

# Reject Inference, “quantization”, interactions, logistic regression trees, and bonuses

Adrien Ehrhardt

Mission Lane, 08/03/2022



# Who am I?

≈ 2016-2019: “CIFRE” PhD student at Inria (consortium of French labs, like CNRS, but specialized in Applied Maths) and Crédit Agricole Consumer Finance (consumer loans).



# Who am I?

≈ 2016-2019: “CIFRE” PhD student at Inria (consortium of French labs, like CNRS, but specialized in Applied Maths) and Crédit Agricole Consumer Finance (consumer loans).



Le périmètre du Groupe Crédit Agricole rassemble Crédit Agricole S.A., l'ensemble des Caisses régionales et des Caisses locales, ainsi que leurs filiales.

## PUBLIC

30,9%  
INVESTISSEURS INSTITUTIONNELS

8,0%  
ACTIONNAIRES INDIVIDUELS

5,8%  
SALARIÉS VIA L'ÉPARGNE SALARIALE

NS<sup>(2)</sup>  
AUTOCONTRÔLE

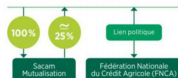


## CAISSES RÉGIONALES

10,9 m  
DE SOCIÉTAIRES  
détenant les parts sociales de  
**2 410**  
CAISSES LOCALES

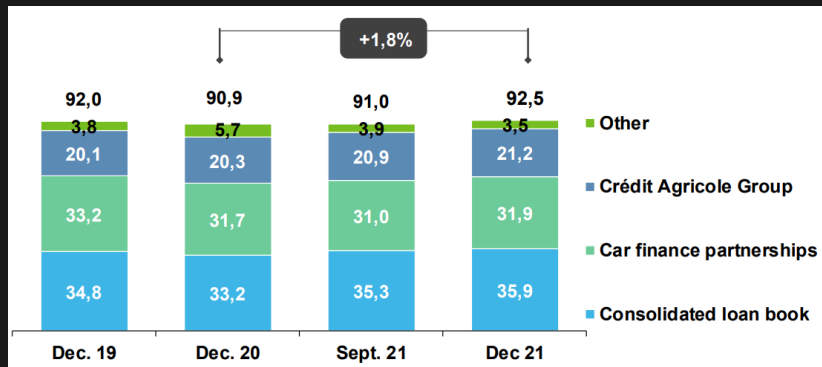
**39**  
CAISSES RÉGIONALES

détenant ensemble la majorité du capital  
de CRÉDIT AGRICOLE S.A.  
via la SAS Rue La Boétie<sup>(3)</sup>



# Who am I?

≈ 2016-2019: “CIFRE” PhD student at Inria (consortium of French labs, like CNRS, but specialized in Applied Maths) and Crédit Agricole Consumer Finance (consumer loans).





# Who am I?

≈ 2016-2019: “CIFRE” PhD student at Inria (consortium of French labs, like CNRS, but specialized in Applied Maths) and Crédit Agricole Consumer Finance (consumer loans).



≈ 2020-now: Machine Learning Engineer at Crédit Agricole S.A. & Associate Professor at École Polytechnique.



# Collaborators



Christophe Biernacki



Vincent Vandewalle



Philippe Heinrich



Elise Bayraktar



Xuwen Liu



Minh Tuan Nguyen



Cléa Laouar

# Context and notations: industrial setting

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

**Table:** Dataset with outliers and missing values.

# Context and notations: industrial setting

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

**Table:** Dataset with outliers and missing values.

1. Discarding not financed applicants
2. Feature selection
3. Discretization / grouping
4. Interaction screening
5. Segmentation
6. Logistic regression fitting

# Context and notations: industrial setting

Job	Home	Time in job	Family status	Wages		Repayment
Craftsman	Owner	20	Widower	2000		1
?	Renter	10	Common-law	1700		0
Engineer	Starter	5	Divorced	4000		1
Executive	By work	8	Married	2700		0
<del>Office employee</del>	<del>Renter</del>	<del>12</del>	<del>Married</del>	<del>1400</del>		NA
<del>Worker</del>	<del>By family</del>	<del>7</del>	<del>7</del>	<del>1200</del>		NA

**Table:** Dataset with outliers and missing values.

1. **Discarding not financed applicants**
2. Feature selection
3. Discretization / grouping
4. Interaction screening
5. Segmentation
6. Logistic regression fitting

# Context and notations: industrial setting

Job			Family status	Wages		Repayment
Craftsman			Widower	2000		1
?			Common-law	1700		0
Engineer			Divorced	4000		1
Executive			Married	2700		0
<del>Office employee</del>	<del>Renter</del>	<del>12</del>	<del>Married</del>	<del>1400</del>		NA
<del>Worker</del>	<del>By family</del>	<del>7</del>	<del>7</del>	<del>1200</del>		NA

**Table:** Dataset with outliers and missing values.

1. Discarding not financed applicants
2. **Feature selection**
3. Discretization / grouping
4. Interaction screening
5. Segmentation
6. Logistic regression fitting

# Context and notations: industrial setting

Job			Family status	Wages		Repayment
Craftsman			Widower	]1500;2000]		1
?			Common-law	]1500;2000]		0
Engineer			Divorced	]2000; $\infty$ [		1
Executive			Married	]2000; $\infty$ [		0
<del>Office employee</del>	<del>Renter</del>	<del>12</del>	<del>Married</del>	<del>1400</del>		NA
<del>Worker</del>	<del>By family</del>	<del>7</del>	<del>7</del>	<del>1200</del>		NA

**Table:** Dataset with outliers and missing values.

1. Discarding not financed applicants
2. Feature selection
3. **Discretization** / grouping
4. Interaction screening
5. Segmentation
6. Logistic regression fitting

# Context and notations: industrial setting

Job			Family status	Wages		Repayment
?+Low-qualified			?+Alone	]1500;2000]		1
?+Low-qualified			Union	]1500;2000]		0
High-qualified			?+Alone	]2000; $\infty$ [		1
High-qualified			Union	]2000; $\infty$ [		0
<del>Office employee</del>	<del>Renter</del>	<del>12</del>	<del>Married</del>	<del>1400</del>		NA
<del>Worker</del>	<del>By family</del>	<del>7</del>	<del>7</del>	<del>1200</del>		NA

**Table:** Dataset with outliers and missing values.

1. Discarding not financed applicants
2. Feature selection
3. Discretization / **grouping**
4. Interaction screening
5. Segmentation
6. Logistic regression fitting



# Context and notations: industrial setting

Job			Family status x Wages		Repayment
?+Low-qualified			?+Alone x ]1500;2000]		1
?+Low-qualified			Union x ]1500;2000]		0
High-qualified			?+Alone x ]2000;∞[		1
High-qualified			Union x ]2000;∞[		0
<del>Office employee</del>	<del>Renter</del>	<del>12</del>	<del>Married</del> 1400		NA
<del>Worker</del>	<del>By family</del>	<del>2</del>	<del>1</del> 1200		NA

**Table:** Dataset with outliers and missing values.

1. Discarding not financed applicants
2. Feature selection
3. Discretization / grouping
4. **Interaction screening**
5. Segmentation
6. Logistic regression fitting

# Context and notations: industrial setting

Job			Family status × Wages		Repayment
?+Low-qualified			?+Alone × ]1500;2000]		1
?+Low-qualified			Union × ]1500;2000]		0
High-qualified			?+Alone × ]2000;∞[		1
High-qualified			Union × ]2000;∞[		0
<del>Office employee</del>	<del>Renter</del>	<del>12</del>	<del>Married 1400</del>		NA
<del>Worker</del>	<del>By family</del>	<del>2</del>	<del>1200</del>		NA

**Table:** Dataset with outliers and missing values.

1. Discarding not financed applicants
2. Feature selection
3. Discretization / grouping
4. Interaction screening
5. **Segmentation**
6. Logistic regression fitting

# Context and notations: industrial setting

Job			Family status × Wages	Score	Repayment
?+Low-qualified			?+Alone × ]1500;2000]	225	1
?+Low-qualified			Union × ]1500;2000]	190	0
High-qualified			?+Alone × ]2000;∞[	218	1
High-qualified			Union × ]2000;∞[	202	0
<del>Office employee</del>	<del>Renter</del>	<del>12</del>	<del>Married 1400</del>	NA	NA
<del>Worker</del>	<del>By family</del>	<del>2</del>	<del>1200</del>	NA	NA

**Table:** Dataset with outliers and missing values.

1. Discarding not financed applicants
2. Feature selection
3. Discretization / grouping
4. Interaction screening
5. Segmentation
6. **Logistic regression fitting**

## Context and notations: available data

Random variables:  $\mathbf{X}, Y, Z$ .

# Context and notations: available data

Random variables:  $\mathbf{X}, Y, Z$ .

## Observations

$\mathbf{x} = (x_1, \dots, x_d)$	characteristics,
$x_j \in \mathbb{R} \text{ or } \{1, \dots, l_j\}$	e.g. rent amount, job, ... ,
$y \in \{0, 1\}$	good or bad,
$z \in \{f, nf\}$	financed or not financed.

# Context and notations: available data

Random variables:  $\mathbf{X}, Y, Z$ .

## Observations

$\mathbf{x} = (x_1, \dots, x_d)$	characteristics,
$x_j \in \mathbb{R} \text{ or } \{1, \dots, l_j\}$	e.g. rent amount, job, ... ,
$y \in \{0, 1\}$	good or bad,
$z \in \{f, nf\}$	financed or not financed.

## Samples

$\mathcal{T}_f = (\mathbf{x}_f, y_f, z_f)$	$n$ -sample of financed clients,
$\mathcal{T}_{nf} = (\mathbf{x}_{nf}, z_{nf})$	$n'$ -sample of not-financed clients,
$\mathcal{T} = \mathcal{T}_f \cup \mathcal{T}_{nf}$	observed sample,
$\mathcal{T}_c = \mathcal{T} \cup \mathbf{y}_{nf}$	complete sample.

# Context and notations: available data

The observed data are the following:

$$\mathcal{T} = \begin{pmatrix} \mathcal{T}_f \\ \mathcal{T}_{nf} \end{pmatrix} = \begin{pmatrix} \begin{pmatrix} x_f \\ \begin{matrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{matrix} \end{pmatrix} \\ \begin{pmatrix} x_{nf} \\ \begin{matrix} x_{n+1,1} & \cdots & x_{n+1,d} \\ \vdots & & \vdots \\ x_{n+n',1} & \cdots & x_{n+n',d} \end{matrix} \end{pmatrix} \end{pmatrix} \begin{pmatrix} \begin{pmatrix} y_f \\ \begin{matrix} y_1 \\ \vdots \\ y_n \end{matrix} \end{pmatrix} \\ \begin{pmatrix} y_{nf} \\ \begin{matrix} NA \\ \vdots \\ NA \end{matrix} \end{pmatrix} \end{pmatrix} \begin{pmatrix} \begin{pmatrix} z_f \\ \begin{matrix} f \\ \vdots \\ f \end{matrix} \end{pmatrix} \\ \begin{pmatrix} z_{nf} \\ \begin{matrix} nf \\ \vdots \\ nf \end{matrix} \end{pmatrix} \end{pmatrix}.$$

## Context and notations: available data

The observed data are the following:

$$\mathcal{T} = \begin{pmatrix} \mathcal{T}_f \\ \mathcal{T}_{nf} \end{pmatrix} = \begin{pmatrix} \begin{pmatrix} x_f \\ \vdots \\ x_n \end{pmatrix} \begin{pmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \vdots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{pmatrix} \begin{pmatrix} y_f \\ \vdots \\ y_n \end{pmatrix} \begin{pmatrix} z_f \\ \vdots \\ f \end{pmatrix} \\ \begin{pmatrix} x_{nf} \\ \vdots \\ x_{n+n',1} \end{pmatrix} \begin{pmatrix} x_{n+1,1} & \cdots & x_{n+1,d} \\ \vdots & \vdots & \vdots \\ x_{n+n',1} & \cdots & x_{n+n',d} \end{pmatrix} \begin{pmatrix} NA \\ \vdots \\ NA \end{pmatrix} \begin{pmatrix} NA \\ \vdots \\ nf \end{pmatrix} \end{pmatrix}.$$

*Credit Scoring* aims at **estimating**  $p(y|\mathbf{x})$  in the form of a simple **parametric** model  $p_\theta(y|\mathbf{x})$  such as logistic regression:



## Context and notations: available data

The observed data are the following:

$$\mathcal{T} = \begin{pmatrix} \mathcal{T}_f \\ \mathcal{T}_{nf} \end{pmatrix} = \begin{pmatrix} \begin{pmatrix} x_f & \begin{pmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \vdots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{pmatrix} & \begin{pmatrix} y_f & \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} & \begin{pmatrix} z_f & \begin{pmatrix} f \\ \vdots \\ f \end{pmatrix} \end{pmatrix} \\ \begin{pmatrix} x_{nf} & \begin{pmatrix} x_{n+1,1} & \cdots & x_{n+1,d} \\ \vdots & \vdots & \vdots \\ x_{n+n',1} & \cdots & x_{n+n',d} \end{pmatrix} & \begin{pmatrix} y_{nf} & \begin{pmatrix} NA \\ \vdots \\ NA \end{pmatrix} & \begin{pmatrix} z_{nf} & \begin{pmatrix} nf \\ \vdots \\ nf \end{pmatrix} \end{pmatrix} \end{pmatrix}.$$

*Credit Scoring* aims at **estimating**  $p(y|\mathbf{x})$  in the form of a simple **parametric** model  $p_\theta(y|\mathbf{x})$  such as logistic regression:

$$\ln \frac{p_\theta(1|\mathbf{x})}{1 - p_\theta(1|\mathbf{x})} = (1, \mathbf{x})' \theta.$$

# Table of Contents

Reject Inference

Feature quantization

Segmentation: logistic regression trees

Missing data imputation

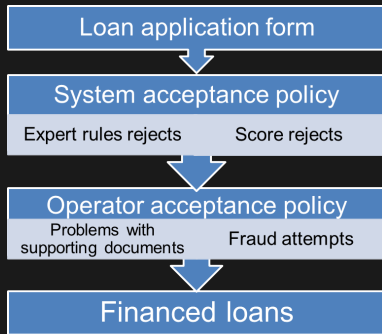
Carbon risk

NLP for extra-financial reports

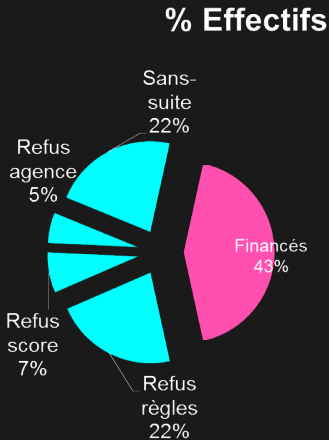
Conclusion and future work

Reject Inference

# Reject Inference: industrial setting



**Figure:** Simplified financing mechanism at Crédit Agricole Consumer Finance



**Figure:** Proportion of “final” lending decisions for CACF France

# Reject Inference: industrial setting

The industry traditionally fits a logistic regression using only  
modelling constraint

financed clients (**fixed parameter space**  $\Theta$ ):

convenience and lack  
of better procedure

$$\hat{\theta}_f = \operatorname{argmax}_{\theta} \ell(\theta; \mathcal{T}_f) = \sum_{i=1}^n \ln p_{\theta}(y_i | \mathbf{x}_i),$$

which asymptotically approximates:

$$\theta_f^* = \operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{X}}[\mathrm{KL}(p || p_{\theta}) | \mathbf{Z} = \mathbf{f}].$$

# Reject Inference: industrial setting

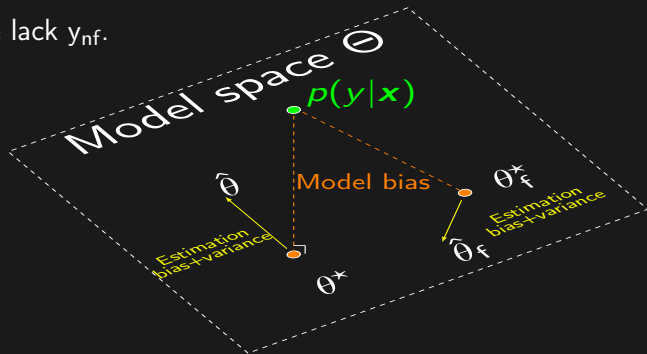
Oracle to be approximated:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{x}} [\text{KL}(p || p_{\theta})] \\ &= \operatorname{argmax}_{\theta} \mathbb{E}_{\mathbf{x}, y \sim p} [\ln p_{\theta}(y | \mathbf{x})],\end{aligned}$$

which standard estimator would be:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \ell(\theta; \mathcal{T}_c),$$

but we lack  $y_{\text{nf}}$ .



## Estimators :

1. “Oracle”:  $\sqrt{n+n'}(\hat{\theta} - \theta_{\text{opt}}) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}})$
2. Current methodology:  $\sqrt{n}(\hat{\theta}^{\text{f}} - \theta_{\text{opt}}^{\text{f}}) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}^{\text{f}}}^{\text{f}})$

---

<sup>1</sup>Zadrozny, “Learning and evaluating classifiers under sample selection bias”.

## Estimators :

1. “Oracle”:  $\sqrt{n + n'}(\hat{\theta} - \theta_{\text{opt}}) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}})$
2. Current methodology:  $\sqrt{n}(\hat{\theta}^{\text{f}} - \theta_{\text{opt}}^{\text{f}}) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}^{\text{f}}}^{\text{f}})$



## Estimators :

1. “Oracle”:  $\sqrt{n+n'}(\hat{\theta} - \theta_{\text{opt}}) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}})$
2. Current methodology:  $\sqrt{n}(\hat{\theta}^f - \theta_{\text{opt}}^f) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}^f}^f)$

---

<sup>1</sup>Zadrozny, “Learning and evaluating classifiers under sample selection bias”.

# Reject Inference: Asymptotics

## Estimators :

1. “Oracle”:  $\sqrt{n + n'}(\hat{\theta} - \theta_{\text{opt}}) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}})$
2. Current methodology:  $\sqrt{n}(\hat{\theta}^f - \theta_{\text{opt}}^f) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}^f}^f)$

What follows will only hold for “local” model which output depends asymptotically only on  $p(y|x)$ , such as logistic regression<sup>1</sup>.

---

<sup>1</sup>Zadrozny, “Learning and evaluating classifiers under sample selection bias”.

# Reject Inference: Asymptotics

## Estimators :

1. “Oracle”:  $\sqrt{n + n'}(\hat{\theta} - \theta_{\text{opt}}) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}})$
2. Current methodology:  $\sqrt{n}(\hat{\theta}^f - \theta_{\text{opt}}^f) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}^f}^f)$

What follows will only hold for “local” model which output depends asymptotically only on  $p(y|x)$ , such as logistic regression<sup>1</sup>.

It can be shown that Bayesian classifiers, SVMs, decision trees are “global” learners<sup>1</sup>.

---

<sup>1</sup>Zadrozny, “Learning and evaluating classifiers under sample selection bias”.

# Reject Inference: modelling the financing mechanism

Due to the financing mechanism, labels  $y$  are not MCAR.

Let  $\{p_\phi(z|\mathbf{x}, y)\}_{\phi \in \Phi}$  denote this hidden financing mechanism (as a parametrized family).

# Reject Inference: modelling the financing mechanism

Due to the financing mechanism, labels  $y$  are not MCAR.

Let  $\{p_\phi(z|\mathbf{x}, y)\}_{\phi \in \Phi}$  denote this hidden financing mechanism (as a parametrized family).

Combining financing and credit-worthiness probability distributions:

$$p_\gamma(y, z|\mathbf{x}) = \underbrace{p_{\theta(\gamma)}(y|\mathbf{x})}_{\text{GCA}} \underbrace{p_{\phi(\gamma)}(z|\mathbf{x}, y)}_{?}.$$

# Reject Inference: modelling the financing mechanism

Due to the financing mechanism, labels  $y$  are not MCAR.

Let  $\{p_\phi(z|\mathbf{x}, y)\}_{\phi \in \Phi}$  denote this hidden financing mechanism (as a parametrized family).

Combining financing and credit-worthiness probability distributions:

$$p_\gamma(y, z|\mathbf{x}) = \underbrace{p_{\theta(\gamma)}(y|\mathbf{x})}_{\text{GCA}} \underbrace{p_{\phi(\gamma)}(z|\mathbf{x}, y)}_{?}.$$

To estimate  $\gamma$ , we could rely on Maximum Likelihood theory:

# Reject Inference: modelling the financing mechanism

Due to the financing mechanism, labels  $y$  are not MCAR.

Let  $\{p_\phi(z|\mathbf{x}, y)\}_{\phi \in \Phi}$  denote this hidden financing mechanism (as a parametrized family).

Combining financing and credit-worthiness probability distributions:

$$p_\gamma(y, z|\mathbf{x}) = \underbrace{p_{\theta(\gamma)}(y|\mathbf{x})}_{\text{GCA}} \underbrace{p_{\phi(\gamma)}(z|\mathbf{x}, y)}_{?}.$$

To estimate  $\gamma$ , we could rely on Maximum Likelihood theory:

$$\ell(\gamma; \mathcal{T}) = \sum_{i=1}^n \ln p_\gamma(y_i, f|\mathbf{x}_i) + \sum_{i=n+1}^{n+n'} \ln \sum_{y \in \{0,1\}} p_\gamma(y, nf|\mathbf{x}_i).$$

# Reject Inference: flawed model selection

No free lunch: **financial or statistical investment to make.**

Because no test-sample  $\mathcal{T}^{\text{test}}$  is available from  $p(\mathbf{x}, y)$ ,  
we cannot resort to error-rate criteria:

$$\text{Error}(\mathcal{T}^{\text{test}}) = \frac{1}{|\mathcal{T}^{\text{test}}|} \sum_{i \in \mathcal{T}^{\text{test}}} \mathbb{I}(\hat{y}_i \neq y_i).$$

~~funding bad clients~~  
~~at a loss~~



# Reject Inference: flawed model selection

No free lunch: **financial or statistical investment to make.**

Because no test-sample  $\mathcal{T}^{\text{test}}$  is available from  $p(\mathbf{x}, y)$ ,  
we cannot resort to error-rate criteria: ~~funding bad clients~~

$$\text{Error}(\mathcal{T}^{\text{test}}) = \frac{1}{|\mathcal{T}^{\text{test}}|} \sum_{i \in \mathcal{T}^{\text{test}}} \mathbb{I}(\hat{y}_i \neq y_i). \quad \text{at a loss}$$

We should use information criteria on the observed data  $\mathcal{T}$  such as:

$$\text{BIC}(\hat{\gamma}; \mathcal{T}) = -2\ell(\hat{\gamma}; \mathcal{T}) + \dim(\mathbf{\Gamma}) \ln n,$$

where  $\hat{\gamma} = \operatorname{argmax}_{\gamma} \ell(\gamma; \mathcal{T})$ , to compare models.

# Reject Inference: flawed model selection

No free lunch: **financial or statistical investment to make.**

Because no test-sample  $\mathcal{T}^{\text{test}}$  is available from  $p(\mathbf{x}, y)$ ,  
we cannot resort to error-rate criteria: ~~funding bad clients~~

$$\text{Error}(\mathcal{T}^{\text{test}}) = \frac{1}{|\mathcal{T}^{\text{test}}|} \sum_{i \in \mathcal{T}^{\text{test}}} \mathbb{I}(\hat{y}_i \neq y_i). \quad \text{at a loss}$$

We should use information criteria on the observed data  $\mathcal{T}$  such as:

$$\text{BIC}(\hat{\gamma}; \mathcal{T}) = -2\ell(\hat{\gamma}; \mathcal{T}) + \dim(\Gamma) \ln n,$$

where  $\hat{\gamma} = \operatorname{argmax}_{\gamma} \ell(\gamma; \mathcal{T})$ , to compare models.

It requires to precisely state the models  $\{p_{\gamma}(y, z | \mathbf{x})\}_{\Gamma}$  that compete and their underlying assumptions.

# Reject Inference: strategies

We gathered 6 so-called Reject Inference methods from the literature that aim at re-injecting  $x_{nf}$  into the estimation procedure of  $\theta$ .

They usually resemble EM-like algorithms:

$$\mathcal{T}_c^{(1)} = \left( \begin{pmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \vdots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \\ x_{n+1,1} & \cdots & x_{n+1,d} \\ \vdots & \vdots & \vdots \\ x_{n+n',1} & \cdots & x_{n+n',d} \end{pmatrix}, \begin{pmatrix} y_1 \\ \vdots \\ y_n \\ \hat{y}_{n+1}^{(1)} \\ \vdots \\ \hat{y}_{n+n'}^{(1)} \end{pmatrix}, \begin{pmatrix} f \\ \vdots \\ f \\ nf \\ \vdots \\ nf \end{pmatrix} \right)$$

Can we **reinterpret these empirical methods** in the missing data and information criterion frameworks and / or expose their **implicit modelling** steps?

## Reject Inference: example of Fuzzy Augmentation<sup>2</sup>

Estimate  $\hat{\theta}_f = \operatorname{argmax}_{\theta} \ell(\theta; \mathcal{T}_f)$ , infer for  $n + 1 \leq i \leq n + n'$ :

$$\hat{y}_i = p_{\hat{\theta}_f}(1|\mathbf{x}_i),$$

---

<sup>2</sup>Nguyen, [Reject inference in application scorecards.](#)

## Reject Inference: example of Fuzzy Augmentation<sup>2</sup>

Estimate  $\hat{\theta}_f = \operatorname{argmax}_{\theta} \ell(\theta; \mathcal{T}_f)$ , infer for  $n + 1 \leq i \leq n + n'$ :

$$\hat{y}_i = p_{\hat{\theta}_f}(1|\mathbf{x}_i),$$

and re-estimate  $\theta$  using the resulting  $\mathcal{T}_c$ . For  $1 \leq j \leq d$ :

$$\frac{\partial \sum_{i=n+1}^{n'+n} \sum_{y_i=0}^1 p_{\hat{\theta}_f}(y_i|\mathbf{x}_i) \ln(p_{\theta}(y_i|\mathbf{x}_i))}{\partial \theta_j} = 0 \Leftrightarrow \theta = \hat{\theta}_f,$$

---

<sup>2</sup>Nguyen, [Reject inference in application scorecards.](#)

## Reject Inference: example of Fuzzy Augmentation<sup>2</sup>

Estimate  $\hat{\theta}_f = \operatorname{argmax}_{\theta} \ell(\theta; \mathcal{T}_f)$ , infer for  $n+1 \leq i \leq n+n'$ :

$$\hat{y}_i = p_{\hat{\theta}_f}(1|\mathbf{x}_i),$$

and re-estimate  $\theta$  using the resulting  $\mathcal{T}_c$ . For  $1 \leq j \leq d$ :

$$\frac{\partial \sum_{i=n+1}^{n'+n} \sum_{y_i=0}^1 p_{\hat{\theta}_f}(y_i|\mathbf{x}_i) \ln(p_{\theta}(y_i|\mathbf{x}_i))}{\partial \theta_j} = 0 \Leftrightarrow \theta = \hat{\theta}_f,$$

such that:

$$\operatorname{argmax}_{\theta \in \Theta} \sum_{i=n+1}^{n'+n} \sum_{y_i=0}^1 p_{\hat{\theta}_f}(y_i|\mathbf{x}_i) \ln(p_{\theta}(y_i|\mathbf{x}_i)) = \hat{\theta}_f.$$

---

<sup>2</sup>Nguyen, [Reject inference in application scorecards.](#)

## Reject Inference: example of Fuzzy Augmentation<sup>2</sup>

Estimate  $\hat{\theta}_f = \operatorname{argmax}_{\theta} \ell(\theta; \mathcal{T}_f)$ , infer for  $n+1 \leq i \leq n+n'$ :

$$\hat{y}_i = p_{\hat{\theta}_f}(1|\mathbf{x}_i),$$

and re-estimate  $\theta$  using the resulting  $\mathcal{T}_c$ . For  $1 \leq j \leq d$ :

$$\frac{\partial \sum_{i=n+1}^{n'+n} \sum_{y_i=0}^1 p_{\hat{\theta}_f}(y_i|\mathbf{x}_i) \ln(p_{\theta}(y_i|\mathbf{x}_i))}{\partial \theta_j} = 0 \Leftrightarrow \theta = \hat{\theta}_f,$$

such that:

$$\operatorname{argmax}_{\theta \in \Theta} \sum_{i=n+1}^{n'+n} \sum_{y_i=0}^1 p_{\hat{\theta}_f}(y_i|\mathbf{x}_i) \ln(p_{\theta}(y_i|\mathbf{x}_i)) = \hat{\theta}_f.$$

Finally:

$$\operatorname{argmax}_{\theta \in \Theta} \ell(\theta; \mathcal{T}_c) = \operatorname{argmax}_{\theta \in \Theta} \ell(\theta; \mathcal{T}_f) = \hat{\theta}_f.$$

---

<sup>2</sup>Nguyen, [Reject inference in application scorecards](#).

## Reject Inference: missingness mechanism

- ▶ **MAR**<sup>3</sup>:  $\forall \mathbf{x}, y, z, p(z|\mathbf{x}, y) = p(z|\mathbf{x})$   
→ Financing is determined by an old score:  $Z = \mathbb{1}_{\{(1, \mathbf{x})' \boldsymbol{\theta} > \text{cut}\}}$ .

---

<sup>3</sup>Little and Rubin, Statistical analysis with missing data.

<sup>4</sup>Molenberghs et al., "Every missingness not at random model has a missingness at random counterpart with equal fit".



# Reject Inference: missingness mechanism

- ▶ **MAR**<sup>3</sup>:  $\forall \mathbf{x}, y, z, p(z|\mathbf{x}, y) = p(z|\mathbf{x})$   
→ Financing is determined by an old score:  $Z = \mathbb{1}_{\{(1, \mathbf{x})' \boldsymbol{\theta} > \text{cut}\}}$ .
- ▶ **MNAR**<sup>3</sup>:  $\exists \mathbf{x}, y, z, p(z|\mathbf{x}, y) \neq p(z|\mathbf{x})$   
→ Operators' hidden "feeling"  $\tilde{\mathbf{X}}$  influence the financing.  
→ Expert rules based on both present and hidden features  $\mathbf{X}$  and  $\tilde{\mathbf{X}}$  resp. where  $\tilde{\mathbf{X}}$  cannot be totally explained by  $\mathbf{X}$ .  
→ Cannot be tested<sup>4</sup>.

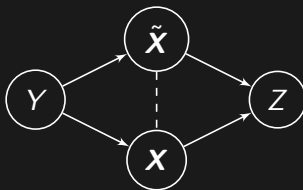
---

<sup>3</sup>Little and Rubin, [Statistical analysis with missing data](#).

<sup>4</sup>Molenberghs et al., "Every missingness not at random model has a missingness at random counterpart with equal fit".

# Reject Inference: missingness mechanism

- ▶ **MAR**<sup>3</sup>:  $\forall \mathbf{x}, y, z, p(z|\mathbf{x}, y) = p(z|\mathbf{x})$   
→ Financing is determined by an old score:  $Z = \mathbb{1}_{\{(1, \mathbf{x})' \theta > \text{cut}\}}$ .
- ▶ **MNAR**<sup>3</sup>:  $\exists \mathbf{x}, y, z, p(z|\mathbf{x}, y) \neq p(z|\mathbf{x})$   
→ Operators' hidden "feeling"  $\tilde{\mathbf{X}}$  influence the financing.  
→ Expert rules based on both present and hidden features  $\mathbf{X}$  and  $\tilde{\mathbf{X}}$  resp. where  $\tilde{\mathbf{X}}$  cannot be totally explained by  $\mathbf{X}$ .  
→ Cannot be tested<sup>4</sup>.



<sup>3</sup>Little and Rubin, [Statistical analysis with missing data](#).

<sup>4</sup>Molenberghs et al., "Every missingness not at random model has a missingness at random counterpart with equal fit".

# Reject Inference: research contribution

Fuzzy Augmentation and Twins produce **the same coefficient  $\hat{\theta}_f$** .

Reclassification<sup>5,6,7</sup> is equivalent to a Classification-EM algorithm, thus introducing a **bias** in the estimation of  $\theta$ .

	MAR	MNAR
Well-specified model	$\hat{\theta}_f$ is unbiased.	$\hat{\theta}_f$ is biased.
Misspecified model	$\hat{\theta}_f$ is biased: Augmentation <sup>2,5,6,7</sup> could be suitable but introduces a <b>new estimation</b> procedure <sup>8</sup> (which requires $\forall \mathbf{x}, p(f \mathbf{x}) > 0$ ).	Any correction relies on <i>a priori</i> <b>unverifiable assumptions</b> about $p_\phi(\mathbf{z} \mathbf{x}, y)$ , e.g. the Parcelling <sup>5,6,7</sup> method.

<sup>5</sup>Guizani et al., “Une Comparaison de quatre Techniques d'Inférence des Refusés dans le Processus d'Octroi de Crédit”.

<sup>6</sup>Soulié and Viennet, “Le Traitement des Refusés dans le Risque Crédit”.

<sup>7</sup>Banasik and Crook, “Reject inference, augmentation, and sample selection”.

<sup>8</sup>Zadrozny, “Learning and evaluating classifiers under sample selection bias”.

## Reject Inference: augmentation

For “local” misspecified models and “global” models:

$$\begin{aligned}\mathbb{E}_{\mathbf{x},y}[\ln[p_{\theta}(y|\mathbf{x})]] &= \sum_{y=0}^1 \int_{\mathcal{X}} \ln p_{\theta}(y|\mathbf{x}) p(y|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= \sum_{y=0}^1 \int_{\mathcal{X}} p(f) \ln p_{\theta}(y|\mathbf{x}) \frac{p(\mathbf{x}|f)}{p(f|\mathbf{x})} p(y|\mathbf{x}) d\mathbf{x} \\ &= \sum_{y=0}^1 \int_{\mathcal{X}} p(f) \frac{\ln p_{\theta}(y|\mathbf{x})}{p(f|\mathbf{x})} p(\mathbf{x}, y|f) d\mathbf{x} \\ &\approx \frac{1}{n} \sum_{i \in \mathcal{T}_f} \frac{p(f)}{p(f|\mathbf{x}_i)} \ln p_{\theta}(y_i|\mathbf{x}_i).\end{aligned}$$

# Reject Inference: augmentation

For “local” misspecified models and “global” models:

$$\begin{aligned}\mathbb{E}_{\mathbf{x},y}[\ln[p_{\theta}(y|\mathbf{x})]] &= \sum_{y=0}^1 \int_{\mathcal{X}} \ln p_{\theta}(y|\mathbf{x}) p(y|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= \sum_{y=0}^1 \int_{\mathcal{X}} p(f) \ln p_{\theta}(y|\mathbf{x}) \frac{p(\mathbf{x}|f)}{p(f|\mathbf{x})} p(y|\mathbf{x}) d\mathbf{x} \\ &= \sum_{y=0}^1 \int_{\mathcal{X}} p(f) \frac{\ln p_{\theta}(y|\mathbf{x})}{p(f|\mathbf{x})} p(\mathbf{x}, y|f) d\mathbf{x} \\ &\approx \frac{1}{n} \sum_{i \in \mathcal{T}_f} \frac{p(f)}{p(f|\mathbf{x}_i)} \ln p_{\theta}(y_i|\mathbf{x}_i).\end{aligned}$$

This assumes  $p(f|\mathbf{x}) > 0 \forall \mathbf{x}$ , which is wrong.

# Reject Inference: augmentation

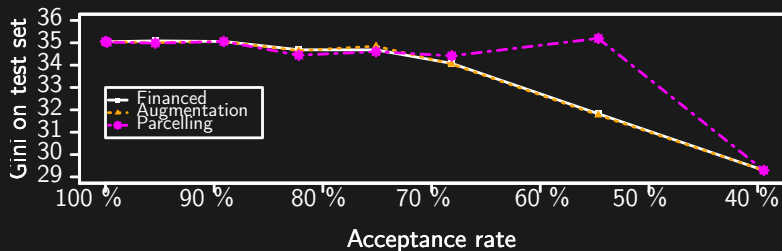
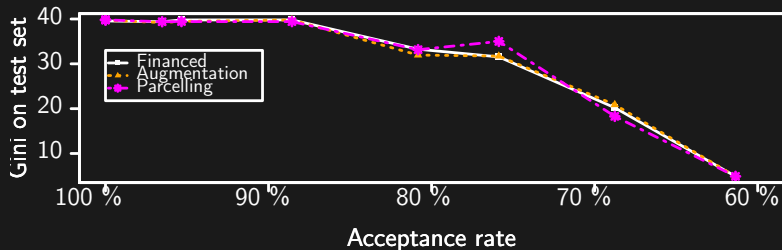
For “local” misspecified models and “global” models:

$$\begin{aligned}\mathbb{E}_{\mathbf{x},y}[\ln[p_{\theta}(y|\mathbf{x})]] &= \sum_{y=0}^1 \int_{\mathcal{X}} \ln p_{\theta}(y|\mathbf{x}) p(y|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= \sum_{y=0}^1 \int_{\mathcal{X}} p(f) \ln p_{\theta}(y|\mathbf{x}) \frac{p(\mathbf{x}|f)}{p(f|\mathbf{x})} p(y|\mathbf{x}) d\mathbf{x} \\ &= \sum_{y=0}^1 \int_{\mathcal{X}} p(f) \frac{\ln p_{\theta}(y|\mathbf{x})}{p(f|\mathbf{x})} p(\mathbf{x}, y|f) d\mathbf{x} \\ &\approx \frac{1}{n} \sum_{i \in \mathcal{T}_f} \frac{p(f)}{p(f|\mathbf{x}_i)} \ln p_{\theta}(y_i|\mathbf{x}_i).\end{aligned}$$

This assumes  $p(f|\mathbf{x}) > 0 \forall \mathbf{x}$ , which is wrong.

Further, one needs to specify / model  $p(f|\mathbf{x})$ .

# Reject Inference: industry contribution



## Feature quantization



## Feature quantization: by an example

For theoretical reasons: bias-variance tradeoff.

# Feature quantization: some more notations I

For practical reasons: interpretability, outliers...  
... at the expense of the statistician's time.

## Quantized data

$$\mathbf{q}(\mathbf{x}) = (\mathbf{q}_1(x_1), \dots, \mathbf{q}_d(x_d))$$

$$\mathbf{q}_j(x_j) = (q_{j,h}(x_j))_1^{m_j} \text{ (one-hot encoding)}$$

$$q_{j,h}(\cdot) = \mathbb{1}(x_j \in C_{j,h}), 1 \leq h \leq m_j$$

# Feature quantization: some more notations II

Quantization is model selection (illustrated here with BIC).

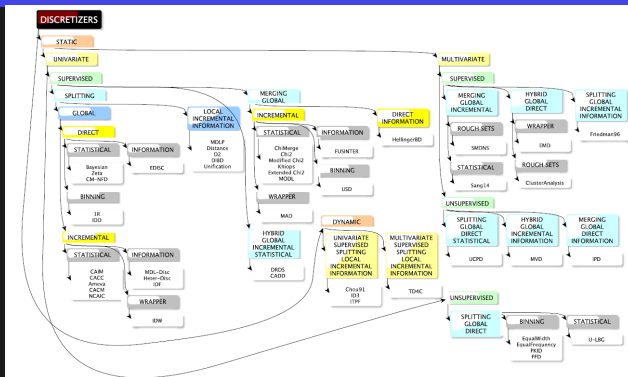
Oracle

$$\begin{aligned}\boldsymbol{\theta}^*, \mathbf{q}^* &= \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta_{\mathbf{q}}, \mathbf{q} \in \mathcal{Q}} \mathbb{E}_{\mathbf{x}, y} [\ln p_{\boldsymbol{\theta}}(y | \mathbf{q}(\mathbf{x}))], \\ \hat{\boldsymbol{\theta}}^{\text{BIC}}, \hat{\mathbf{q}}^{\text{BIC}} &= \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta_{\mathbf{q}}, \mathbf{q} \in \mathcal{Q}} \text{BIC}(\hat{\boldsymbol{\theta}}_{\mathbf{q}}; y_{\text{f}}, \mathbf{q}(\mathbf{x}_{\text{f}})), \\ &\text{where } \hat{\boldsymbol{\theta}}_{\mathbf{q}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta_{\mathbf{q}}} \ell(\boldsymbol{\theta}; y_{\text{f}}, \mathbf{q}(\mathbf{x}_{\text{f}})).\end{aligned}$$

Implicitly assumes quantizations are “well” separated.

Quantization becomes an algorithmic problem.

# Feature quantization: existing approaches



These approaches<sup>9</sup> maximize an “intermediary” criterion, e.g.:

$$\hat{\mathbf{q}}_j^{\chi^2} = \operatorname{argmax}_{\mathbf{q}_j} \chi^2(\mathbf{q}_j(\mathbf{x}_f), y_f) \stackrel{?}{\approx} \mathbf{q}_j^*,$$

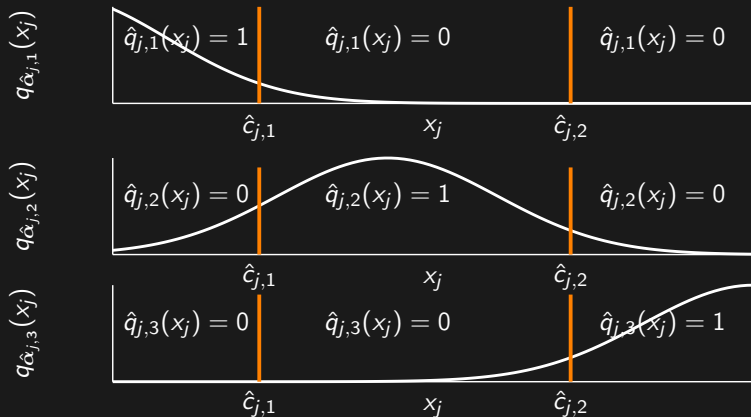
and we hope that it’s aligned with our original goal s.t.:

$$\hat{\theta}^{\chi^2} = \operatorname{argmax}_{\theta} \ell(\theta; y_f, \hat{\mathbf{q}}^{\chi^2}(\mathbf{x}_f)) \stackrel{?}{\approx} \theta^*.$$

<sup>9</sup>Ramirez-Gallego et al., “Data Discretization: Taxonomy and Big Data Challenge”.

# Feature quantization: MAP estimation

$$\hat{q}_{j,h}(x_j) = 1 \text{ if } h = \operatorname{argmax}_{1 \leq h' \leq m_j} q_{\hat{\alpha}_{j,h'}}, 0 \text{ otherwise}^{10,11}.$$



<sup>10</sup>Chamroukhi et al., "A regression model with a hidden logistic process for feature extraction from time series".

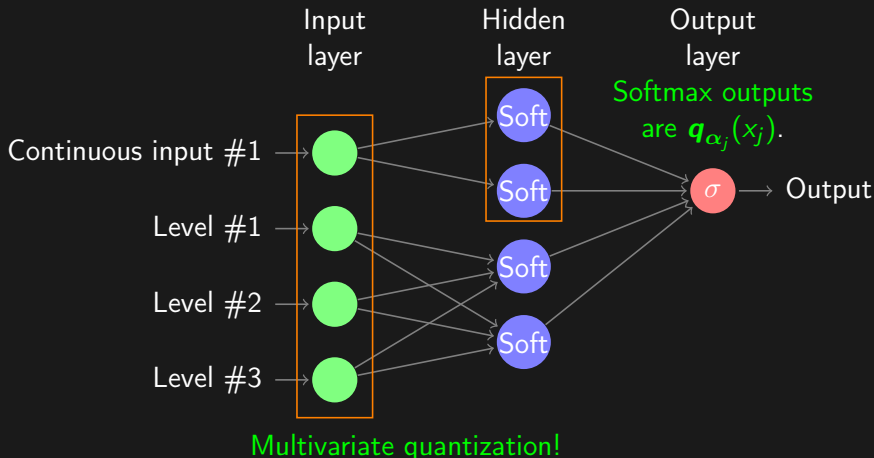
<sup>11</sup>Samé et al., "Model-based clustering and segmentation of time series with changes in regime".

# Feature quantization: neural networks

Very simple neural network.

Very fast implementations available, e.g. TensorFlow.

No guarantee of global optimum (but works well in practice).



# Feature quantization: neural networks

# Feature quantization: results

## Simulated data

**Table:** For different sample sizes  $n$ , (A) CI of  $\hat{c}_{j,2}$  for  $c_{j,2} = 2/3$ . (B) CI of  $\hat{m}$  for  $m_1 = 3$ . (C) CI of  $\hat{m}_3$  for  $m_3 = 1$ .

$n$	(A) $\hat{c}_{j,2}$	(B) $\hat{m}_1$	(C) $\hat{m}_3$
1,000	[0.656, 0.666]	1	60
		90	32
		9	8
10,000	[0.666, 0.666]	0	88
		100	12
		0	0



# Feature quantization: results

## CACF data

**Table:** Gini indices (the greater the value, the better the performance) of our proposed quantization algorithm *glmdisc*, the two baselines and the current scorecard.

Portfolio	ALLR	Current performance	<i>ad hoc</i> methods	Our proposal: <i>glmdisc</i> -NN	Our proposal: <i>glmdisc</i> -SEM	<i>glmdisc</i> -SEM w. interactions
Automobile	59.3 (3.1)	55.6 (3.4)	59.3 (3.0)	58.9 (2.6)	57.8 (2.9)	<b>64.8</b> (2.0)
Renovation	52.3 (5.5)	50.9 (5.6)	54.0 (5.1)	<b>56.7</b> (4.8)	55.5 (5.2)	55.5 (5.2)
Standard	39.7 (3.3)	37.1 (3.8)	45.3 (3.1)	43.8 (3.2)	36.7 (3.7)	<b>47.2</b> (2.8)
Revolving	62.7 (2.8)	58.5 (3.2)	63.2 (2.8)	62.3 (2.8)	60.7 (2.8)	<b>67.2</b> (2.5)
Mass retail	52.8 (5.3)	48.7 (6.0)	61.4 (4.7)	<b>61.8</b> (4.6)	61.0 (4.7)	60.3 (4.8)
Electronics	52.9 (11.9)	55.8 (10.8)	56.3 (10.2)	<b>72.6</b> (7.4)	62.0 (9.5)	63.7 (9.0)

Segmentation: logistic regression trees

# Segmentation: logistic regression trees

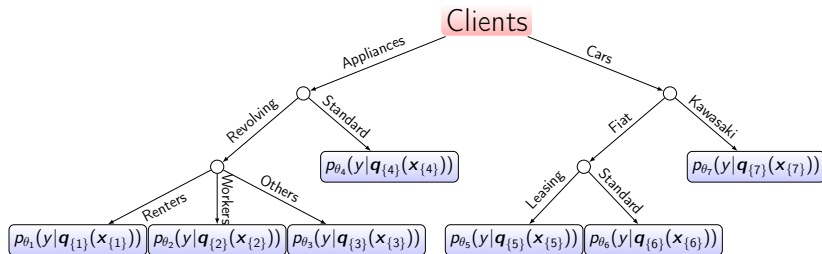


Figure: Scorecards tree structure in acceptance system.

## Segmentation: logistic regression trees

Current procedure(s):

# Segmentation: logistic regression trees

Current procedure(s):

- ▶ Promise a new partner their own score to maximize acceptance;

# Segmentation: logistic regression trees

Current procedure(s):

- ▶ Promise a new partner their own score to maximize acceptance;
- ▶ Merge existing “close” branches that show similar performance;

# Segmentation: logistic regression trees

Current procedure(s):

- ▶ Promise a new partner their own score to maximize acceptance;
- ▶ Merge existing “close” branches that show similar performance;
- ▶ Try basic “clustering” techniques, e.g. visual separation of the data and / or levels on the two first MCA axes.

# Segmentation: logistic regression trees

Current procedure(s):

- ▶ Promise a new partner their own score to maximize acceptance;
- ▶ Merge existing “close” branches that show similar performance;
- ▶ Try basic “clustering” techniques, e.g. visual separation of the data and / or levels on the two first MCA axes.

Problem(s):



# Segmentation: logistic regression trees

## Current procedure(s):

- ▶ Promise a new partner their own score to maximize acceptance;
- ▶ Merge existing “close” branches that show similar performance;
- ▶ Try basic “clustering” techniques, e.g. visual separation of the data and / or levels on the two first MCA axes.

## Problem(s):

- ▶ This structure is not the result of optimization and is probably suboptimal (by how much?);

# Segmentation: logistic regression trees

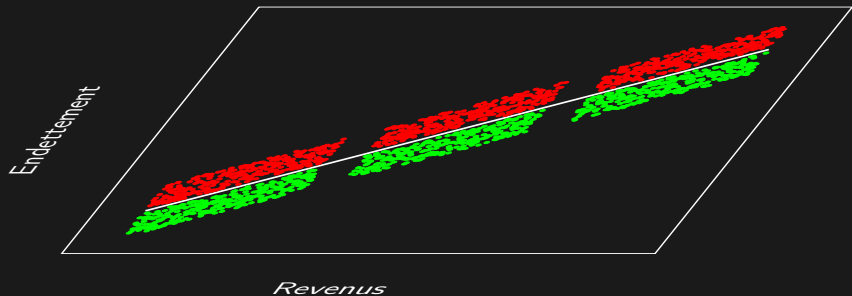
## Current procedure(s):

- ▶ Promise a new partner their own score to maximize acceptance;
- ▶ Merge existing “close” branches that show similar performance;
- ▶ Try basic “clustering” techniques, e.g. visual separation of the data and / or levels on the two first MCA axes.

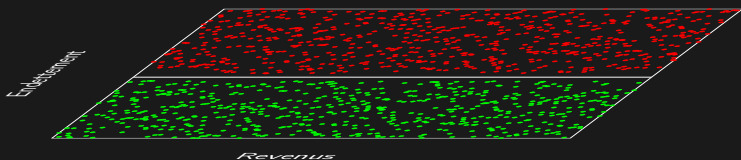
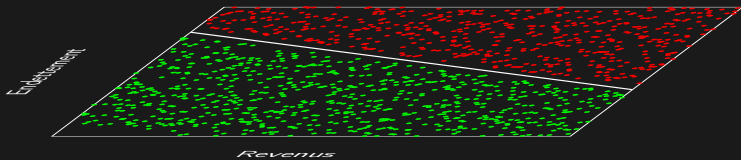
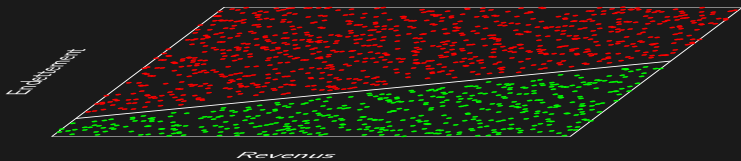
## Problem(s):

- ▶ This structure is not the result of optimization and is probably suboptimal (by how much?);
- ▶ There are situations in which it severely fails.

## Segmentation: logistic regression trees



# Segmentation: logistic regression trees



## Segmentation: logistic regression trees: contribution

Similarly to the quantization proposal: **ability to be in several segments at a time.**

## Segmentation: logistic regression trees: contribution

Similarly to the quantization proposal: **ability to be in several segments at a time.**

$$p(y|\mathbf{x}) = \sum_{c=1}^K \underbrace{p_{\theta}(y|\mathbf{x}; c)}_{\text{"optimized" GCA constraint}} \underbrace{p_{\beta}(c|\mathbf{x})}_{\text{"unoptimized" relaxed CACF constraint}},$$

where  $p_{\beta}(c|\mathbf{x})$  is given by the classification tree as the proportion of training samples in each leaf (**not** majority vote).

## Segmentation: logistic regression trees: contribution

Similarly to the quantization proposal: **ability to be in several segments at a time.**

$$p(y|\mathbf{x}) = \sum_{c=1}^K \underbrace{p_{\theta}(y|\mathbf{x}; c)}_{\text{"optimized" GCA constraint}} \underbrace{p_{\beta}(c|\mathbf{x})}_{\text{"unoptimized" relaxed CACF constraint}},$$

where  $p_{\beta}(c|\mathbf{x})$  is given by the classification tree as the proportion of training samples in each leaf (**not** majority vote).

$$c_i^{(s)} \sim p_{\theta \cdot (s-1)}(y_i | \mathbf{x}_i; \cdot) p_{\beta(s-1)}(\cdot | \mathbf{x}_i).$$

# Segmentation: logistic regression trees: contribution

Similarly to the quantization proposal: **ability to be in several segments at a time.**

$$p(y|\mathbf{x}) = \sum_{c=1}^K \underbrace{p_{\theta}(y|\mathbf{x}; c)}_{\text{"optimized" GCA constraint}} \underbrace{p_{\beta}(c|\mathbf{x})}_{\text{"unoptimized" relaxed CACF constraint}},$$

where  $p_{\beta}(c|\mathbf{x})$  is given by the classification tree as the proportion of training samples in each leaf (**not** majority vote).

$$c_i^{(s)} \sim p_{\theta^{(s-1)}}(y_i|\mathbf{x}_i; \cdot) p_{\beta^{(s-1)}}(\cdot|\mathbf{x}_i).$$

$$\theta^{c(s)} = \operatorname{argmax}_{\theta^c} \sum_{i=1}^n \mathbb{1}_c(c_i^{(s)}) \ln p_{\theta^c}(y_i|\mathbf{x}_i; c_i).$$



# Segmentation: logistic regression trees: contribution

Similarly to the quantization proposal: **ability to be in several segments at a time.**

$$p(y|\mathbf{x}) = \sum_{c=1}^K \underbrace{p_{\theta}(y|\mathbf{x}; c)}_{\text{"optimized" GCA constraint}} \underbrace{p_{\beta}(c|\mathbf{x})}_{\text{"unoptimized" relaxed CACF constraint}},$$

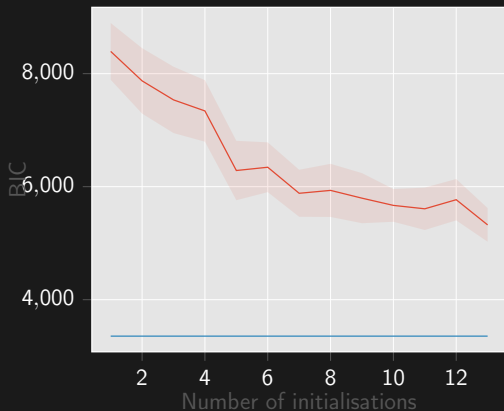
where  $p_{\beta}(c|\mathbf{x})$  is given by the classification tree as the proportion of training samples in each leaf (**not** majority vote).

$$c_i^{(s)} \sim p_{\theta^{(s-1)}}(y_i|\mathbf{x}_i; \cdot) p_{\beta^{(s-1)}}(\cdot|\mathbf{x}_i).$$

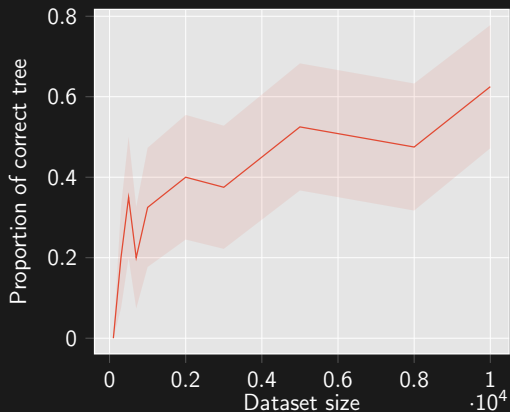
$$\theta^{c(s)} = \operatorname{argmax}_{\theta^c} \sum_{i=1}^n \mathbb{1}_c(c_i^{(s)}) \ln p_{\theta^c}(y_i|\mathbf{x}_i; c_i).$$

$$\beta^{(s)} = \text{C4.5}(c^{(s)}, \mathbf{x}).$$

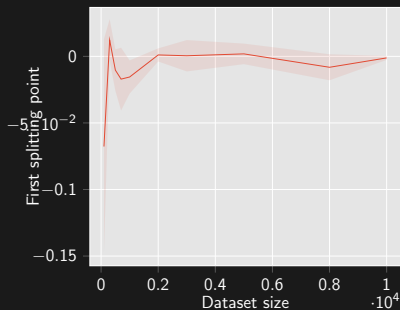
## Segmentation: logistic regression trees: some results



## Segmentation: logistic regression trees: some results



## Segmentation: logistic regression trees: some results



## Segmentation: logistic regression trees: some results

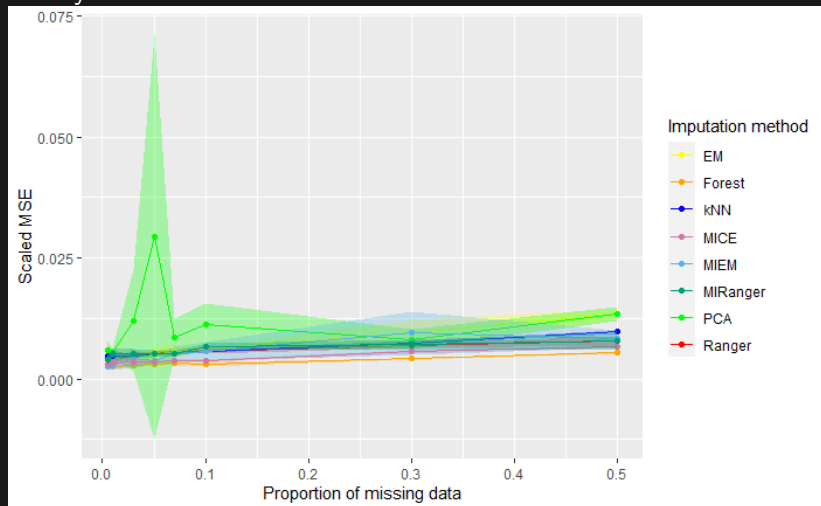
	Logistic regression	Decision Tree	SEM	Gradient Boosting
AUC ( $\pm$ vs current method)	-3,02	-2,66	-1,78	-0,17

	SEM	LMT	MOB
# segment (current: 9)	2	11	1
AUC ( $\pm$ vs current method)	-1,52	-7,70	-5,21

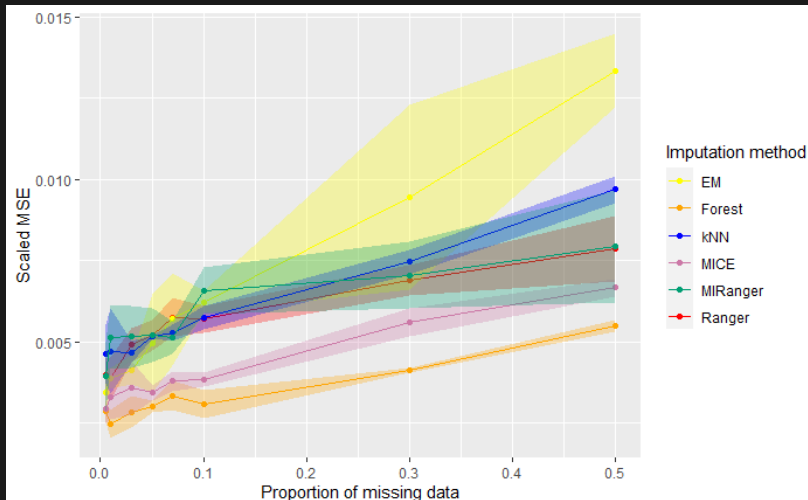
## Missing data imputation

# Missing data imputation: some results I

Research internship: comparing missing data imputation methods, mostly in MAR situations.

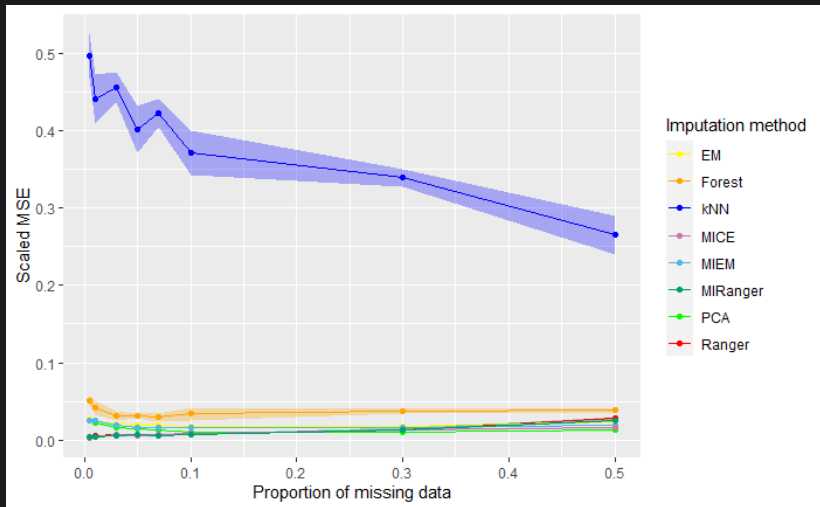


# Missing data imputation: some results II

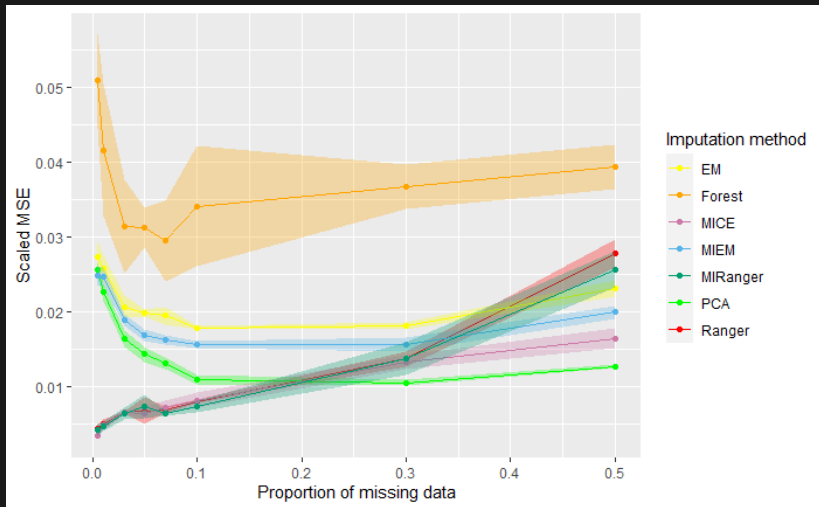




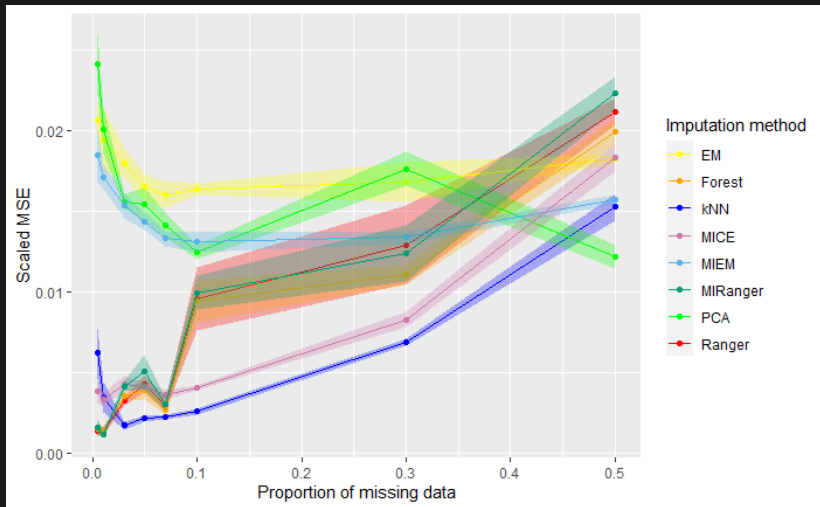
# Missing data imputation: some results III



# Missing data imputation: some results IV



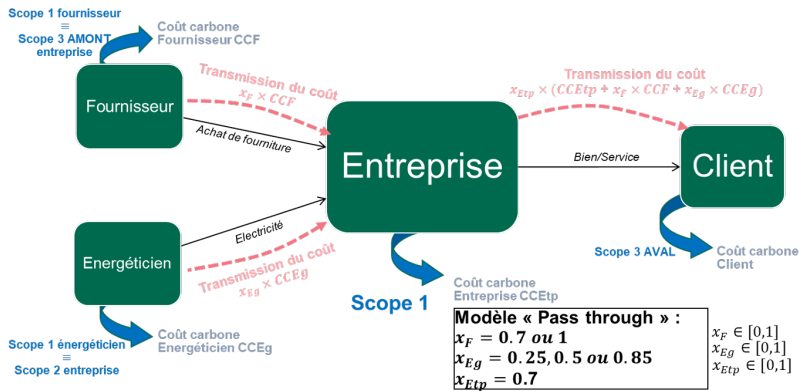
# Missing data imputation: some results V



## Carbon risk

# Carbon risk: some results

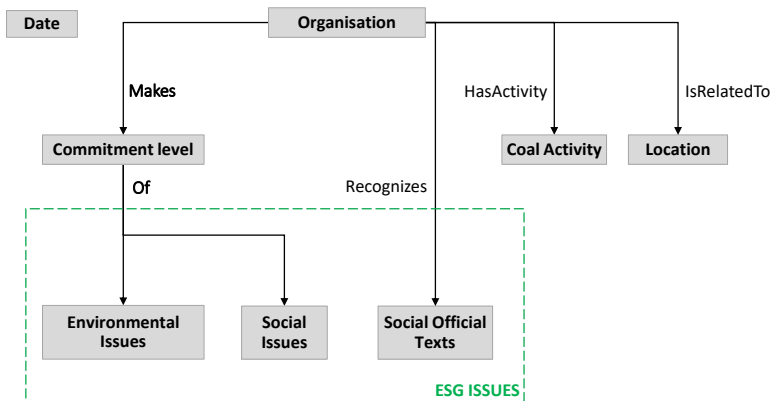
Research internship: use carbon price scenarios to impact the earnings of big corporations and adjust their default probability accordingly.



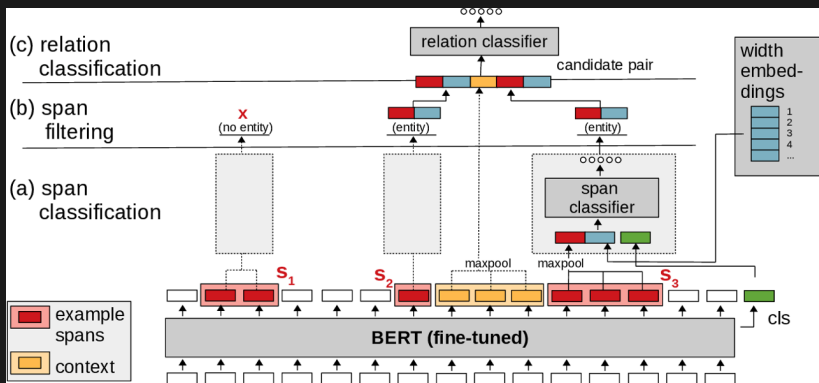
## NLP for extra-financial reports

# NLP for extra-financial reports: some results I

Research internship: build joint NER and RE models to automatically read through extra-financial reports.



# NLP for extra-financial reports: some results II





## Conclusion and future work

# Conclusions from my PhD

This PhD tackled three main issues of “traditional” Credit Scoring:

1. Reject inference: impact of tossing away not-financed clients,

# Conclusions from my PhD

This PhD tackled three main issues of “traditional” Credit Scoring:

1. Reject inference: impact of tossing away not-financed clients,

Conclusion: sound problem reformulation, no method recommended, `scoringTools` R package.

# Conclusions from my PhD

This PhD tackled three main issues of “traditional” Credit Scoring:

1. Reject inference: impact of tossing away not-financed clients,
2. “Constrained” representation learning: discretization, grouping, interaction screening,

# Conclusions from my PhD

This PhD tackled three main issues of “traditional” Credit Scoring:

1. Reject inference: impact of tossing away not-financed clients,
2. “Constrained” representation learning: discretization, grouping, interaction screening,

Conclusion: better performance, less time-consuming, glmdisc R and Python packages.

# Conclusions from my PhD

This PhD tackled three main issues of “traditional” Credit Scoring:

1. Reject inference: impact of tossing away not-financed clients,
2. “Constrained” representation learning: discretization, grouping, interaction screening,
3. Predictive segmentation: logistic regression trees,

# Conclusions from my PhD

This PhD tackled three main issues of “traditional” Credit Scoring:

1. Reject inference: impact of tossing away not-financed clients,
2. “Constrained” representation learning: discretization, grouping, interaction screening,
3. Predictive segmentation: logistic regression trees,

Conclusion: first experiments on simulated and real data are encouraging, glmtree R package.

# Future work as presented for my PhD - might be helpful?

There remains a lot of open questions:

1. Credit Scoring for profit: swap “ $p(2 \text{ unpaid instalments})$ ” for  $p(\text{profit} > 0)$  or  $\mathbb{E}[\text{profit}]$ ,



# Future work as presented for my PhD - might be helpful?

There remains a lot of open questions:

1. Credit Scoring for profit: swap “ $p(2 \text{ unpaid instalments})$ ” for  $p(\text{profit} > 0)$  or  $\mathbb{E}[\text{profit}]$ ,

Perspective: experiment observation-wise misclassification costs.

# Future work as presented for my PhD - might be helpful?

There remains a lot of open questions:

1. Credit Scoring for profit: swap “ $p(2 \text{ unpaid instalments})$ ” for  $p(\text{profit} > 0)$  or  $\mathbb{E}[\text{profit}]$ ,
2. Representation learning for fine-grained unstructured data,

# Future work as presented for my PhD - might be helpful?

There remains a lot of open questions:

1. Credit Scoring for profit: swap “ $p(2 \text{ unpaid instalments})$ ” for  $p(\text{profit} > 0)$  or  $\mathbb{E}[\text{profit}]$ ,
2. Representation learning for fine-grained unstructured data,  
Perspective: provide statistically sound methods to aggregate “behavioural” data, e.g. web visitation patterns.

Thanks!

- [1] John Banasik and Jonathan Crook. “Reject inference, augmentation, and sample selection”. In: European Journal of Operational Research 183.3 (2007), pp. 1582–1594. url: <http://www.sciencedirect.com/science/article/pii/S0377221706011969> (visited on 08/25/2016).
- [2] Faicel Chamroukhi et al. “A regression model with a hidden logistic process for feature extraction from time series”. In: International Joint Conference on Neural Networks, 2009. IJCNN 2009. IEEE. 2009, pp. 489–496.
- [3] Asma Guizani et al. “Une Comparaison de quatre Techniques d’Inférence des Refusés dans le Processus d’Octroi de Crédit”. In: 45 èmes Journées de statistique. 2013. url: [http://cedric.cnam.fr/fichiers/art\\_2753.pdf](http://cedric.cnam.fr/fichiers/art_2753.pdf) (visited on 08/25/2016).
- [4] Roderick JA Little and Donald B Rubin. Statistical analysis with missing data. John Wiley & Sons, 2014.

- [5] Geert Molenberghs et al. “Every missingness not at random model has a missingness at random counterpart with equal fit”. In: Journal of the Royal Statistical Society: Series B 7.2 (2008), pp. 371–388.
- [6] Ha Thu Nguyen. Reject inference in application scorecards: evidence from France. Tech. rep. University of Paris West-Nanterre la Défense, EconomiX, 2016. url: [http://economix.fr/pdf/dt/2016/WP\\_EcoX\\_2016-10.pdf](http://economix.fr/pdf/dt/2016/WP_EcoX_2016-10.pdf) (visited on 08/25/2016).
- [7] Sergio Ramirez-Gallego et al. “Data Discretization: Taxonomy and Big Data Challenge”. In: Wiley Int. Rev. Data Min. and Knowl. Disc. 6.1 (Jan. 2016), pp. 5–21. issn: 1942-4787. doi: [10.1002/widm.1173](https://doi.org/10.1002/widm.1173). url: <http://dx.doi.org/10.1002/widm.1173>.

- [8] Allou Samé et al. “Model-based clustering and segmentation of time series with changes in regime”. In: Advances in Data Analysis and Classification 5.4 (2011), pp. 301–321.
- [9] Françoise Fogelman Soulié and Emmanuel Viennet. “Le Traitement des Refusés dans le Risque Crédit”. In: Revue des Nouvelles Technologies de l’Information Data Mining et Apprentissage Statistique : application en assurance, banque et marketing, RNTI-A-1 (2007), pp. 22–44.
- [10] Bianca Zadrozny. “Learning and evaluating classifiers under sample selection bias”. In: Proceedings of the twenty-first ICML. ACM. 2004, p. 114.

# Quantization



“Soft” approximation:

$$\mathbf{q}_{\alpha_j}(\cdot) = \left( q_{\alpha_{j,h}}(\cdot) \right)_{h=1}^{m_j} \text{ with } \begin{cases} \sum_{h=1}^{m_j} q_{\alpha_{j,h}}(\cdot) = 1, \\ 0 \leq q_{\alpha_{j,h}}(\cdot) \leq 1, \end{cases}$$

“Soft” approximation:

$$\mathbf{q}_{\alpha_j}(\cdot) = \left( q_{\alpha_{j,h}}(\cdot) \right)_{h=1}^{m_j} \quad \text{with} \quad \begin{cases} \sum_{h=1}^{m_j} q_{\alpha_{j,h}}(\cdot) = 1, \\ 0 \leq q_{\alpha_{j,h}}(\cdot) \leq 1, \end{cases}$$

For continuous features, we set for  $\alpha_{j,h} = (\alpha_{j,h}^0, \alpha_{j,h}^1) \in \mathbb{R}^2$

$$q_{\alpha_{j,h}}(\cdot) = \frac{\exp(\alpha_{j,h}^0 + \alpha_{j,h}^1 \cdot)}{\sum_{g=1}^{m_j} \exp(\alpha_{j,g}^0 + \alpha_{j,g}^1 \cdot)}.$$

“Soft” approximation:

$$q_{\alpha_j}(\cdot) = \left( q_{\alpha_{j,h}}(\cdot) \right)_{h=1}^{m_j} \text{ with } \begin{cases} \sum_{h=1}^{m_j} q_{\alpha_{j,h}}(\cdot) = 1, \\ 0 \leq q_{\alpha_{j,h}}(\cdot) \leq 1, \end{cases}$$

For continuous features, we set for  $\alpha_{j,h} = (\alpha_{j,h}^0, \alpha_{j,h}^1) \in \mathbb{R}^2$

$$q_{\alpha_{j,h}}(\cdot) = \frac{\exp(\alpha_{j,h}^0 + \alpha_{j,h}^1 \cdot)}{\sum_{g=1}^{m_j} \exp(\alpha_{j,g}^0 + \alpha_{j,g}^1 \cdot)}.$$

For categorical features, we set for

$$\alpha_{j,h} = (\alpha_{j,h}(1), \dots, \alpha_{j,h}(l_j)) \in \mathbb{R}^{l_j}$$

$$q_{\alpha_{j,h}}(\cdot) = \frac{\exp(\alpha_{j,h}(\cdot))}{\sum_{g=1}^{m_j} \exp(\alpha_{j,g}(\cdot))}.$$

# Quantization: research contribution

We wish to maximize the following likelihood:

$$(\hat{\theta}, \hat{\alpha}) = \operatorname{argmax}_{\theta, \alpha} \ell(\theta, \alpha; \mathbf{x}_f, \mathbf{y}_f) = \operatorname{argmax}_{\theta, \alpha} \sum_{i=1}^n \ln p_{\theta}(y_i | \mathbf{q}_{\alpha}(\mathbf{x}_i)).$$

" $\alpha^* = \lim_{n \rightarrow \infty} \hat{\alpha}$ " should be such that  $\mathbf{q}_{\alpha^*} = \mathbf{q}^*$ .

We wish to maximize the following likelihood:

$$(\hat{\theta}, \hat{\alpha}) = \operatorname{argmax}_{\theta, \alpha} \ell(\theta, \alpha; \mathbf{x}_f, \mathbf{y}_f) = \operatorname{argmax}_{\theta, \alpha} \sum_{i=1}^n \ln p_{\theta}(y_i | \mathbf{q}_{\alpha}(\mathbf{x}_i)).$$

“ $\alpha^* = \lim_{n \rightarrow \infty} \hat{\alpha}$ ” should be such that  $\mathbf{q}_{\alpha^*} = \mathbf{q}^*$ .

Problem:  $\hat{\alpha}$  has to **diverge**, the MLE is at the border of the parameter space which could hinder its properties.

# Quantization: research contribution

We wish to maximize the following likelihood:

$$(\hat{\theta}, \hat{\alpha}) = \operatorname{argmax}_{\theta, \alpha} \ell(\theta, \alpha; \mathbf{x}_f, \mathbf{y}_f) = \operatorname{argmax}_{\theta, \alpha} \sum_{i=1}^n \ln p_{\theta}(y_i | \mathbf{q}_{\alpha}(\mathbf{x}_i)).$$

“ $\alpha^* = \lim_{n \rightarrow \infty} \hat{\alpha}$ ” should be such that  $\mathbf{q}_{\alpha^*} = \mathbf{q}^*$ .

Problem:  $\hat{\alpha}$  has to **diverge**, the MLE is at the border of the parameter space which could hinder its properties.

Anyway, or more generally if there is no true quantization  $\mathbf{q}^*$ ,  $\hat{\mathbf{q}}$  is used instead as a **quantization candidate**.

# Quantization: research contribution

We wish to maximize the following likelihood:

$$(\hat{\theta}, \hat{\alpha}) = \operatorname{argmax}_{\theta, \alpha} \ell(\theta, \alpha; \mathbf{x}_f, \mathbf{y}_f) = \operatorname{argmax}_{\theta, \alpha} \sum_{i=1}^n \ln p_{\theta}(y_i | \mathbf{q}_{\alpha}(\mathbf{x}_i)).$$

“ $\alpha^* = \lim_{n \rightarrow \infty} \hat{\alpha}$ ” should be such that  $\mathbf{q}_{\alpha^*} = \mathbf{q}^*$ .

Problem:  $\hat{\alpha}$  has to **diverge**, the MLE is at the border of the parameter space which could hinder its properties.

Anyway, or more generally if there is no true quantization  $\mathbf{q}^*$ ,  $\hat{\mathbf{q}}$  is used instead as a **quantization candidate**.

**Problem:**  $\ell(\theta, \alpha; \mathbf{x}_f, \mathbf{y}_f)$  cannot be directly maximized.

**Solution:** Resort to (stochastic) gradient descent which **each step (s) will yield**  $\hat{\alpha}^{(s)}$  and **quantization candidate**  $\hat{\mathbf{q}}^{(s)}$ .

# Quantization: model = quantization selection

Quantization provider to original selection criterion

We have **drastically restricted the search space** to *iter* well-chosen candidates resulting from the the gradient descent steps.

$$s^* = \operatorname{argmin}_{s=1,\dots,iter} \operatorname{BIC}(\hat{\theta}_{\hat{\mathbf{q}}^{(s)}})$$



# Quantization: model = quantization selection

Quantization provider to original selection criterion

We have **drastically restricted the search space** to *iter* well-chosen candidates resulting from the the gradient descent steps.

$$s^* = \operatorname{argmin}_{s=1,\dots,iter} \operatorname{BIC}(\hat{\theta}_{\hat{q}^{(s)}})$$

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

# Quantization: model = quantization selection

Quantization provider to original selection criterion

We have **drastically restricted the search space** to *iter* well-chosen candidates resulting from the the gradient descent steps.

$$s^* = \operatorname{argmin}_{s=1,\dots,iter} \operatorname{BIC}(\hat{\theta}_{\hat{q}^{(s)}})$$

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

In practice if  $\forall i, q_{\alpha_{j,h}}(x_j) \ll 1$ , then level *h* disappears while performing the argmax.

# Quantization: model = quantization selection

Quantization provider to original selection criterion

We have **drastically restricted the search space** to *iter* well-chosen candidates resulting from the the gradient descent steps.

$$s^* = \operatorname{argmin}_{s=1,\dots,iter} \operatorname{BIC}(\hat{\theta}_{\hat{\mathbf{q}}^{(s)}})$$

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

In practice if  $\forall i, q_{\alpha_{j,h}}(x_j) \ll 1$ , then level *h* disappears while performing the argmax.

Start with ***m*** =  $(m_{\max})_1^d$  and “wait” ...

## Bivariate interactions

## Bivariate interactions: notations

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

$$\text{logit}(p_{\theta}(1|\mathbf{q}(\mathbf{x}))) = \theta_0 + \sum_{j=1}^d \theta_j^{\mathbf{q}_j(x_j)} + \sum_{1 \leq k < \ell \leq d} \delta_{k,\ell} \theta_{k,\ell}^{\mathbf{q}_k(x_k) \mathbf{q}_\ell(x_\ell)}.$$

## Bivariate interactions: notations

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

$$\text{logit}(p_{\theta}(1|\mathbf{q}(\mathbf{x}))) = \theta_0 + \sum_{j=1}^d \theta_j^{\mathbf{q}_j(x_j)} + \sum_{1 \leq k < \ell \leq d} \delta_{k,\ell} \theta_{k,\ell}^{\mathbf{q}_k(x_k) \mathbf{q}_\ell(x_\ell)}.$$

Imagine for now that the discretization  $\mathbf{q}(\mathbf{x})$  is fixed. The criterion becomes:

$$(\boldsymbol{\theta}^*, \boldsymbol{\delta}^*) = \underset{\boldsymbol{\theta}, \boldsymbol{\delta} \in \{0,1\}^{\frac{d(d-1)}{2}}}{\text{argmin}} \quad \text{BIC}(\hat{\boldsymbol{\theta}}_{\boldsymbol{\delta}}; \mathcal{T}_{\text{f}}).$$

## Bivariate interactions: notations

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

$$\text{logit}(p_{\theta}(1|\mathbf{q}(\mathbf{x}))) = \theta_0 + \sum_{j=1}^d \theta_j^{\mathbf{q}_j(x_j)} + \sum_{1 \leq k < \ell \leq d} \delta_{k,\ell} \theta_{k,\ell}^{\mathbf{q}_k(x_k) \mathbf{q}_\ell(x_\ell)}.$$

Imagine for now that the discretization  $\mathbf{q}(\mathbf{x})$  is fixed. The criterion becomes:

$$(\boldsymbol{\theta}^*, \boldsymbol{\delta}^*) = \underset{\boldsymbol{\theta}, \boldsymbol{\delta} \in \{0,1\}^{\frac{d(d-1)}{2}}}{\text{argmin}} \quad \text{BIC}(\hat{\boldsymbol{\theta}}_{\boldsymbol{\delta}}; \mathcal{T}_f).$$

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

## Bivariate interactions: model proposal

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



## Bivariate interactions: model proposal

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

Strategy used here: Metropolis-Hastings sampling algorithm.

**Idea:** propose well-chosen interactions and accept / reject them based on the BIC criterion of the resulting logistic regression.

# Bivariate interactions: model proposal

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

Strategy used here: Metropolis-Hastings sampling algorithm.

**Idea:** propose well-chosen interactions and accept / reject them based on the BIC criterion of the resulting logistic regression.

$$p(y|\mathbf{q}) = \sum_{\delta \in \{0,1\}^{\frac{d(d-1)}{2}}} p(y|\mathbf{q}, \delta) p(\delta)$$
$$p(\delta|\mathbf{q}, y) \propto \exp(-\text{BIC}[\delta]/2) p(\delta)$$

## Bivariate interactions: model proposal

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

Strategy used here: Metropolis-Hastings sampling algorithm.

**Idea:** propose well-chosen interactions and accept / reject them based on the BIC criterion of the resulting logistic regression.

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

$$p(\delta|\mathbf{q}, y) \propto \exp(-\text{BIC}[\delta]/2) p(\delta) \quad p(\delta_{p,q}) = \frac{1}{2}$$

# Bivariate interactions: model proposal

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

Strategy used here: Metropolis-Hastings sampling algorithm.

**Idea:** propose well-chosen interactions and accept / reject them based on the BIC criterion of the resulting logistic regression.

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

$$p(\delta|\mathbf{q}, y) \propto \exp(-\text{BIC}[\delta]/2) p(\delta) \quad p(\delta_{p,q}) = \frac{1}{2}$$

Which transition proposal  $T : (\{0,1\}^{\frac{d(d-1)}{2}}, \{0,1\}^{\frac{d(d-1)}{2}}) \mapsto [0; 1]$ ?

## Bivariate interactions: model proposal

$2^{d(d-1)}$  probabilities to calculate. . .

# Bivariate interactions: model proposal

$2^{d(d-1)}$  probabilities to calculate. . .

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

# Bivariate interactions: model proposal

$2^{d(d-1)}$  probabilities to calculate...

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

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

# Bivariate interactions: model proposal

$2^{d(d-1)}$  probabilities to calculate...

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

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

If the interaction between two features is meaningful when only these two features are considered, there is a good chance that it will be in the full multivariate model.



# Bivariate interactions: model proposal

$2^{d(d-1)}$  probabilities to calculate...

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

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

If the interaction between two features is meaningful when only these two features are considered, there is a good chance that it will be in the full multivariate model.

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

## SEM-Gibbs quantization

# SEM-Gibbs quantization

Originally (and as implemented in the R package `glmddisc`), the optimization was a bit different:

Originally (and as implemented in the R package `glmdisc`), the optimization was a bit different:

- ▶  $\mathbf{q}$  is considered a latent (unobserved) feature  $\mathbf{q}$ ;

Originally (and as implemented in the R package `glmdisc`), the optimization was a bit different:

- ▶  $\mathbf{q}$  is considered a latent (unobserved) feature  $\mathbf{q}$ ;
- ▶ A classical EM algorithm is intractable since it requires an Expectation step over all possible quantizations;

Originally (and as implemented in the R package `glmdisc`), the optimization was a bit different:

- ▶  $\mathbf{q}$  is considered a latent (unobserved) feature  $\mathbf{q}$ ;
- ▶ A classical EM algorithm is intractable since it requires an Expectation step over all possible quantizations;
- ▶ Solution: random draw  $\approx$  Bayesian statistics.

# SEM-Gibbs quantization: estimation

“Classical” estimation strategy with latent variables: EM algorithm.

## SEM-Gibbs quantization: estimation

“Classical” estimation strategy with latent variables: EM algorithm.

There would still be a sum over  $\mathcal{Q}_m$ :

$$p(y|\mathbf{x}, \theta, \alpha) = \sum_{\mathbf{q} \in \mathcal{Q}_m} p_{\theta}(y|\mathbf{q}) \prod_{j=1}^d p_{\alpha_j}(\mathbf{q}_j|x_j)$$



## SEM-Gibbs quantization: estimation

“Classical” estimation strategy with latent variables: EM algorithm.

There would still be a sum over  $\mathcal{Q}_m$ :

$$p(y|\mathbf{x}, \theta, \alpha) = \sum_{\mathbf{q} \in \mathcal{Q}_m} p_{\theta}(y|\mathbf{q}) \prod_{j=1}^d p_{\alpha_j}(\mathbf{q}_j|x_j)$$

Use a Stochastic-EM! Draw  $\mathbf{q}$  knowing that:

# SEM-Gibbs quantization: estimation

“Classical” estimation strategy with latent variables: EM algorithm.

There would still be a sum over  $\mathcal{Q}_m$ :

$$p(y|\mathbf{x}, \theta, \alpha) = \sum_{\mathbf{q} \in \mathcal{Q}_m} p_{\theta}(y|\mathbf{q}) \prod_{j=1}^d p_{\alpha_j}(\mathbf{q}_j|x_j)$$

Use a Stochastic-EM! Draw  $\mathbf{q}$  knowing that:

$$p(\mathbf{q}|\mathbf{x}, y) = \frac{p_{\theta}(y|\mathbf{q}) \prod_{j=1}^d p_{\alpha_j}(\mathbf{q}_j|x_j)}{\underbrace{\sum_{\mathbf{q} \in \mathcal{Q}_m} p_{\theta}(y|\mathbf{q}) \prod_{j=1}^d p_{\alpha_j}(\mathbf{q}_j|x_j)}_{\text{still difficult to calculate}}}$$

# SEM-Gibbs quantization: estimation

“Classical” estimation strategy with latent variables: EM algorithm.

There would still be a sum over  $\mathcal{Q}_m$ :

$$p(y|\mathbf{x}, \theta, \alpha) = \sum_{\mathbf{q} \in \mathcal{Q}_m} p_{\theta}(y|\mathbf{q}) \prod_{j=1}^d p_{\alpha_j}(\mathbf{q}_j|x_j)$$

Use a Stochastic-EM! Draw  $\mathbf{q}$  knowing that:

$$p(\mathbf{q}|\mathbf{x}, y) = \frac{p_{\theta}(y|\mathbf{q}) \prod_{j=1}^d p_{\alpha_j}(\mathbf{q}_j|x_j)}{\underbrace{\sum_{\mathbf{q} \in \mathcal{Q}_m} p_{\theta}(y|\mathbf{q}) \prod_{j=1}^d p_{\alpha_j}(\mathbf{q}_j|x_j)}_{\text{still difficult to calculate}}}$$

Gibbs-sampling step:

$$p(\mathbf{q}_j|\mathbf{x}, y, \mathbf{q}_{\{-j\}}) \propto p_{\theta}(y|\mathbf{q}) p_{\alpha_j}(\mathbf{q}_j|x_j)$$

# SEM-Gibbs quantization: algorithm

## Initialization

$$\begin{pmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \vdots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{pmatrix} \xRightarrow{\text{at random}} \begin{pmatrix} q_{1,1} & \cdots & q_{1,d} \\ \vdots & \vdots & \vdots \\ q_{n,1} & \cdots & q_{n,d} \end{pmatrix}$$

## Loop

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \xRightarrow{\text{logistic regression}} \begin{pmatrix} q_{1,1} & \cdots & q_{1,d} \\ \vdots & \vdots & \vdots \\ q_{n,1} & \cdots & q_{n,d} \end{pmatrix} \xRightarrow{\text{polytomous regression}} \begin{pmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \vdots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{pmatrix}$$

## Updating $q$

$$\begin{pmatrix} p(y_1, q_{1,j} = k | x_i) \\ \vdots \\ p(y_n, q_{n,j} = k | x_i) \end{pmatrix} \xRightarrow{\text{random sampling}} \begin{pmatrix} q_{1,j} \\ \vdots \\ q_{n,j} \end{pmatrix}$$

## Calculating $q^{\text{MAP}}$

$$\begin{pmatrix} q^{\text{MAP},1,j} \\ \vdots \\ q^{\text{MAP},n,j} \end{pmatrix} \xRightarrow[\text{estimate}]{\text{MAP}} \begin{pmatrix} \operatorname{argmax}_{q_j} p_{\alpha_j}(q_j | x_{1,j}) \\ \vdots \\ \operatorname{argmax}_{q_j} p_{\alpha_j}(q_j | x_{n,j}) \end{pmatrix}$$

# SEM-Gibbs quantization: simulations