Discrétisation, regroupement de modalités et introduction d'interactions en régression logistique

Adrien Ehrhardt^{1,2}

Christophe Biernacki² Philippe Heinrich³ Vincent Vandewalle^{2,4}

¹Crédit Agricole Consumer Finance
²Inria Lille - Nord-Europe
³Université de Lille, Paul Painlevé
⁴Université de Lille, EA2694

03/10/2018





<□▶ <□▶ < 글▶ < 글▶ < 글▶ 글 < 의 < ℃ 1/38

Context and basic notations

Supervised multivariate discretization and factor levels grouping

<□ ▶ < □ ▶ < 三 ▶ < 三 ▶ 三 りへで 2/38

Selecting interactions in logistic regression

Conclusion and future work

Context and basic notations

୬୯୯ 3/38

3

.≣ ▶

(日) (同) (三) (

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

4/38

.≣ ►

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

	Family status	Wages	Repayment
Craftsman	Widower	2000	0
	Common-law	1700	0
	Divorced	4000	1
	Single	2700	1
Office employee	Married	1400	0
	?	1200	0

Table: Dataset with outliers and missing values.

- 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]-∞ ; 1500]	0
Worker]-∞ ; 1500]	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
]- ∞ ; 1500]	0
]- ∞ ; 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 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]- ∞ ; 1500]	0
?+Low-qualified		?+Alone x]- ∞ ; 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; ∞ [1
High-qualified		?+Alone x]2000; ∞ [1
?+Low-qualified		Union x]- ∞ ; 1500]		0
?+Low-qualified		?+Alone x]- ∞ ; 1500]		0

Table: Dataset with outliers and missing values.

< □ > < 同

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

Feature	Level	Points
	18-25	10
Age	25-45	20
	45- $+\infty$	30
	<u>−∞-1000</u>	15
Wages	1000-2000	25
	2000-+∞	35
Glucose level		

Table: Final scorecard.

The whole process can be decomposed into two steps:

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

< □ > < □ > < □ > < □ > < □ > < □ > < □ >

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

The whole process can be decomposed into two steps:

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

▲□▶▲□▶▲壹▶▲壹▶ 壹 ∽९ペ 6/38

Selected features: $\mathbf{x} = (x_j)_1^d$ (continuous or categorical).

The whole process can be decomposed into two steps:

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

▲□▶ ▲□▶ ▲壹▶ ▲壹▶ 壹 少�♡ 6/38

Selected features: $\mathbf{x} = (x_j)_1^d$ (continuous or categorical). \mathbf{f} is "component-wise", *i.e.* $\mathbf{f}(\mathbf{x}) = (f_j(x_j))_1^d$. We restrict to discretization and grouping of factor levels.

<□▶ <□▷ <□▷ < ⊇▶ < ⊇▶ < ⊇▶ < ⊇ < つへで 7/38

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

Discretization

Into *m* intervals with associated cutpoints $c = (c_0 = -\infty, c_1, \dots, c_{m-1}, c_m = +\infty).$

Discretization function

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

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







Grouping

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

Grouping function

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

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

SAC.

$$f_j \in \mathcal{M}_j = \{f_j(\cdot; \boldsymbol{c}_j, m_j) | m_j \in \mathbb{N}, c_{j,1} < \ldots < c_{j,m_j-1}\}$$



$$f_j \in \mathcal{M}_j = \{f_j(\cdot; \boldsymbol{c}_j, m_j) | m_j \in \mathbb{N}, c_{j,1} < \ldots < c_{j,m_j-1}\}$$

 \mathcal{M}_j is seemingly continuous but with a finite sample, a countable Feature Space can be recovered by remarking:

<□ ▶ < □ ▶ < 三 ▶ < 三 ▶ 三 りへで 9/38

$$f_j \in \mathcal{M}_j = \{f_j(\cdot; \boldsymbol{c}_j, m_j) | m_j \in \mathbb{N}, c_{j,1} < \ldots < c_{j,m_j-1}\}$$

 \mathcal{M}_j is seemingly continuous but with a finite sample, a countable Feature Space can be recovered by remarking:

< □ > < □ > < □ > < □ > < □ > < □ >

596

$$f_j \in \mathcal{M}_j = \{f_j(\cdot; \boldsymbol{c}_j, m_j) | m_j \in \mathbb{N}, c_{j,1} < \ldots < c_{j,m_j-1}\}$$

 \mathcal{M}_j is seemingly continuous but with a finite sample, a countable Feature Space can be recovered by remarking:

< □ > < □ > < □ > < □ > < □ > < □ >

596

$$f_j \in \mathcal{M}_j = \{f_j(\cdot; \boldsymbol{c}_j, m_j) | m_j \in \mathbb{N}, c_{j,1} < \ldots < c_{j,m_j-1}\}$$

 \mathcal{M}_j is seemingly continuous but with a finite sample, a countable Feature Space can be recovered by remarking:

$$f_j \in \mathcal{M}_j = \{f_j(\cdot; \boldsymbol{c}_j, m_j) | m_j \in \mathbb{N}, c_{j,1} < \ldots < c_{j,m_j-1}\}$$

 \mathcal{M}_j is seemingly continuous but with a finite sample, a countable Feature Space can be recovered by remarking:

<□ ▶ < □ ▶ < 三 ▶ < 三 ▶ 三 りへで 9/38

Example (n = 20, d = 10): $\approx 10^{57}$ models in \mathcal{M}_i^d .

Grouping

 $f_j \in \mathcal{M}_j = \{ \text{Partitions from } \{1, \dots, o_j \} \text{ to } \{1, \dots, m_j \}; m_j \leq o_j \}.$

<u>↓□ ▶ ↓ @ ▶ ↓ ! ↓ ↓ ! ♪ ↓ ! ♪ ↓ 0 ↓ 10/38</u>

Grouping

 $f_j \in \mathcal{M}_j = \{ \text{Partitions from } \{1, \dots, o_j \} \text{ to } \{1, \dots, m_j \}; m_j \leq o_j \}.$

Its cardinality is given by the Stirling number of the second kind: $|\mathcal{M}_j| = \sum_{m_j=1}^{o_j} \frac{1}{m_j!} \sum_{i=0}^{m_j} (-1)^{m_j-i} {m_j \choose i} i^{o_j}.$

▲□▶ ▲□▶ ▲ 글▶ ▲ 글▶ 글 ∽ ९ ℃ 10/38

Grouping

 $f_j \in \mathcal{M}_j = \{ \text{Partitions from } \{1, \dots, o_j \} \text{ to } \{1, \dots, m_j\}; m_j \leq o_j \}.$

Its cardinality is given by the Stirling number of the second kind: $|\mathcal{M}_j| = \sum_{m_j=1}^{o_j} \frac{1}{m_j!} \sum_{i=0}^{m_j} (-1)^{m_j-i} {m_j \choose i} i^{o_j}.$

▲□▶ ▲□▶ ▲ 글▶ ▲ 글▶ 글 ∽ ९ ℃ 10/38

Exhaustive search is untractable.

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

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

On "raw" data, logistic regression yields:

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三 ▶ ● ⑦ � () 11/38

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

On "raw" data, logistic regression yields:

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

On discretized / grouped data, logistic regression yields:

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

◆□ ▶ ◆□ ▶ ◆ 三 ▶ ▲ 三 ▶ ● ♡ � ① 11/38

Mathematical reinterpretation: Objective



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

< □ >

4

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

13/38

Э

Logistic regression on "raw" data:

 $\mathsf{logit}(p_{{m{ heta}}_\mathsf{raw}}(1|{m{x}})) = heta_0 + {m{ heta}}_1 {m{x}}_1$



Figure: Linear logistic regression fit

< □ ▶

▲□ ► ▲ □ ► ▲

E

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

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

< □ >

JAC.

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

$$\mathsf{logit}(p_{m{ heta}_f}(1|m{f}(m{x}))) = heta_0 + \underbrace{ heta_1^{f_1(m{x}_1)}}_{ heta_1^{1},\ldots, heta_1^{1}}$$



Figure: Good (bias/variance tradeoff) discretization

▲□▶ ▲ 同

⊒ ▶ 3

∍ ►

JAC.

 θ can be estimated for each discretization f and f^* can be chosen through our favorite model choice criterion *e.g.* BIC.
Criterion

 θ can be estimated for each discretization f and f^* can be chosen through our favorite model choice criterion *e.g.* BIC.

A model selection problem

$$(\boldsymbol{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmin}_{\boldsymbol{f} \in \mathcal{F}, \boldsymbol{\theta} \in \Theta_{\boldsymbol{f}}} - 2 \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_i | \boldsymbol{f}(\boldsymbol{x}_i)) + |\boldsymbol{\theta}| \times \ln(n),$$

≣ ��� 14/38

4 □ > 4 同 > 4 Ξ > 4 Ξ >

where θ is classicaly estimated via MLE.

Criterion

 θ can be estimated for each discretization f and f^* can be chosen through our favorite model choice criterion *e.g.* BIC.

A model selection problem

$$(\boldsymbol{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmin}_{\boldsymbol{f} \in \mathcal{F}, \boldsymbol{\theta} \in \Theta_{\boldsymbol{f}}} - 2 \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_i | \boldsymbol{f}(\boldsymbol{x}_i)) + |\boldsymbol{\theta}| \times \ln(n),$$

where θ is classicaly estimated via MLE.

Compromise between (over-)fitting the data and model complexity (and explainability in a sense!).

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三 ∽९ペ 14/38

Criterion

 θ can be estimated for each discretization f and f^* can be chosen through our favorite model choice criterion *e.g.* BIC.

A model selection problem

$$(\boldsymbol{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmin}_{\boldsymbol{f} \in \mathcal{F}, \boldsymbol{\theta} \in \Theta_{\boldsymbol{f}}} - 2 \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_i | \boldsymbol{f}(\boldsymbol{x}_i)) + |\boldsymbol{\theta}| \times \ln(n),$$

where θ is classicaly estimated via MLE.

Compromise between (over-)fitting the data and model complexity (and explainability in a sense!).

 ${\mathcal F}$ is discrete and combinatorial: how can we get around this problem?

Current academic methods:

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



Quick example of χ^2 :

Category	# samples	# cases	p-value
18-20	10	5	0.3
20-22	10	6	0.3
22-24	10	4	0.2

Supervised multivariate discretization and factor levels grouping

17/38

<u>・ロト</u> ・ (日) ト ・ (三) ト ・ (

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

 $e_j = f_j(x_j).$

<□ ▶ < @ ▶ < \ > ↓ \ \ = ▶ \ = りへで 18/38

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

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□

かくで 18/38

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

↓□▶ <□▶ < => < => < </p>

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.

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

First set of hypotheses

H1: implicit hypothesis of every discretization:

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

<□ ▶ < @ ▶ < \ = ▶ \ = りへで 19/38

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 $x_k, k \neq j$.

<□ ▶ < @ ▶ < \ > ↓ \ \ \ = ▶ \ \ = り \ \ 0 \ \ 19/38

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

▲□▶ ▲□▶ ▲豆▶ ▲豆▶ 豆 りへで 19/38

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:

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

<□ ▶ < @ ▶ < \ > ▶ < \ > > \ > \ \ 20/38

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:

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

<□▶ <♬▶ < \= ▶ < \= ♪ < \= ♡ < \C 20/38

Intuitions about how it works: model proposal

p(

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

<□▶

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

୬ ବ୍ ତି 21/38

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" $p(y|\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\alpha})$.

$$\max_{\substack{\boldsymbol{\theta}, \boldsymbol{e}}} p_{\boldsymbol{\theta}}(y | \boldsymbol{e}) \simeq \max_{\substack{\boldsymbol{\theta}, \boldsymbol{\alpha}}} p(y | \boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{\alpha})$$

 $\mathcal{O} \diamond \mathcal{O}$

21/38

Go back to "hard" thresholding: MAP estimation



Two very different estimation strategies

Two very different estimation strategies

1. In the statistics community: latent feature = EM-like algorithm. We try to get $\max_{\theta,\alpha} p(y|\mathbf{x}; \theta, \alpha)$ through SEM algorithm + Gibbs sampling step that explicitly draws \mathbf{e} .

↓□▶ ↓ @ ▶ ↓ ! ● ▶ ↓ ! ● ♡ \ ○ 23/38

Two very different estimation strategies

1. In the statistics community: latent feature = EM-like algorithm. We try to get $\max_{\theta,\alpha} p(y|\mathbf{x}; \theta, \alpha)$ through SEM algorithm + Gibbs sampling step that explicitly draws \mathbf{e} .

2. Machine Learning: neural networks natively learn representations of the data.

A 1-hidden layer neural network with softmax activation function that *via* Stochastic Gradient Descent tries to maximize the likelihood of $p_{\theta}(y|\tilde{e} = (p_{\alpha_j}(1|x_j), \dots, p_{\alpha_j}(m_j|x_j))_1^d)$.

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

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

・□ ▶ < □ ▶ < 三 ▶ < 三 ▶ 三 り < ℃ 24/38</p>

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.

<□ ▶ < □ ▶ < 三 ▶ < 三 ▶ < 三 ♪ ○ ○ ○ 24/38

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

< □ ▶ < □ ▶ < 三 ▶ < 三 ▶ ミ の Q @ 24/38

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

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

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

Gibbs-sampling step:

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

<□ ▶ < □ ▶ < 三 ▶ < 三 ▶ < 三 ♪ ○ ○ 24/38

Algorithm

Initialization

(×1,1	$x_{1,d}$			($e_{1,1}$	$e_{1,d}$	
				at random				
\langle	$x_{n,1}$	× _{n,d})		\langle	$e_{n,1}$	e _{n,d})

Loop

/ Y1 \		(e 1 ,1	e1,d	١	/ ×1,1	×1,d \
	logistic			polytomous		
	regression			regression		
· ·	regression			regression		
	\Rightarrow			\Rightarrow		
y _n /		e _{n,1}	en,d)		×n,1	×n,d /

Updating e

$$\left(\begin{array}{c}p(\mathbf{y_1}, \mathbf{e_{1,j}} = k | \mathbf{x}_i)\\\vdots\\p(\mathbf{y}_n, \mathbf{e}_{n,j} = k | \mathbf{x}_i)\end{array}\right) \quad \begin{array}{c}\text{random}\\\text{sampling}\\\Rightarrow\\ \end{array} \left(\begin{array}{c}\mathbf{e_{1,j}}\\\vdots\\\mathbf{e}_{n,j}\end{array}\right)$$

Calculating eMAP

$$\left(\begin{array}{c} \hat{f}_{j}(\mathbf{x}_{1,j}) \\ \vdots \\ \hat{f}_{j}(\mathbf{x}_{n,j}) \end{array}\right) \begin{array}{c} \mathsf{MAP} \\ \mathsf{estimate} \\ = \end{array} \left(\begin{array}{c} \operatorname{argmax}_{e_{j}} p_{\mathbf{\alpha}_{j}}(e_{j} | \mathbf{x}_{1,j}) \\ \vdots \\ \operatorname{argmax}_{e_{j}} p_{\mathbf{\alpha}_{j}}(e_{j} | \mathbf{x}_{n,j}) \end{array}\right)$$

Estimation via neural nets



< □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □

We have drastically restricted the search space to provably clever candidates $\hat{f}^{(1)}, \ldots, \hat{f}^{(\text{iter})}$ resulting either from the Gibbs sampling or the neural network and MAP estimation.

$$(\mathbf{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmin}_{\hat{\mathbf{f}} \in \{\hat{\mathbf{f}}^{(1)}, \dots, \hat{\mathbf{f}}^{(\mathrm{iter})}\}, \boldsymbol{\theta} \in \Theta_{m}} -2\sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_{i}|\hat{\mathbf{f}}(\mathbf{x}_{i})) + (m_{1} \times \dots \times m_{d} - d) \times \ln(n)$$

▲□▶ ▲□▶ ▲ 壹▶ ▲ 壹 ▶ 壹 → り ९ (~ 27/38

We have drastically restricted the search space to provably clever candidates $\hat{f}^{(1)}, \ldots, \hat{f}^{(\text{iter})}$ resulting either from the Gibbs sampling or the neural network and MAP estimation.

$$(\mathbf{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmin}_{\hat{\mathbf{f}} \in \{\hat{\mathbf{f}}^{(1)}, \dots, \hat{\mathbf{f}}^{(\mathrm{itor})}\}, \boldsymbol{\theta} \in \Theta_{m}} - 2\sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_{i}|\hat{\mathbf{f}}(\mathbf{x}_{i})) + (m_{1} \times \dots \times m_{d} - d) \times \ln(n)$$

<□ ▶ < □ ▶ < 亘 ▶ < 亘 ▶ < 亘 ▶ ○ ♀ ♀ 27/38

We would still need to loop over candidates m!

We have drastically restricted the search space to provably clever candidates $\hat{f}^{(1)}, \ldots, \hat{f}^{(\text{iter})}$ resulting either from the Gibbs sampling or the neural network and MAP estimation.

$$(\mathbf{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmin}_{\mathbf{\hat{f}} \in \{\mathbf{\hat{f}}^{(1)}, \dots, \mathbf{\hat{f}}^{(\mathrm{iter})}\}, \boldsymbol{\theta} \in \Theta_{m}} - 2\sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_{i} | \mathbf{\hat{f}}(\mathbf{x}_{i})) + (m_{1} \times \dots \times m_{d} - d) \times \ln(n)$$

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...

<□▶ <□▶ < 壹▶ < 壹▶ < 壹▶ 壹 の 𝔅 27/38

We have drastically restricted the search space to provably clever candidates $\hat{f}^{(1)}, \ldots, \hat{f}^{(\text{iter})}$ resulting either from the Gibbs sampling or the neural network and MAP estimation.

$$(\mathbf{f}^{\star}, \boldsymbol{\theta}^{\star}) = \operatorname*{argmin}_{\hat{\mathbf{f}} \in \{\hat{\mathbf{f}}^{(1)}, \dots, \hat{\mathbf{f}}^{(\mathrm{ter})}\}, \boldsymbol{\theta} \in \Theta_{m}} - 2\sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_{i} | \hat{\mathbf{f}}(\mathbf{x}_{i})) + (m_{1} \times \dots \times m_{d} - d) \times \ln(n)$$

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 $\boldsymbol{m} = 1$.

Selecting interactions in logistic regression
Notations

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

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

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

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

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

$$(\boldsymbol{\theta}^{\star}, \boldsymbol{\delta}^{\star}) = \underset{\boldsymbol{\theta}, \boldsymbol{\delta} \in \{0, 1\}}{\operatorname{argmin}} \underbrace{-2 \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_i | \boldsymbol{f}(\boldsymbol{x}_i), \boldsymbol{\delta}) + |\boldsymbol{\theta}| \ln(n)}_{\operatorname{BIC}[\boldsymbol{\delta}]}$$

4 □ ▶ 4 @ ▶ 4 E ▶ 4 E ▶

1

かくで 29/38

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

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

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

$$(\boldsymbol{\theta}^{\star}, \boldsymbol{\delta}^{\star}) = \underset{\boldsymbol{\theta}, \boldsymbol{\delta} \in \{0, 1\}}{\operatorname{argmin}} \underbrace{-2 \sum_{i=1}^{n} \ln p_{\boldsymbol{\theta}}(y_i | \boldsymbol{f}(\boldsymbol{x}_i), \boldsymbol{\delta}) + |\boldsymbol{\theta}| \ln(n)}_{\operatorname{BIC}[\boldsymbol{\delta}]}$$

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

δ is latent and hard to optimize over: use a stochastic algorithm!

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

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

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

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

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

4 □ ▶ 4 @ ▶ 4 E ▶ 4 E ▶

E

୬ ୧୯ 30/38

Wh

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

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

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

$$\approx \exp(-\text{BIC}[\boldsymbol{\delta}]/2) p(\boldsymbol{\delta}) \qquad p(\boldsymbol{\delta}_{p,q}) = \frac{1}{2}$$
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.

▲□▶ ▲□▶ ▲ 글▶ ▲ 글▶ 글 ∽ ९ ペ 31/38

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.

▲□▶ ▲□▶ ▲ 글▶ ▲ 글▶ 글 ∽ ९ ° 31/38

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

Gini	Current performance	glmdisc	Basic glm
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.

	Pima	Breast	Heart	Birthwt
Naïve LR	0.73	0.94	0.78	0.34
Naïve LR w. interactions	0.60	0.51	0.47	0.15
glmdisc	0.57	0.93	0.82	0.18
glmdisc w. interactions	0.62	0.95	0.67	0.54

Conclusion and future work

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

<□▶ </p>

 Interpretability + good empirical results and statistical guarantees (to some extent...),

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

▲□▶ ▲□▶ ▲臺▶ ▲臺▶ 喜 ∽�� 35/38

- Interpretability + 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,

- Interpretability + 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 익 (0 35/38)

► Big gain for statisticians relying on logistic regression.

- Interpretability + 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 relying on logistic regression.

Perspectives

Tested for logistic regression and polytomous logistic links: can be adapted to other models p_{θ} and p_{α} !

Thanks!

↓□ → ↓ □ → ↓ Ξ → ↓ Ξ → ↓ Ξ → ↓ 0 ↓ 0 ↓ 36/38

 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.

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

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

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

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

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

<□ ▶ < □ ▶ < □ ▶ < 三 ▶ < 三 ▶ = りへで 38/38