Supervised multivariate discretization and levels merging for logistic regression

Adrien Ehrhardt<sup>1,2</sup>

Christophe Biernacki<sup>2</sup> Philippe Heinrich<sup>3,4</sup> Vincent Vandewalle<sup>2,3</sup>

<sup>1</sup>Crédit Agricole Consumer Finance <sup>2</sup>Inria Lille - Nord-Europe <sup>3</sup>Université de Lille <sup>4</sup>CNRS

31/08/2018





▲□▶ ▲⑦

クヘペ 1/34

Context and basic notations

Supervised multivariate discretization and factor levels grouping

▲□▶ ▲□▶ ▲ 三▶ ▲ 三 ▶ ● 三 ∽ ९ ℃ 2/34

Interactions in logistic regression

Conclusion and future work

## Context and basic notations

E

.≣ ►

(日) (同) (三) (

৩৭৫ 3/34

Home	Time in job	Family status	Wages		Repayment
Owner	20	Widower	2000		0
Renter	10	Common-law	1700		0
Starter	5	Divorced	4000		1
By work	8	Single	2700		1
Renter	12	Married	1400		0
By family	2		1200		0
	Home Owner Renter Starter By work Renter By family	HomeTime in jobOwner20Renter10Starter5By work8Renter12By family2	HomeTime in jobFamily statusOwner20WidowerRenter10Common-lawStarter5DivorcedBy work8SingleRenter12MarriedBy family2?	HomeTime in jobFamily statusWagesOwner20Widower2000Renter10Common-law1700Starter5Divorced4000By work8Single2700Renter12Married1400By family2?1200	HomeTime in jobFamily statusWagesOwner20Widower2000Renter10Common-law1700Starter5Divorced4000By work8Single2700Renter12Married1400By family2?1200

Table: Dataset with outliers and missing values.

< □ > < 同

▶ ◀ ≣ ▶ ◀

E

4/34

Job	Home	Time in job	Family status	Wages	Repayment
Craftsman	Owner	20	Widower	2000	0
?	Renter	10	Common-law	1700	0
Licensed profes- sional	Starter	5	Divorced	4000	1
Executive	By work	8	Single	2700	1
Office employee	Renter	12	Married	1400	0
Worker	By family	2	?	1200	0

Table: Dataset with outliers and missing values.

< □ ▶ < 同

- 1. Feature selection
- 2. Discretization / grouping
- 3. Interaction screening
- 4. Logistic regression fitting

		Wages	Repayment
Craftsman	Widower	2000	0
		1700	0
	Divorced	4000	1
		2700	1
Office employee		1400	0
		1200	0

Table: Dataset with outliers and missing values.

< □ > < 同

- 1. Feature selection
- 2. Discretization / grouping
- 3. Interaction screening
- 4. Logistic regression fitting

Job	Family status	Wages	Repayment
Craftsman	Widower	]1500;2000]	0
	Common-law	]1500;2000]	0
Licensed profes- sional	Divorced	]2000;∞[	1
Executive	Single	]2000;∞[	1
Office employee	Married	]-∞ ; <b>1500</b> ]	0
Worker		]-∞ ; <b>1500</b> ]	0

Table: Dataset with outliers and missing values.

< □ > < 同

- 1. Feature selection
- 2. Discretization / grouping
- 3. Interaction screening
- 4. Logistic regression fitting

		Wages	Repayment
?+Low-qualified	?+Alone	]1500;2000]	0
		]1500;2000]	0
		]2000;∞[	1
		]2000;∞[	1
		]- $\infty$ ; 1500]	0
		]- $\infty$ ; 1500]	o

Table: Dataset with outliers and missing values.

< □ > < 同

- 1. Feature selection
- 2. Discretization / grouping
- 3. Interaction screening
- 4. Logistic regression fitting

Job		Family status x Wages	Repayment
?+Low-qualified		?+Alone × ]1500;2000]	0
?+Low-qualified		Union × ]1500;2000]	0
High-qualified		?+Alone × ]2000;∞[	1
High-qualified		?+Alone x ]2000;∞[	1
?+Low-qualified		Union x ]- $\infty$ ; 1500]	0
?+Low-qualified		?+Alone x ]- $\infty$ ; 1500]	0

Table: Dataset with outliers and missing values.

< □ > < 同

- 1. Feature selection
- 2. Discretization / grouping
- 3. Interaction screening
- 4. Logistic regression fitting

Job		Family status × Wages		Repayment
?+Low-qualified		?+Alone × ]1500;2000]	225	0
?+Low-qualified		Union × ]1500;2000]		0
High-qualified		?+Alone x ]2000; $\infty$ [		1
High-qualified		?+Alone x ]2000; $\infty$ [		1
?+Low-qualified		Union x ]- $\infty$ ; 1500]		0
?+Low-qualified		?+Alone x ]- $\infty$ ; 1500]		0

Table: Dataset with outliers and missing values.

< □ > < 同

- 1. Feature selection
- 2. Discretization / grouping
- 3. Interaction screening
- 4. Logistic regression fitting

The whole process can be decomposed into two steps:

$$egin{aligned} \mathcal{X} & o \mathcal{E} & o \mathcal{Y} \ \mathbf{x} &\mapsto \mathbf{e} = \mathbf{f}(\mathbf{x}) \mapsto y \end{aligned}$$

The whole process can be decomposed into two steps:

$$egin{aligned} \mathcal{X} & o \mathcal{E} & o \mathcal{Y} \ \mathbf{x} &\mapsto \mathbf{e} = \mathbf{f}(\mathbf{x}) \mapsto y \end{aligned}$$

• □ ▶ • 一冊 ▶ •

.≣ ▶ ∢

JAC.

Selected features: 
$$\mathbf{x} = (\underbrace{(x_j)_1^{d_1}}_{\in \mathbb{R}}, \underbrace{(x_j)_{d_1+1}^d}_{\in \{1, \dots, o_i\}}).$$

The whole process can be decomposed into two steps:

$$egin{aligned} \mathcal{X} & o \mathcal{E} & o \mathcal{Y} \ \mathbf{x} &\mapsto \mathbf{e} = \mathbf{f}(\mathbf{x}) \mapsto y \end{aligned}$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三 のへで

5/34

Selected features: 
$$\mathbf{x} = (\underbrace{(x_j)_1^{d_1}}_{\in \mathbb{R}}, \underbrace{(x_j)_{d_1+1}^d}_{\in \{1, \dots, o_j\}}).$$

f must be "simple" and "component-wise", *i.e.*  $f = (f_j)_1^d$ . We restrict to discretization and grouping of factor levels.

< □ ▶ < @ ▶ < \overline ₽ ♡ < \cap 6/34

$$f_j(x_j) = 1 \qquad f_j(x_j) = 2 \qquad f_j(x_j) = 3 \qquad \Rightarrow x_j$$

Discretization  $(1 \le j \le d_1)$ 

Into *m* intervals with associated cutpoints  $\boldsymbol{c} = (c_1, \ldots, c_{m-1})$ .

#### Discretization function

$$f_{j}(\cdot; \boldsymbol{c}, m) \colon \mathbb{R} \to \{1, \dots, m\}$$

$$x \mapsto \mathbb{1}_{]-\infty; c_{1}]}(x) + \sum_{k=1}^{m-2} (k+1) \,\mathbb{1}_{]c_{k}; c_{k+1}]}(x)$$

$$+ m \,\mathbb{1}_{]c_{m-1}, \infty[}(x)$$

596







Grouping  $(d_1 < j \le d)$ 

Grouping o values into  $m, m \leq o$ .

#### Grouping function

$$f_j \colon \{1,\ldots,o\} o \{1,\ldots,m\}$$

 $f_j$  surjective: it defines a partition of  $\{1, \ldots, o\}$  in *m* elements.

୬ ୯୯ 7/34

Target feature  $y \in \{0, 1\}$  must be predicted given engineered features  $f(x) = (f_j(x_j))_1^d$ .

↓□▶ < □▶ < 三▶ < 三▶ < 三 < つへで 8/34</p>

Target feature  $y \in \{0, 1\}$  must be predicted given engineered features  $f(x) = (f_j(x_j))_1^d$ .

▲□▶ ▲□▶ ▲壹▶ ▲壹▶ 壹 ∽९ペ 8/34

We restrict to binary logistic regression.

Target feature  $y \in \{0, 1\}$  must be predicted given engineered features  $f(x) = (f_j(x_j))_1^d$ .

We restrict to binary logistic regression.

On "raw" data, logistic regression yields:

$$\mathsf{logit}(p_{m{ heta}_{\mathsf{raw}}}(1|m{x})) = heta_0 + \sum_{j=1}^{d} heta_j x_j + \sum_{j=d_1+1}^{d} heta_j^{x_j}$$

Target feature  $y \in \{0, 1\}$  must be predicted given engineered features  $f(x) = (f_j(x_j))_1^d$ .

We restrict to binary logistic regression.

On "raw" data, logistic regression yields:

$$\mathsf{logit}(p_{m{ heta}_{\mathsf{raw}}}(1|m{x})) = heta_0 + \sum_{j=1}^{d_1} heta_j x_j + \sum_{j=d_1+1}^{d} heta_j^{x_j}$$

On discretized / grouped data, logistic regression yields:

$$\mathsf{logit}(p_{m{ heta}_f}(1|m{f}(m{x}))) = heta_0 + \sum_{j=1}^d heta_j^{f_j(x_j)}$$

 $\mathcal{O} \mathcal{Q} \mathcal{O}$ 

True data

$$\mathsf{logit}(p_{\mathsf{true}}(1|\boldsymbol{x})) = \mathsf{ln}\left(\frac{p_{\mathsf{true}}(1|\boldsymbol{x})}{1 - p_{\mathsf{true}}(1|\boldsymbol{x})}\right) = \mathsf{sin}((x_1 - 0.7) \times 7)$$



Figure: True relationship between predictor and outcome

< □

Logistic regression on "raw" data:

 $\operatorname{logit}(p_{\overline{ heta_{raw}}}(1|m{x})) = heta_0 + \overline{ heta_1 x_1}$ 



Figure: Linear logistic regression fit

< D >

 $\mathcal{O}\mathcal{A}\mathcal{O}$ 

**Logistic regression on discretized data:** If *f* is not carefully chosen ...

$$\operatorname{logit}(p_{\theta_f}(1|\boldsymbol{f}(\boldsymbol{x}))) = \theta_0 + \underbrace{\theta_1^{f_1(\boldsymbol{x}_1)}}_{\theta_1^1, \dots, \theta_1^{f_0}}$$



Figure: Bad (high variance) discretization

< □

#### Logistic regression on discretized data: If *f* is carefully chosen ...

$$\mathsf{logit}(p_{m{ heta}_f}(1|m{f}(m{x}))) = heta_0 + \underbrace{m{ heta}_1^{m{f}_1(m{x}_1)}}_{ heta_1^1,\dots,m{ heta}_1^1}$$



Figure: Good (bias/variance tradeoff) discretization

< D >

# $\theta$ can be estimated for each discretization f and $f^*$ can be chosen through our favorite model choice criterion: BIC, AIC, ...

# $\theta$ can be estimated for each discretization f and $f^*$ can be chosen through our favorite model choice criterion: BIC, AIC, ...

### A model selection problem

$$(\boldsymbol{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmax}_{\boldsymbol{f} \in \mathcal{F}, \boldsymbol{\theta} \in \Theta_{\boldsymbol{f}}} \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_i | \boldsymbol{f}(\boldsymbol{x}_i)) - \operatorname{penalty}(n; \boldsymbol{\theta})$$

# $\theta$ can be estimated for each discretization f and $f^*$ can be chosen through our favorite model choice criterion: BIC, AIC, ...

### A model selection problem

$$(\boldsymbol{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmax}_{\boldsymbol{f} \in \mathcal{F}, \boldsymbol{\theta} \in \Theta_{\boldsymbol{f}}} \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_i | \boldsymbol{f}(\boldsymbol{x}_i)) - \operatorname{penalty}(n; \boldsymbol{\theta})$$

How to efficiently explore  $\mathcal{F}$ ?



"Functional" space  $\mathcal{F}$  where f lives is continuous:



#### "Functional" space $\mathcal{F}$ where f lives is continuous:

$$f_j(x_j) = 1 \qquad f_j(x_j) = 2 \qquad f_j(x_j) = m_j$$

<□ > < 母 > < 壹 > < 亘 > < 重 > う < ♡ 11/34

"Functional" space  $\mathcal{F}$  where f lives is continuous:

$$f_j(x_j) = 1 \qquad f_j(x_j) = 2 \qquad f_j(x_j) = m_j$$

However, for a fixed design  $\mathbf{x} = (\mathbf{x}_i)_1^n$  there is a countable space  $\tilde{\mathcal{F}}$  in which  $\mathbf{f}\mathcal{R}\mathbf{g} \Leftrightarrow \forall i, j, f_j(x_i) = g_j(x_i)$ 

◆□ ▶ ▲□ ▶ ▲ 三 ▶ ▲ 三 ♪ り へ ( 11/34

# Exploring $\mathcal{F}$

#### Example of discretization

"Functional" space  $\mathcal{F}$  where f lives is continuous:

$$f_j(x_j) = 1 \qquad f_j(x_j) = 2 \qquad f_j(x_j) = m_j$$

However, for a fixed design  $\mathbf{x} = (\mathbf{x}_i)_1^n$  there is a countable space  $\tilde{\mathcal{F}}$  in which  $\mathbf{f}\mathcal{R}\mathbf{g} \Leftrightarrow \forall i, j, f_j(x_i) = g_j(x_i)$ 

$$f_j(x_j) = 1 \qquad f_j(x_j) = 2$$

# Exploring $\mathcal{F}$

#### Example of discretization

"Functional" space  $\mathcal{F}$  where f lives is continuous:

$$f_j(x_j) = 1 \qquad f_j(x_j) = 2 \qquad f_j(x_j) = m_j$$

However, for a fixed design  $\mathbf{x} = (\mathbf{x}_i)_1^n$  there is a countable space  $\tilde{\mathcal{F}}$  in which  $\mathbf{f}\mathcal{R}\mathbf{g} \Leftrightarrow \forall i, j, f_j(x_i) = g_j(x_i)$ 

$$g_j(x_j) = 1 \qquad g_j(x_j) = 2$$

# Exploring $\mathcal{F}$

#### Example of discretization

"Functional" space  $\mathcal{F}$  where f lives is continuous:

$$f_j(x_j) = 1 \qquad f_j(x_j) = 2 \qquad f_j(x_j) = m_j$$

However, for a fixed design  $\mathbf{x} = (\mathbf{x}_i)_1^n$  there is a countable space  $\tilde{\mathcal{F}}$  in which  $\mathbf{f}\mathcal{R}\mathbf{g} \Leftrightarrow \forall i, j, f_j(x_i) = g_j(x_i)$ 

"Functional" space  $\mathcal{F}$  where f lives is continuous:

$$f_j(x_j) = 1 \qquad f_j(x_j) = 2 \qquad f_j(x_j) = m_j$$

However, for a fixed design  $\mathbf{x} = (\mathbf{x}_i)_1^n$  there is a countable space  $\tilde{\mathcal{F}}$  in which  $\mathbf{f}\mathcal{R}\mathbf{g} \Leftrightarrow \forall i, j, f_j(x_i) = g_j(x_i)$ 

$$(\boldsymbol{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmax}_{\boldsymbol{f} \in \tilde{\mathcal{F}}, \boldsymbol{\theta} \in \Theta_{\boldsymbol{f}}} \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_i | \boldsymbol{f}(\boldsymbol{x}_i)) - \operatorname{penalty}(n; \boldsymbol{\theta})$$

#### Current academic methods:

A lot of existing heuristics, see [Ramírez-Gallego et al., 2016]:


Most of these methods are:

<□><</li>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□>
<□></li

► Univariate,

Most of these methods are:

- Univariate,
- Test statistics more or less justified ( $\chi^2$ -based).

<□ > < @ > < \ > < \ > < \ > < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? <\ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? < \ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ? <\ ?

# Supervised multivariate discretization and factor levels grouping

<u> イロト</u> イタト イミト イ

⊒ ►

∃ つへで 14/34

Discretized / grouped  $x_j$  denoted by  $e_j$  has been seen up to now as the result of a function of  $x_j$ :

 $e_j = f_j(x_j).$ 

▲□▶ ▲□▶ ▲ 글▶ ▲ 글▶ 글 ∽ ९ ℃ 15/34

Discretized / grouped  $x_j$  denoted by  $e_j$  has been seen up to now as the result of a function of  $x_j$ :

 $e_j = f_j(x_j).$ 

Discretization / grouping  $e_j$  can be seen as a latent random variable for which

 $p(e_j|x_j) = \mathbb{1}_{e_j}(f_j(x_j)).$ 



Discretized / grouped  $x_j$  denoted by  $e_j$  has been seen up to now as the result of a function of  $x_j$ :

 $e_j = f_j(x_j).$ 

Discretization / grouping  $e_j$  can be seen as a latent random variable for which

$$p(e_j|x_j) = \underbrace{\mathbb{1}_{e_j}(f_j(x_j))}$$

Heaviside-like function difficult to optimize

<□▶ < □▶ < 亘▶ < 亘▶ < 亘▶ Ξ りへで 15/34

Discretized / grouped  $x_j$  denoted by  $e_j$  has been seen up to now as the result of a function of  $x_j$ :

 $e_j = f_j(x_j).$ 

Discretization / grouping  $e_j$  can be seen as a latent random variable for which

$$p(e_j|x_j) = \underbrace{\mathbb{1}_{e_j}(f_j(x_j))}$$

Heaviside-like function difficult to optimize

▲□▶ ▲□▶ ▲壹▶ ▲壹▶ 壹 - ∽�や 15/34

Suppose for now that  $\boldsymbol{m} = (m_j)_1^d$  is fixed.

Discretized / grouped  $x_j$  denoted by  $e_j$  has been seen up to now as the result of a function of  $x_j$ :

 $e_j = f_j(x_j).$ 

Discretization / grouping  $e_j$  can be seen as a latent random variable for which

Heaviside-like function difficult to optimize

Suppose for now that  $\boldsymbol{m} = (m_j)_1^d$  is fixed.

$$\boldsymbol{e} \in \boldsymbol{\mathcal{E}}_{\boldsymbol{m}} = \{1, \ldots, m_1\} \times \ldots \times \ldots \times \{1, \ldots, m_d\}.$$

#### Model selection criterion

We want the "best" model  $p_{\theta^*}(y|e^*)$  where  $\theta^*$  is the maximum likelihood estimator and  $e^*$  is determined by AIC, BIC...

$$(\boldsymbol{e}^{\star}, \boldsymbol{ heta}^{\star}) = \operatorname*{argmax}_{\boldsymbol{e} \in \boldsymbol{\mathcal{E}}_{\boldsymbol{m}}, \boldsymbol{ heta} \in \boldsymbol{\Theta}_{\boldsymbol{m}}} \sum_{i=1}^{n} \ln p_{\boldsymbol{ heta}}(y_i | \boldsymbol{e}_i) - \operatorname{penalty}(n; \boldsymbol{ heta})$$

16/34

#### Model selection criterion

We want the "best" model  $p_{\theta^*}(y|e^*)$  where  $\theta^*$  is the maximum likelihood estimator and  $e^*$  is determined by AIC, BIC...

$$(\boldsymbol{e}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmax}_{\boldsymbol{e} \in \boldsymbol{\mathcal{E}}_{m}, \boldsymbol{\theta} \in \Theta_{m}} \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_{i} | \boldsymbol{e}_{i}) - \operatorname{penalty}(n; \boldsymbol{\theta})$$

୬ ୯୦ 16/34

 $\mathcal{E}_m$  is still too big, so there is a need for a "path" in  $\mathcal{E}_m$ .

# First set of hypotheses

#### H1: implicit hypothesis of every discretization:

Predictive information about y in x is "squeezed" in e, i.e.  $p_{\text{true}}(y|x, e) = p_{\text{true}}(y|e)$ .

# First set of hypotheses

#### H1: implicit hypothesis of every discretization:

Predictive information about y in x is "squeezed" in e, i.e.  $p_{\text{true}}(y|x, e) = p_{\text{true}}(y|e)$ .

H2: conditional independence:

Conditional independence of  $e_j|x_j$  with other features  $x_k, k \neq j$ .

# First set of hypotheses

#### H1: implicit hypothesis of every discretization:

Predictive information about y in x is "squeezed" in e, i.e.  $p_{\text{true}}(y|x, e) = p_{\text{true}}(y|e)$ .

H2: conditional independence:

Conditional independence of  $e_j | x_j$  with other features  $x_k, k \neq j$ .



Figure: Dependance structure between  $x_i, e_i$  and y

#### Proposal: continuous relaxation

H3: link between  $x_j$  and  $e_j$ :

H3: link between  $x_j$  and  $e_j$ : Continuous relaxation of a discrete problem (cf neural nets)

Continuous features: relaxation of the "hard" discretization

Link between  $e_j$  and  $x_j$  is supposed to be polytomous logistic:

 $p_{\alpha_j}(e_j|x_j).$ 

<□ ▶ < □ ▶ < 三 ▶ < 三 ▶ 三 りへで 18/34

H3: link between  $x_j$  and  $e_j$ : Continuous relaxation of a discrete problem (cf neural nets)

Continuous features: relaxation of the "hard" discretization

Link between  $e_j$  and  $x_j$  is supposed to be polytomous logistic:

 $p_{\alpha_j}(e_j|x_j).$ 

Categorical features: relaxation of the grouping problem

A simple contingency table is used:

$$p_{\alpha_j}(e_j = k | x_j = \ell) = \alpha_j^{k,\ell}.$$

<□> <⊡> <⊡> < 클> < 클> < 클> · 클 · 𝔍 𝔅 18/34

#### Intuitions about how it works: model proposal

p(

$$egin{aligned} p(oldsymbol{x},oldsymbol{ heta},oldsymbol{lpha}) &= \sum_{oldsymbol{e}\in\mathcal{E}_{oldsymbol{m}}} p(y|oldsymbol{e}) \prod_{j=1}^{d} p(e_j|x_j) \ &= \sum_{oldsymbol{e}\in\mathcal{E}_{oldsymbol{m}}} p_{oldsymbol{ heta}e}(y|oldsymbol{e}) \prod_{j=1}^{d} p_{oldsymbol{lpha}_j}(e_j|x_j) \ &= \sum_{oldsymbol{e}\in\mathcal{E}_{oldsymbol{m}}} p_{oldsymbol{ heta}e}(y|oldsymbol{e}) \prod_{j=1}^{d} p_{oldsymbol{lpha}_j}(e_j|x_j) \ &= \sum_{oldsymbol{e}\in\mathcal{E}_{oldsymbol{m}}} p_{oldsymbol{ heta}e}(y|oldsymbol{e}) \prod_{j=1}^{d} p_{oldsymbol{lpha}_j}(e_j|x_j) \ &= p_{oldsymbol{ heta}}(y|oldsymbol{e}^*) \end{aligned}$$

<u>▲□▶ ▲□▶</u> ▲ ె ▶ ▲ 王 → 의 < ℃ 19/34

#### Intuitions about how it works: model proposal

р

$$y|\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\alpha}) = \sum_{\boldsymbol{e} \in \mathcal{E}_{m}} p(y|\mathbf{x}, \boldsymbol{e}) p(\boldsymbol{e}|\mathbf{x})$$
$$= \sum_{\boldsymbol{e} \in \mathcal{E}_{m}} p(y|\boldsymbol{e}) \prod_{j=1}^{d} p(e_{j}|x_{j})$$
$$= \sum_{\boldsymbol{e} \in \mathcal{E}_{m}} \underbrace{p_{\boldsymbol{\theta} \boldsymbol{e}}(y|\boldsymbol{e})}_{\text{logistic}} \prod_{j=1}^{d} \underbrace{p_{\boldsymbol{\alpha}_{j}}(e_{j}|x_{j})}_{\text{logistic or table}}$$
$$\approx p_{\boldsymbol{\theta}^{\star}}(y|\boldsymbol{e}^{\star})$$

▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
▲ □ ▶
<li

クマ 19/34

E

Subsequently, it is equivalent to "optimize"  $p(y|\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\alpha})$ .

#### Intuitions about how it works: model proposal

р

$$egin{aligned} & egin{aligned} & egin{aligned} & egin{aligned} & egin{aligned} & eta &$$

Subsequently, it is equivalent to "optimize"  $p(y|\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\alpha})$ .

$$\max_{oldsymbol{ heta},oldsymbol{e}} p_{oldsymbol{ heta}}(y|oldsymbol{e}) \simeq \max_{oldsymbol{ heta},oldsymbol{lpha}} p(y|oldsymbol{x},oldsymbol{ heta},oldsymbol{lpha})$$

うくで

19/34

"Classical" estimation strategy with latent variables: EM algorithm.

"Classical" estimation strategy with latent variables: EM algorithm.

There would still be a sum over  $\mathcal{E}_{\boldsymbol{m}}$ :  $p(y|\boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{\alpha}) = \sum_{\boldsymbol{e} \in \mathcal{E}_{\boldsymbol{m}}} p_{\boldsymbol{\theta}}(y|\boldsymbol{e}) \prod_{j=1}^{d} p_{\alpha_j}(e_j|x_j)$ 

"Classical" estimation strategy with latent variables: EM algorithm.

▲□▶ ▲□▶ ▲壹▶ ▲壹▶ 壹 ∽ ♀♀ 20/34

There would still be a sum over  $\mathcal{E}_{m}$ :  $p(y|\mathbf{x}, \theta, \alpha) = \sum_{\mathbf{e} \in \mathcal{E}_{m}} p_{\theta}(y|\mathbf{e}) \prod_{j=1}^{d} p_{\alpha_{j}}(e_{j}|x_{j})$ 

Use a Stochastic-EM! Draw *e* knowing that:

"Classical" estimation strategy with latent variables: EM algorithm.

There would still be a sum over  $\mathcal{E}_{m}$ :  $p(y|\mathbf{x}, \theta, \alpha) = \sum_{\mathbf{e} \in \mathcal{E}_{m}} p_{\theta}(y|\mathbf{e}) \prod_{j=1}^{d} p_{\alpha_{j}}(e_{j}|x_{j})$ 

Use a Stochastic-EM! Draw *e* knowing that:

$$p(\boldsymbol{e}|\boldsymbol{x}, y) = \frac{p_{\boldsymbol{\theta}}(y|\boldsymbol{e}) \prod_{j=1}^{d} p_{\boldsymbol{\alpha}_{j}}(e_{j}|x_{j})}{\sum_{\boldsymbol{e} \in \boldsymbol{\mathcal{E}}_{m}} p_{\boldsymbol{\theta}}(y|\boldsymbol{e}) \prod_{j=1}^{d} p_{\boldsymbol{\alpha}_{j}}(e_{j}|x_{j})}$$
still difficult to calculate

"Classical" estimation strategy with latent variables: EM algorithm.

There would still be a sum over  $\mathcal{E}_{m}$ :  $p(y|\mathbf{x}, \theta, \alpha) = \sum_{\mathbf{e} \in \mathcal{E}_{m}} p_{\theta}(y|\mathbf{e}) \prod_{j=1}^{d} p_{\alpha_{j}}(e_{j}|x_{j})$ 

Use a Stochastic-EM! Draw *e* knowing that:

$$p(\boldsymbol{e}|\boldsymbol{x}, y) = \frac{p_{\boldsymbol{\theta}}(y|\boldsymbol{e}) \prod_{j=1}^{d} p_{\alpha_j}(e_j|x_j)}{\sum_{\boldsymbol{e} \in \boldsymbol{\mathcal{E}}_m} p_{\boldsymbol{\theta}}(y|\boldsymbol{e}) \prod_{j=1}^{d} p_{\alpha_j}(e_j|x_j)}$$
still difficult to calculate

Gibbs-sampling step:

$$p(e_j|m{x},y,m{e}_{\{-j\}}) \propto p_{m{ heta}}(y|m{e})p_{m{lpha}_j}(e_j|x_j)$$

↓ □ ▶ ↓ □ ▶ ↓ Ξ ▶ ↓ Ξ ▶ ↓ Ξ ♪ ♀ ♀ ♀ 20/34

# Algorithm

#### Initialization

1	$x_{1,1}$	$x_{1,d}$			(	$e_{1,1}$	$e_{1,d}$	
				at random				
$\langle$	$x_{n,1}$	× <sub>n,d</sub>	)		$\left( \right)$	$e_{n,1}$	e <sub>n,d</sub>	)

Loop

/ Y1 \		( e <b>1</b> ,1	e1,d	١	/ ×1,1	×1,d \
	logistic			polytomous		
	regression			regression		
· ·	regression			regression		
	$\Rightarrow$			$\Rightarrow$		
y <sub>n</sub> /		e <sub>n,1</sub>	en,d )		×n,1	×n,d /

<□▶ < □▶ < 壹▶ < 壹▶ < 壹▶ 壹 りへで 21/34

Updating e

$$\left(\begin{array}{c}p(y_{1}, e_{1,j} = k | x_{j})\\\vdots\\p(y_{n}, e_{n,j} = k | x_{j})\end{array}\right) \xrightarrow{\text{random}}_{\substack{\text{sampling}\\\Rightarrow\\}} \left(\begin{array}{c}e_{1,j}\\\vdots\\e_{n,j}\end{array}\right)$$

Calculating eMAP

$$\left( \begin{array}{c} {}^{e_{\textbf{MAP}, \mathbf{1}, j}} \\ \vdots \\ {}^{e_{\textbf{MAP}, n, j}} \end{array} \right) \begin{array}{c} \mathsf{MAP} \\ \text{estimate} \\ = \end{array} \left( \begin{array}{c} \operatorname{argmax}_{e_{j}} \rho_{\boldsymbol{\alpha}_{j}}(e_{j} | \mathbf{x}_{\mathbf{1}, j}) \\ \vdots \\ \operatorname{argmax}_{e_{j}} \rho_{\boldsymbol{\alpha}_{j}}(e_{j} | \mathbf{x}_{n, j}) \end{array} \right)$$

# Go back to "hard" thresholding: MAP estimation



We have drastically restricted the search space to provably clever candidates  $\boldsymbol{e}_{MAP}^{(1)}, \ldots, \boldsymbol{e}_{MAP}^{(iter)}$  resulting from the Gibbs sampling and MAP estimation.

$$(\boldsymbol{e}^{\star}, \boldsymbol{\theta}^{\star}) = \underset{\boldsymbol{e} \in \{\boldsymbol{e}_{\mathsf{MAP}}^{(1)}, \dots, \boldsymbol{e}_{\mathsf{MAP}}^{(\mathsf{ter})}\}, \boldsymbol{\theta} \in \Theta_{m}}{\operatorname{argmax}} \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}_{\boldsymbol{e}}}(y_{i} | \boldsymbol{e}_{i}) - \mathsf{penalty}(n; \boldsymbol{\theta})$$

↓ □ ▶ ↓ □ ▶ ↓ Ξ ▶ ↓ Ξ ▶ ↓ Ξ → ♡ 𝔅 ♡ 𝔅 23/34

We have drastically restricted the search space to provably clever candidates  $\boldsymbol{e}_{MAP}^{(1)}, \ldots, \boldsymbol{e}_{MAP}^{(iter)}$  resulting from the Gibbs sampling and MAP estimation.

$$(\boldsymbol{e}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmax}_{\boldsymbol{e} \in \{\boldsymbol{e}_{\mathsf{MAP}}^{(1)}, \dots, \boldsymbol{e}_{\mathsf{MAP}}^{(\mathsf{iter})}\}, \boldsymbol{\theta} \in \Theta_{m}} \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}_{\boldsymbol{e}}}(y_{i} | \boldsymbol{e}_{i}) - \mathsf{penalty}(n; \boldsymbol{\theta})$$

<□ ▶ < □ ▶ < 三 ▶ < 三 ▶ < 三 ♪ ○ ○ 23/34

We would still need to loop over candidates m!

We have drastically restricted the search space to provably clever candidates  $\boldsymbol{e}_{MAP}^{(1)}, \ldots, \boldsymbol{e}_{MAP}^{(iter)}$  resulting from the Gibbs sampling and MAP estimation.

$$(\boldsymbol{e}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmax}_{\boldsymbol{e} \in \{\boldsymbol{e}_{\mathsf{MAP}}^{(1)}, \dots, \boldsymbol{e}_{\mathsf{MAP}}^{(\mathsf{iter})}\}, \boldsymbol{\theta} \in \Theta_{m}} \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}_{\boldsymbol{e}}}(y_{i} | \boldsymbol{e}_{i}) - \mathsf{penalty}(n; \boldsymbol{\theta})$$

We would still need to loop over candidates m!

In practice if  $\forall i, \ p(e_{i,j} = 1 | x_{i,j}, y_i) \ll 1$ , then  $e_j = 1$  disappears...

<□▶ < @▶ < \= ▶ < \= ♪ > \= ♡ \ \ C 23/34

We have drastically restricted the search space to provably clever candidates  $\boldsymbol{e}_{MAP}^{(1)}, \ldots, \boldsymbol{e}_{MAP}^{(iter)}$  resulting from the Gibbs sampling and MAP estimation.

$$(\boldsymbol{e}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmax}_{\boldsymbol{e} \in \{\boldsymbol{e}_{\mathsf{MAP}}^{(1)}, \dots, \boldsymbol{e}_{\mathsf{MAP}}^{(\mathsf{iter})}\}, \boldsymbol{\theta} \in \Theta_{m}} \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}_{\boldsymbol{e}}}(y_{i} | \boldsymbol{e}_{i}) - \mathsf{penalty}(n; \boldsymbol{\theta})$$

We would still need to loop over candidates *m*!

In practice if  $\forall i$ ,  $p(e_{i,j} = 1 | x_{i,j}, y_i) \ll 1$ , then  $e_j = 1$  disappears... Start with  $\boldsymbol{m} = (m_{\max})_1^d$  and "wait" ... eventually until  $\boldsymbol{m} = \mathbf{1}$ .

# Interactions in logistic regression

<□ ▶ < □ ▶ < ■ ▶ < ■ ▶ < ■ ▶ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■

Upper triangular matrix with  $\delta_{k,\ell} = 1$  if  $k < \ell$  and features p and q "interact" in the logistic regression.

$$\mathsf{logit}(p_{\boldsymbol{\theta}_{\boldsymbol{f}}}(1|\boldsymbol{f}(\boldsymbol{x}))) = \theta_0 + \sum_{j=1}^{d} \theta_j^{f_j(x_j)} + \sum_{1 \leq k < \ell \leq d} \delta_{k,\ell} \theta_{k,\ell}^{f_k(x_k)f_\ell(x_\ell)}$$

Upper triangular matrix with  $\delta_{k,\ell} = 1$  if  $k < \ell$  and features p and q "interact" in the logistic regression.

$$\mathsf{logit}(p_{\boldsymbol{\theta}_{\boldsymbol{f}}}(1|\boldsymbol{f}(\boldsymbol{x}))) = \theta_0 + \sum_{j=1}^{d} \theta_j^{f_j(x_j)} + \sum_{1 \leq k < \ell \leq d} \delta_{k,\ell} \theta_{k,\ell}^{f_k(x_k)f_\ell(x_\ell)}$$

Imagine for now that the discretization e = f(x) is fixed. The criterion becomes:

$$egin{aligned} & (oldsymbol{ heta}^{\star}, oldsymbol{\delta}^{\star}) = rgmax \ & ext{argmax} \ & ext{$ heta$}_{oldsymbol{ heta}, oldsymbol{\delta} = \{oldsymbol{0}, 1\}^{rac{d(d-1)}{2}} \sum_{i=1}^n \ln p_{oldsymbol{ heta}}(y_i | oldsymbol{e}_i, oldsymbol{\delta}) - ext{penalty}(n; oldsymbol{ heta}) \end{aligned}$$

▲□▶ ▲□▶ ▲臺▶ ▲臺▶ 喜 り९℃ 25/34

Upper triangular matrix with  $\delta_{k,\ell} = 1$  if  $k < \ell$  and features p and q "interact" in the logistic regression.

$$\mathsf{logit}(p_{\boldsymbol{\theta}_{\boldsymbol{f}}}(1|\boldsymbol{f}(\boldsymbol{x}))) = \theta_0 + \sum_{j=1}^{d} \theta_j^{f_j(x_j)} + \sum_{1 \leq k < \ell \leq d} \delta_{k,\ell} \theta_{k,\ell}^{f_k(x_k)f_\ell(x_\ell)}$$

Imagine for now that the discretization e = f(x) is fixed. The criterion becomes:

$$egin{aligned} & (m{ heta}^{\star}, m{\delta}^{\star}) = lpha^{}_{m{ heta}, m{\delta} \in \{0,1\}} \sum_{i=1}^{n} \ln p_{m{ heta}}(y_i | m{e}_i, m{\delta}) - ext{penalty}(n; m{ heta}) \end{aligned}$$

▲□▶ ▲圖▶ ▲ 볼▶ ▲ 볼▶ 볼 - 의 역 <sup>(0)</sup> 25/34

Analogous to previous problem:  $2^{\frac{d(d-1)}{2}}$  models.

 $\delta$  is latent and hard to optimize over: use a stochastic algorithm!

 $\delta$  is latent and hard to optimize over: use a stochastic algorithm!

Strategy used here: Metropolis-Hastings algorithm.


## Model proposal

 $\delta$  is latent and hard to optimize over: use a stochastic algorithm! Strategy used here: Metropolis-Hastings algorithm.

$$p(y|oldsymbol{e}) = \sum_{oldsymbol{\delta} \in \{0,1\}^{rac{d(d-1)}{2}}} p(y|oldsymbol{e}, \delta) p(\delta) 
onumber \ p(\delta|oldsymbol{e}, y) \propto p(y|oldsymbol{e}, \delta) p(\delta) 
onumber \ pprox \exp(-\mathsf{BIC}[\delta]/2) p(\delta)$$

 $\delta$  is latent and hard to optimize over: use a stochastic algorithm! Strategy used here: Metropolis-Hastings algorithm.

$$p(y|\boldsymbol{e}) = \sum_{\boldsymbol{\delta} \in \{0,1\}} p(y|\boldsymbol{e}, \boldsymbol{\delta}) p(\boldsymbol{\delta})$$
$$p(\boldsymbol{\delta}|\boldsymbol{e}, y) \propto p(y|\boldsymbol{e}, \boldsymbol{\delta}) p(\boldsymbol{\delta})$$
$$\approx \exp(-\mathsf{BIC}[\boldsymbol{\delta}]/2) p(\boldsymbol{\delta}) \qquad p(\delta_{p,q}) = \frac{1}{2}$$

Whi

 $\delta$  is latent and hard to optimize over: use a stochastic algorithm!

Strategy used here: Metropolis-Hastings algorithm.

$$p(y|\boldsymbol{e}) = \sum_{\delta \in \{0,1\}^{\frac{d(d-1)}{2}}} p(y|\boldsymbol{e},\delta)p(\delta)$$

$$p(\delta|\boldsymbol{e},y) \propto p(y|\boldsymbol{e},\delta)p(\delta)$$

$$\approx \exp(-\mathsf{BIC}[\delta]/2)p(\delta) \qquad p(\delta_{p,q}) = \frac{1}{2}$$
ch transition proposal  $q : (\{0,1\}^{\frac{d(d-1)}{2}}, \{0,1\}^{\frac{d(d-1)}{2}}) \mapsto [0;1]?$ 



We restrict changes to only one entry  $\delta_{k,\ell}$ .



We restrict changes to only one entry  $\delta_{k,\ell}$ .

**Proposal:** gain/loss in BIC between **bivariate** models with / without the interaction.

<□▶ <□▶ < 壹▶ < 壹▶ ≤ うへで 27/34

We restrict changes to only one entry  $\delta_{k,\ell}$ .

**Proposal:** gain/loss in BIC between **bivariate** models with / without the interaction.

**Trick:** alternate one discretization / grouping step and one "interaction" step.

<□▶ <□▶ < 壹▶ < 壹▶ ≤ うへで 27/34

#### Performance asserted on simulated data. Good performance on real data:

Gini	Current performance	glmdisc	Basic <b>glm</b>
Auto (n=50,000 ; d=15)	57.9	64.84	58
Revolving (n=48,000 ; d=9)	58.57	67.15	53.5
Prospects (n=5,000 ; d=25)	35.6	47.18	32.7
Electronics (n=140,000 ; d=8)	57.5	58	-10
Young (n=5,000 ; d=25)	pprox 15	30	12.2
Basel II (n=70,000 ; d=13)	70	71.3	19

Relatively fast computing time: between 2 hours and a day on a laptop according to number of observations, features, ...

"Inexisting" human time.

# Conclusion and future work

<□ ▶ < □ ▶ < ■ ▶ < ■ ▶ < ■ ▶ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■ ♪ < ■



#### Conclusion

<□ ▶ < 酉 ▶ < 重 ▶ < 重 ▶ Ξ ♪ 𝔅 𝔅 30/34

#### Conclusion

Reinterpretation as a latent variable problem,

#### Conclusion

- Reinterpretation as a latent variable problem,
- Resolution proposal relying on MCMC and "soft" discretization,

▲□▶▲□▶▲壹▶▲壹▶ 壹 ∽९℃ 30/34

#### Conclusion

- Reinterpretation as a latent variable problem,
- ► Resolution proposal relying on MCMC and "soft" discretization,

▲□▶▲□▶▲壹▶▲壹▶ 壹 ∽९℃ 30/34

 Good empirical results and statistical guarantees (to some extent...),

#### Conclusion

- Reinterpretation as a latent variable problem,
- ► Resolution proposal relying on MCMC and "soft" discretization,

- Good empirical results and statistical guarantees (to some extent...),
- R implementation of glmdisc available on Github, to be submitted to CRAN,

#### Conclusion

- Reinterpretation as a latent variable problem,
- ► Resolution proposal relying on MCMC and "soft" discretization,
- Good empirical results and statistical guarantees (to some extent...),
- R implementation of glmdisc available on Github, to be submitted to CRAN,
- Python implementation of glmdisc available on Github and PyPi,

▲□▶ ▲圖▶ ▲ 볼▶ ▲ 볼▶ 볼 · 의 역 @ 30/34

#### Conclusion

- Reinterpretation as a latent variable problem,
- Resolution proposal relying on MCMC and "soft" discretization,
- Good empirical results and statistical guarantees (to some extent...),
- R implementation of glmdisc available on Github, to be submitted to CRAN,
- Python implementation of glmdisc available on Github and PyPi,

▲□▶ ▲圖▶ ▲ 볼▶ ▲ 볼▶ 볼 · 의 역 @ 30/34

▶ Big gain for statisticians in the field of *Credit Scoring*.

#### Conclusion

- Reinterpretation as a latent variable problem,
- ► Resolution proposal relying on MCMC and "soft" discretization,
- Good empirical results and statistical guarantees (to some extent...),
- R implementation of glmdisc available on Github, to be submitted to CRAN,
- Python implementation of glmdisc available on Github and PyPi,

▲□▶ ▲圖▶ ▲ 볼▶ ▲ 볼▶ 볼 · 의 역 @ 30/34

▶ Big gain for statisticians in the field of *Credit Scoring*.

Perspectives

#### Conclusion

- Reinterpretation as a latent variable problem,
- ► Resolution proposal relying on MCMC and "soft" discretization,
- Good empirical results and statistical guarantees (to some extent...),
- R implementation of glmdisc available on Github, to be submitted to CRAN,
- Python implementation of glmdisc available on Github and PyPi,
- Big gain for statisticians in the field of Credit Scoring.

#### Perspectives

Tested for logistic regression and polytomous logistic links: can be adapted to other models p<sub>θ</sub> and p<sub>α</sub>!

#### Conclusion

- Reinterpretation as a latent variable problem,
- ► Resolution proposal relying on MCMC and "soft" discretization,
- Good empirical results and statistical guarantees (to some extent...),
- R implementation of glmdisc available on Github, to be submitted to CRAN,
- Python implementation of glmdisc available on Github and PyPi,
- Big gain for statisticians in the field of Credit Scoring.

#### Perspectives

- Tested for logistic regression and polytomous logistic links: can be adapted to other models p<sub>θ</sub> and p<sub>α</sub>!
- ► The same model can be estimated with shallow neural networks.

# Shallow neural nets as a substitute estimation procedure



# Shallow neural nets as a substitute estimation procedure



# Thanks!

↓□ → ↓ □ → ↓ Ξ → ↓ Ξ → ↓ Ξ → ↓ □ → ↓ □ → ↓ □ → ↓ Ξ → ↓ Ξ → ↓ Ξ → ↓ Ξ → ↓ Ξ → ↓ □ → ↓

 Ramírez-Gallego, S., García, S., Mouriño-Talín, H., Martínez-Rego, D., Bolón-Canedo, V., Alonso-Betanzos, A., Benítez, J. M., and Herrera, F. (2016).
 Data discretization: taxonomy and big data challenge.
 Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 6(1):5–21.

$$p(\delta_{k,\ell} = 1 | e_k, e_\ell, y) = g(\mathsf{BIC}[\delta_{k,\ell} = 1] - \mathsf{BIC}[\delta_{k,\ell} = 0])$$
  

$$\approx \exp\left(\frac{1}{2}(\mathsf{BIC}[p_\theta(y|e_k, e_\ell, \delta_{k,\ell} = 0)] - \mathsf{BIC}[p_\theta(y|e_k, e_\ell, \delta_{k,\ell} = 1)])\right)$$
  

$$q(\delta, \delta') = |\delta_{k,\ell} - p_{k,\ell}| \text{ for the unique couple } (k,\ell) \text{ s.t. } \delta_{k,\ell}^{(s)} \neq \delta'_{k,\ell}$$
  

$$\alpha = \min\left(1, \frac{p(\delta'|e,y)}{p(\delta|e,y)} \frac{1-q(\delta,\delta')}{q(\delta,\delta')}\right)$$
  

$$\approx \min\left(1, \exp\left(\frac{1}{2}(\mathsf{BIC}[p_\theta(y|e,\delta)] - \mathsf{BIC}[p_\theta(y|e,\delta')])\right) \frac{1-q(\delta,\delta')}{q(\delta,\delta')}\right)$$

◆□ ▶ < 昼 ▶ < 重 ▶ < 重 ▶ 至 り < で 34/34</p>