesian net-  
 1413 works via factorized conditional log-likelihood. *Journal of Machine Learning Research* 12 (2011), 2181–  
 1414 2210.
- 1415 J. Cerquides and R. López de Mántaras. 2005a. Robust Bayesian linear classifier ensembles. In *Proceedings  
 1416 of the 16th European Conference on Machine Learning (ECML-2005)*. *Lecture Notes in Computer Science*,  
 1417 Vol. 3720. Springer, 72–83.
- 1418 J. Cerquides and R. López de Mántaras. 2005b. TAN classifiers based on decomposable distributions. *Machine  
 1419 Learning* 59, 3 (2005), 323–354.
- 1420 B. Cestnik. 1990. Estimating probabilities: A crucial task in machine learning. In *Proceedings of the European  
 1421 Conference in Artificial Intelligence*. 147–149.
- 1422 X. Chai, L. Deng, Q. Yang, and C. X. Ling. 2004. Test-cost sensitive naive Bayes classification. In *Proceed-  
 1423 ings of the 4th IEEE International Conference on Data Mining (ICDM-2004)*. IEEE Computer Society,  
 1424 51–58.
- 1425 J. Cheng and R. Greiner. 1999. Comparing Bayesian network classifiers. In *Proceedings of the 15th Conference  
 1426 on Uncertainty in Artificial Intelligence (UAI-1999)*. Morgan Kaufmann Publishers, 101–108.
- 1427 J. Cheng and R. Greiner. 2001. Learning Bayesian belief networks classifiers: Algorithms and system.  
 1428 In *Proceedings of the 14th Biennial Conference of the Canadian Society for Computational Studies of  
 1429 Intelligence (CSCSI-2001)*, Vol. 2056. Springer, 141–151.
- 1430 D. M. Chickering. 1995. A transformational characterization of equivalent Bayesian network structures.  
 1431 In *Proceedings of the 11th Conference on Uncertainty in Artificial Intelligence (UAI-1995)*. Morgan  
 1432 Kaufmann, 87–98.
- 1433 D. M. Chickering, D. Heckerman, and C. Meek. 2004. Large-sample learning of Bayesian networks is NP-  
 1434 hard. *Journal of Machine Learning Research* 5 (2004), 1287–1330.
- 1435 C. Chow and C. Liu. 1968. Approximating discrete probability distributions with dependency trees. *IEEE  
 1436 Transactions on Information Theory* 14 (1968), 462–467.
- 1437 G. F. Cooper and E. Herskovits. 1992. A Bayesian method for the induction of probabilistic networks from  
 1438 data. *Machine Learning* 9 (1992), 309–347.
- 1439 D. Dash and G. F. Cooper. 2004. Model averaging for prediction with discrete Bayesian networks. *Journal of  
 1440 Machine Learning Research* 5 (2004), 1177–1203.

## Discrete Bayesian Network Classifiers: A Survey

60:37

- D. Dash and G. F. Cooper. 2002. Exact model averaging with naïve Bayesian classifiers. In *Proceedings of the 19th International Conference on Machine Learning (ICML-2002)*. 91–98. 1441  
1442
- A. P. Dempster, N. M. Laird, and D. B. Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B* 39, 1 (1977), 1–38. 1443  
1444
- P. Domingos and M. Pazzani. 1997. On the optimality of the simple Bayesian classifier under zero-one loss. *Machine Learning* 29 (1997), 103–130. 1445  
1446
- E. B. dos Santos, E. R. Hruschka Jr., E. R. Hruschka, and N. F. F. Ebecken. 2011. Bayesian network classifiers: Beyond classification accuracy. *Intelligent Data Analysis* 15, 3 (2011), 279–298. 1447  
1448
- M. M. Drugan and M. A. Wiering. 2010. Feature selection for Bayesian network classifiers using the MDL-FS score. *International Journal of Approximate Reasoning* 51 (2010), 695–717. 1449  
1450
- R. Duda, P. Hart, and D. G. Stork. 2001. *Pattern Classification*. John Wiley and Sons. 1451
- D. Edwards and S. L. Lauritzen. 2001. The TM algorithm for maximising a conditional likelihood function. *Biometrika* 88 (2001), 961–972. 1452  
1453
- M. Ekdahl and T. Koski. 2006. Bounds for the loss in probability of correct classification under model based approximation. *Journal of Machine Learning Research* 7 (2006), 2449–2480. 1454  
1455
- S. Eyheramendy, D. D. Lewis, and D. Madigan. 2002. On the naive Bayes model for text categorization. In *Proceedings of the 9th International Workshop on Artificial Intelligence and Statistics (AISTATS-2002)*. 1456  
1457
- K. J. Ezawa and S. W. Norton. 1996. Constructing Bayesian networks to predict uncollectible telecommunications accounts. *IEEE Expert* 11, 5 (1996), 45–51. 1458  
1459
- A. J. Feelders and J. Ivanovs. 2006. Discriminative scoring of Bayesian network classifiers: A comparative study. In *Proceedings of the 3rd European Workshop on Probabilistic Graphical Models (PGM-2006)*. 75–82. 1460  
1461  
1462
- Q. Feng, F. Tian, and H. Huang. 2007. A discriminative learning method of TAN classifier. In *Proceedings of the 9th European Conference in Symbolic and Quantitative Approaches to Reasoning with Uncertainty (ECSQARU-2007)*. *Lecture Notes in Artificial Intelligence*, Vol. 4724. Springer, 443–452. 1463  
1464  
1465
- J. Flores, J. A. Gámez, and A. M. Martínez. 2012. Supervised classification with Bayesian networks: A review on models and applications. In *Intelligent Data Analysis for Real World Applications. Theory and Practice*. IGI Global, 72–102. 1466  
1467  
1468
- M. J. Flores, J. A. Gámez, A. M. Martínez, and J. M. Puerta. 2009. HODE: Hidden one-dependence estimator. In *Proceedings of the 10th European Conference in Symbolic and Quantitative Approaches to Reasoning with Uncertainty (ECSQARU-2009)*. *Lecture Notes in Artificial Intelligence*, Vol. 5590. Springer, 481–492. 1469  
1470  
1471  
1472
- O. François and P. Leray. 2006. Learning the tree augmented naive Bayes classifier from incomplete datasets. In *Proceedings of the 3rd European Workshop on Probabilistic Graphical Models (PGM-2006)*. 91–98. 1473  
1474
- E. Frank, M. Hall, and B. Pfahringer. 2003. Locally weighted naive Bayes. In *Proceedings of the 19th Conference in Uncertainty in Artificial Intelligence (UAI-2003)*. Morgan Kaufmann, 249–256. 1475  
1476
- M. L. Fredman and R. E. Tarjan. 1987. Fibonacci heaps and their uses in improved network optimization algorithms. *Journal ACM* 34, 3 (1987), 596–615. 1477  
1478
- N. Friedman. 1997. Learning belief networks in the presence of missing values and hidden variables. In *Proceedings of the 14th International Conference on Machine Learning (ICML-1997)*. Morgan Kaufmann, 125–133. 1479  
1480  
1481
- N. Friedman, D. Geiger, and M. Goldszmidt. 1997. Bayesian network classifiers. *Machine Learning* 29 (1997), 131–163. 1482  
1483
- N. Friedman, M. Goldszmidt, and A. Wyner. 1999. Data analysis with Bayesian networks: A bootstrap approach. In *Proceedings of the 15th Conference on Uncertainty in Artificial Intelligence (UAI-1999)*. Morgan Kaufmann, 196–205. 1484  
1485  
1486
- S. Fu and M. Desmarais. 2007. Local learning algorithm for Markov blanket discovery. In *Proceedings of the 20th Australian Joint Conference on Artificial Intelligence (AI-2007)*. *Lecture Notes in Computer Science*, Vol. 4830. Springer, 68–79. 1487  
1488  
1489
- A. Fujino, N. Ueda, and K. Saito. 2007. A hybrid generative/discriminative approach to text classification with additional information. *Information Processing and Management* 43, 2 (2007), 379–392. 1490  
1491
- J. Gama. 1999. Iterative naive Bayes. *Theoretical Computer Science* 292, 2 (1999), 417–430. 1492
- D. Geiger and D. Heckerman. 1996. Knowledge representation and inference in similarity networks and Bayesian multinets. *Artificial Intelligence* 82 (1996), 45–74. 1493  
1494
- M. Goldszmidt. 2010. Bayesian network classifiers. In *Wiley Encyclopedia of Operations Research and Management Science*. John Wiley & Sons, 1–10. 1495  
1496
- I. J. Good. 1965. *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*. The MIT Press. 1497

- 1498 R. Greiner, X. Su, B. Shen, and W. Zhou. 2005. Structural extension to logistic regression: Discriminative  
1499 parameter learning of belief net classifiers. *Machine Learning* 59, 3 (2005), 297–322.
- 1500 R. Greiner and W. Zhou. 2002. Structural extension to logistic regression: Discriminative parameter learning  
1501 of belief net classifiers. In *Proceedings of the 18th National Conference on Artificial Intelligence (AAAI-  
1502 2002)*. AAAI Press/MIT Press, 167–173.
- 1503 D. Grossman and P. Domingos. 2004. Learning Bayesian network classifiers by maximizing conditional  
1504 likelihood. In *Proceedings of the 21st International Conference on Machine Learning (ICML-2004)*. 361–  
1505 368.
- 1506 Y. Guo and R. Greiner. 2005. Discriminative model selection for belief net structures. In *Proceedings of  
1507 the 20th National Conference on Artificial Intelligence (AAAI-2005)*. AAAI Press / The MIT Press, 770–  
1508 776.
- 1509 Y. Guo, D. F. Wilkinson, and D. Schuurmans. 2005. Maximum margin Bayesian networks. In *Proceedings of  
1510 the 21st Conference in Uncertainty in Artificial Intelligence (UAI-2005)*. AUAI Press, 233–242.
- 1511 Y. Gurwicz and B. Lerner. 2006. Bayesian class-matched multinet classifier. In *Proceedings of the 2006 Joint  
1512 IAPR international Conference on Structural, Syntactic, and Statistical Pattern Recognition (SSPR-  
1513 2006/SPR-2006)*. Lecture Notes in Computer Science, Vol. 4109. Springer, 145–153.
- 1514 M. A. Hall. 1999. *Correlation-Based Feature Selection for Machine Learning*. Ph.D. Dissertation. Department  
1515 of Computer Science, University of Waikato.
- 1516 M. Hall. 2007. A decision tree-based attribute weighting filter for naive Bayes. *Knowledge-Based Systems*  
1517 20, 2 (2007), 120–126.
- 1518 M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. 2009. The WEKA data mining  
1519 software: An update. *SIGKDD Explorations* 11, 1 (2009), 10–18.
- 1520 D. J. Hand and K. Yu. 2001. Idiot’s Bayes - not so stupid after all? *International Statistical Review* 69, 3  
1521 (2001), 385–398.
- 1522 D. Heckerman, D. Geiger, and D. Chickering. 1995. Learning Bayesian networks: The combination of knowl-  
1523 edge and statistical data. *Machine Learning* 20 (1995), 197–243.
- 1524 J. Hilden and B. Bjerregaard. 1976. Computer-aided diagnosis and the atypical case. In *Decision Making  
1525 and Medical Care. Can Information Science Help?* 365–378.
- 1526 E. R. Hruschka and N. F. F. Ebecken. 2007. Towards efficient variables ordering for Bayesian network  
1527 classifiers. *Data and Knowledge Engineering* 63 (2007), 258–269.
- 1528 H. Huang and C. Hsu. 2002. Bayesian classification for data from the same unknown class. *IEEE Transactions  
1529 on Systems, Man, and Cybernetics Part B* 32, 2 (2002), 137–145.
- 1530 K. Huang, I. King, and M. R. Lyu. 2003. Discriminative training of Bayesian Chow-Liu multinet classifiers.  
1531 In *Proceedings of the International Joint Conference on Neural Networks (IJCNN-2003)*, Vol. 1. 484–  
1532 488.
- 1533 A. Hussein and E. Santos. 2004. Exploring case-based Bayesian networks and Bayesian multi-nets for  
1534 classification. In *Proceedings of the 17th Conference of the Canadian Society for Computational Studies  
1535 of Intelligence (CSCSI-2004)*. Lecture Notes in Computer Science, Vol. 3060. Springer, 485–492.
- 1536 K.-B. Hwang and B. T. Zhang. 2005. Bayesian model averaging of Bayesian network classifiers over multiple  
1537 node-orders: Application to sparse datasets. *IEEE Transactions on Systems, Man, and Cybernetics. Part  
1538 B: Cybernetics* 35, 6 (2005), 1302–1310.
- 1539 A. Ibáñez, P. Larrañaga, and C. Bielza. 2014. Cost-sensitive selective naive Bayes classifiers for predicting  
1540 the increase of the h-index for scientific journals. *Neurocomputing* in press (2014).
- 1541 I. Inza, P. Larrañaga, R. Blanco, and A. J. Cerrolaza. 2004. Filter versus wrapper gene selection approaches  
1542 in DNA microarray domains. *Artificial Intelligence in Medicine* 31, 2 (2004), 91–103.
- 1543 I. Inza, P. Larrañaga, R. Etxeberria, and B. Sierra. 2000. Feature subset selection by Bayesian network-based  
1544 optimization. *Artificial Intelligence* 123, 1–2 (2000), 157–184.
- 1545 A. G. Ivakhnenko. 1970. Heuristic self-organization in problems of engineering cybernetics. *Automatica* 6, 2  
1546 (1970), 207–219.
- 1547 N. Japkowicz and S. Mohak. 2011. *Evaluating Learning Algorithms. A Classification Perspective*. Cambridge  
1548 University Press.
- 1549 T. Jebara. 2004. *Machine Learning: Discriminative and Generative*. Kluwer Academic Publishers.
- 1550 L. Jiang, Z. Cai, D. Wang, and H. Zhang. 2012. Improving tree augmented Naive Bayes for class probability  
1551 estimation. *Knowledge-Based Systems* 26 (2012), 239–245.
- 1552 L. Jiang and H. Zhang. 2006. Lazy averaged one-dependence estimators. In *Proceedings of the 19th Cana-  
1553 dian Conference on AI (Canadian AI-2006)*. Lecture Notes in Computer Science, Vol. 4013. Springer,  
1554 515–525.

## Discrete Bayesian Network Classifiers: A Survey

60:39

- L. Jiang, H. Zhang, and Z. Cai. 2009. A novel Bayes model: Hidden naive Bayes. *IEEE Transactions on Knowledge and Data Engineering* 21, 10 (2009), 1361–1371. 1555
- L. Jiang, H. Zhang, Z. Cai, and D. Wang. 2012. Weighted average of one-dependence estimators. *Journal of Experimental and Theoretical Artificial Intelligence* 24, 2 (2012), 219–230. 1557
- Y. Jing, V. Pavlovic, and J. M. Rehg. 2008. Boosted Bayesian network classifiers. *Machine Learning* 73 (2008), 155–184. 1559
- C. Kang and J. Tian. 2006. A Hybrid generative/discriminative Bayesian classifier. In *Proceedings of the 19th International Florida Artificial Intelligence Research Society Conference (FLAIRS-2006)*. AAAI Press, 562–567. 1561
- E. J. Keogh and M. J. Pazzani. 2002. Learning the structure of augmented Bayesian classifiers. *International Journal on Artificial Intelligence Tools* 11, 4 (2002), 587–601. 1564
- M. A. Klopotek. 2005. Very large Bayesian multinets for text classification. *Future Generation Computer Systems* 21, 7 (2005), 1068–1082. 1566
- R. Kohavi. 1996. Scaling up the accuracy of naive-Bayes classifiers: A decision-tree hybrid. In *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining (KDD-1996)*. 202–207. 1568
- R. Kohavi, B. Becker, and D. Sommerfield. 1997. *Improving Simple Bayes*. Technical Report. Data Mining and Visualization Group, Silicon Graphics. 1570
- R. Kohavi and G. H. John. 1997. Wrappers for feature subset selection. *Artificial Intelligence* 97, 1 (1997), 273–324. 1571
- D. Koller and M. Sahami. 1996. Toward optimal feature selection. In *Proceedings of the 13th International Conference on Machine Learning (ICML-1996)*. 284–292. 1574
- I. Kononenko. 1993. Successive naive Bayesian classifier. *Informatica (Slovenia)* 17, 2 (1993), 167–174. 1576
- P. Kontkanen, P. Myllymäki, T. Silander, and H. Tirri. 1998. BAYDA: Software for Bayesian classification and feature selection. In *Proceedings of the 4th International Conference on Knowledge Discovery and Data Mining (KDD-1998)*. AAAI Press, 254–258. 1577
- P. Kontkanen, P. Myllymäki, and H. Tirri. 1996. *Constructing Bayesian Finite Mixture Models by the EM Algorithm*. Technical Report C-1996-9. Department of Computer Science, University of Helsinki. 1581
- J. B. Kruskal. 1956. On the shortest spanning subtree of a graph and the traveling salesman problem. *Proceedings of the American Mathematical Society* 7 (1956), 48–50. 1582
- C. K. Kwoh and D. Gillies. 1996. Using hidden nodes in Bayesian networks. *Artificial Intelligence* 88 (1996), 1–38. 1584
- P. Langley. 1993. Induction of recursive Bayesian classifiers. In *Proceedings of the 8th European Conference on Machine Learning (ECML-1993)*. 153–164. 1586
- P. Langley and S. Sage. 1994. Induction of selective Bayesian classifiers. In *Proceedings of the 10th Conference on Uncertainty in Artificial Intelligence (UAI-1994)*. Morgan Kaufmann, 399–406. 1588
- H. Langseth and T. D. Nielsen. 2006. Classification using hierarchical naïve Bayes models. *Machine Learning* 63, 2 (2006), 135–159. 1590
- J. Li, C. Zhang, T. Wang, and Y. Zhang. 2007. Generalized additive Bayesian network classifiers. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI-2007)*. 913–918. 1592
- J. N. K. Liu, N. L. Li, and T. S. Dillon. 2001. An improved naïve Bayes classifier technique coupled with a novel input solution method. *IEEE Transactions on Systems, Man, and Cybernetics* 31 (2001), 249–256. 1594
- D. J. Lizotte, O. Madani, and R. Greiner. 2003. Budgeted learning of naive-Bayes classifiers. In *Proceedings of the 19th Conference in Uncertainty in Artificial Intelligence (UAI-2003)*. Morgan Kaufmann, 378–385. 1596
- F. Louzada and A. Ara. 2012. Bagging k-dependence probabilistic networks: An alternative powerful fraud detection tool. *Expert Systems with Applications* 39, 14 (2012), 11583–11592. 1598
- P. Lucas. 2004. Restricted Bayesian network structure learning. In *Advances in Bayesian Networks*. Springer, 217–232. 1600
- S.-C. Ma and H.-B. Shi. 2004. Tree-augmented naive Bayes ensembles. In *Proceedings of the 3rd International Conference on Machine Learning and Cybernetics*. IEEE, 1497–1502. 1602
- M. G. Madden. 2009. On the classification performance of TAN and general Bayesian networks. *Knowledge-Based Systems* 22, 7 (2009), 489–495. 1604
- M. G. Madden. 2002. A new Bayesian network structure for classification tasks. In *Proceedings of the 13th Irish Conference on Artificial Intelligence and Cognitive Science*. 203–208. 1606
- D. Margaritis and S. Thrun. 2000. Bayesian network induction via local neighborhoods. In *Advances in Neural Information Processing Systems 12 (NIPS-1999)*. MIT Press, 505–511. 1608
- M. Maron and J. Kuhns. 1960. On relevance, probabilistic indexing, and information retrieval. *Journal of the Association for Computing Machinery* 7 (1960), 216–244. 1610

- 1612 W. J. McGill. 1954. Multivariate information transmission. *Psychometrika* 19 (1954), 97–116.
- 1613 R. S. Michalski, I. Mozetic, J. Hong, and N Lavrac. 1986. The multi-purpose incremental learning system  
1614 AQ15 and its testing application to three medical domains. In *Proceedings of the 5th National Conference*  
1615 *on Artificial Intelligence*. Morgan Kaufman, 1041–1045.
- 1616 M. Minsky. 1961. Steps toward artificial intelligence. *Transactions on Institute of Radio Engineers* 49 (1961),  
1617 8–30.
- 1618 D. Mladenic and M. Grobelnik. 1999. Feature selection for unbalanced class distribution and naive Bayes. In  
1619 *Proceedings of the 16th International Conference on Machine Learning (ICML-1999)*. Morgan Kaufmann,  
1620 258–267.
- 1621 S. Monti and G. F. Cooper. 1999. A Bayesian network classifier that combines a finite mixture model and  
1622 a naive Bayes model. In *Proceedings of the 15th Conference on Uncertainty in Artificial Intelligence*  
1623 *(UAI-1999)*. 447–456.
- 1624 M. Možina, J. Demšar, M. Kattan, and B. Zupan. 2004. Nomograms for visualization of naive Bayesian clas-  
1625 sifier. In *Proceedings of the 8th European Conference on Principles and Practice of Knowledge Discovery*  
1626 *in Databases (PKDD-2004)*. 337–348.
- 1627 J. F. Murray, G. F. Hughes, and K. Kreutz-Delgado. 2005. Machine learning methods for predicting failures in  
1628 hard drives: A multiple-instance application. *Journal of Machine Learning Research* 6 (2005), 783–816.
- 1629 M. Narasimhan and J. A. Bilmes. 2005. A submodular-supermodular procedure with applications to discrim-  
1630 inative structure learning. In *Proceedings of the 21st Conference in Uncertainty in Artificial Intelligence*  
1631 *(UAI-2005)*. AUAI Press, 404–412.
- 1632 A. Ng and M. Jordan. 2001. On discriminative vs. generative classifiers: A comparison of logistic regression  
1633 and naive Bayes. In *Advances in Neural Information Processing Systems 14 (NIPS-2001)*. MIT Press,  
1634 841–848.
- 1635 G. N. Norén and R. Orre. 2005. Case based imprecision estimates for Bayes classifiers with the Bayesian  
1636 bootstrap. *Machine Learning* 58, 1 (2005), 79–94.
- 1637 M. Pazzani. 1996. Constructive induction of Cartesian product attributes. In *Proceedings of the Information,*  
1638 *Statistics and Induction in Science Conference (ISIS-1996)*. 66–77.
- 1639 M. Pazzani and D. Billsus. 1997. Learning and revising user profiles: the identification of interesting web  
1640 sites. *Machine Learning* 27 (1997), 313–331.
- 1641 J. Pearl. 1988. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann, Palo Alto, CA.
- 1642 J. M. Peña, R. Nilsson, J. Björkegren, and J. Tegnér. 2007. Towards scalable and data efficient learning of  
1643 Markov boundaries. *International Journal of Approximate Reasoning* 45, 2 (2007), 211–232.
- 1644 F. Pernkopf. 2005. Bayesian network classifiers versus selective  $k$ -NN classifier. *Pattern Recognition* 38  
1645 (2005), 1–10.
- 1646 F. Pernkopf and J. A. Bilmes. 2005. Discriminative versus generative parameter and structure learning of  
1647 Bayesian network classifiers. In *Proceedings of the 22nd International Conference on Machine Learning*  
1648 *(ICML-2005)*. ACM, 657–664.
- 1649 F. Pernkopf and J. A. Bilmes. 2010. Efficient heuristics for discriminative structure learning of Bayesian  
1650 network classifiers. *Journal of Machine Learning Research* 11 (2010), 2323–2360.
- 1651 F. Pernkopf and P. O’Leary. 2003. Floating search algorithm for structure learning of Bayesian network  
1652 classifiers. *Pattern Recognition Letters* 24 (2003), 2839–2848.
- 1653 F. Pernkopf and M. Wohlmayr. 2009. On discriminative parameter learning of Bayesian network classifiers.  
1654 In *Proceedings of the 20th European Conference on Machine Learning (ECML-2009)*. *Lecture Notes in*  
1655 *Computer Science*, Vol. 5782. Springer, 221–237.
- 1656 F. Pernkopf and M. Wohlmayr. 2013. Stochastic margin-based structure learning of Bayesian network clas-  
1657 sifiers. *Pattern Recognition* 46, 2 (2013), 464–471.
- 1658 F. Pernkopf, M. Wohlmayr, and S. Tschiatschek. 2012. Maximum margin Bayesian network classifiers. *IEEE*  
1659 *Transactions on Pattern Analysis and Machine Intelligence* 34, 3 (2012), 521–532.
- 1660 T. V. Pham, M. Worring, and A. W. M. Smeulders. 2002. Face detection by aggregated Bayesian network  
1661 classifiers. *Pattern Recognition Letters* 23, 4 (2002), 451–461.
- 1662 B. Poulin, R. Eisner, D. Szafron, Paul Lu, R. Greiner, D. S. Wishart, A. Fyshe, B. Percy, C. MacDonell, and J.  
1663 Anvik. 2006. Visual explanation of evidence with additive classifiers. In *Proceedings of the 21th National*  
1664 *Conference on Artificial Intelligence (AAAI-2006)*. AAAI Press/MIT Press, 1822–1829.
- 1665 A. Prinzie and D. Van den Poel. 2007. Random multiclass classification: Generalizing random forests to  
1666 random MNL and random NB. In *Proceedings of the Database and Expert Systems Applications. Lecture*  
1667 *Notes in Computer Science*. Vol. 4653. Springer, 349–358.
- 1668 G. M. Provan and M. Singh. 1995. Learning Bayesian networks using feature selection. In *Proceedings of*  
1669 *the 5th International Workshop on Artificial Intelligence and Statistics (AISTATS-1995)*. 450–456.

## Discrete Bayesian Network Classifiers: A Survey

60:41

- R. Raina, Y. Shen, A. Y. Ng, and A. McCallum. 2004. Classification with hybrid generative/discriminative models. In *Advances in Neural Information Processing Systems 16 (NIPS-2003)*. The MIT Press. 1670–1671
- M. Ramoni and P. Sebastiani. 2001a. Robust Bayes classifiers. *Artificial Intelligence* 125 (2001), 209–226. 1672–1673
- M. Ramoni and P. Sebastiani. 2001b. Robust learning with missing data. *Machine Learning* 45, 2 (2001), 147–170. 1674–1675
- C. A. Ratanamahatana and D. Gunopulos. 2003. Feature selection for the naive Bayesian classifier using decision trees. *Applied Artificial Intelligence* 17, 5–6 (2003), 475–487. 1676–1677
- S. Renooij and L. C. van der Gaag. 2008. Evidence and scenario sensitivities in naive Bayesian classifiers. *International Journal of Approximate Reasoning* 49, 2 (2008), 398–416. 1678–1679
- G. Ridgeway, D. Madigan, and T. Richardson. 1998. Interpretable boosted naive Bayes classification. In *Proceedings of the 4th International Conference on Knowledge Discovery and Data Mining (KDD-1998)*. 101–104. 1680–1681–1682
- V. Robles, P. Larrañaga, J. M. Peña, E. Menasalvas, and M. S. Pérez. 2003. Interval estimation naive Bayes. In *Proceedings of the 5th International Symposium on Intelligent Data Analysis (IDA-2003). Lecture Notes in Computer Science*, Vol. 2810. Springer, 143–154. 1683–1684–1685
- V. Robles, P. Larrañaga, J. M. Peña, E. Menasalvas, M. S. Pérez, and V. Herves. 2004. Bayesian networks as consensed voting system in the construction of a multi-classifier for protein secondary structure prediction. *Artificial Intelligence in Medicine* 31 (2004), 117–136. 1686–1687–1688
- V. Robles, P. Larrañaga, J. M. Peña, M. S. Pérez, E. Menasalvas, and V. Herves. 2003. Learning semi naive Bayes structures by estimation of distribution algorithms. In *Proceedings of the 11th Portuguese Conference on Artificial Intelligence (EPLA-2003). Lecture Notes in Computer Science*. 244–258. 1689–1690–1691
- S. Rodrigues de Moraes and A. Aussem. 2010. A novel Markov boundary based feature subset selection algorithm. *Neurocomputing* 73, 4–6 (2010), 578–584. 1692–1693
- J. J. Rodríguez and L. I. Kuncheva. 2007. Naive Bayes ensembles with a random oracle. In *Proceedings of the 7th International Workshop on Multiple Classifier Systems (MCS-2007). Lecture Notes in Computer Science*, Vol. 4472. Springer, 450–458. 1694–1695–1696
- T. Roos, H. Wettig, P. Grünwald, P. Myllymäki, and H. Tirri. 2005. On discriminative Bayesian network classifiers and logistic regression. *Machine Learning* 59, 3 (2005), 267–296. 1697–1698
- G. A. Ruz and D. T. Pham. 2009. Building Bayesian networks classifiers through a Bayesian monitoring system. *Proceedings of the Institution of Mechanical Engineers, Part C: Journal of Mechanical Engineering Science* 223 (2009), 743–755. 1699–1700–1701
- Y. Saeyns, I. Inza, and P. Larrañaga. 2007. A review of feature selection techniques in bioinformatics. *Bioinformatics* 23, 19 (2007), 2507–2517. 1702–1703
- M. Sahami. 1996. Learning limited dependence Bayesian classifiers. In *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining (KDD-1996)*. 335–338. 1704–1705
- G. Santafé, J. A. Lozano, and P. Larrañaga. 2005. Discriminative learning of Bayesian network classifiers via the TM algorithm. In *Proceedings of the 8th European Conference in Symbolic and Quantitative Approaches to Reasoning with Uncertainty (ECSQARU-2005). Lecture Notes in Artificial Intelligence*, Vol. 3571. Springer, 148–160. 1706–1707–1708–1709
- B. Sierra and P. Larrañaga. 1998. Predicting the survival in malignant skin melanoma using Bayesian networks automatically induced by genetic algorithms. An empirical comparison between different approaches. *Artificial Intelligence in Medicine* 14 (1998), 215–230. 1710–1711–1712
- B. Sierra, E. Lazkano, E. Jauregi, and I. Irigoien. 2009. Histogram distance-based Bayesian network structure learning: A supervised classification specific approach. *Decision Support Systems* 48, 1 (2009), 180–190. 1713–1714–1715
- B. Sierra, N. Serrano, P. Larrañaga, E. J. Plasencia, I. Inza, J. J. Jiménez, P. Revuelta, and M. L. Mora. 2001. Using Bayesian networks in the construction of a bi-level multi-classifier. A case study using intensive care unit patient data. *Artificial Intelligence in Medicine* 22 (2001), 233–248. 1716–1717–1718
- M. Singh and G. Provan. 1996. Efficient learning of selective Bayesian network classifiers. In *Proceedings of the 13th International Conference on Machine Learning (ICML-1996)*. 453–461. 1719–1720
- M. Singh and M. Valtorta. 1995. Construction of Bayesian network structures from data: A brief survey and an efficient algorithm. *International Journal of Approximate Reasoning* 12, 2 (1995), 111–131. 1721–1722
- P. Spirtes, C. Glymour, and R. Scheines. 1993. *Causation, Prediction, and Search*. 1723
- J. Su, H. Zhang, C. X. Ling, and S. Matwin. 2008. Discriminative parameter learning for Bayesian networks. In *Proceedings of the 25th International Conference on Machine Learning (ICML-2008)*, Vol. 307. ACM, 1016–1023. 1724–1725–1726

- 1727 J.-N. Sulzmann, J. Fürnkranz, and E. Hüllermeier. 2007. On pairwise naive Bayes classifiers. In *Proceedings*  
1728 *of the 18th European Conference on Machine Learning (ECML-2007)*. *Lecture Notes in Computer Science*,  
1729 Vol. 4701. Springer, 371–381.
- 1730 D. M. Titterington, G. D. Murray, L. S. Spiegelhalter, A. M. Skene, J. D. F. Habbema, and G. J. Gelpke. 1981.  
1731 Comparison of discrimination techniques applied to a complex data set of head injured patients (with  
1732 discussion). *Journal of the Royal Statistical Society Series A* 144, 2 (1981), 145–175.
- 1733 I. Tsamardinos and C. F. Aliferis. 2003. Towards principled feature selection: Relevancy, filters and wrappers.  
1734 In *Proceedings of the 9th International Workshop on Artificial Intelligence and Statistics (AISTATS-*  
1735 *2003)*.
- 1736 I. Tsamardinos, C. F. Aliferis, and A. R. Statnikov. 2003a. Algorithms for large scale Markov blanket discov-  
1737 ery. In *Proceedings of the 16th International Florida Artificial Intelligence Research Society Conference*  
1738 *(FLAIRS-2003)*. AAAI Press, 376–381.
- 1739 I. Tsamardinos, C. F. Aliferis, and A. R. Statnikov. 2003b. Time and sample efficient discovery of Markov  
1740 blankets and direct causal relations. In *Proceedings of the 9th ACM SIGKDD International Conference*  
1741 *on Knowledge Discovery and Data Mining (KDD-2003)*. 673–678.
- 1742 M. van Gerven and P. J. F. Lucas. 2004. Employing maximum mutual information for Bayesian classification.  
1743 In *Proceedings of the 5th International Symposium on Biological and Medical Data Analysis (ISBMDA-*  
1744 *2004)*. *Lecture Notes in Computer Science*, Vol. 3337. Springer, 188–199.
- 1745 T. Verma and J. Pearl. 1990. Equivalence and synthesis of causal models. In *Proceedings of the 6th Conference*  
1746 *on Uncertainty in Artificial Intelligence (UAI-1990)*. Elsevier, 255–270.
- 1747 D. Vidaurre, C. Bielza, and P. Larrañaga. 2012. Forward stagewise naive Bayes. *Progress in Artificial Intel-*  
1748 *ligence* 1 (2012), 57–69.
- 1749 R. Vilalta and I. Rish. 2003. A decomposition of classes via clustering to explain and improve naive Bayes.  
1750 In *Proceedings of the 14th European Conference on Machine Learning (ECML-2003)*. *Lecture Notes in*  
1751 *Computer Science*, Vol. 2837. Springer, 444–455.
- 1752 G. I. Webb, J. Boughton, and Z. Wang. 2005. Not so naive Bayes: Aggregating one-dependence estimators.  
1753 *Machine Learning* 58 (2005), 5–24.
- 1754 G. I. Webb and M. J. Pazzani. 1998. Adjusted probability naive Bayesian induction. In *Proceedings of the*  
1755 *11th Australian Joint Conference on Artificial Intelligence (AI-1998)*. *Lecture Notes in Computer Science*,  
1756 Vol. 1502. Springer.
- 1757 T.-T. Wong. 2009. Alternative prior assumptions for improving the performance of naive Bayesian classifiers.  
1758 *Data Mining and Knowledge Discovery* 18, 2 (2009), 183–213.
- 1759 J. Xiao, C. He, and X. Jiang. 2009. Structure identification of Bayesian classifiers based on GMDH. *Knowledge-*  
1760 *Based Systems* 22 (2009), 461–470.
- 1761 J.-H. Xue and D. M. Titterington. 2010. Joint discriminative-generative modelling based on statistical tests  
1762 for classification. *Pattern Recognition Letters* 31, 9 (2010), 1048–1055.
- 1763 Y. Yang, K. B. Korb, K. M. Ting, and G. I. Webb. 2005. Ensemble selection for superparent-one-dependence  
1764 estimators. In *Proceedings of the 18th Australian Conference on Artificial Intelligence*. 102–112.
- 1765 Y. Yang, G. I. Webb, J. Cerquides, K. B. Korb, J. Boughton, and K. M. Ting. 2007. To select or to weigh:  
1766 A comparative study of linear combination schemes for superparent-one-dependence estimators. *IEEE*  
1767 *Transactions on Knowledge and Data Engineering* 19 (2007), 1652–1665.
- 1768 S. Yaramakala and D. Margaritis. 2005. Speculative Markov blanket discovery for optimal feature selection.  
1769 In *Proceedings of the 5th IEEE International Conference on Data Mining (ICDM-2005)*. IEEE Computer  
1770 Society, 809–812.
- 1771 M. Zaffalon. 2002. The naive credal classifier. *Journal of Statistical Planning and Inference* 105, 1 (2002),  
1772 5–21.
- 1773 M. Zaffalon and E. Fagioli. 2003. Tree-based credal networks for classification. *Reliable Computing* 9, 6  
1774 (2003), 487–509.
- 1775 H. Zhang and S. Sheng. 2004. Learning weighted naive Bayes with accurate ranking. In *Proceedings of the*  
1776 *5th IEEE International Conference on Data Mining (ICDM-2005)*. IEEE Computer Society, 567–570.
- 1777 H. Zhang and J. Su. 2008. Naive Bayes for optimal ranking. *Journal of Experimental & Theoretical Artificial*  
1778 *Intelligence* 20, 2 (2008), 79–93.
- 1779 N. L. Zhang, T. D. Nielsen, and F. V. Jensen. 2004. Latent variable discovery in classification models. *Artificial*  
1780 *Intelligence in Medicine* 30, 3 (2004), 283–299.
- 1781 F. Zheng and G. I. Webb. 2006. Efficient lazy elimination for averaged one-dependence estimators. In  
1782 *Proceedings of the 23rd International Conference on Machine Learning (ICML-2006)*, Vol. 148. ACM,  
1783 1113–1120.

## Discrete Bayesian Network Classifiers: A Survey

60:43

- Z. Zheng. 1998. Naïve Bayesian classifier committees. In *Proceedings of the 10th European Conference on Machine Learning (ECML-1998). Lecture Notes in Computer Science*, Vol. 1398. Springer, 196–207. 1784  
1785
- Z. Zheng and G. I. Webb. 2000. Lazy learning of Bayesian rules. *Machine Learning* 41 (2000), 53–84. 1786
- B. Ziebart, A. K. Dey, and J. A. Bagnell. 2007. Learning selectively conditioned forest structures with applications to DBNs and classification. In *Proceedings of the 23rd Conference Annual Conference on Uncertainty in Artificial Intelligence (UAI-2007)*. AUAI Press, 458–465. 1787  
1788  
1789

**Q2** Received January 2013; revised October 2013; accepted xxx

## QUERIES

**Q1:** AU: Please provide full author address for correspondence.

**Q2:** AU: Please provide accepted date.