

# Generalized Bayesian Network Classifiers

Master's Thesis Presentation

Yang Yang

Supervisors: Cassio de Campos, Vu-Linh Nguyen

November 5, 2022

1. Introduction
2. Background
3. Generalized Bayesian Network Classifiers (GBNCs)
4. Experiments
5. Conclusion

# Introduction

In a nutshell

## Multi-Dimensional Classification

- ▶ Each instance is characterized by **multiple** class variables  $\mathbf{Y} = \{Y_1, \dots, Y_d\}$  ( $d \geq 2$ ).
- ▶ Each class variable can take **multiple** ( $\geq 2$ ) values ( $r_j = |\mathcal{Y}_j| \geq 2, \forall 1 \leq j \leq d$ ).

In a nutshell

## Multi-Dimensional Classification

- ▶ Each instance is characterized by **multiple** class variables  $\mathbf{Y} = \{Y_1, \dots, Y_d\}$  ( $d \geq 2$ ).
- ▶ Each class variable can take **multiple** ( $\geq 2$ ) values ( $r_j = |\mathcal{Y}_j| \geq 2, \forall 1 \leq j \leq d$ ).

## A general MDC task

$$\underbrace{x_1, x_2, \dots, x_m}_{\mathbf{x}} \xrightarrow{h(\mathbf{x})} \underbrace{y_1, y_2, \dots, y_d}_{\mathbf{y}}$$

In a nutshell

## Multi-Dimensional Classification

- ▶ Each instance is characterized by **multiple** class variables  $\mathbf{Y} = \{Y_1, \dots, Y_d\}$  ( $d \geq 2$ ).
- ▶ Each class variable can take **multiple** ( $\geq 2$ ) values ( $r_j = |\mathcal{Y}_j| \geq 2, \forall 1 \leq j \leq d$ ).

## A general MDC task

$$\underbrace{x_1, x_2, \dots, x_m}_{\mathbf{x}} \xrightarrow{h(\mathbf{x})} \underbrace{y_1, y_2, \dots, y_d}_{\mathbf{y}}$$

In a nutshell

## Multi-Dimensional Classification

- ▶ Each instance is characterized by **multiple** class variables  $\mathbf{Y} = \{Y_1, \dots, Y_d\}$  ( $d \geq 2$ ).
- ▶ Each class variable can take **multiple** ( $\geq 2$ ) values ( $r_j = |\mathcal{Y}_j| \geq 2, \forall 1 \leq j \leq d$ ).

## A general MDC task

$$\underbrace{x_1, x_2, \dots, x_m}_{\mathbf{x}} \xrightarrow{h(\mathbf{x})} \underbrace{y_1, y_2, \dots, y_d}_{\mathbf{y}}$$

Feature Space	Class Space		
	Color	Brand	Type
	White	BMW	Car
	White	Mercedes	Truck
	Yellow	Lamborghini	Supercar
:	:	:	:
	Red/White	Toyota	Offroad
	?	?	?

**Table 1:** Multi-dimensional vehicle image classification.  
Adapted from [3].

## Multi-Dimensional Classification (MDC)

While MDC is hard, it may become harder!

While MDC is hard, it may become harder!

## Various types of features can coexist

- ▶ Numeric values.
  - ▶ 0.6, 2.36, 17.85, ...
- ▶ Binary values.
  - ▶ 0, 1
- ▶ Ordinal values.
  - ▶ Likert scale: Like (0), Like Somewhat (1), Neutral (2), Dislike Somewhat (3), Dislike (4)
- ▶ Non-ordinal discrete signals.
  - ▶ Color: Red, Blue, Yellow, Green, ...
  - ▶ Gender: Male, Female, ...

# Multi-Dimensional Classification (MDC)

While MDC is hard, it may become harder!

## Various types of features can coexist

- ▶ Numeric values.
  - ▶ 0.6, 2.36, 17.85, ...
- ▶ Binary values.
  - ▶ 0, 1
- ▶ Ordinal values.
  - ▶ Likert scale: Like (0), Like Somewhat (1), Neutral (2), Dislike Somewhat (3), Dislike (4)
- ▶ Non-ordinal discrete signals.
  - ▶ Color: Red, Blue, Yellow, Green, ...
  - ▶ Gender: Male, Female, ...

World is **multi-modal!**



**Figure 1:** "A cat is playing with a dog."

## A General MDC Task

$$\underbrace{x_1, x_2, \dots, x_m}_{\mathbf{x}} \xrightarrow{h(\mathbf{x})} \underbrace{y_1, y_2, \dots, y_d}_{\mathbf{y}}$$

## A Probabilistic MDC Task

$$\underbrace{x_1, x_2, \dots, x_m}_{\mathbf{x}}$$

$$\xrightarrow{\text{Modeling}} P(\mathbf{Y} \mid \mathbf{X})$$

$$\xrightarrow{\text{Inference}} \underbrace{y_1, y_2, \dots, y_d}_{\mathbf{y}}$$

# Probabilistic MDC

## A General MDC Task

$$\underbrace{x_1, x_2, \dots, x_m}_{\mathbf{x}} \xrightarrow{h(\mathbf{x})} \underbrace{y_1, y_2, \dots, y_d}_{\mathbf{y}}$$

## A Probabilistic MDC Task

$$\begin{aligned} & \xrightarrow{\quad \quad \quad} \underbrace{x_1, x_2, \dots, x_m}_{\mathbf{x}} \\ \text{Modeling} \xrightarrow{\quad \quad \quad} & P(\mathbf{Y} \mid \mathbf{X}) \\ \text{Inference} \xrightarrow{\quad \quad \quad} & \underbrace{y_1, y_2, \dots, y_d}_{\mathbf{y}} \end{aligned}$$



## A General MDC Task

$$\underbrace{x_1, x_2, \dots, x_m}_{\mathbf{x}} \xrightarrow{h(\mathbf{x})} \underbrace{y_1, y_2, \dots, y_d}_{\mathbf{y}}$$

## A Probabilistic MDC Task

$$\begin{array}{c} \xrightarrow{\text{Modeling}} \quad \underbrace{x_1, x_2, \dots, x_m}_{\mathbf{x}} \\ \qquad \quad P(\mathbf{Y} \mid \mathbf{X}) \\ \xrightarrow{\text{Inference}} \quad \underbrace{y_1, y_2, \dots, y_d}_{\mathbf{y}} \end{array}$$



- ▶ Probabilities associated to predictions.
- ▶ Optimal predictions under different loss functions.
- ▶ How to capture the **probabilistic relationships** between class variables?



# Background

What is a Bayesian network (BN)?

What is a Bayesian network (BN)?

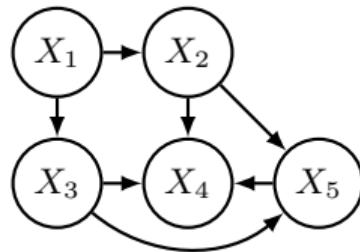
$$\mathcal{B} = (\mathcal{G}, \Theta)$$

- ▶  $\mathcal{G} = (\mathbf{V}, \mathbf{E})$  is a DAG with  $\mathbf{V}$  a collection of nodes associated to random variables (RVs)  $\mathbf{X}$ .
- ▶  $\Theta = \{\Theta_i \mid 1 \leq i \leq m\}$  is a collection of parameters
  - ▶ encodes local conditional probability distributions (CPDs)  $\{P(X_i \mid \mathbf{Pa}(X_i)) \mid 1 \leq i \leq m\}$  of  $\mathbf{X}$ .

What is a Bayesian network (BN)?

$$\mathcal{B} = (\mathcal{G}, \Theta)$$

- ▶  $\mathcal{G} = (\mathbf{V}, \mathbf{E})$  is a DAG with  $\mathbf{V}$  a collection of nodes associated to random variables (RVs)  $\mathbf{X}$ .
- ▶  $\Theta = \{\Theta_i \mid 1 \leq i \leq m\}$  is a collection of parameters
  - ▶ encodes local conditional probability distributions (CPDs)  $\{P(X_i \mid \text{Pa}(X_i)) \mid 1 \leq i \leq m\}$  of  $\mathbf{X}$ .
- ▶  $\mathcal{B}$  represents a joint probability distribution over  $\mathbf{X}$ .



**Figure 2:**  $P_{\mathcal{B}} = P(X_1)P(X_2 \mid X_1)P(X_3 \mid X_1)P(X_4 \mid X_2, X_3, X_5)P(X_5 \mid X_2, X_3)$

Why using a Bayesian network (BN)?

Why using a Bayesian network (BN)?

## Advantages of Bayesian Networks (BNs)

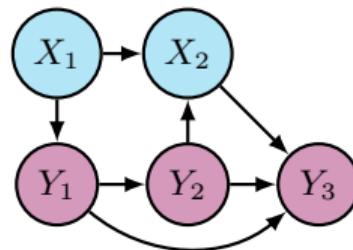
- ▶ Intuitive graphical formalisms of probabilistic relationships.
- ▶ Interpretable representations of uncertainties supported by probability theory.

Why using a Bayesian network (BN)?

## Advantages of Bayesian Networks (BNs)

- ▶ Intuitive graphical formalisms of probabilistic relationships.
- ▶ Interpretable representations of uncertainties supported by probability theory.

(Multi-dimensional) BNCs are simply BNs applied to (multi-dimensional) classification problems.



**Figure 3:** An example (M)BNC over  $\mathbf{X} = \{X_1, X_2\}$  and  $\mathbf{Y} = \{Y_1, Y_2, Y_3\}$ .

## Learning an (M)BNC

Learning an (M)BNC is essentially learning the underlying Bayesian network.

Learning an (M)BNC is essentially learning the underlying Bayesian network.

## Structure Learning: Learning $\mathcal{G}$ from data $\mathcal{D}$

- ▶ Constraint-based
  - ▶ Employ statistical tests to identify independencies between variables.
- ▶ Score-based
  - ▶ Define a structure score metric (such as BD scores and the BIC score).
  - ▶ Find a structure achieving the maximum score from the structure space.

Learning an (M)BNC is essentially learning the underlying Bayesian network.

## Structure Learning: Learning $\mathcal{G}$ from data $\mathcal{D}$

- ▶ Constraint-based
  - ▶ Employ statistical tests to identify independencies between variables.
- ▶ Score-based
  - ▶ Define a structure score metric (such as BD scores and the BIC score).
  - ▶ Find a structure achieving the maximum score from the structure space.

## Parameter Learning: Learning $\Theta$ from data $\mathcal{D}$

- ▶ Bayesian Learning
  - ▶ Treat parameters as random variables and update  $\Theta$  by using Bayes' rule:  $p(\Theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \Theta)p(\Theta)}{p(\mathcal{D})}$ .
- ▶ Maximum Likelihood Estimation (MLE)
  - ▶ Pick parameters that maximize the model's probability of generating  $\mathcal{D}$ .
  - ▶ Given enough data, uncover the real **data-generating** distribution  $P(\mathbf{X}, \mathbf{Y})$ .

## Maximum Likelihood Estimation (MLE)

- ▶ Pick parameters to make the model to be close to the **data-generating** distribution  $P(\mathbf{X}, \mathbf{Y})$ .

## Maximum Likelihood Estimation (MLE)

- ▶ Pick parameters to make the model to be close to the **data-generating** distribution  $P(\mathbf{X}, \mathbf{Y})$ .

However, classification is a **discriminative** or **supervised** task.

In an (M)BNC, what we need is the **conditional distribution**  $P(\mathbf{Y} \mid \mathbf{X})$  of class variables.

## Maximum Likelihood Estimation (MLE)

- ▶ Pick parameters to make the model to be close to the **data-generating** distribution  $P(\mathbf{X}, \mathbf{Y})$ .

However, classification is a **discriminative** or **supervised** task.

In an (M)BNC, what we need is the **conditional distribution**  $P(\mathbf{Y} | \mathbf{X})$  of class variables.

## Maximizing the Conditional Log-Likelihood (CLL)

Given data  $\mathcal{D} = \{(\mathbf{x}^l, \mathbf{y}^l) \mid 1 \leq l \leq N\}$ , the CLL is a **discriminative** objective function:

$$\text{CLL}(\mathcal{G}, \Theta \mid \mathcal{D}) := \log \prod_{l=1}^N p_{\mathcal{B}}(\mathbf{y}^l \mid \mathbf{x}^l) = \log \prod_{l=1}^N \frac{f_{\mathcal{B}}(\mathbf{x}^l, \mathbf{y}^l)}{\sum_{\mathbf{y} \in \mathcal{Y}} f_{\mathcal{B}}(\mathbf{x}^l, \mathbf{y})} \quad (1)$$

## Maximum Likelihood Estimation (MLE)

- ▶ Pick parameters to make the model to be close to the **data-generating** distribution  $P(\mathbf{X}, \mathbf{Y})$ .

However, classification is a **discriminative** or **supervised** task.

In an (M)BNC, what we need is the **conditional distribution**  $P(\mathbf{Y} \mid \mathbf{X})$  of class variables.

## Maximizing the Conditional Log-Likelihood (CLL)

Given data  $\mathcal{D} = \{(\mathbf{x}^l, \mathbf{y}^l) \mid 1 \leq l \leq N\}$ , the CLL is a **discriminative** objective function:

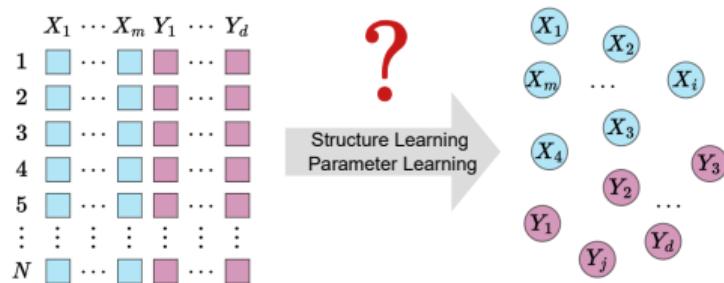
$$\text{CLL}(\mathcal{G}, \boldsymbol{\Theta} \mid \mathcal{D}) := \log \prod_{l=1}^N p_{\mathcal{B}}(\mathbf{y}^l \mid \mathbf{x}^l) = \log \prod_{l=1}^N \frac{f_{\mathcal{B}}(\mathbf{x}^l, \mathbf{y}^l)}{\sum_{\mathbf{y} \in \mathcal{Y}} f_{\mathcal{B}}(\mathbf{x}^l, \mathbf{y})} \quad (1)$$

- ▶ Unfortunately, CLL does not decompose over  $\mathcal{G}$  into a separate term for each variable.
- ▶ There is no closed-form solution for optimizing parameters to maximize the CLL.

# Generalized Bayesian Network Classifiers (GBNCs)

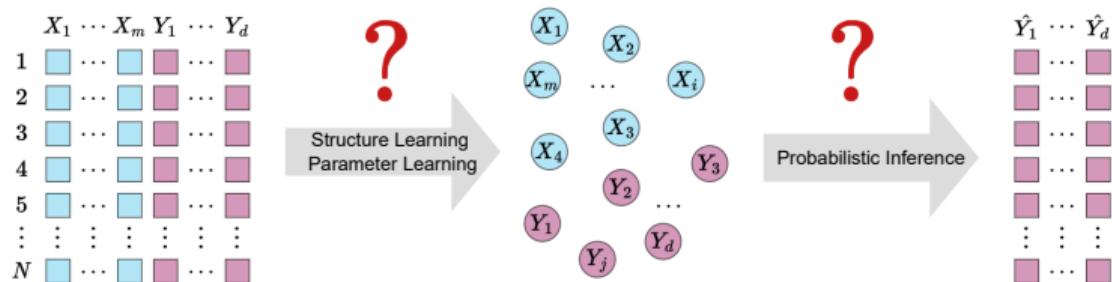
# Research Questions

- ▶ How to **discriminatively** learn the **parameters** and the **structure** of a BNC?



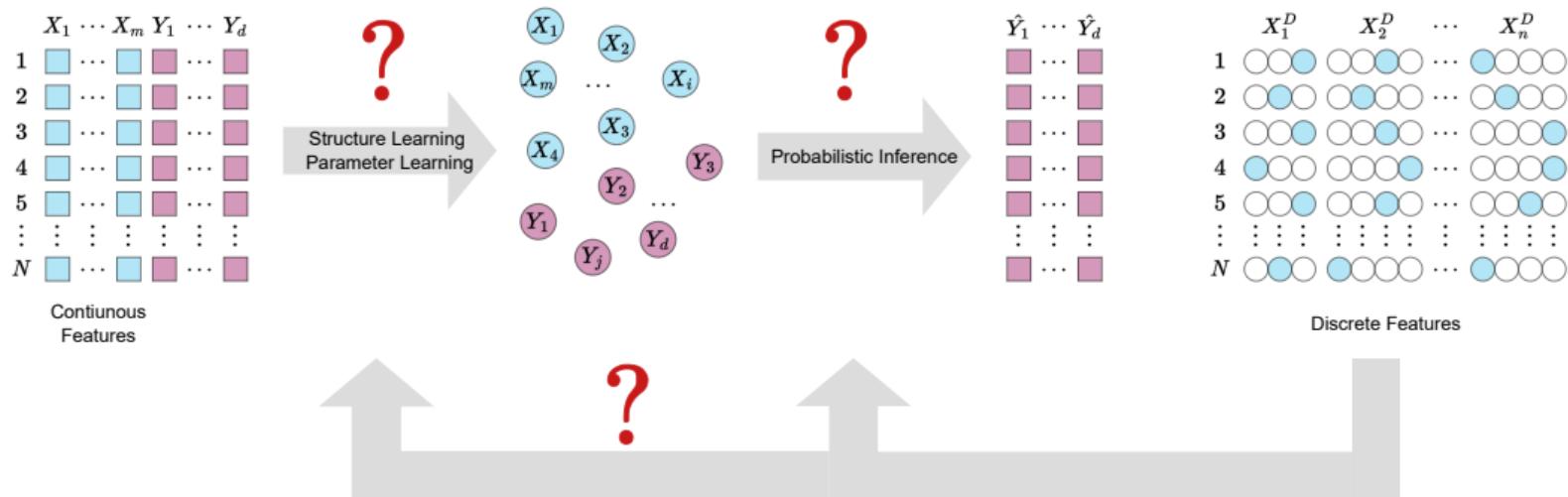
# Research Questions

- ▶ How to **discriminatively** learn the **parameters and the structure** of a BNC?
- ▶ How to perform **probabilistic inference** to compute optimal predictions under **different loss functions**?



# Research Questions

- ▶ How to **discriminatively** learn the **parameters and the structure** of a BNC?
- ▶ How to perform **probabilistic inference** to compute optimal predictions under **different loss functions**?
- ▶ How to handle **mixed data**, i.e., continuous and discrete feature variables coexist?





To enable the CLL decomposability

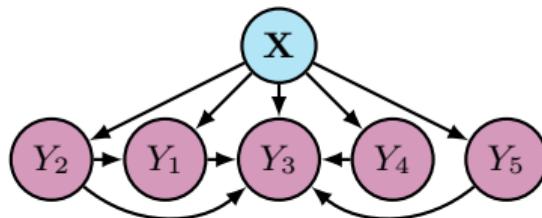
## Structural constraints in GBNCs

- ▶ There is no directed edge from class variables to feature variables.
- ▶ There is a directed edge from any feature variable to any class variable.

To enable the CLL decomposability

## Structural constraints in GBNCs

- ▶ There is no directed edge from class variables to feature variables.
- ▶ There is a directed edge from any feature variable to any class variable.



**Figure 4:** An example GBNC over  $\mathbf{X}$  and  $\mathbf{Y} = \{Y_1, Y_2, Y_3, Y_4, Y_5\}$ .

## Proposition 1: CLL Decomposition

Under our structural constraints, the CLL can be simplified as

$$\text{CLL}(\mathcal{G}, \Theta \mid \mathcal{D}) = \sum_{l=1}^N \sum_{j=1}^d \log P_{\mathcal{B}}(y_j^l \mid \mathbf{pa}(y_j)^l) \quad (2)$$

- ▶  $\Pi_j = \mathbf{Pa}(Y_j) \cap \mathbf{Y}$ , with the number of configurations of  $\Pi_j$  as  $q_j$ .
- ▶  $\Phi_j = \mathbf{Pa}(Y_j) \cap \mathbf{X}$ .

The CLL can be further simplified as

$$\text{CLL}(\mathcal{G}, \Theta \mid \mathcal{D}) = \sum_{l=1}^N \sum_{j=1}^d \sum_{\substack{k=1: \\ \mathbb{1}_{\pi_j^l} = \pi_{jk}}}^{q_j} \log P_j^{\pi_{jk}}(y_j^l \mid \phi_j^l) \quad (3)$$

$$\text{CLL}(\mathcal{G}, \Theta \mid \mathcal{D}) = \sum_{l=1}^N \sum_{j=1}^d \sum_{\substack{k=1: \\ \mathbb{1}_{\pi_j^l = \pi_{jk}}}}^{q_j} \log P_j^{\pi_{jk}}(y_j^l \mid \phi_j^l) \quad (4)$$

- ▶ In words, we will need  $q = q_1 + \dots + q_d$  probabilistic models to represent the distribution  $P_{\mathcal{B}}(\mathbf{Y} \mid \mathbf{X})$ .
  - ▶ One for each  $P_j^{\pi_{jk}}(Y_j \mid \Phi_j)$  ( $1 \leq j \leq d, 1 \leq k \leq q_j$ ).

$$\text{CLL}(\mathcal{G}, \Theta \mid \mathcal{D}) = \sum_{l=1}^N \sum_{j=1}^d \sum_{\substack{k=1: \\ \mathbb{1}_{\pi_j^l = \pi_{jk}}}}^{q_j} \log P_j^{\pi_{jk}}(y_j^l \mid \phi_j^l) \quad (4)$$

- ▶ In words, we will need  $q = q_1 + \dots + q_d$  probabilistic models to represent the distribution  $P_{\mathcal{B}}(\mathbf{Y} \mid \mathbf{X})$ .
  - ▶ One for each  $P_j^{\pi_{jk}}(Y_j \mid \Phi_j)$  ( $1 \leq j \leq d, 1 \leq k \leq q_j$ ).
- ▶ The probabilistic relationships between **feature variables** do not affect the CLL.

$$\text{CLL}(\mathcal{G}, \Theta \mid \mathcal{D}) = \sum_{l=1}^N \sum_{j=1}^d \sum_{\substack{k=1: \\ \mathbb{1}_{\pi_j^l = \pi_{jk}}}}^{q_j} \log P_j^{\pi_{jk}}(y_j^l \mid \phi_j^l) \quad (4)$$

- ▶ In words, we will need  $q = q_1 + \dots + q_d$  probabilistic models to represent the distribution  $P_{\mathcal{B}}(\mathbf{Y} \mid \mathbf{X})$ .
  - ▶ One for each  $P_j^{\pi_{jk}}(Y_j \mid \Phi_j)$  ( $1 \leq j \leq d, 1 \leq k \leq q_j$ ).
- ▶ The probabilistic relationships between **feature variables** do not affect the CLL.
- ▶ Although we assume any  $Y \in \mathbf{Y}$  is connected to all  $X \in \mathbf{X}$

$$\text{CLL}(\mathcal{G}, \Theta \mid \mathcal{D}) = \sum_{l=1}^N \sum_{j=1}^d \sum_{\substack{k=1: \\ \mathbb{1}_{\pi_j^l = \pi_{jk}}}}^{q_j} \log P_j^{\pi_{jk}}(y_j^l \mid \phi_j^l) \quad (4)$$

- ▶ In words, we will need  $q = q_1 + \dots + q_d$  probabilistic models to represent the distribution  $P_{\mathcal{B}}(\mathbf{Y} \mid \mathbf{X})$ .
  - ▶ One for each  $P_j^{\pi_{jk}}(Y_j \mid \Phi_j)$  ( $1 \leq j \leq d, 1 \leq k \leq q_j$ ).
- ▶ The probabilistic relationships between **feature variables** do not affect the CLL.
- ▶ Although we assume any  $Y \in \mathbf{Y}$  is connected to all  $X \in \mathbf{X}$ 
  - ▶ Learning a **base probabilistic model**  $C_j^{\pi_{jk}}$  for each  $P_j^{\pi_{jk}}(y_j \mid \phi_j)$  is essentially a **feature selection!**

# Discriminative Parameter Learning: Input Space Partitioning

How do we learn a base probabilistic model  $C_j^{\pi_{jk}}$ ?

	$X_1$	$\dots$	$X_m$	$Y_1$	$Y_2$	$Y_3$	$Y_4$	$\dots$	$Y_d$
1	1			1	1	2	2		1
2	2			2	3	4	2		1
3	3			3	1	2	2		2
4	4			3	1	2	2		1
5	5			1	2	2	2		3
6	6			4	3	4	1		3
7	7			2	1	3	3		4
8	8			2	1	2	2		1
$\vdots$									
$N$	$N$			1	1	2	2		2

$\mathcal{D}$

# Discriminative Parameter Learning: Input Space Partitioning

How do we learn a base probabilistic model  $C_j^{\pi_{jk}}$ ?

- ▶ Extract  $\mathcal{D}_j^{\pi_{jk}}$  from  $\mathcal{D}$  according to  $\pi_{jk}$ .

	$X_1$	$\dots$	$X_m$	$Y_1$	$Y_2$	$Y_3$	$Y_4$	$\dots$	$Y_d$
1	1	...	1	1	1	2	2	...	1
2	2	...	2	2	3	4	2	...	1
3	3	...	3	3	1	2	2	...	2
4	4	...	4	3	1	2	2	...	1
5	5	...	5	1	2	2	2	...	3
6	6	...	6	4	3	4	1	...	3
7	7	...	7	2	1	3	3	...	4
8	8	...	8	2	1	2	2	...	1
$\vdots$									
$N$	$N$	...	$N$	1	1	2	2	...	2

$\mathcal{D}$

$\Pi_1 = \{Y_2, Y_3, Y_4\}$   
Partitioning  
 $\pi_{1k} = \{1, 2, 2\}$

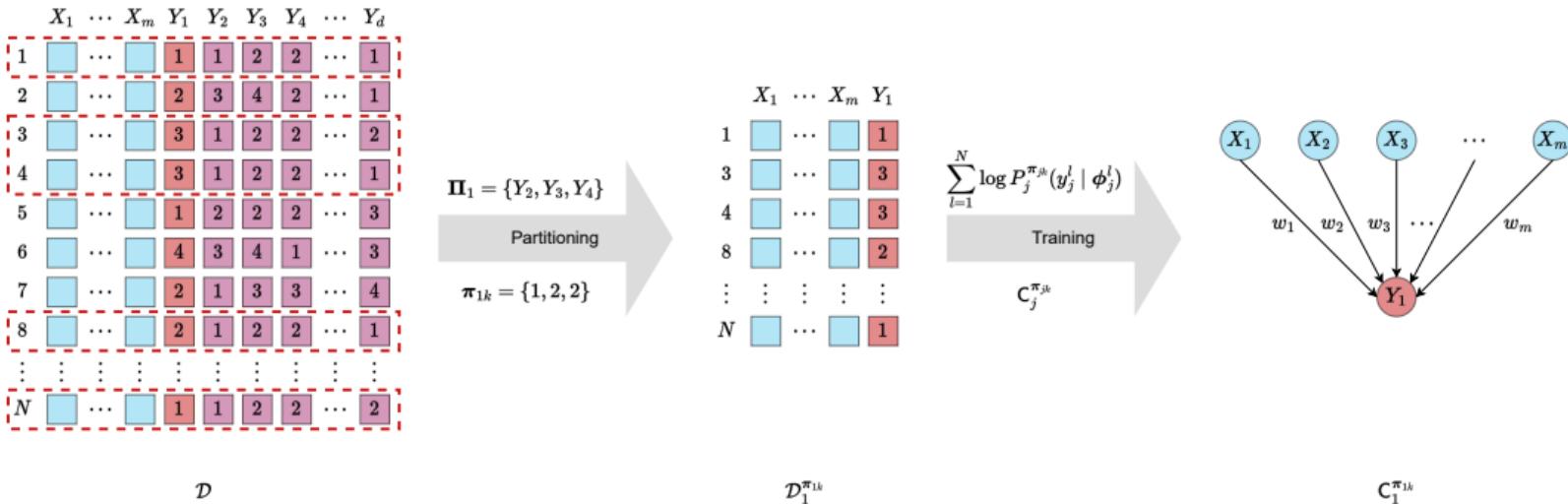
$\mathcal{D}_1^{\pi_{1k}}$

	$X_1$	$\dots$	$X_m$	$Y_1$
1	1	...	1	1
3	3	...	3	3
4	4	...	4	3
8	8	...	8	2
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$N$	$N$	...	$N$	1

# Discriminative Parameter Learning: Input Space Partitioning

How do we learn a base probabilistic model  $C_j^{\pi_{jk}}$ ?

- ▶ Extract  $\mathcal{D}_j^{\pi_{jk}}$  from  $\mathcal{D}$  according to  $\pi_{jk}$ .
- ▶ Train  $C_j^{\pi_{jk}}$  using  $\mathcal{D}_j^{\pi_{jk}}$  by iteratively optimizing  $\sum_{l=1}^N \log P_j^{\pi_{jk}}(y_j^l | \phi_j^l)$ .





The probabilistic relationships between **feature variables** do not affect the CLL.

- ▶ Learning  $\mathcal{G}$  is essentially learning the **class subgraph**  $\mathcal{G}_Y$ .

The probabilistic relationships between **feature variables** do not affect the CLL.

- ▶ Learning  $\mathcal{G}$  is essentially learning the **class subgraph**  $\mathcal{G}_Y$ .

## A score-based structure learning approach

- ▶ Use a penalized CLL score as a (**decomposable**) structure score function:

$$\begin{aligned} S(\mathcal{G}, \Theta \mid \mathcal{D}) &= \text{CLL}(\mathcal{G}, \Theta \mid \mathcal{D}) + \text{PEN}(\mathcal{G}, \Theta \mid \mathcal{D}) \\ &= \sum_{j=1}^d (\text{CLL}(\mathcal{G}_{\Pi_j}, \Theta_j \mid \mathcal{D}) + \text{PEN}(\mathcal{G}_{\Pi_j}, \Theta_j \mid \mathcal{D})) \\ &= \sum_{j=1}^d S(\mathcal{G}_{\Pi_j}, \Theta_j \mid \mathcal{D}) \end{aligned} \tag{5}$$

where  $\text{PEN}(\mathcal{G}_{\Pi_j}, \Theta_j \mid \mathcal{D}) = -\frac{\log N}{2}(r_j - 1)q_j$ .

The probabilistic relationships between **feature variables** do not affect the CLL.

- ▶ Learning  $\mathcal{G}$  is essentially learning the **class subgraph**  $\mathcal{G}_Y$ .

## A score-based structure learning approach

- ▶ Use a penalized CLL score as a (**decomposable**) structure score function:

$$\begin{aligned} S(\mathcal{G}, \Theta \mid \mathcal{D}) &= \text{CLL}(\mathcal{G}, \Theta \mid \mathcal{D}) + \text{PEN}(\mathcal{G}, \Theta \mid \mathcal{D}) \\ &= \sum_{j=1}^d (\text{CLL}(\mathcal{G}_{\Pi_j}, \Theta_j \mid \mathcal{D}) + \text{PEN}(\mathcal{G}_{\Pi_j}, \Theta_j \mid \mathcal{D})) \\ &= \sum_{j=1}^d S(\mathcal{G}_{\Pi_j}, \Theta_j \mid \mathcal{D}) \end{aligned} \tag{5}$$

where  $\text{PEN}(\mathcal{G}_{\Pi_j}, \Theta_j \mid \mathcal{D}) = -\frac{\log N}{2}(r_j - 1)q_j$ .

- ▶ We can further prune the search space of potential  $\Pi_j$  for  $Y_j$  by using rules from [2].

# Mixed Data

Continuous and discrete feature variables coexist.

	$X_1^C$	$X_2^C$	$X_3^C$	$X_4^C$	$\dots$	$X_m^C$	$X_1^D$	$X_2^D$	$X_3^D$	$X_4^D$	$\dots$	$X_n^D$
1	.26	.88	.24	.51	$\dots$	.35	1	2	2	1	$\dots$	2
2	.57	.37	.33	.89	$\dots$	.29	3	4	2	1	$\dots$	5
3	.43	.21	.93	.77	$\dots$	.93	1	2	2	2	$\dots$	9
4	.06	.16	.99	.31	$\dots$	.23	1	2	2	1	$\dots$	8
5	.22	.92	.67	.32	$\dots$	.87	2	2	2	3	$\dots$	2
6	.36	.16	.75	.41	$\dots$	.32	3	4	1	3	$\dots$	6
7	.12	.46	.96	.55	$\dots$	.55	1	3	3	4	$\dots$	7
8	.86	.01	.32	.56	$\dots$	.02	1	2	2	1	$\dots$	1
$\vdots$												
$N$	.06	.21	.90	.72	$\dots$	.49	1	2	2	2	$\dots$	5

Continuous Features

Discrete Features

Labels

	$Y_1$	$Y_2$	$Y_3$	$Y_4$	$\dots$	$Y_d$
1	1	2	3	3	$\dots$	1
2	2	3	1	2	$\dots$	1
3	3	2	3	3	$\dots$	2
3	3	2	3	2	$\dots$	1
1	1	3	2	2	$\dots$	3
4	4	3	4	1	$\dots$	3
2	2	1	3	3	$\dots$	4
2	2	2	1	2	$\dots$	1
$\vdots$						
1	1	2	3	3	$\dots$	2



Continuous and discrete feature variables coexist.

- ▶ How to handle discrete features?

	$X_1^C$	$X_2^C$	$X_3^C$	$X_4^C$	$\dots$	$X_m^C$	$X_1^D$	$X_2^D$	$X_3^D$	$X_4^D$	$\dots$	$X_n^D$
1	.26	.88	.24	.51	$\dots$	.35	1	2	2	1	$\dots$	2
2	.57	.37	.33	.89	$\dots$	.29	3	4	2	1	$\dots$	5
3	.43	.21	.93	.77	$\dots$	.93	1	2	2	2	$\dots$	9
4	.06	.16	.99	.31	$\dots$	.23	1	2	2	1	$\dots$	8
5	.22	.92	.67	.32	$\dots$	.87	2	2	2	3	$\dots$	2
6	.36	.16	.75	.41	$\dots$	.32	3	4	1	3	$\dots$	6
7	.12	.46	.96	.55	$\dots$	.55	1	3	3	4	$\dots$	7
8	.86	.01	.32	.56	$\dots$	.02	1	2	2	1	$\dots$	1
$\vdots$												
$N$	.06	.21	.90	.72	$\dots$	.49	1	2	2	2	$\dots$	5

Continuous Features

Discrete Features

Labels

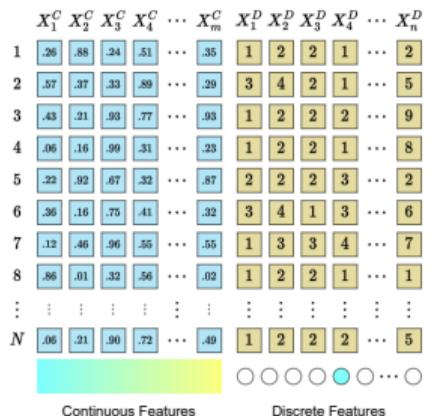
	$Y_1$	$Y_2$	$Y_3$	$Y_4$	$\dots$	$Y_d$
1	1	2	3	3	$\dots$	1
2	2	3	1	2	$\dots$	1
3	3	2	3	3	$\dots$	2
3	3	2	3	2	$\dots$	1
1	1	3	2	2	$\dots$	3
4	4	3	4	1	$\dots$	3
2	2	1	3	3	$\dots$	4
2	2	2	1	2	$\dots$	1
$\vdots$						
1	1	2	3	3	$\dots$	2



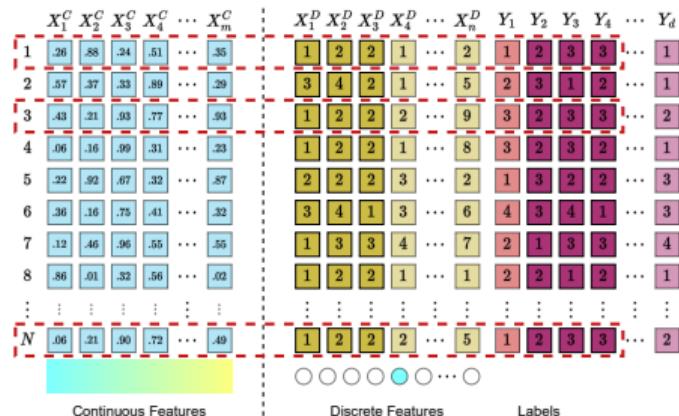
# Mixed Data

Continuous and discrete feature variables coexist.

- ▶ How to handle **discrete features**?
- ▶  $\Lambda_j = \text{Pa}(Y_j) \cap \mathbf{X}^D$ .
- ▶ Treat **discrete feature variables** as "special" class variables.
- ▶ Use **discrete features** to do further input space partitioning!



$$\begin{aligned}\Lambda_1 &= \{X_1^D, X_2^D, X_3^D\} \\ \Pi_1 &= \{Y_2, Y_3, Y_4\} \\ \text{Partitioning} & \\ \lambda_1 &= \{1, 2, 2\} \\ \pi_1 &= \{2, 3, 3\}\end{aligned}$$





Given a loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  and an input  $\mathbf{x} \in \mathcal{X}$ , the Bayes-Optimal Prediction (BOP)  $\hat{\mathbf{y}}$  is

$$\hat{\mathbf{y}} = \arg \min_{\mathbf{y}' \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}, \mathbf{y}') P(\mathbf{y} \mid \mathbf{x}) \quad (6)$$

Given a loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  and an input  $\mathbf{x} \in \mathcal{X}$ , the Bayes-Optimal Prediction (BOP)  $\hat{\mathbf{y}}$  is

$$\hat{\mathbf{y}} = \arg \min_{\mathbf{y}' \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}, \mathbf{y}') P(\mathbf{y} \mid \mathbf{x}) \quad (6)$$

**Subset 0/1 Loss:**  $\ell_{0/1}(\mathbf{y}, \hat{\mathbf{y}}) = \mathbb{1}_{\mathbf{y} \neq \hat{\mathbf{y}}}$

Most Probable Explanation (MPE) Inference for  $\mathbf{Y}$ .

$$\hat{\mathbf{y}} = \arg \min_{\mathbf{y}' \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} \mathbb{1}_{\mathbf{y} \neq \mathbf{y}'} P(\mathbf{y} \mid \mathbf{x}) = \arg \max_{\mathbf{y}' \in \mathcal{Y}} P(\mathbf{y}' \mid \mathbf{x}) \quad (7)$$

**Hamming Loss:**  $\ell_{\text{HL}}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{j=1}^d \mathbb{1}_{y_j \neq \hat{y}_j}$

Maximum A Posteriori (MAP) Inference for each  $Y \in \mathbf{Y}$ .

$$\begin{aligned} \hat{\mathbf{y}} &= \{\hat{y}_j \mid 1 \leq j \leq d\} \\ \text{where, } \hat{y}_j &= \arg \min_{y' \in \mathcal{Y}_j} \sum_{y \in \mathcal{Y}_j} \mathbb{1}_{y \neq y'} P(y \mid \mathbf{x}) = \arg \max_{y' \in \mathcal{Y}_j} P(y' \mid \mathbf{x}) \end{aligned} \quad (8)$$

---

**Algorithm 1** BOP under the subset 0/1 loss

---

**Require:** Test set  $\mathcal{T} = \{\mathbf{x}^l \mid 1 \leq l \leq L\}$ , local probabilistic classifiers  $\mathbf{C}$ , class subgraph  $\mathcal{G}_Y$

**Ensure:** The set of predictions  $\mathcal{T}_Y = \{\hat{\mathbf{y}}^l \mid 1 \leq l \leq L\}$

```

1: Initialize  $\mathcal{T}_Y \leftarrow \emptyset$ 
2: Initialize CPT of  $\mathcal{G}_Y$ :  $P(y_{jt} \mid \pi_{jk}) \leftarrow 0$ , with all  $j \in \{1, 2, \dots, d\}$ ,  $t \in \{1, 2, \dots, r_j\}$  and  $k \in \{1, 2, \dots, q_j\}$ 
3: for  $l = 1, 2, \dots, L$  do
4:   for  $j = 1, 2, \dots, d$  do
5:     for  $t = 1, 2, \dots, r_j$  do
6:       for  $k = 1, 2, \dots, q_j$  do
7:         Update  $P(y_{jt} \mid \pi_{jk}) \leftarrow C_j^{\pi_{jk}}(\mathbf{x}^l)$ 
8:       end for
9:     end for
10:   end for
11:   Perform MPE inference on  $\mathcal{G}_Y$  to compute  $\hat{\mathbf{y}}^l \leftarrow \arg \max_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} \mid \mathbf{x}^l)$ 
12:   Update  $\mathcal{T}_Y \leftarrow \mathcal{T}_Y \cup \{\hat{\mathbf{y}}^l\}$ 
13: end for

```

---

# Experiments

- ▶ 17 data sets containing only continuous feature variables.
- ▶ 3 mixed data sets containing both continuous and discrete feature variables.
- ▶ The number of class variables range from 2 to 16.
- ▶ 5, 6 and 22 discrete feature variables in the 3 mixed data sets, respectively.

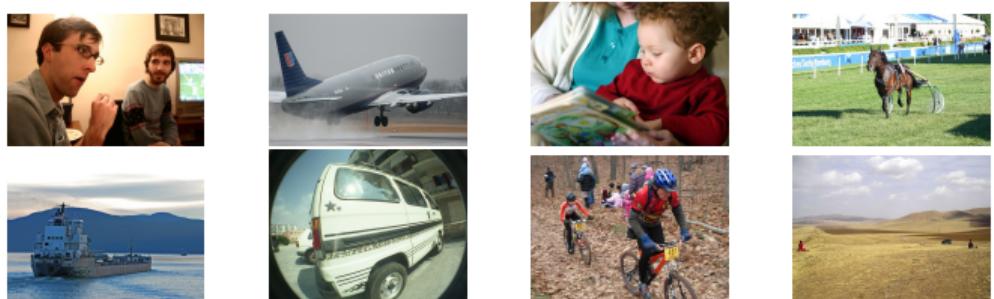
Data Set	#CV	#Samples	#States/CV	#Features
Edm	2	154	3	16n
Jura	2	359	4,5	9n
Enb	2	768	2,4	6n
Voice	2	3136	4,2	19n
Song	3	785	3	98n
Flickr	5	12198	3, 4, 3, 4, 4	1536n
Fera	5	14052	6	136n
WQplants	7	1060	4	16n
WQanimals	7	1060	4	16n
Rf1	8	8987	4, 4, 3, 4, 4, 3, 4, 3	64n
Pain	10	9734	2, 5, 4, 2, 2, 5, 2, 5, 2, 2	136n
Disfa	12	13095	5, 5, 6, 3, 4, 4, 5, 4, 4, 4, 6, 4	136n
WaterQuality	14	1060	4	16n
Oes97	16	334	3	263n
Oes10	16	403	3	298n
Scm20d	16	8966	4	61n
Scm1d	16	9803	4	280n
Adult	4	18419	7, 7, 5, 2	5n, 5x
Default	4	28779	2, 7, 4, 2	14n, 6x
Thyroid	7	9172	5, 5, 3, 2, 4, 4, 3	7n, 22x

**Table 2:** Statistics of the tabular benchmark data sets.

CV	States
Person	no person, person
Animal	no animals, bird, cat, cow, dog, horse, sheep
Vehicle	no vehicles, aeroplane, bicycle, boat, bus, car, motorbike, train
Indoor	no indoor objects, bottle, chair, dining table, potted plant, sofa, tv/monitor

- ▶ **9963** natural images.
  - ▶ **50%** for training and **50%** for test.
- ▶ **4** class variables with different numbers of states.
- ▶ Remove images that take multiple values for a class variable
  - ▶ E.g., an image that contains both cats and dogs.
- ▶ Resize to  $256 \times 256$ .

**Table 3:** Characteristics of the PASCAL VOC 2007 data set in an MDC setting.



**Figure 5:** Example images from the PASCAL VOC 2007 data set.

Given a test set  $\mathcal{T} = \{(\mathbf{x}^l, \mathbf{y}^l) \mid 1 \leq l \leq L\}$  and an MDC classifier  $h : \mathcal{X} \rightarrow \mathcal{Y}$ :

**Exact Match Score (EMS)**, which is equivalent to the subset 0/1 loss  $\ell_{0/1}$ .

$$\begin{aligned} \text{EMS}(h) &= \frac{1}{L} \sum_{l=1}^L \mathbb{1}_{\mathbf{y}^l = \hat{\mathbf{y}}^l} \\ &= \frac{1}{L} \sum_{l=1}^L (1 - \ell_{0/1}(\mathbf{y}^l, \hat{\mathbf{y}}^l)) \end{aligned} \tag{9}$$

**Hamming Score (HS)**, which is equivalent to the Hamming loss  $\ell_{\text{HS}}$ .

$$\begin{aligned} \text{HS}(h) &= \frac{1}{L} \sum_{l=1}^L \frac{1}{d} \sum_{j=1}^d \mathbb{1}_{y_j^l = \hat{y}_j^l} \\ &= \frac{1}{L} \sum_{l=1}^L (1 - \ell_{\text{HL}}(\mathbf{y}^l, \hat{\mathbf{y}}^l)) \end{aligned} \tag{10}$$

## Binary Relevance (BR)

- ▶ Assume all class variables are **independent** of each other.
- ▶ Train a base classifier for each class variable.

## Class Powerset (CP)

- ▶ Transform each unique combination of class values in the training set into a new **meta** class.
- ▶ Train an overall MCC classifier on the transformed meta classes.

## Classifier Chain (CC)

- ▶ Represent the class dependencies using a **chain-like structure**.
- ▶ Train a base classifier for each class variable by incorporating previous predictions in the chain.

# Experimental Results: Tabular Data #1

Data Set	Exact Match Score (EMS)			
	GBNC ( $\ell_{0/1}$ )	BR	CC	CP
Edm	<b>0.570 ± 0.126</b>	0.475 ± 0.129	0.469 ± 0.126	0.450 ± 0.141 ↑
Jura	<b>0.415 ± 0.054</b>	0.409 ± 0.066	0.368 ± 0.114	0.014 ± 0.042 ↑
Enb	0.561 ± 0.061	<b>0.574 ± 0.090</b>	0.525 ± 0.041	0.371 ± 0.065 ↑
Voice	<b>0.858 ± 0.024</b>	0.836 ± 0.031 ↑	0.837 ± 0.033 ↑	0.187 ± 0.024 ↑
Song	<b>0.438 ± 0.049</b>	0.422 ± 0.041	0.392 ± 0.061 ↑	0.068 ± 0.035 ↑
Flickr	0.289 ± 0.013	0.324 ± 0.012 ↓	<b>0.325 ± 0.013 ↓</b>	0.042 ± 0.006 ↑
Fera	<b>0.196 ± 0.012</b>	0.187 ± 0.014	0.193 ± 0.011	0.164 ± 0.015 ↑
WQplants	0.093 ± 0.034	0.093 ± 0.026	<b>0.095 ± 0.029</b>	0.075 ± 0.033 ↑
WQanimals	0.046 ± 0.008	0.047 ± 0.020	<b>0.057 ± 0.015 ↓</b>	0.023 ± 0.016 ↑
Rf1	<b>0.532 ± 0.018</b>	0.283 ± 0.020 ↑	0.291 ± 0.018 ↑	0.062 ± 0.009 ↑
Pain	<b>0.758 ± 0.018</b>	0.751 ± 0.015 ↑	0.754 ± 0.017 ↑	0.751 ± 0.015 ↑
Disfa	<b>0.401 ± 0.011</b>	0.393 ± 0.012 ↑	0.394 ± 0.013 ↑	0.371 ± 0.011 ↑
WaterQuality	0.005 ± 0.005	<b>0.007 ± 0.006</b>	<b>0.007 ± 0.006</b>	0.006 ± 0.006
Oes97	<b>0.057 ± 0.044</b>	0.051 ± 0.046	0.030 ± 0.014 ↑	0.000 ↑
Oes10	0.087 ± 0.041	0.089 ± 0.032	<b>0.092 ± 0.034</b>	0.005 ± 0.010 ↑
Scm20d	<b>0.124 ± 0.014</b>	0.045 ± 0.008 ↑	0.062 ± 0.011 ↑	0.076 ± 0.012 ↑
Scm1d	<b>0.189 ± 0.010</b>	0.098 ± 0.009 ↑	0.109 ± 0.008 ↑	0.091 ± 0.009 ↑
No. of Wins	11	2	5	0

**Table 4:** Exact match scores (mean ± std.) of each MDC approach (base classifier: *logisitic regression*). GBNC performs inference by optimizing  $\ell_{0/1}$ . The best performance is highlighted in bold, and ↑ /↓ indicates whether GBNC is significantly superior/inferior to other approaches on each data set by using a Wilcoxon signed-rank test.

## Experimental Results: Tabular Data #2

Data Set	Hamming Score (HS)			
	GBNC ( $\ell_{HL}$ )	BR	CC	CP
Edm	<b>0.736 ± 0.091</b>	0.696 ± 0.093	0.685 ± 0.099	0.725 ± 0.071
Jura	<b>0.634 ± 0.029</b>	0.625 ± 0.043	0.607 ± 0.075	0.317 ± 0.051 ↑
Enb	0.781 ± 0.031	<b>0.787 ± 0.045</b>	0.762 ± 0.020	0.685 ± 0.033 ↑
Voice	<b>0.927 ± 0.013</b>	0.915 ± 0.015 ↑	0.916 ± 0.017 ↑	0.584 ± 0.013 ↑
Song	<b>0.766 ± 0.028</b>	0.752 ± 0.023	0.738 ± 0.039 ↑	0.507 ± 0.040 ↑
Flickr	0.784 ± 0.006	<b>0.797 ± 0.006 ↓</b>	0.796 ± 0.006 ↓	0.506 ± 0.004 ↑
Fera	<b>0.624 ± 0.010</b>	0.616 ± 0.007 ↑	0.605 ± 0.009 ↑	0.475 ± 0.018 ↑
WQplants	<b>0.658 ± 0.013</b>	0.655 ± 0.010	0.650 ± 0.014	0.611 ± 0.024 ↑
WQanimals	<b>0.630 ± 0.018</b>	0.626 ± 0.019	0.624 ± 0.021 ↑	0.579 ± 0.024 ↑
Rf1	<b>0.902 ± 0.008</b>	0.836 ± 0.004 ↑	0.835 ± 0.006 ↑	0.635 ± 0.007 ↑
Pain	<b>0.953 ± 0.003</b>	<b>0.953 ± 0.003</b>	0.951 ± 0.003 ↑	0.948 ± 0.003 ↑
Disfa	<b>0.897 ± 0.002</b>	0.894 ± 0.003 ↑	0.894 ± 0.002 ↑	0.871 ± 0.003 ↑
WaterQuality	<b>0.642 ± 0.018</b>	0.637 ± 0.011	0.639 ± 0.017	0.597 ± 0.018 ↑
Oes97	<b>0.731 ± 0.023</b>	0.716 ± 0.018 ↑	0.706 ± 0.030 ↑	0.521 ± 0.032 ↑
Oes10	<b>0.809 ± 0.014</b>	0.801 ± 0.019	0.791 ± 0.021 ↑	0.605 ± 0.046 ↑
Scm20d	<b>0.685 ± 0.007</b>	0.640 ± 0.008 ↑	0.613 ± 0.011 ↑	0.424 ± 0.012 ↑
Scm1d	<b>0.815 ± 0.003</b>	0.763 ± 0.006 ↑	0.748 ± 0.005 ↑	0.444 ± 0.011 ↑
No. of Wins	15	3	0	0

**Table 5:** Hamming scores (mean ± std.) of each MDC approach (base classifier: *logistic regression*). GBNC performs inference by optimizing  $\ell_{HL}$ . The best performance is highlighted in bold, and ↑ /↓ indicates whether GBNC is significantly superior/inferior to other approaches on each data set by using a *Wilcoxon signed-rank test*.

# Experimental Results: Mixed Data

Data Set	Exact Match Score (EMS)			
	GBNC ( $\ell_{0/1}$ )	BR	CC	CP
Adult	$0.245 \pm 0.006$	<b><math>0.274 \pm 0.011 \downarrow</math></b>	<b><math>0.274 \pm 0.013 \downarrow</math></b>	$0.134 \pm 0.007 \uparrow$
Default	<b><math>0.187 \pm 0.005</math></b>	$0.177 \pm 0.009 \uparrow$	$0.177 \pm 0.012 \uparrow$	$0.060 \pm 0.004 \uparrow$
Thyroid	<b><math>0.784 \pm 0.018</math></b>	$0.774 \pm 0.014 \uparrow$	$0.768 \pm 0.016 \uparrow$	$0.751 \pm 0.016 \uparrow$

**Table 6:** Exact match scores (mean  $\pm$  std.) of each MDC approach (base classifier: *logistic regression*) on mixed data. GBNC performs inference by optimizing  $\ell_{0/1}$ .

Data Set	Hamming Score (HS)			
	GBNC ( $\ell_{HL}$ )	BR	CC	CP
Adult	$0.676 \pm 0.004$	<b><math>0.718 \pm 0.005 \downarrow</math></b>	$0.715 \pm 0.005 \downarrow$	$0.585 \pm 0.006 \uparrow$
Default	<b><math>0.666 \pm 0.004</math></b>	$0.664 \pm 0.007$	$0.663 \pm 0.009$	$0.561 \pm 0.004 \uparrow$
Thyroid	<b><math>0.966 \pm 0.003</math></b>	$0.965 \pm 0.002 \uparrow$	$0.964 \pm 0.003 \uparrow$	$0.962 \pm 0.003 \uparrow$

**Table 7:** Hamming scores (mean  $\pm$  std.) of each MDC approach (base classifier: *logistic regression*) on mixed data. GBNC performs inference by optimizing  $\ell_{HL}$ .

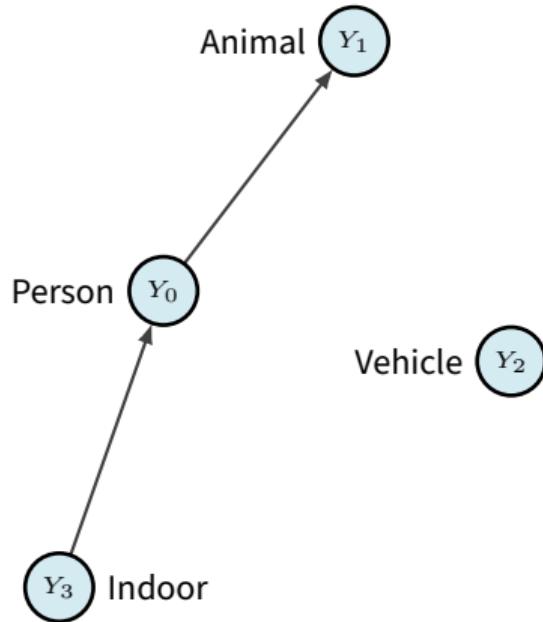
- ▶ 5 continuous and 5 discrete feature variables in the Adult data set.

Metrics	Methods			
	GBNC ( $\ell_{0/1}$ )	GBNC ( $\ell_{HL}$ )	BR	CP
HS	<b>0.897</b>	<b>0.897</b>	0.891	0.887
EMS	<b>0.662</b>	<b>0.662</b>	0.604	0.655

**Table 8:** Results of each MDC approach (base classifier: *ResNet-18*) for the PASCAL VOC 2007 data set.

- ▶ Sparse BN structure since there are only 4 class variables.
- ▶ Extreme probability estimates due to uncalibrated neural networks.

# Analysis of Learned BN Structures



**Figure 6:** Learned BN structure of GBNC on the PASCAL VOC 2007 data set. The base classifier is ResNet-18 with weights pre-trained on ImageNet.



**Figure 7:** Example images from the PASCAL VOC 2007 data set.

- ▶ The Vehicle dimension appear to be independent.
- ▶ GBNC tends to consider the mutual relationships between Indoor, Person and Animal.

## Conclusion

Proposed GBNC, a generalized framework for solving probabilistic MDC tasks with complex types of input.

- ▶ The **first** approach that exactly optimizes the CLL function in learning MBNCs.
- ▶ Introduced a new structural constraint for learning (multi-dimensional) BNCs, which enables the **decomposability** of the CLL function over a BN structure.
- ▶ Proposed a input space partitioning algorithm to **discriminatively** learn BNC structures and parameters simultaneously, in which the CLL function is optimized.
- ▶ GBNC converts the prediction problem into an inference problem in the learned class BN structure, which allows computing the **Bayes-Optimal Prediction** under different loss functions.

Proposed GBNC, a generalized framework for solving probabilistic MDC tasks with complex types of input.

- ▶ The **first** approach that exactly optimizes the CLL function in learning MBNCs.
- ▶ Introduced a new structural constraint for learning (multi-dimensional) BNCs, which enables the **decomposability** of the CLL function over a BN structure.
- ▶ Proposed a input space partitioning algorithm to **discriminatively** learn BNC structures and parameters simultaneously, in which the CLL function is optimized.
- ▶ GBNC converts the prediction problem into an inference problem in the learned class BN structure, which allows computing the **Bayes-Optimal Prediction** under different loss functions.
- ▶ By employing the pruning rules from [2] and using GOBNILP [1], GBNC is able to **exactly and efficiently** learn class BN structures with the **penalized CLL** as a structure score function.
- ▶ By using the same partitioning idea, GBNC is able to handle **mixed data**.
- ▶ GBNC achieves **leading performance** among other discriminative MDC approaches.

- ▶ Structure complexity from feature variables.
- ▶ Compute Bayes-optimal prediction under more loss functions, such as the Brier score and the AUC score.
- ▶ Further utilize the discrete features other than input space partitioning.
- ▶ More efficient inference algorithms.

Thank You!  
Questions?

---

**Algorithm 2** Extract training data

---

**Require:** Training data  $\mathcal{D} = \{(\mathbf{x}^l, \mathbf{y}^l) \mid 1 \leq l \leq N\}$ , class variable  $Y_j$  ( $1 \leq j \leq d$ ), parent set  $\Pi_j$ , parent configuration  $\pi_{jk}$ , parent set  $\Phi_j$

**Ensure:**  $\mathcal{D}_j^{\pi_{jk}}$

- 1: Initialize  $\mathcal{D}_j^{\pi_{jk}} \leftarrow \emptyset$
  - 2: **for**  $l = 1, 2, \dots, N$  **do**
  - 3:   **if**  $\pi_j^l = \pi_{jk}$  **then**
  - 4:     Update  $\mathcal{D}_j^{\pi_{jk}} \leftarrow \mathcal{D}_j^{\pi_{jk}} \cup \{(\phi_j^l, y_j)\}$
  - 5:   **end if**
  - 6: **end for**
-

---

**Algorithm 3** Train local probabilistic model

---

**Require:** Training data  $\mathcal{D} = \{(\mathbf{x}^l, \mathbf{y}^l) \mid 1 \leq l \leq N\}$ , class variable  $Y_j$  ( $1 \leq j \leq d$ ), parent set  $\Pi$ , parent configuration  $\pi$ , parent set  $\Phi$  ( $\Phi = \mathbf{X}$  by default if not specified)

**Ensure:** Local classifier  $C_j^\pi$

- 1: Extract  $\mathcal{D}_j^\pi$  using 2 with input  $\mathcal{D}, Y_j, \Pi, \pi$  and  $\Phi$
  - 2: Train  $C_j^\pi$  using  $\mathcal{D}_j^\pi$  by maximizing the CLL
-

---

**Algorithm 4** Parameter learning in  $\mathcal{G}$ 

**Require:** Training data  $\mathcal{D} = \{(\mathbf{x}^l, \mathbf{y}^l) \mid 1 \leq l \leq N\}$ , BN structure  $\mathcal{G}$

**Ensure:** Local classifiers  $\mathbf{C}$

- 1: Initialize  $\mathbf{C} \leftarrow \emptyset$
  - 2: **for**  $j = 1, 2, \dots, d$  **do**
  - 3:   Extract  $\boldsymbol{\Pi}_j$  and  $\boldsymbol{\Phi}_j$  from  $\mathcal{G}$
  - 4:   **for**  $k = 1, 2, \dots, q_j$  **do**
  - 5:     Train local classifier  $C_j^{\pi_{jk}}$  using 3 with input  $\mathcal{D}, Y_j, \boldsymbol{\Pi}_j, \pi_{jk}$  and  $\boldsymbol{\Phi}_j$
  - 6:     Update  $\mathbf{C} \leftarrow \mathbf{C} \cup C_j^{\pi_{jk}}$
  - 7:   **end for**
  - 8: **end for**
-

# Structure Learning Algorithm

---

**Algorithm 5** Structure learning of  $\mathcal{G}_Y$ 

```
Require: Training data  $\mathcal{D} = \{(\mathbf{x}^l, \mathbf{y}^l) \mid 1 \leq l \leq N\}$ 
Ensure: Class subgraph  $\mathcal{G}_Y$ 
1: Initialize a score dict  $\mathbf{S}$   $\triangleright$  where  $\mathbf{S}[j][\Pi]$  stores the local score for  $Y_j$  given the parent set  $\Pi$ 
2: for  $j = 1, 2, \dots, d$  do
3:   Initialize candidate parent sets  $\mathbf{K}_j$  of  $Y_j$  as the powerset of  $\mathbf{Y} \setminus Y_j$ 
4:   for  $u = 1, 2, \dots, |\mathbf{K}_j|$  do
5:     Compute  $\text{PEN}(\mathcal{G}_{\mathbf{K}_{ju}}) \leftarrow -\frac{\log N}{2}(r_j - 1)q_j$   $\triangleright$  here  $q_j = \prod_{c=1: Y_c \in \mathbf{K}_{ju}}^d r_c$ 
6:     if  $\mathbf{K}_{ju}$  can be pruned then  $\triangleright$  using Lemma 2 in the thesis
7:       Remove all supersets of  $\mathbf{K}_{ju}$  from  $\mathbf{K}$ 
8:       continue
9:     end if
10:    Initialize  $\mathbf{S} \leftarrow 0$ 
11:    for  $k = 1, 2, \dots, q_j$  do
12:      Extract  $\mathcal{D}_j^{juk}$  using 2 with input  $\mathcal{D}, Y_j, \Pi_j = \mathbf{K}_{ju}, \pi_j = \mathbf{k}_{juk}$  and  $\Phi = \mathbf{X}$ 
13:      Train local classifier  $C_j^{juk}$  using 3 with input  $\mathcal{D}, Y_j, \Pi_j = \mathbf{K}_{ju}, \pi_j = \mathbf{k}_{juk}$  and  $\Phi_j = \mathbf{X}$ 
14:      Update  $\mathbf{S} \leftarrow \mathbf{S} + C_j^{juk}(\mathcal{D}_j^{juk})$ 
15:    end for
16:    if  $\mathbf{K}_{ju}$  can be pruned then  $\triangleright$  using Lemma 1 in the thesis
17:      continue
18:    end if
19:     $\mathbf{S}[j][\mathbf{K}_{ju}] \leftarrow \mathbf{S}$ 
20:  end for
21: end for
22: Learn  $\mathcal{G}_Y$  using GOBNILP with input  $\mathbf{S}$ 
```

---

-  J. Cussens.  
Gobnilp: Learning bayesian network structure with integer programming.  
In *International Conference on Probabilistic Graphical Models*, pages 605–608. PMLR, 2020.
-  C. P. de Campos, M. Scanagatta, G. Corani, and M. Zaffalon.  
Entropy-based pruning for learning bayesian networks using BIC.  
*Artif. Intell.*, 260:42–50, 2018.
-  B.-B. Jia and M.-L. Zhang.  
Decomposition-based classifier chains for multi-dimensional classification.  
*IEEE Transactions on Artificial Intelligence*, 3(2):176–191, 2021.