

---

# Probabilistic Multi-Dimensional Classification

---

Vu-Linh Nguyen<sup>\*1</sup>

Yang Yang<sup>\*2</sup>

Cassio de Campos<sup>3</sup>

<sup>1</sup>Heudiasyc Laboratory, University of Technology of Compiègne, France

<sup>2</sup>Department of Computer Science, KU Leuven, Belgium

<sup>3</sup>Eindhoven University of Technology, The Netherlands

## Abstract

Multi-dimensional classification (MDC) can be employed in a range of applications where one needs to predict multiple class variables for each given instance. Many existing MDC methods suffer from at least one of inaccuracy, scalability, limited use to certain types of data, hardness of interpretation or lack of probabilistic (uncertainty) estimations. This paper is an attempt to address all these disadvantages simultaneously. We propose a formal framework for probabilistic MDC in which learning an optimal multi-dimensional classifier can be decomposed, without loss of generality, into learning a set of (smaller) single-variable multi-class probabilistic classifiers and a directed acyclic graph. Current and future developments of both probabilistic classification and graphical model learning can directly enhance our framework, which is flexible and provably optimal. A collection of experiments is conducted to highlight the usefulness of this MDC framework.

## 1 INTRODUCTION

In (multi-class) classification, a predictive system makes use of a training data set (consisting of input-output pairs which specify individuals of a population) and a hypothesis space (consisting of the possible classifiers), and seeks for a classifier that optimizes its chance of making accurate predictions with respect to some given evaluation criterion (such as a loss function or an accuracy measure). Numerous studies on classification have been devoted to learning probabilistic classifiers which predict, for each observation of the input space, a univariate probability distribution over the output space. The intention of probabilistic classification is not only to provide the end user with all necessary information about

the optimal predictions of different loss functions [Elkan, 2001, Mortier et al., 2021], but also information about the uncertainty associated with the possible predictions.

To overcome the assumption that the output space must be fully characterized by a single class variable, MDC has been proposed in which the output space is characterized by multiple class variables which can be correlated. MDC appears in important applications. An example of MDC is predicting subtypes/stages of diseases associated with each patient given his/her medical image and/or demographic information. Few other examples of MDC tasks are classification of biomedical text [Shatkay et al., 2008], vehicle classification [Jia and Zhang, 2021] and beyond [Gil-Begue et al., 2021, Jia and Zhang, 2022].

Existing multi-dimensional classifiers are non-probabilistic [Jia and Zhang, 2021], relatively inaccurate [Jia and Zhang, 2021][Section II & III], or unscalable [Gil-Begue et al., 2021]. To the best of our knowledge, no existing method specific for MDC is capable of directly handling mixed data, i.e., continuous and discrete features coexisting (without preprocessing or other external manipulations). Problem transformation methods [Jia and Zhang, 2021] which transform the original MDC problem into either a huge multi-class classification (MCC) problem, for example using the class powerset (CP) classifier, or a set of independent MCC problems, for example the Binary relevance (BR) classifier, can be combined with deep multimodal learning [Kline et al., 2022, Xu et al., 2021] to handle mixed data and other complex types of input. They suffer from the aforementioned issues and are arguably hard to interpret. The set of marginal probability distributions provided by BR can be associated to (infinitely) many joint distributions over the class variables<sup>1</sup> and does not inform much about the (true) joint distribution, while the joint distribution provided by CP contains an exponential number of masses and is not

---

<sup>\*</sup>These authors contributed equally to this work.

---

<sup>1</sup>Thus, BR can be seen as a credal classifier and would be useful when targeting reliable set-valued predictions [Augustin et al., 2014, Jansen et al., 2022, Troffaes, 2007].

easily interpretable for end users.

We present a framework to learn probabilistic multi-dimensional classifiers addressing those issues. This formal framework allows us to learn an optimal multi-dimensional classifier, without loss of generality/optimalty, by decomposing the task into learning a set of probabilistic MCC models plus a directed acyclic graph (DAG). Notably, the framework inherits the interpretability of Bayesian networks (BNs) [Atienza et al., 2022, Kitson et al., 2023, Koller and Friedman, 2009], which is a compact representation of quantitative and qualitative probabilistic relationships among class variables, and the scalability and flexibility of deep (multimodal) learning [Kline et al., 2022, LeCun et al., 2015, Xu et al., 2021], i.e., handling complex types of data. Moreover, the probabilistic nature allows the framework, among other characteristics, to optimize different loss functions by only learning a single probabilistic model. We prove that the probabilistic model learned by this framework is universal and the learning procedure is globally optimal whenever MCC is universal and can be solved optimally too. We formalize the probabilistic MDC problem in Section 2, present formal results on the optimality of the framework in Section 3, followed by a practical algorithm and properties of the learning framework in Sections 3.1 to 3.4. Section 4 discusses the inference task, and Section 5 further motivates the framework by presenting a collection of experiments indicating the advantages of the framework against existing MDC approaches. Section 6 concludes this paper. All formal results in this paper (propositions) are stated without proofs, which are deferred to Appendix B. Some technical details and experiments were also given in [Yang, 2022].

## 2 PROBABILISTIC MDC

Let  $\mathbf{X} = \{X^1, \dots, X^Q\}$  be a finite set of features, let  $\mathcal{X} := \mathcal{X}^1 \times \dots \times \mathcal{X}^Q$  denote an input space, and let  $\mathbf{Y} = \{Y^1, \dots, Y^K\}$  be a finite set of class variables. Let  $\mathcal{Y}^k = \{y^{k,1}, \dots, y^{k,M_k}\}$  be the set of  $M_k$  possible outcomes for the  $k^{\text{th}}$  class variable  $Y^k$ ,  $k \in [K] := \{1, \dots, K\}$ . We define  $\mathbf{Z} := \mathbf{Y} \cup \mathbf{X}$ . We denote by  $\mathbf{X}_d$  and  $\mathbf{X}_c$  the discrete feature set and continuous feature set, respectively. We also define  $\mathbf{Z}_d := \mathbf{Y} \cup \mathbf{X}_d$ . For each instance  $\mathbf{x} \in \mathcal{X}$ , we say it is associated with a (vector)class  $\mathbf{y} \in \mathcal{Y} = \mathcal{Y}^1 \times \dots \times \mathcal{Y}^K$ .

We assume observations to be realizations of independently and identically distributed (i.i.d.) random variables generated according to a probability distribution on  $\mathcal{X} \times \mathcal{Y}$ , i.e., an observation  $\mathbf{y} = (y^1, \dots, y^K)$  is the realization of a corresponding random vector  $\mathbf{Y} = (Y^1, \dots, Y^K)$ . Let  $p(\mathbf{X}, \mathbf{Y})$  be a (mixed) joint density function. We denote by  $p(\mathbf{Y} | \mathbf{x})$  the conditional joint distribution of  $\mathbf{Y}$  given  $\mathbf{X} = \mathbf{x}$ , whose probability mass function is given by

$$p(\mathbf{y} | \mathbf{x}) := \frac{p(\mathbf{x}, \mathbf{y})}{\sum_{\mathbf{y}' \in \mathcal{Y}} p(\mathbf{x}, \mathbf{y}')}, \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}. \quad (1)$$

We assume the denominator to be non-zero whenever needed. We denote by  $p(Y^k | \mathbf{x})$ ,  $k \in [K]$ , the marginal distribution of  $Y^k$ , whose probability mass function is

$$p(y^k | \mathbf{x}) := \sum_{\mathbf{y} \in \mathcal{Y}: Y^k = y^k} p(\mathbf{y} | \mathbf{x}), \forall y^k \in \mathcal{Y}^k. \quad (2)$$

Given training data in the form of a finite set of observations  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N \subset \mathcal{X} \times \mathcal{Y}$  drawn independently from a distribution, MDC aims to learn a predictive classifier model  $\mathbf{h} : \mathcal{X} \rightarrow \mathcal{Y}$  assigning  $\mathbf{y} \in \mathcal{Y}$  to each  $\mathbf{x} \in \mathcal{X}$ . The output of  $\mathbf{h}$  is a vector

$$\hat{\mathbf{y}} := \mathbf{h}(\mathbf{x}) = (h^1(\mathbf{x}), \dots, h^K(\mathbf{x})) \in \mathcal{Y}. \quad (3)$$

In a probabilistic setting, a classification task can be viewed as a two-stage problem, in which a mapping  $\mathbf{h} : \mathcal{X} \rightarrow \mathcal{Y}$  is not learned directly, but in a more indirect way. Roughly speaking, one can split a probabilistic classification into two tasks: learning a function  $p : \mathcal{X} \rightarrow p(\mathcal{Y} | \mathcal{X})$  (with abuse of notation) and constructing an efficient inference operator  $o : p(\mathcal{Y} | \mathcal{X}) \rightarrow \mathcal{Y}$  (we will deal with  $o$  in Section 4).

Motivated by the observations that discriminative models can perform better than generative models in many classification tasks [Bouchard and Triggs, 2004, Carvalho et al., 2011, Ng and Jordan, 2001, Ulusoy and Bishop, 2006], and by the fact that in M-open cases [Bernardo and Smith, 2000], maximizing the (log) likelihood function may not converge to a best possible distribution as maximizing the conditional (log) likelihood function does [Roos et al., 2005], we learn a multi-dimensional classifier encoding  $p$  which maximizes the conditional log likelihood (CLL) function  $C(p | \mathcal{D})$ :

$$C(p | \mathcal{D}) := \log \prod_{n=1}^N p(\mathbf{y}_n | \mathbf{x}_n). \quad (4)$$

This idea has been mentioned before [Benjumbeda et al., 2018], but, to the best of our knowledge, it has been left open until now. Let  $\mathcal{P}^0$  be a hypothesis space for  $p$ . The learning problem can be defined as finding

$$p^* \in \arg \max_{p \in \mathcal{P}^0} C(p | \mathcal{D}). \quad (5)$$

To avoid overfitting, the CLL function is often augmented by a regularization term. We will discuss it later.

## 3 A LEARNING FRAMEWORK

The optimization problem (5) is very generic and its complexity highly depends on the given hypothesis space  $\mathcal{P}^0$ . We present reformulations of this problem that are more suitable to be optimized based on some assumptions about the hypothesis space. We proceed under the assumption that the features  $X^q$ ,  $q \in [Q]$ , are always made available. This

means we neither admit missing values at the training time [Nguyen et al., 2021] nor admit missing features at the prediction time [Saar-Tsechansky and Provost, 2007]. This is not a limitation of the approach and missing data can be tackled using a variation of structure EM [Friedman, 1998, Rancoita et al., 2016], but the discussion goes beyond the scope of this paper (see Appendix D for a quick discussion).

Throughout, we assume the chain rule of probability [Koller and Friedman, 2009][Section 2.1.3.4] holds <sup>2</sup>. Using the concept of conditional independence, we can assume without loss of generality that any  $p(\mathbf{X}, \mathbf{Y})$  can be fully encoded by a DAG  $G$  and a parameter set  $\theta$  inducing the factorization

$$p_{\theta}^G(\mathbf{x}, \mathbf{y}) = \prod_{X \in \mathbf{X}_c} p_{\theta}(x | \pi_x) \prod_{Z \in \mathbf{Z}_d} p_{\theta}(z | \pi_z), \quad (6)$$

where  $\pi_x$  and  $\pi_z$  are (with abuse of notation) called configurations (compatible with  $(\mathbf{x}, \mathbf{y})$ ) of the parent sets  $\Delta_G^X$  and  $\Delta_G^Z$  (for easiness, we assume that discrete parts of configurations are dictionaries with pairs (variable, value), and continuous parts are given via the appropriate functionals). The complexity of this factorization depends on  $G$ .

Therefore, the hypothesis space of any probabilistic MDC can be defined as  $\mathcal{P} := \mathcal{G} \times \Theta$ , where  $\mathcal{G}$  and  $\Theta$  are respectively the set of possible DAGs and the set of possible parameter sets, and the problem (5) becomes

$$p_{\theta^*}^{G^*} : (G^*, \theta^*) \in \arg \max_{(G, \theta) \in \mathcal{P}} C(p_{\theta}^G | \mathcal{D}). \quad (7)$$

A learning procedure is optimal if it can find an optimal pair  $(G^*, \theta^*)$ . Parameter learning is optimally solved if we can find  $\theta^*$  in (7) for a given  $G \in \mathcal{G}$ . In the following, we show that the factorization in (6) can lead to a great simplification of the learning problem (7).

**Proposition 3.1.** *Assume the parameter learning problem is optimally solved. We have*

$$\max_{\mathbf{p} \in \mathcal{P}^0} C(\mathbf{p} | \mathcal{D}) = \max_{(G, \theta) \in \mathcal{P}} C(p_{\theta}^G | \mathcal{D}) = \max_{(G, \theta) \in \mathcal{P}^1} C(p_{\theta}^G | \mathcal{D}), \quad (8)$$

where  $\mathcal{P}^1 := \mathcal{G}^1 \times \Theta$  and  $\mathcal{G}^1 \subsetneq \mathcal{G}$  is the set of DAGs which contain no edge of the form<sup>3</sup>  $Y \rightarrow X$ .

We assume in this document that parameter learning can be optimally solved. In general, this is a strong assumption. However, we often deal with factorizations of  $\mathbf{p}$  where each factor involves a small number of variables. In these cases, we hope one can learn the parameters well (certainly

much better than in a global model). This is a condition we expect from local models in the factorization in order to prove the optimality of the framework. Note that the cardinality  $|\mathcal{G}^1| = R(K)2^{KQ}R(Q)$  can be much smaller than  $|\mathcal{G}| = R(K+Q)$ , where  $R(\cdot)$  is Robinson’s formula [Bielza et al., 2011]. Thus, looking for the best  $(G, \theta)$  over  $\mathcal{P}^1$  can be much more practical than doing so over  $\mathcal{P}$ . The next proposition shows that finding an optimal pair  $(G, \theta) \in \mathcal{P}^1$  is equivalent to finding an optimal pair whose  $G$  contains no edge between features.

**Proposition 3.2.** *For any  $G \in \mathcal{G}^1$ , the joint conditional distribution (1) can be factorized (according to  $G$ ):*

$$p_{\theta}^G(\mathbf{y} | \mathbf{x}) = \prod_{Y \in \mathbf{Y}} p_{\theta}(y | \pi_y), \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}, \quad (9)$$

where  $\pi_y$  is the configuration for the parents of  $Y$  (according to  $G$ ) that is compatible with  $(\mathbf{x}, \mathbf{y})$ . Moreover, the following relation holds:

$$\max_{(G, \theta) \in \mathcal{P}^1} C(p_{\theta}^G | \mathcal{D}) = \max_{(G, \theta) \in \mathcal{P}^2} C(p_{\theta}^G | \mathcal{D}), \quad (10)$$

where  $\mathcal{P}^2 := \mathcal{G}^2 \times \Theta$  and  $\mathcal{G}^2 \subsetneq \mathcal{G}^1$  consists of  $R(K)2^{KQ}$  DAGs with no edges between any two elements of  $\mathbf{X}$ .

Thus, we formulate the new optimization problem:

$$p_{\theta^*}^{G^*} : (G^*, \theta^*) \in \arg \max_{(G, \theta) \in \mathcal{P}^2} \log \prod_{n=1}^N \prod_{Y \in \mathbf{Y}} p_{\theta}(y_n | \pi_{y_n}). \quad (11)$$

It is clear that solving (11) may lead to sub-optimal solutions, compared to solving (7) if the assumption that the parameter learning problem is optimally solved does not hold, and in that case the relation  $\mathcal{G}^2 \subsetneq \mathcal{G}$  implies that the best CLL score attained over  $\mathcal{G}^2$  is at best the one attained over  $\mathcal{G}$ . However, there are strong motivations for why one should solve (11) in practice, instead of (7).

First, the optimality of (11) can be reachable under milder conditions, while the optimality of (7) is often unreachable. In fact, solving (7) is often impractical because optimizing the CLL function can be impractical even if  $G \in \mathcal{G}$  is given [Friedman et al., 1997]. However, one can be much more optimistic about solving (11). As will be shown in Section 3.1, solving (11) is possible as long as one can learn a set of (independent) probabilistic classifiers, plus learning an optimal DAG over the class variables. So one can use all current/future developments of both probabilistic classification and graphical model learning towards solving (11).

Second, as will be shown in Section 3.1,  $\forall G \in \mathcal{G}^2$  and  $\forall \mathbf{x} \in \mathcal{X}$ ,  $p_{\theta}^G(\mathbf{Y} | \mathbf{x})$  can be factorized as a product of conditional probability distributions whose conditional part is always specified by a multivariate continuous variable. This

<sup>2</sup>An intensive study on the conditions under which the chain rule of probability is (in)valid is beyond the scope of this paper.

<sup>3</sup>To the best of our knowledge, we are the first who extend/adapt the setting suggested in [Lerner et al., 2001] to do probabilistic multi-dimensional classification when targeting the (regularized) joint conditional likelihood function.

provides us with a rich representational capacity as discussed in Section 3.2. In particular, any probabilistic classifier can be directly employed to model conditional probability distributions without requiring any data preprocessing transformation, leading to a rich framework for the employment of sophisticated techniques. The representational capacity would be much weaker if one had to parameterize  $G \in \mathcal{G} \setminus \mathcal{G}^1 \supsetneq \mathcal{G} \setminus \mathcal{G}^2$  because it would be needed to find some parametric model to encode all conditional density functions  $p_\theta^G(z | \pi_z)$  whose conditional part would be specified by a mixture of discrete and continuous variables. This would be a challenging problem by itself, especially if one does not want to use any data preprocessing transformation either before or during the training phase.

Our final simplification of the optimization problem while keeping optimality is to realize that we can seek for an optimal  $G$  where all continuous variables are parents of every class variable, that is,  $\mathbf{X}_c \subset \Delta_G^Y$ ,  $\forall Y \in \mathbf{Y}$ . Besides being non-restrictive (we are forcing arcs to stay put, hence we can always fit any “simpler” distribution which would have dropped some connections by the appropriate parameter learning), this condition has also a positive consequence, as it allows us to use methods which are not able to handle mixed setups of continuous and discrete variables.

Therefore, we introduce an updated version of (11) in which we only force the global learning algorithm to explicitly handle the discrete features, while assuming all continuous ones are passed on to the learning of local models. More formally, let  $\mathcal{G}^3 \subseteq \mathcal{G}^2$  be the set of  $R(K)2^{K|\mathbf{X}_d|}$  DAGs such that,  $\forall G \in \mathcal{G}^3$  and  $\forall Y \in \mathbf{Y}$ , we have  $\mathbf{X}_c \subset \Delta_G^Y$ . We formulate the optimization problem as:

$$p_{\theta^*}^{G^*} : (G^*, \theta^*) \in \arg \max_{(G, \theta) \in \mathcal{P}^3} \log \prod_{n=1}^N \prod_{Y \in \mathbf{Y}} p_\theta(y_n | \pi_{y_n}), \quad (12)$$

where  $\mathcal{P}^3 = \mathcal{G}^3 \times \Theta$ .

**Proposition 3.3.** *Assume the parameter learning problem is optimally solved. The following relation holds*

$$\max_{(G, \theta) \in \mathcal{P}^2} C(p_\theta^G | \mathcal{D}) = \max_{(G, \theta) \in \mathcal{P}^3} C(p_\theta^G | \mathcal{D}). \quad (13)$$

The conclusion here is that we can have a globally optimal probabilistic MDC whose optimization is done via (12), potentially saving significant time and data requirements for training. One needs “only” to learn the local conditional models (factors) of the expression, so long as we have an efficient solver to find the DAG  $G$  inducing a good factorization. Moreover, we hope for a valid (in terms of being an I-map for the true distribution [Bouckaert, 1994, Koller and Friedman, 2009]) yet simple  $G$ . Hence, in the next section, we show that solving (12) can be optimally decomposed into learning a set of probabilistic classifiers and learning an optimal DAG.

### 3.1 ALGORITHMIC SOLUTION

In order to solve (12), we first need to model the local conditional probability distributions:

$$p_\theta(Y | \Delta_G^Y), \forall G \in \mathcal{G}^3, \forall Y \in \mathbf{Y}. \quad (14)$$

Given  $G$ , for any  $Y \in \mathbf{Y}$ , let  $\Delta_d^Y = \Delta_G^Y \setminus \mathbf{X}_c$  be the set of all discrete variables in  $\Delta_G^Y$ . Let  $\Pi_d^Y$  be the set of all configurations of  $\Delta_d^Y$ . Hence, each local distribution (14) is represented by  $|\Pi_d^Y|$  distributions

$$p_\theta(Y | \pi, \mathbf{X}_c), \forall \pi \in \Pi_d^Y. \quad (15)$$

Thus, the optimization problem (12) becomes

$$(G^*, \theta^*) \in \arg \max_{(G, \theta) \in \mathcal{P}^3} \sum_{Y \in \mathbf{Y}} \sum_{\pi \in \Pi_d^Y} \log \prod_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_\pi} p_\theta(y | \pi, \mathbf{x}^c),$$

with  $\mathcal{D}_\pi := \{(\mathbf{x}, \mathbf{y}) \in \mathcal{D} | \pi_y^d = \pi\}$ . A key point is the separation of discrete conditionals  $\pi$  and continuous conditionals  $\mathbf{x}^c$ . Such separations were used in learning BNs optimizing the likelihood function [Atienza et al., 2022]. Moreover, we have  $\max_{(G, \theta) \in \mathcal{P}^3} C(p_\theta^G | \mathcal{D})$

$$= \max_{G \in \mathcal{G}^3} \sum_{Y \in \mathbf{Y}} \sum_{\pi \in \Pi_d^Y} \max_{\theta \in \Theta} C(p_\theta | Y, \pi, \mathcal{D}), \quad (16)$$

where  $C(p_\theta | Y, \pi, \mathcal{D}) = \log \prod_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_\pi} p_\theta(y | \pi, \mathbf{x}^c)$ . This means that we can reformulate the optimization problem (12) as a two-phase optimization problem: (P1) for any tuple  $(Y, \pi) \in \mathbf{Y} \times \Pi_d^Y$  (considering the possible  $\Delta_d^Y$ ), learn the optimal parameter set  $\theta^*$  of each distribution (15) which optimizes the local CLL function, i.e.,

$$\theta_{Y, \pi}^* \in \arg \max_{\theta \in \Theta} C(p_\theta | Y, \pi, \mathcal{D}), \quad (17)$$

and then (P2) learn the best DAG  $G^* \in \mathcal{G}^3$  which maximizes the CLL function:  $G^* = \arg \max_G C(p_{\theta^*}^G | \mathcal{D})$  and

$$C(p_{\theta^*}^G | \mathcal{D}) = \sum_{Y \in \mathbf{Y}} \sum_{\pi \in \Pi_d^Y} C(p_{\theta_{Y, \pi}^*} | Y, \pi, \mathcal{D}). \quad (18)$$

Problem (P1) can be solved for each tuple  $(Y, \pi) \in \mathbf{Y} \times \Pi_d^Y$ , for each possible  $\Delta_d^Y \in \mathcal{F}^Y$  independently (where  $\mathcal{F}^Y$  is a set of candidate parent sets for  $Y$ ). (P2) can be cast as the structure learning for BNs, so we can leverage the research on that topic [Kitson et al., 2023]. The elephant in the room here is the size of  $\mathcal{F}^Y$  (for each  $Y$ ), which will be discussed in Section 3.4.

In this paper, we solve (18) using GOBNILP [Bartlett and Cussens, 2017, Cussens et al., 2017] which is a state-of-the-art anytime globally optimal algorithm and can be easily adapted to handle regularized variants of CLL function as presented in Section 3.4. Intuitively, GOBNILP, which was designed for generative learning of Bayesian networks, can

be instead used to reformulate the problem (P2) as learning a collection of parent sets  $\{\Delta_d^Y : Y \in \mathbf{Y}\}$  which optimizes the CLL function (18) and together satisfy the DAG properties. It uses the local scores:  $\forall Y \in \mathbf{Y}, \forall \Delta_d^Y \in \mathcal{F}^Y$ :

$$C(Y, \Delta_d^Y) = \sum_{\pi \in \Pi_d^Y} C(p_{\theta_{Y,\pi}}^* | Y, \pi, \mathcal{D}), \quad (19)$$

where we simplified the notation by removing  $\theta$  and  $\mathcal{D}$ , since parameters have been already learned via (17) and data are fixed. Problem (P2) can be expressed as an Integer Programming (IP) problem:

$$\begin{aligned} & \text{Maximize } \sum_{Y \in \mathbf{Y}} \sum_{\Delta_d^Y \in \mathcal{F}^Y} \gamma(\Delta_d^Y) \cdot C(Y, \Delta_d^Y), \quad (20) \\ & \text{Subject to } \sum_{\Delta_d^Y \in \mathcal{F}^Y} \gamma(\Delta_d^Y) = 1, \forall Y \in \mathbf{Y}, \\ & \sum_{Y \in \mathbf{Y}'} \sum_{\substack{\Delta_d^Y \in \mathcal{F}^Y \\ \Delta_d^Y \cap \mathbf{Y}' = \emptyset}} \gamma(\Delta_d^Y) > 1, \forall \mathbf{Y}' \subseteq \mathbf{Y}, |\mathbf{Y}'| > 1, \\ & \gamma(\Delta_d^Y) \in \{0, 1\}, \forall Y \in \mathbf{Y}, \forall \Delta_d^Y \in \mathcal{F}^Y. \end{aligned}$$

The implementation is given in Algorithm 1, which returns a  $(G^*, \theta^*) \in \mathcal{P}^3$  of (12). We call this type of model defined by  $(G^*, \theta^*)$  a generalized Bayesian Network classifier (GBNC). Note that the loops starting in lines 2 and 3 can be easily parallelized since the local distributions (15) can be learned independently.

---

**Algorithm 1** Learning a GBNC of (12)

---

- 1: **Input:** Data  $\mathcal{D}$ , Probabilistic hypothesis spaces.
  - 2: **for**  $Y \in \mathbf{Y}$  **do**
  - 3:   **for**  $\Delta_d^Y \in \mathcal{F}^Y$  **do**
  - 4:     **for**  $\pi \in \Pi_d^Y$  **do**
  - 5:       Solve (17) and store it in a proper data structure
  - 6:     **end for**
  - 7:     Compute  $C(Y, \Delta_d^Y)$  by (19) using stored values
  - 8:   **end for**
  - 9: **end for**
  - 10: Find a best collection  $\{\Delta_d^Y : Y \in \mathbf{Y}\}$  which optimizes (20) using GOBNILP
  - 11: **Output:** A GBNC  $(G^*, \theta^*) \in \mathcal{P}^3$  of (12)
- 

The optimality of the proposed framework can be derived as a consequence of Proposition 3.1–3.3.

**Corollary 3.4.** *Assume the chain rule of probability holds. Assume the parameter learning problem is optimally solved. The procedure to learn a classifier  $(G^*, \theta^*)$  by Algorithm 1 is universal (for distributions in  $\mathcal{P}^0$ ).*

### 3.2 REPRESENTATIONAL CAPACITY

To represent the joint conditional probability distribution  $p(\mathbf{Y} | \mathbf{X})$ , we need a set of probabilistic classifiers  $p'$ :

$\mathcal{X}_c \rightarrow \mathcal{Y}^k$  to estimate the local conditional probability distributions (15). Local models  $p'$  are trained with what we call base learners. Note that discrete variables are not included in the input for  $p'$  (they are dealt with through the DAG optimization), which also facilitates learning and representational capacity.

First, it allows us to represent the distribution  $p(\mathbf{Y} | \mathbf{X})$  where  $\mathbf{X}$  can contain both continuous features and discrete features without requiring any preprocessing transformation either before or during the training phase. We never face the problem of representing qualitative data for use as input as deep learning does [Hancock and Khoshgoftaar, 2020]. Besides, representing qualitative data for use as input is arguably the most critical obstacle for generalizing Classifier Chains (CCs) [Dembczyński et al., 2010, Read et al., 2021], which is a state-of-the-art multi-label classification framework, to cope with MDC. Moreover, we naturally overcome a bottleneck in the development of Multi-dimensional Bayesian network classifiers (MDBNCs) [Gil-Begue et al., 2021] that is a shortage of classifiers for the cases of continuous features, and mixed features.

Second, the probabilistic classifier inducing  $p'$  can be freely chosen according to our needs. It can be as intuitive as  $k$ -NN classifiers [Cover and Hart, 1967] and can be as counter-intuitive as ensembles of deep networks [Ganaie et al., 2022]. This allows us to employ sophisticated probabilistic classifiers to encode complex probabilistic relationships within  $p'_{Y,\pi} := p_\theta(Y | \pi, \mathbf{X}_c)$ ,  $\forall \pi \in \Pi_d^Y$ . For example, when each image is encoded using an  $\mathbf{x}$ , a convolutional network [LeCun et al., 2015] can be employed to encode  $p'_{Y,\pi}$ . If one seeks for more accurate GBNCs, there should be no restriction on the use of ensemble learning methods, except the availability of computational resources. This flexibility of the framework is remarkably different from existing probabilistic MDC approaches [Gil-Begue et al., 2021, Jia and Zhang, 2022]. Roughly speaking, so long as you train good local models  $p'_{Y,\pi} : \mathcal{X}_c \rightarrow \mathcal{Y}^k$  (for which you can use all toolsets available in the literature for “standard” single-class-variable classification), the framework in this paper does the rest to combine them optimally into an MDC solution.

### 3.3 INTERPRETABILITY

GBNCs are interpretable at both the population and individual levels. At the population level, the structure  $G$  provides a compact representation of the qualitative probabilistic relationships among feature and class variables. This graph representation is easy to interpret to end users when compared to an exponential number of masses provided by CP [Jia and Zhang, 2021] and the (infinitely) many joint conditional distributions associated with the set of marginal probability distributions provided by BR [Jia and Zhang, 2021]. At the individual level, the structure  $G$  and its parameters specified by  $\theta$  under the particular value of an individual  $\mathbf{x}$  form a

compact representation of the qualitative and quantitative probabilistic relationships within  $p(\mathcal{Y} | \mathbf{x})$ , which can be seen as a BN over the class variables.

As an example, we provide in Figure 1 a DAG over class variables learned from the PASCAL VOC 2007 data set whose description is given in Section 5.

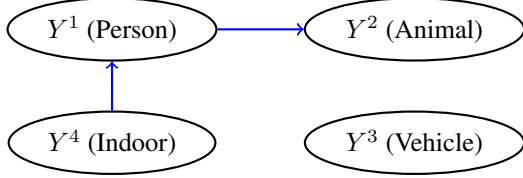


Figure 1: A DAG over class variables learned from the PASCAL VOC 2007 data set.

### 3.4 REGULARIZATION

While Algorithm 1 helps to find an optimal GBNC which maximizes the CLL function, the next proposition suggests that this best GBNC may not always be the one we want, especially with regard to overfitting.

**Proposition 3.5.** *Assume local models have parameters optimally learned. Then  $\forall Y \in \mathbf{Y}$  and  $\forall \Delta, \Delta' \in \mathcal{F}^Y$  such that  $\Delta_d \subset \Delta'_d$ , we have*

$$C(Y, \Delta_d) \leq C(Y, \Delta'_d). \quad (21)$$

Therefore, at least one optimal solution of the Algorithm 1 is a fully connected DAG  $G$ .

Over-complex DAGs can happen frequently, especially when the local classifiers are learned without enforcing regularization terms. To seek for a better generalization, we propose a regularized variant of the CLL function:

$$S(\mathbf{p}_\theta^G | \mathcal{D}) = C(\mathbf{p}_\theta^G | \mathcal{D}) - \sum_{Y \in \mathbf{Y}} \text{pen}(|\Delta_d^Y|, |\mathcal{D}|), \quad (22)$$

where  $\text{pen}(\Delta_d^Y, |\mathcal{D}|)$  can be the penalty term of any decomposable scoring function [Liu et al., 2012]. Even a mild penalty can already help to reduce model complexity, but we leave this study to future work.

Algorithm 1 can be revised to learn GBNCs of regularized variants (22) as presented in Appendix C.1 and C.2. Moreover, as shown in Appendix C.2, pruning rules [de Campos et al., 2018] can be employed to find GBNCs which optimize regularized variants (22) without losing any optimality. This helps to greatly reduce the learning time because for each  $Y \in \mathbf{Y}$ , large candidate parent sets  $\Delta_d^Y \in \mathcal{F}^Y$  are often pruned due to high penalties [de Campos et al., 2018]. Finally, for a very large number of class variables, it is not unreasonable to expect the treewidth of the true distribution to be limited, so that one can bound the size of  $\mathcal{F}^Y$

and use the scalability of (approximate) bounded-treewidth learning [Scanagatta et al., 2016].

## 4 INFERENCE

The learned function  $\mathbf{p}$  (defined via  $G$  and  $\theta$ ) provides, given an  $\mathbf{x} \in \mathcal{X}$ , a conditional joint probability distribution  $p(\mathcal{Y} | \mathbf{x})$  which is used to find the Bayes-optimal prediction (BOP)  $\hat{\mathbf{y}}$  w.r.t a target loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ :

$$\hat{\mathbf{y}} := o(\mathbf{p}(\mathcal{Y} | \mathbf{x})) \in \argmin_{\bar{\mathbf{y}} \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}, \bar{\mathbf{y}}) p(\mathbf{y} | \mathbf{x}). \quad (23)$$

Yet, different loss functions may call for different BOPs (23) [Dembczyński et al., 2012, Gil-Begue et al., 2021, Nguyen and Hüllermeier, 2021, Waegeman et al., 2014]. Knowledge about the probability distribution  $p(\mathcal{Y} | \mathbf{x})$  is necessary for finding BOP (23) of any loss function. The complexity of finding BOP can greatly depend on the nature of the chosen loss function. This problem has been studied rarely in the MDC setting. An exception is [Bielza et al., 2011, Gil-Begue et al., 2021]. Notably, in these works, finding BOP (23) of some commonly used loss functions is shown to be equivalent to computing the most probable explanations (MPEs) of class variables when the classifier is an MDBNC. This is an interesting finding because it implies that the complexity of finding BOP (23) depends on the nature of both the chosen loss function and the classifier. While this finding allows us to directly employ any current/future developments on exact/approximate MPE inference [Gil-Begue et al., 2021] to find BOP (23) of some loss functions, one cannot get rid of the computational burden introduced by large numbers of features when working with MDBNCs.

In our framework, we can also show that finding BOP (23) of some loss functions is computing the MPEs of class variables. In the following, we describe the problem of finding BOP (23) of two commonly used loss functions<sup>4</sup> which are the *Hamming loss* (24) and the *subset 0/1 loss* (25):

$$\ell_H(\mathbf{y}, \hat{\mathbf{y}}) := \frac{1}{K} \sum_{k=1}^K \mathbb{I}[y^k \neq \hat{y}^k], \quad (24)$$

$$\ell_S(\mathbf{y}, \hat{\mathbf{y}}) := \mathbb{I}[\mathbf{y} \neq \hat{\mathbf{y}}]. \quad (25)$$

The indicator  $\mathbb{I}[A]$  equals 1 if the  $A$  is true and 0 otherwise. Thus, both losses generalize the standard 0/1 loss in binary classification. As noted in [Bielza et al., 2011], finding a BOP of  $\ell_H$  and  $\ell_S$  are respectively equivalent to finding  $K$  marginals (26) and equivalent to finding one MPE (27):

$$\hat{y}^k \in \argmax_{\bar{y}^k \in \mathcal{Y}^k} p(\bar{y}^k | \mathbf{x}), \forall k \in [K], \quad (26)$$

$$\hat{\mathbf{y}} \in \argmax_{\bar{\mathbf{y}} \in \mathcal{Y}} p(\bar{\mathbf{y}} | \mathbf{x}). \quad (27)$$

<sup>4</sup>We defer intensive studies on finding BOPs of other loss functions [Gil-Begue et al., 2021][Section 4] to future work.

Hence, the model does not require retraining to allow for different BOP. Exact MPE and marginal inferences are NP-hard problems [de Campos, 2020, Roth, 1996, Shimony, 1994]. However, in our framework, the complexity of MPE and marginal inferences only depend on the number of class variables. Thus, we do not encounter the computational burden introduced by large numbers of features, making the framework usable in practice in spite of that. Moreover, one can control the graph complexity among class variables by employing bounded-treewidth learning [Nie et al., 2017].

## 5 EXPERIMENTS

This section presents a set of experiments to assess the usefulness of our proposal.

### 5.1 EXPERIMENTAL SETTING

We compare two instantiations of GBNCs (GBNC-S which optimizes (22) and produces BOP (27) of  $\ell_S$ , and GBNC-H which optimizes (22) and produces BOP (26) of  $\ell_H$ ) with three probabilistic competitors found in the literature on 20 tabular data sets [Jia and Zhang, 2021] and one image data set [Everingham et al., 2010]. The number of instances varies from 154 to 28779, the number of features varies from 10 to 1536, and the number of class variables varies from 2 to 16. It also contains 3 data sets with mixed discrete and continuous features.

We utilize an MDC version of the PASCAL VOC 2007 data set [Everingham et al., 2010]. We encode the objects found in that data set using 4 class variables: Person (Yes and No), Animal (No animal, Bird, Cat, Cow, Dog, Horse and Sheep), Vehicle (No vehicle, Aeroplane, Bicycle, Boat, Bus, Car, Motorbike, Train) and Indoor (No indoor object, Bottle, Chair, Dining table, Potted plant, Sofa, TV/Monitor).

For tabular data sets, we compare GBNCs with BR and PC [Jia and Zhang, 2021][Section II], and CC [Jia and Zhang, 2021][Section III]. Because of the limitations of competitors to deal with mixed data, we follow the suggestion of [Jia and Zhang, 2021] and convert discrete features/variables into continuous variables using one-hot encoding whenever they appear as parts of input of local classifiers of BR, PC and CC. Because we are not aware of any refinement of CC which can handle image data sets, we eliminate it from our comparison on the PASCAL VOC 2007. For tabular data sets, we use logistic regression (LR) [Menard, 2002] and Naive Bayes (NB) classifiers [Domingos and Pazzani, 1996] to estimate the local distributions (15) (one can use more complex models, but as we see in the remainder, these choices already yield state-of-the-art results, so we decided that further tuning would go beyond our scope). For the image data set, distributions (15) are estimated using ResNet-18 [He et al., 2016] with the weights pre-trained on ImageNet

[Deng et al., 2009], which are calibrated using temperature scaling [Guo et al., 2017]. Following the suggestion of [Zhang et al., 2017], we also employ *mixup* to improve the generalization of ResNet-18.

In our experiments,  $\text{pen}(|\Delta_d^Y|, |\mathcal{D}|)$  is the penalty term of the Bayesian Information Criterion (BIC) [Schwarz, 1978]. The experimental setting is detailed in Appendix E.1. The source code has been made public at <https://github.com/yangyang-pro/probabilistic-mdc>.

### 5.2 RESULTS

Overall, the results suggest the superiority of our framework against existing probabilistic MDC frameworks (See Table 1–3, and Figure 2). On the image data set, GBNCs indeed provide the most promising  $\ell_H$  and  $\ell_S$  (See Table 1).

Table 1: Results (mean  $\pm$  std.) on the image data set.

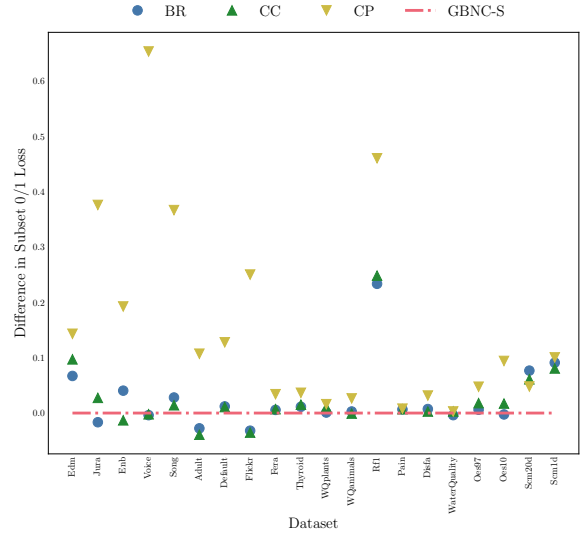
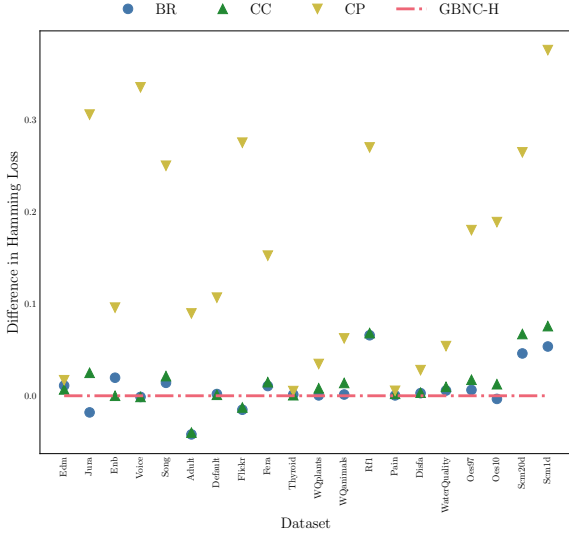
Hamming loss ( $\ell_H$ )		
GBNC-H	BR	CP
<b>11.41 <math>\pm</math> 0.35</b>	12.51 $\pm$ 1.71	21.81 $\pm$ 7.62
Subset 0/1 loss ( $\ell_S$ )		
GBNC-S	BR	CP
<b>37.31 <math>\pm</math> 0.84</b>	41.57 $\pm$ 5.16	56.57 $\pm$ 13.28

GBNCs yield the best average ranks over the 20 tabular data sets, both for  $\ell_H$  and  $\ell_S$ . Furthermore, Friedman tests [Demšar, 2006] on the ranks yield small p-values, and strongly suggest performance differences between the classifiers. We also conduct Nemenyi post-hoc test [Nemenyi, 1963] and Conover post-hoc test [Conover, 1999, Conover and Iman, 1979] (see Table 3) to see if there are significant differences between pairs of classifiers. For each combination (among the 12 combinations) of competitor, loss and local models, we find at least one test where GBNCs is significantly better than that competitor in almost all cases.

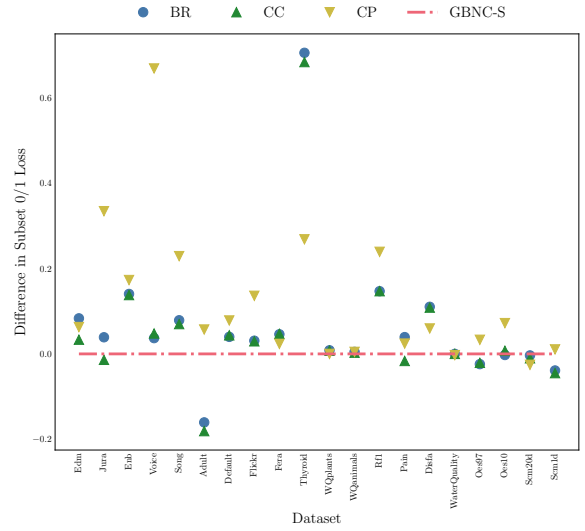
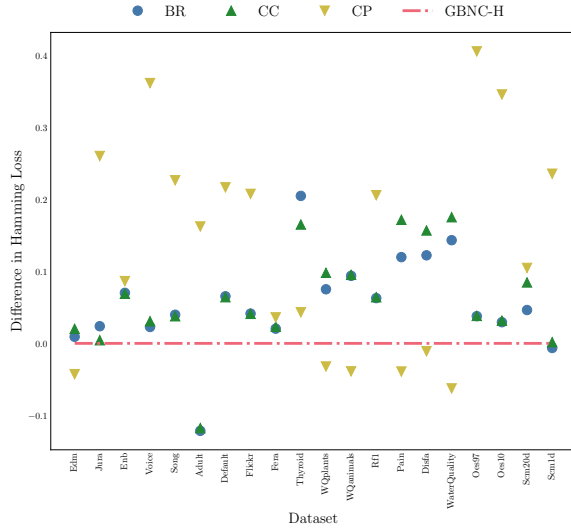
Table 2: Average ranks and p-values of Friedman tests.

The cases of Hamming loss ( $\ell_H$ )					
Learner	GBNC-H	BR	CC	CP	p-value
LR	<b>1.43</b>	1.98	2.60	4.00	<b>1.1e-09</b>
NB	<b>1.40</b>	2.70	2.95	2.95	<b>1.8e-04</b>
The cases of Subset 0/1 loss ( $\ell_S$ )					
Learner	GBNC-S	BR	CC	CP	p-value
LR	<b>1.55</b>	2.20	2.38	3.88	<b>1.2e-07</b>
NB	<b>1.73</b>	2.80	2.28	3.20	<b>1.6e-03</b>

Even if the Nemenyi post-hoc test may be too conservative, has low power, and may not detect existing differences when Friedman’s test rejects the null hypothesis (as elaborated in [Ulař et al., 2012] and also elsewhere), it already informs



(a) Base learner: Logistic Regression.



(b) Base learner: Naive Bayes.

Figure 2: Tabular data sets: Performance differences to GBNCs (negative means better than GBNCs). Data sets (x-axis) are ordered by number of class variables.

significant differences. Table 3 suggests that the use of both LR and NB as local models (i.e. base learners) leads to improvements with respect to other approaches. Actually, LR performs better with more class variables, while NB with fewer (these differences can be appreciated in the Appendices). Yet, it is not the goal of this work to answer this question. The experiments with two different local models (LR and NB) have the purpose of demonstrating the capabilities of the overall idea.

Our experimental results are in agreement with the results found in literature. First, CC can hardly be a state-of-the-art MDC approach [Jia and Zhang, 2021]. Second, BR may provide competitive performance, especially when the number of class variables is not large [Wu and Zhu, 2020]. On

the other hand, our experiments suggest a very interesting result that GBNC-H which estimates the joint conditional distribution and extracts marginal distributions using Definitions (2) often outperforms BR which directly estimates the marginal distributions. This suggests that capturing the dependency relationships can lead to more accurate estimates of the marginal probability distributions.

Although comparing ranks [Demšar, 2006] of classifiers is a common practice when one seeks short summaries of the performances, there is no golden rule about how the classifiers should be ranked. In this case, ranking the losses can not tell us whether there is any visible gain/loss. To gain more insights into the differences between classifiers, we make scatter plots for the losses provided by pairs of



Table 3: Post-hoc tests: p-values.

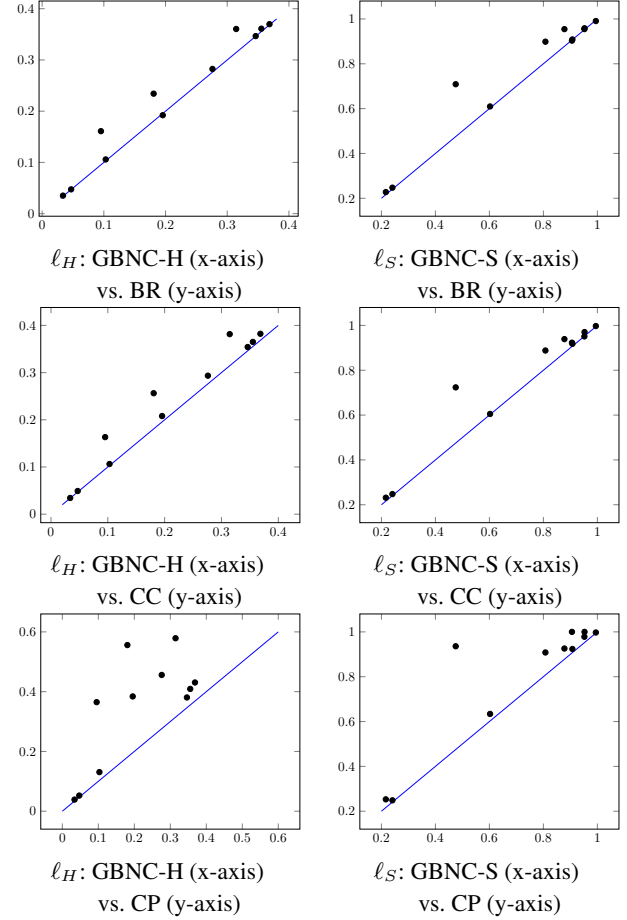
The cases of $\ell_H$ : p-values < 0.05 are given in bold				
$H_0$	Nemenyi		Conover	
	LR	NB	LR	NB
GBNC-H = BR	0.529	<b>0.008</b>	0.184	<b>0.002</b>
GBNC-H = CC	<b>0.021</b>	<b>0.001</b>	<b>0.006</b>	<b>3.9e-04</b>
GBNC-H = CP	<b>0.001</b>	<b>0.001</b>	<b>4.6e-08</b>	<b>3.9e-04</b>
BR = CP	<b>0.001</b>	0.9	<b>7.0e-07</b>	0.545
CC = CP	<b>0.003</b>	0.9	<b>0.001</b>	1
BR = CC	0.42	0.9	0.132	0.545
The cases of $\ell_S$ : p-values < 0.05 are given in bold				
$H_0$	Nemenyi		Conover	
	LR	NB	LR	NB
GBNC-S = BR	0.384	<b>0.042</b>	0.118	<b>0.01</b>
GBNC-S = CC	0.180	0.528	<b>0.049</b>	0.178
GBNC-S = CP	<b>0.001</b>	<b>0.002</b>	<b>4.9e-07</b>	<b>5.6e-04</b>
BR = CP	<b>0.001</b>	0.735	<b>1.4e-04</b>	0.326
CC = CP	<b>0.001</b>	0.106	<b>5.5e-04</b>	<b>0.026</b>
BR = CC	0.9	0.563	0.671	0.198

classifiers (See Figure 4–7 in Appendix E.2). In all cases, GBNC-H and GBNC-S are rarely worse than others with visible differences, and visible gains of GBNC-H and GBNC-S are observed in all cases. Again, those figures suggest that GBNC-H and GBNC-S can consistently provide promising performance. In practice, we would expect to see approaches which take into account dependencies among the class variables brings more advantages when the number of class variables  $K$  increases and the base learner is accurate. To show this ability of GBNCs, we make scatter plots for the losses provided by pairs of classifiers on 11 data sets with  $K \geq 7$  with LR as the base learner (which is often more accurate than NB on these data sets). Figure 3 confirms that GBNCs indeed provide visible gains on these data sets.

Finally, we acknowledge that one can devise creative ideas to tackle MDC indirectly via other approaches, so one might ask to which extent our experiments yield state-of-the-art performance in a broader sense. We emphasize that our goal is to improve on probabilistic MDC itself and to demonstrate the usefulness of this framework which has proven optimality properties and is very flexible to work with many other (off-the-shelf) classifiers as internal local models (i.e. base learners). If one embraces the framework and chooses strong local models, this is likely (based on the theoretical results) to perform very well for MDC.

## 6 CONCLUSION

We propose a formal framework for probabilistic multi-dimensional classification (MDC) in which learning an optimal multi-dimensional classifier can be decomposed into learning a set of probabilistic classifiers and learning

Figure 3:  $\ell_H$  and  $\ell_S$  with  $K \geq 7$  (base learner: LR)

an optimal Bayesian network (BN) structure. We discuss how single-class-variable probabilistic classification and BN learning can be directly integrated into the framework with respect to optimality, representational capacity and scalability. We present algorithmic solutions for the learning and inference problems and discuss on their complexity. Finally, a set of experiments highlights the usefulness of the MDC framework. We hope that this paper can open doors for further research on all these strongly related topics.

## Acknowledgements

This work was initiated when all authors were at the TU Eindhoven. Vu-Linh Nguyen has been funded by the Junior Professor Chair in Trustworthy AI (Ref. ANR-R311CHD). Yang Yang has been funded by the Research Foundation – Flanders (FWO, G097720N). This work was partially funded/supported by the EU European Defence Fund Project KOIOS (EDF-2021-DIGIT-R-FL-KOIOS) and Dutch NWO Perspectief 2022 Project PersOn (P21-03).

## References

- David Atienza, Pedro Larrañaga, and Concha Bielza. Hybrid semiparametric Bayesian networks. *TEST*, 31(2):299–327, 2022.
- Thomas Augustin, Frank PA Coolen, Gert De Cooman, and Matthias CM Troffaes. *Introduction to imprecise probabilities*. John Wiley & Sons, 2014.
- Mark Bartlett and James Cussens. Integer linear programming for the Bayesian network structure learning problem. *Artificial Intelligence*, 244:258–271, 2017.
- Marco Benjumbeda, Concha Bielza, and Pedro Larrañaga. Tractability of most probable explanations in multidimensional Bayesian network classifiers. *International Journal of Approximate Reasoning*, 93:74–87, 2018.
- J. M. Bernardo and A. F. M. Smith. *Bayesian Theory*. Chichester: Wiley, 1st edition, 2000.
- Concha Bielza, Guangdi Li, and Pedro Larranaga. Multidimensional classification with Bayesian networks. *International Journal of Approximate Reasoning*, 52(6):705–727, 2011.
- Guillaume Bouchard and Bill Triggs. The tradeoff between generative and discriminative classifiers. In *Proceedings of the 16th IASC International Symposium on Computational Statistics (COMPSTAT)*, pages 721–728, 2004.
- Remco R Bouckaert. Properties of Bayesian belief network learning algorithms. In *Proceedings of the Tenth International Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 102–109, 1994.
- Alexandra M Carvalho, Teemu Roos, Arlindo L Oliveira, and Petri Myllymäki. Discriminative learning of Bayesian networks via factorized conditional log-likelihood. *Journal of Machine Learning Research*, 12(7), 2011.
- William Jay Conover. *Practical nonparametric statistics*, volume 350. John Wiley & Sons, 1999.
- William Jay Conover and Ronald L Iman. On multiple-comparisons procedures. *Los Alamos Sci. Lab. Tech. Rep. LA-7677-MS*, 1:14, 1979.
- Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13(1):21–27, 1967.
- James Cussens, Matti Järvisalo, Janne H Korhonen, and Mark Bartlett. Bayesian network structure learning with integer programming: Polytopes, facets and complexity. *Journal of Artificial Intelligence Research*, 58:185–229, 2017.
- Cassio P de Campos. Almost No News on the Complexity of MAP in Bayesian Networks. In *Proceedings of the 10th International Conference on Probabilistic Graphical Models (PGM)*, pages 149–160. PMLR, 2020.
- Cassio P de Campos, Mauro Scanagatta, Giorgio Corani, and Marco Zaffalon. Entropy-based pruning for learning Bayesian networks using BIC. *Artificial Intelligence*, 260:42–50, 2018.
- Krzysztof Dembczyński, Weiwei Cheng, and Eyke Hüllermeier. Bayes optimal multilabel classification via probabilistic classifier chains. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pages 279–286, 2010.
- Krzysztof Dembczyński, Willem Waegeman, Weiwei Cheng, and Eyke Hüllermeier. On label dependence and loss minimization in multi-label classification. *Machine Learning*, 88(1):5–45, 2012.
- Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine learning research*, 7:1–30, 2006.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the 2009 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255. Ieee, 2009.
- Pedro M Domingos and Michael J Pazzani. Beyond independence: conditions for the optimality of the simple Bayesian classifier. In *Proceedings of the Thirteenth International Conference on Machine Learning (ICML)*, pages 105–112, 1996.
- Charles Elkan. The foundations of cost-sensitive learning. In *Proceedings of the 17th International Joint Conference on Artificial intelligence (IJCAI)*, pages 973–978, 2001.
- Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The Pascal visual object classes (VOC) challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010.
- Nir Friedman. The Bayesian structural EM algorithm. In *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 129–138, 1998.
- Nir Friedman, Dan Geiger, and Moises Goldszmidt. Bayesian network classifiers. *Machine Learning*, 29(2):131–163, 1997.
- Mudasir A Ganaie, Minghui Hu, AK Malik, M Tanveer, and PN Suganthan. Ensemble deep learning: A review. *Engineering Applications of Artificial Intelligence*, 115:105151, 2022.

- Santiago Gil-Begue, Concha Bielza, and Pedro Larrañaga. Multi-dimensional Bayesian network classifiers: A survey. *Artificial Intelligence Review*, 54(1):519–559, 2021.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pages 1321–1330, 2017.
- John T Hancock and Taghi M Khoshgoftaar. Survey on categorical data for neural networks. *Journal of Big Data*, 7(1):1–41, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the 29th IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- Christoph Jansen, Georg Schollmeyer, and Thomas Augustin. Quantifying degrees of e-admissibility in decision making with imprecise probabilities. In *Reflections on the Foundations of Probability and Statistics: Essays in Honor of Teddy Seidenfeld*, pages 319–346. Springer, 2022.
- Bin-Bin Jia and Min-Ling Zhang. Decomposition-based classifier chains for multi-dimensional classification. *IEEE Transactions on Artificial Intelligence*, 3(2):176–191, 2021.
- Bin-Bin Jia and Min-Ling Zhang. Multi-dimensional classification via selective feature augmentation. *Machine Intelligence Research*, 19(1):38–51, 2022.
- Neville Kenneth Kitson, Anthony C Constantinou, Zhigao Guo, Yang Liu, and Kiattikun Chobtham. A survey of Bayesian network structure learning. *Artificial Intelligence Review*, pages 1–94, 2023.
- Adrienne Kline, Hanyin Wang, Yikuan Li, Saya Dennis, Meghan Hutch, Zhenxing Xu, Fei Wang, Feixiong Cheng, and Yuan Luo. Multimodal machine learning in precision health: A scoping review. *NPJ Digital Medicine*, 5(1):1–14, 2022.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT Press, 2009.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- Uri Lerner, Eran Segal, and Daphne Koller. Exact inference in networks with discrete children of continuous parents. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 319–328, 2001.
- Zhifa Liu, Brandon Malone, and Changhe Yuan. Empirical evaluation of scoring functions for Bayesian network model selection. In *BMC Bioinformatics*, volume 13, pages 1–16. Springer, 2012.
- Scott Menard. *Applied logistic regression analysis*. Sage, 2002.
- Thomas Mortier, Marek Wydmuch, Krzysztof Dembczyński, Eyke Hüllermeier, and Willem Waegeman. Efficient set-valued prediction in multi-class classification. *Data Mining and Knowledge Discovery*, 35(4):1435–1469, 2021.
- Peter Bjorn Nemenyi. Distribution-free multiple comparisons. Master’s thesis, Princeton University, United States, 1963.
- Andrew Y Ng and Michael I Jordan. On discriminative vs. generative classifiers: a comparison of logistic regression and naive Bayes. In *Proceedings of the 14th International Conference on Neural Information Processing Systems (NIPS)*, pages 841–848, 2001.
- Vu-Linh Nguyen and Eyke Hüllermeier. Multilabel classification with partial abstention: Bayes-optimal prediction under label independence. *Journal of Artificial Intelligence Research*, 72:613–665, 2021.
- Vu-Linh Nguyen, Sébastien Destercke, Marie-Hélène Masson, and Rashad Ghassani. Racing trees to query partial data. *Soft Computing*, 25(14):9285–9305, 2021.
- Siqi Nie, Cassio P. de Campos, and Qiang Ji. Efficient learning of Bayesian networks with bounded tree-width. *International Journal of Approximate Reasoning*, 80:412–427, 2017.
- Paola M.V. Rancoita, Marco Zaffalon, Emanuele Zucca, Francesco Bertoni, and Cassio P. de Campos. Bayesian network data imputation with application to survival tree analysis. *Computational Statistics & Data Analysis*, 93:373–387, 2016.
- Jesse Read, Bernhard Pfahringer, Geoffrey Holmes, and Eibe Frank. Classifier chains: a review and perspectives. *Journal of Artificial Intelligence Research*, 70:683–718, 2021.
- Teemu Roos, Hannes Wettig, Peter Grünwald, Petri Myllymäki, and Henry Tirri. On discriminative Bayesian network classifiers and logistic regression. *Machine Learning*, 59(3):267–296, 2005.
- Dan Roth. On the hardness of approximate reasoning. *Artificial Intelligence*, 82(1):273–302, 1996.
- Maytal Saar-Tsechansky and Foster Provost. Handling missing values when applying classification models. *Journal of Machine Learning Research*, 8:1625–1657, 2007.
- Mauro Scanagatta, Giorgio Corani, Cassio P de Campos, and Marco Zaffalon. Learning treewidth-bounded Bayesian networks with thousands of variables. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.

- Gideon Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, pages 461–464, 1978.
- Hagit Shatkay, Fengxia Pan, Andrey Rzhetsky, and W John Wilbur. Multi-dimensional classification of biomedical text: Toward automated, practical provision of high-utility text to diverse users. *Bioinformatics*, 24(18):2086–2093, 2008.
- Solomon Eyal Shimony. Finding MAPs for belief networks is NP-hard. *Artificial Intelligence*, 68(2):399–410, 1994.
- Matthias CM Troffaes. Decision making under uncertainty using imprecise probabilities. *International Journal of Approximate Reasoning*, 45(1):17–29, 2007.
- Aydın Ulaş, Olcay Taner Yıldız, and Ethem Alpaydın. Cost-conscious comparison of supervised learning algorithms over multiple data sets. *Pattern Recognition*, 45(4):1772–1781, 2012.
- Ilkay Ulusoy and Christopher M Bishop. Comparison of generative and discriminative techniques for object detection and classification. In *Toward Category-Level Object Recognition*, pages 173–195. Springer, 2006.
- Willem Waegeman, Krzysztof Dembczyński, Arkadiusz Jachnik, Weiwei Cheng, and Eyke Hüllermeier. On the Bayes-optimality of F-measure maximizers. *Journal of Machine Learning Research*, 15:3333–3388, 2014.
- Guoqiang Wu and Jun Zhu. Multi-label classification: do Hamming loss and subset accuracy really conflict with each other? In *Proceedings of the 34th International Conference on Neural Information Processing Systems (NeurIPS)*, pages 3130–3140, 2020.
- Zhen Xu, David R So, and Andrew M Dai. Mufasa: Multimodal fusion architecture search for electronic health records. In *Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI)*, volume 35, pages 10532–10540, 2021.
- Yang Yang. Generalized Bayesian network classifiers. Master’s thesis, Eindhoven University of Technology, The Netherlands, 2022.
- Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *Proceedings of the 5th International Conference on Learning Representations (ICLR)*, 2017.