

Some thoughts about current Credit Scoring practices

Adrien Ehrhardt
AGOS Machine Learning Day

20/06/2019



Table of Contents

Context and notations

Reject Inference

Feature quantization

Bivariate interactions

Segmentation: logistic regression trees

Bonus

Context and notations

Context and notations: Industrial setting

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

Table: Dataset with outliers and missing values.

1. Discarding rejected 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			0
?	Renter	10	Common-law	1700			1
Licensed professional	Starter	5	Divorced	4000			0
Executive	By work	8	Single	2700			1
<u>Office employee</u>	<u>Renter</u>	<u>✓</u>	<u>Married</u>	<u>1400</u>			NA
<u>Worker</u>	<u>By family</u>	<u>?</u>	<u>?</u>	<u>1200</u>			NA

Table: Dataset with outliers and missing values.

1. Discarding rejected 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			0
?			Common-law	1700			1
Licensed professional			Divorced	4000			0
Executive			Single	2700			1
<u>Office employee</u>	Renter	12	Married	1400			NA
Worker	By family	?	?	1200			NA

Table: Dataset with outliers and missing values.

1. Discarding rejected 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] $1500;2000]$			0
?			Common-law] $1500;2000]$			1
Licensed professional			Divorced] $2000;\infty[$			0
Executive			Single] $2000;\infty[$			1
<u>Office employee</u>	Renter	$\cancel{12}$	<u>Married</u>	<u>1400</u>			NA
<u>Worker</u>	<u>By family</u>	$\cancel{\frac{1}{2}}$	$\cancel{\frac{1}{2}}$	<u>1200</u>			NA

Table: Dataset with outliers and missing values.

1. Discarding rejected 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]		0
?+Low-qualified			Union]1500;2000]		1
High-qualified			?+Alone]2000;∞[0
High-qualified			?+Alone]2000;∞[1
<u>Office employee</u>	Renter	✓	<u>Married</u>	<u>1400</u>		NA
<u>Worker</u>	By family	✗	✗	<u>1200</u>		NA

Table: Dataset with outliers and missing values.

1. Discarding rejected 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]		0
?+Low-qualified			Union x]1500;2000]		1
High-qualified			?+Alone x]2000;∞[0
High-qualified			?+Alone x]2000;∞[1
<u>Office employee</u>	Renter	12	<u>Married</u> <u>1400</u>		NA
<u>Worker</u>	By family	?	† <u>1200</u>		NA

Table: Dataset with outliers and missing values.

1. Discarding rejected 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]		0
?+Low-qualified			Union x]1500;2000]		1
High-qualified			?+Alone x]2000;∞[0
High-qualified			?+Alone x]2000;∞[1
<u>Office employee</u>	Renter	12	<u>Married</u> <u>1400</u>		NA
<u>Worker</u>	By family	?	† <u>1200</u>		NA

Table: Dataset with outliers and missing values.

1. Discarding rejected 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	Score	Repayment
?+Low-qualified			?+Alone x]1500;2000]	225	0
?+Low-qualified			Union x]1500;2000]	190	1
High-qualified			?+Alone x]2000;∞[218	0
High-qualified			?+Alone x]2000;∞[202	1
<u>Office employee</u>	Renter	12	<u>Married</u> <u>1400</u>	NA	NA
<u>Worker</u>	<u>By family</u>	?	?	1200	NA

Table: Dataset with outliers and missing values.

1. Discarding rejected 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.

$x_j \in \mathbb{R}$ or $\{1, \dots, l_j\}$: e.g. rent amount, job, ...

$y \in \{0, 1\}$: good or bad.

$z \in \{\text{f, nf}\}$: financed or not financed.

Context and notations: Available data

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

Observations:

$\mathbf{x} = (x_1, \dots, x_d)$: characteristics.

$x_j \in \mathbb{R}$ or $\{1, \dots, l_j\}$: e.g. rent amount, job, ...

$y \in \{0, 1\}$: good or bad.

$z \in \{\text{f, nf}\}$: financed or not financed.

True distribution of good and bad clients: $p(y|\mathbf{x})$

Context and notations: Available data

Need for a **computable model** that resembles p , often in the form of a **parametric model** $p_{\theta}(y|\mathbf{x})$, which we can calculate for a new client.

Context and notations: Available data

Need for a **computable model** that resembles p , often in the form of a **parametric model** $p_{\theta}(y|\mathbf{x})$, which we can calculate for a new client.

Example: logistic regression

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

Context and notations: Available data

Need for a **computable model** that resembles p , often in the form of a **parametric model** $p_{\theta}(y|\mathbf{x})$, which we can calculate for a new client.

Example: logistic regression

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

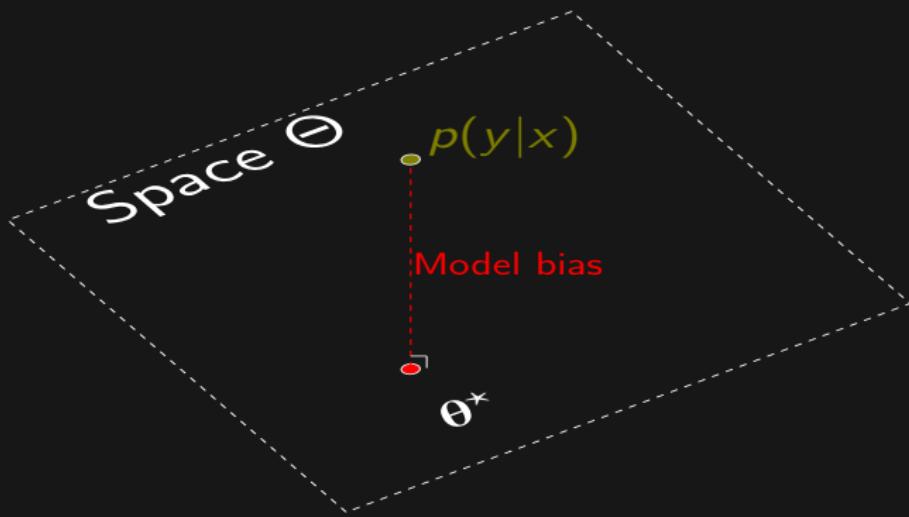
There is $\boldsymbol{\theta}^*$ that makes p_{θ^*} “close” to p .

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x}}[\text{KL}(p||p_{\boldsymbol{\theta}})] = \int_{\mathcal{X}} \sum_{y \in \{0,1\}} p(y|\mathbf{x}) \ln \frac{p(y|\mathbf{x})}{p_{\boldsymbol{\theta}}(y|\mathbf{x})}.$$

Well-specified model assumption

$$\mathbb{E}_{\mathbf{x}}[\text{KL}(p||p_{\boldsymbol{\theta}^*})] = 0,$$
$$p_{\boldsymbol{\theta}^*}(y|\mathbf{x}) = p(y|\mathbf{x}).$$

Context and notations: Available data



Context and notations: Available data

p is unknown: access to an i.i.d. $n + n'$ -sample
 $\mathcal{T} = (\mathbf{x}_i, y_i, z_i)_{1}^{n+n'} \sim p.$

Context and notations: Available data

p is unknown: access to an i.i.d. $n + n'$ -sample
 $\mathcal{T} = (\mathbf{x}_i, y_i, z_i)_{1}^{n+n'} \sim p$.

We can deduce from the KL divergence the (log-)likelihood:

$$\ell(\boldsymbol{\theta}; \mathcal{T}) = \sum_{i=1}^{n+n'} \ln p_{\boldsymbol{\theta}}(y_i | \mathbf{x}_i).$$

Context and notations: Available data

p is unknown: access to an i.i.d. $n + n'$ -sample
 $\mathcal{T} = (\mathbf{x}_i, y_i, z_i)_{1}^{n+n'} \sim p$.

We can deduce from the KL divergence the (log-)likelihood:

$$\ell(\boldsymbol{\theta}; \mathcal{T}) = \sum_{i=1}^{n+n'} \ln p_{\boldsymbol{\theta}}(y_i | \mathbf{x}_i).$$

The MLE $\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{T})$ is a good approximation of $\boldsymbol{\theta}^*$.

Context and notations: Available data

p is unknown: access to an i.i.d. $n + n'$ -sample
 $\mathcal{T} = (\mathbf{x}_i, y_i, z_i)_{1}^{n+n'} \sim p$.

We can deduce from the KL divergence the (log-)likelihood:

$$\ell(\boldsymbol{\theta}; \mathcal{T}) = \sum_{i=1}^{n+n'} \ln p_{\boldsymbol{\theta}}(y_i | \mathbf{x}_i).$$

The MLE $\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{T})$ is a good approximation of $\boldsymbol{\theta}^*$.

Unfortunately, $\hat{\boldsymbol{\theta}}$ is not directly computable (no closed form solution).

Context and notations: Available data

p is unknown: access to an i.i.d. $n + n'$ -sample
 $\mathcal{T} = (\mathbf{x}_i, y_i, z_i)_{1}^{n+n'} \sim p.$

We can deduce from the KL divergence the (log-)likelihood:

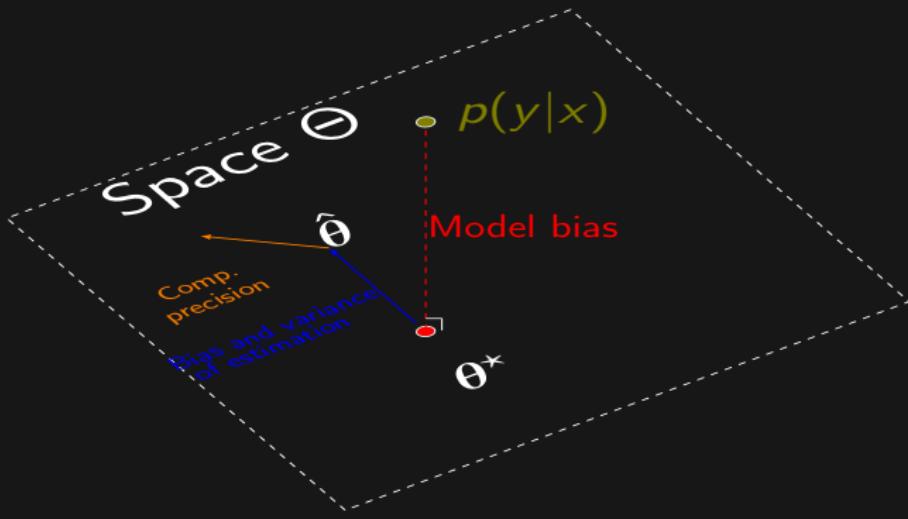
$$\ell(\boldsymbol{\theta}; \mathcal{T}) = \sum_{i=1}^{n+n'} \ln p_{\boldsymbol{\theta}}(y_i | \mathbf{x}_i).$$

The MLE $\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{T})$ is a good approximation of $\boldsymbol{\theta}^*$.

Unfortunately, $\hat{\boldsymbol{\theta}}$ is not directly computable (no closed form solution).

$$\tilde{\boldsymbol{\theta}} = \text{Newton-Raphson}(\ell(\boldsymbol{\theta}; \mathcal{T})) \neq \hat{\boldsymbol{\theta}}.$$

Context and notations: Available data



All of this is “hidden” in your favourite statistical language / package / library but is essential to understanding Reject Inference.

Context and notations: Feature / model selection

Up to now, we assumed a parameter space Θ **fixed**.

Context and notations: Feature / model selection

Up to now, we assumed a parameter space Θ **fixed**.

Comparing models = different parameter spaces $\Theta^1, \Theta^2, \dots$

Corresponding to feature subsets, different discretizations,
interactions, . . . since **we don't know which parameter space is
closest to the "truth" p** .

Context and notations: Feature / model selection

Up to now, we assumed a parameter space Θ **fixed**.

Comparing models = different parameter spaces $\Theta^1, \Theta^2, \dots$

Corresponding to feature subsets, different discretizations,
interactions, . . . since **we don't know which parameter space is closest to the "truth" p** .

Model selection tools

$$\hat{\theta}^{\text{best}} = \operatorname{argmin}_{\hat{\theta}^k \in \Theta^k} \text{BIC}(\hat{\theta}^k) = -2\ell(\hat{\theta}^k, \mathcal{T}) + \dim(\Theta^k) \ln n.$$

Context and notations: Feature / model selection

Up to now, we assumed a parameter space Θ **fixed**.

Comparing models = different parameter spaces $\Theta^1, \Theta^2, \dots$

Corresponding to feature subsets, different discretizations, interactions, . . . since **we don't know which parameter space is closest to the “truth” p** .

Model selection tools

$$\hat{\theta}^{\text{best}} = \underset{\hat{\theta}^k \in \Theta^k}{\operatorname{argmin}} \text{BIC}(\hat{\theta}^k) = -2\ell(\hat{\theta}^k, \mathcal{T}) + \dim(\Theta^k) \ln n.$$

BIC has nice statistical properties (consistency) but can be swapped in the entire presentation with your favourite model selection tool like Gini on $\mathcal{T}^{\text{test}}$.

Reject Inference

Reject Inference: Industrial setting

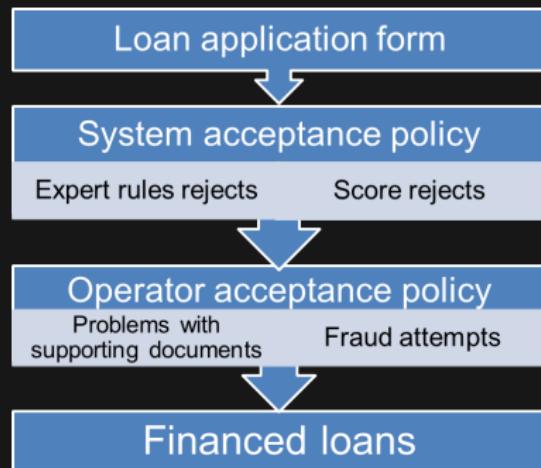


Figure: Simplified Acceptance mechanism in Crédit Agricole Consumer Finance

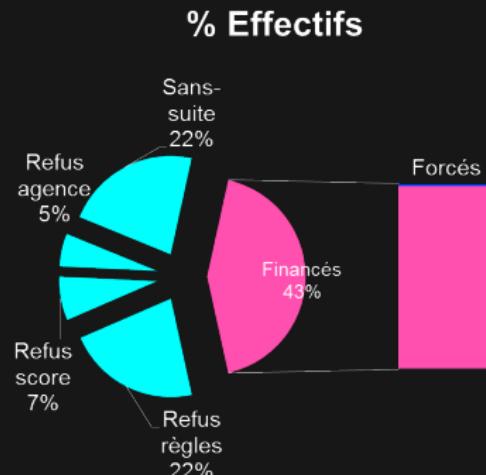


Figure: Proportion of “final” lending decisions for CACF France

Reject Inference: Industrial setting

The observed data are the following:

$$\mathcal{T} = \mathcal{T}_f \cup \mathcal{T}_{nf}$$
$$\mathcal{T}_f = \left(\begin{array}{c} \mathbf{x}_f \\ \boxed{\begin{array}{ccc} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \vdots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \end{array}} \\ \mathbf{y}_f \\ \mathbf{z}_f \\ \mathbf{f} \end{array} \right).$$
$$\mathcal{T}_{nf} = \left(\begin{array}{c} \mathbf{x}_{nf} \\ \boxed{\begin{array}{ccc} x_{n+1,1} & \cdots & x_{n+1,d} \\ \vdots & \vdots & \vdots \\ x_{n+n',1} & \cdots & x_{n+n',d} \end{array}} \\ \mathbf{y}_{nf} \\ \mathbf{z}_{nf} \\ \mathbf{nf} \end{array} \right).$$

We traditionally build a logistic regression using only financed clients (**fixed parameter space** Θ):

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

which asymptotically approximates:

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

Reject Inference: Industrial setting

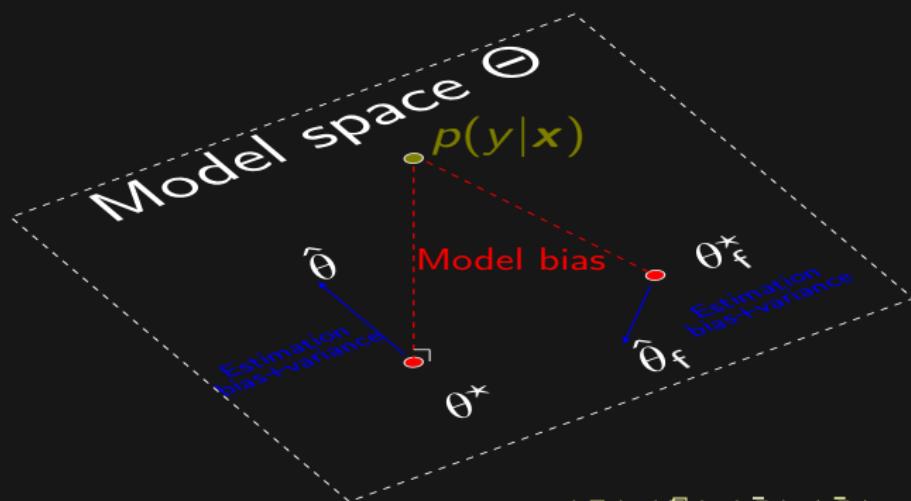
We wish we had:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \ell(\theta; \mathbf{x}, \mathbf{y}),$$

which asymptotically approximates:

$$\theta^* = \operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{x}} [\text{KL}(p || p_{\theta})].$$

But we lack \mathbf{y}_{nf} .



Reject Inference: What is at stake?

Estimators :

1. "Oracle": $\sqrt{n + n'}(\hat{\theta} - \theta^*) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta^*})$
2. Current methodology: $\sqrt{n}(\hat{\theta}_f - \theta_f^*) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{f, \theta_f^*})$

Reject Inference: What is at stake?

Estimators :

1. "Oracle": $\sqrt{n + n'}(\hat{\theta} - \theta^*) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta^*})$
2. Current methodology: $\sqrt{n}(\hat{\theta}_f - \theta_f^*) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{f, \theta_f^*})$

Reject Inference: What is at stake?

Estimators :

1. "Oracle": $\sqrt{n + n'}(\hat{\theta} - \theta^*) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta^*})$
2. Current methodology: $\sqrt{n}(\hat{\theta}_f - \theta_f^*) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{f, \theta_f^*})$

Reject Inference: What is at stake?

Estimators :

1. "Oracle": $\sqrt{n + n'}(\hat{\theta} - \theta^*) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta^*})$
2. Current methodology: $\sqrt{n}(\hat{\theta}_f - \theta_f^*) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{f, \theta_f^*})$

Reject Inference: What is at stake?

Estimators :

1. "Oracle": $\sqrt{n+n'}(\hat{\theta} - \theta^*) \xrightarrow[n, n' \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta^*})$
2. Current methodology: $\sqrt{n}(\hat{\theta}_f - \theta_f^*) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{f, \theta_f^*})$

Question 1 : asymptotics of the estimators

$$\boxed{(Q1) \theta^* \stackrel{?}{=} \theta_f^*}$$
$$\boxed{(Q2) \Sigma_{\theta^*} \stackrel{?}{=} \Sigma_{f, \theta_f^*}}$$

Reject Inference: Missingness mechanism

- ▶ **MAR** : $\forall x, y, z, p(z|x, y) = p(z|x)$
→ Acceptance is determined by an old score: $Z = \mathbb{1}_{\{\theta'X > \text{cut}\}}$.
- ▶ **MNAR** : $\exists x, y, z, p(z|x, y) \neq p(z|x)$
→ Operators' "feeling" $\tilde{\mathbf{X}}$ influence the acceptance.
→ Expert rules based on features $\tilde{\mathbf{X}}$ not in \mathbf{X} .

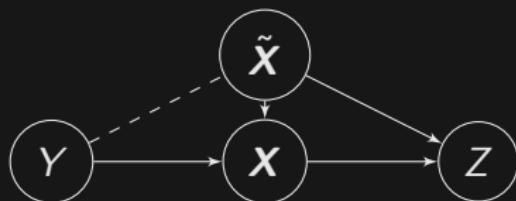


Figure: Dependencies between random variables Y , $\tilde{\mathbf{X}}$, \mathbf{X} and Z

Reject Inference: Model specification

- ▶ **Well-specified model** : $p(y|\mathbf{x}) = p_{\theta^*}(y|\mathbf{x})$.
→ With real data \Rightarrow hypothesis unlikely to be true.
- ▶ **Misspecified model** : θ^* is the “best” in the Θ family.
→ Logistic regression commonly used for its robustness to misspecification (no assumption about $p(\mathbf{x})$).

$p_{\theta}(y \mathbf{x})$	$p(z \mathbf{x}, y)$	MAR	MNAR
Well specified	$\theta_f^* = \theta^*$ $\Sigma_{f,\theta_f^*} \neq \Sigma_{\theta^*}$	$\theta_f^* \neq \theta^*$	$\Sigma_{f,\theta_f^*} \neq \Sigma_{\theta^*}$
Misspecified	$\theta_f^* \neq \theta^*$ $\Sigma_{f,\theta_f^*} \neq \Sigma_{\theta^*}$		

Table: (Q1) and (Q2) w.r.t. model specification and missingness mechanism

Reject Inference: How to use \mathbf{x}_{nf} ?

Question 2: How to construct a better estimator than $\hat{\theta}_f$?

Reject Inference: How to use \mathbf{x}_{nf} ?

Question 2: How to construct a better estimator than $\hat{\theta}_f$?

Scope for action:

- ▶ Change model space Θ ,

Reject Inference: How to use \mathbf{x}_{nf} ?

Question 2: How to construct a better estimator than $\hat{\theta}_f$?

Scope for action:

- ▶ Change model space Θ ,
- ▶ Model acceptance/rejection process (i.e. $p_\gamma(z|\mathbf{x}, y)$),

Reject Inference: How to use \mathbf{x}_{nf} ?

Question 2: How to construct a better estimator than $\hat{\theta}_f$?

Scope for action:

- ▶ Change model space Θ ,
- ▶ Model acceptance/rejection process (i.e. $p_\gamma(z|\mathbf{x}, y)$),
- ▶ Use \mathbf{x}_{nf} .

Reject Inference: How to use \mathbf{x}_{nf} ?

Question 2: How to construct a better estimator than $\hat{\theta}_f$?

Scope for action:

- ▶ Change model space Θ ,
- ▶ Model acceptance/rejection process (i.e. $p_\gamma(z|\mathbf{x}, y)$),
- ▶ Use \mathbf{x}_{nf} .

Natural way to achieve all three: generative approach

$$p_{\alpha}(\mathbf{x}, y, z) = p_{\beta_{\alpha}}(\mathbf{x})p_{\theta_{\alpha}}(y|\mathbf{x})p_{\gamma_{\alpha}}(z|\mathbf{x}, y).$$

$$\begin{aligned} (\boxed{\hat{\theta}_{\alpha}}, \hat{\beta}_{\alpha}, \hat{\gamma}_{\alpha}) &= \underset{\theta_{\alpha}, \beta_{\alpha}, \gamma_{\alpha}}{\operatorname{argmax}} \ell(\alpha; \mathbf{x}, \mathbf{y}_f) = \underset{\theta_{\alpha}, \beta_{\alpha}, \gamma_{\alpha}}{\operatorname{argmax}} \sum_{i=1}^n \ln(p_{\theta_{\alpha}}(y_i|\mathbf{x}_i)) \\ &\quad + \sum_{i=1}^{n+n'} \ln(p_{\beta_{\alpha}}(\mathbf{x}_i)) \left(+ \sum_{i=1}^n \ln(p_{\gamma_{\alpha}}(z_i|\mathbf{x}_i, y_i)) \right). \end{aligned}$$

Reject Inference: How to use \mathbf{x}_{nf} ?

Question 2: How to construct a better estimator than $\hat{\theta}_f$?

Scope for action:

- ▶ Change model space Θ logistic regression,
- ▶ Model acceptance/rejection process (i.e. $p_\gamma(z|\mathbf{x}, y)$),
- ▶ Use \mathbf{x}_{nf} .

Natural way to achieve all three: generative approach

$$p_{\alpha}(\mathbf{x}, y, z) = p_{\beta_{\alpha}}(\mathbf{x})p_{\theta_{\alpha}}(y|\mathbf{x})p_{\gamma_{\alpha}}(z|\mathbf{x}, y).$$

$$(\boxed{\hat{\theta}_{\alpha}}, \hat{\beta}_{\alpha}, \hat{\gamma}_{\alpha}) = \underset{\alpha}{\operatorname{argmax}} \ell(\alpha; \mathbf{x}, \mathbf{y}_f) = \underset{\theta_{\alpha}, \beta_{\alpha}, \gamma_{\alpha}}{\operatorname{argmax}} \sum_{i=1}^n \ln(p_{\theta_{\alpha}}(y_i|x_i))$$

$$+ \sum_{i=1}^{n+n'} \ln(p_{\beta_{\alpha}}(\mathbf{x}_i)) \left(+ \sum_{i=1}^n \ln(p_{\gamma_{\alpha}}(z_i|x_i, y_i)) \right).$$

Reject Inference: How to use \mathbf{x}_{nf} ?

Question 2: How to construct a better estimator than $\hat{\theta}_f$?

Scope for action:

- ▶ Change model space Θ logistic regression,
- ▶ Model acceptance/rejection process (i.e. $p_\gamma(z|\mathbf{x}, y)$)
 γ cannot be estimated,
- ▶ Use \mathbf{x}_{nf} .

Reject Inference: How to use \mathbf{x}_{nf} ?

Question 2: How to construct a better estimator than $\hat{\theta}_f$?

Scope for action:

- ▶ Change model space Θ logistic regression,
- ▶ Model acceptance/rejection process (i.e. $p_\gamma(z|\mathbf{x}, y)$)
 γ cannot be estimated,
- ▶ Use \mathbf{x}_{nf} .

Remember that \mathcal{T}^{OOT} also comes from $p(y|\mathbf{x}, f)$ such that applying a *Reject Inference* method and getting a higher Gini is no guarantee that it would on the Through-the-Door population (on the contrary!).

Reject Inference: How to use \mathbf{x}_{nf} ?

For logistic regression, *Reject Inference* methods amount to:

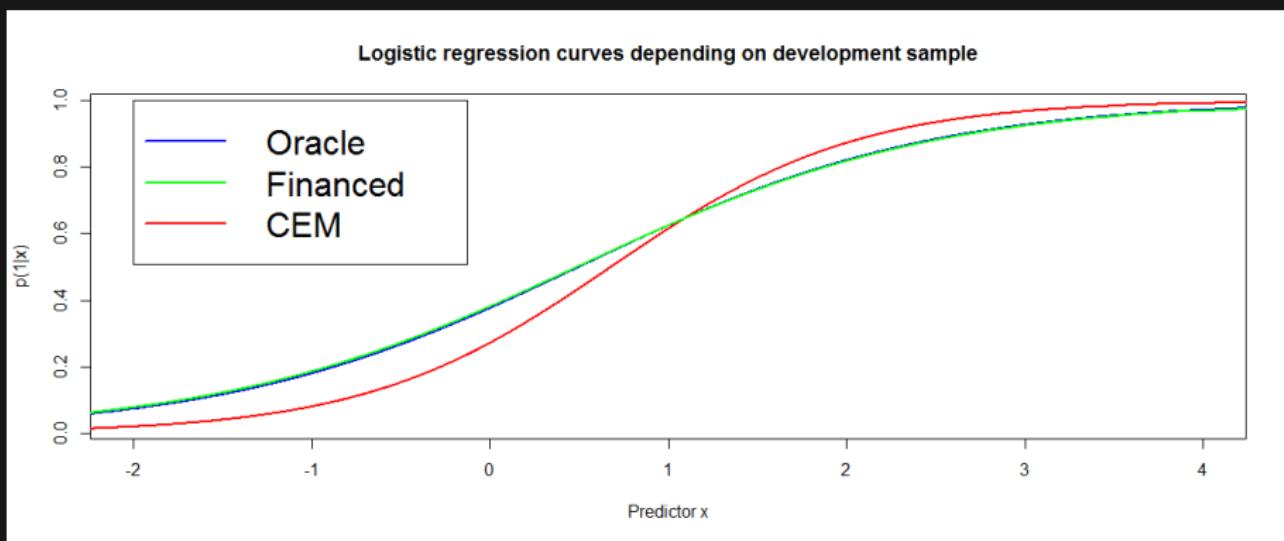
$$\mathcal{T}_c^{(1)} = \begin{matrix} & \mathbf{x}_f \\ \mathbf{x}_{nf} & \end{matrix} \begin{pmatrix} x_{1,1} & \cdots & x_{1,d} \\ \vdots & \vdots & \vdots \\ x_{n,1} & \cdots & x_{n,d} \\ x_{n+1,1} & \cdots & x_{n+1,d} \\ \vdots & \vdots & \vdots \\ x_{n+n',1} & \cdots & x_{n+n',d} \end{pmatrix}, \quad \begin{matrix} & \mathbf{y}_f \\ \mathbf{y}_{nf} & \end{matrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \\ \hat{y}_{n+1}^{(1)} \\ \vdots \\ \hat{y}_{n+n'}^{(1)} \end{pmatrix}, \quad \begin{matrix} & \mathbf{z}_f \\ \mathbf{z}_{nf} & \end{matrix} \begin{pmatrix} f \\ \vdots \\ f \\ nf \\ \vdots \\ nf \end{pmatrix}$$

Reject Inference: How to use \mathbf{x}_{nf} ?

Reclassification¹ :

$$(\hat{\theta}^{\text{CEM}}, \hat{\mathbf{y}}^{\text{nf}}) = \underset{\theta, \mathbf{y}^{\text{nf}}}{\operatorname{argmax}} \ell(\theta; \mathcal{T}_c^{(1)}) \text{ where } \hat{y}_i = \underset{y_i}{\operatorname{argmax}} p_{\hat{\theta}_f}(y_i | \mathbf{x}_i).$$

Problem: inconsistent estimator.



¹[4, 1, 2]

Reject Inference: How to use \mathbf{x}_{nf} ?

Augmentation²: MAR / misspecified model.

$$\ell_{\text{Aug}}(\boldsymbol{\theta}; \mathcal{T}_f) = \sum_{i=1}^n \frac{1}{p(f|\mathbf{x}_i)} \ln(p_{\boldsymbol{\theta}}(y_i|\mathbf{x}_i)).$$

Problem: estimation of $p(f|x_i)$ + assumes $p(f|x_i) > 0$ (clearly not true).

Parcelling³:

$$\ell(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y}_f, \mathbf{\hat{y}}_{\text{nf}}) \text{ where } \hat{y}_i = \begin{cases} 1 \text{ w.p. } \alpha_i p_{\hat{\theta}_f}(1|\mathbf{x}_i, f) \\ 0 \text{ w.p. } 1 - \alpha_i p_{\hat{\theta}_f}(1|\mathbf{x}_i, f) \end{cases}.$$

Problem: MNAR assumptions hidden in $\mathbf{\hat{y}}_{\text{nf}}$ (α_i) impossible to test.

²[4, 1, 2, 3]

³[4, 1, 2]

Reject Inference: Additional remarks

All this stands for logistic regression and all “local” methods [5].

All “global” methods (explicit or implicit modelling of $p(\mathbf{x})$) will produce biased estimates under MAR.

We might have:

Gini	Logistic regression	Decision trees
Financed	40	45
Through-the-door	40	35

Feature quantization

Feature quantization: By an example

Some more notations I

Raw data

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

$x_j \in \mathbb{R}$ (continuous case)

$x_j \in \{1, \dots, l_j\}$ (categorical case)

$y \in \{0, 1\}$ (target)

Quantized data

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

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

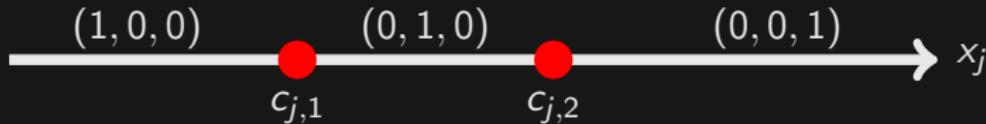
$$q_{j,h}(\cdot) = 1 \text{ if } x_j \in C_{j,h}, 0 \text{ otherwise, } 1 \leq h \leq m_j$$

Some more notations II

Discretization

$$C_{j,h} = (c_{j,h-1}, c_{j,h}]$$

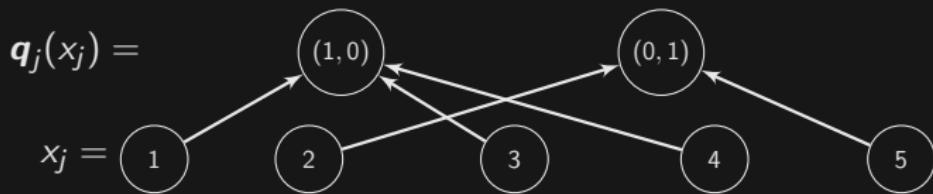
where $c_{j,1}, \dots, c_{j,m_j-1}$ are increasing numbers called cutpoints,
 $c_{j,0} = -\infty$ and $c_{j,m_j} = +\infty$.



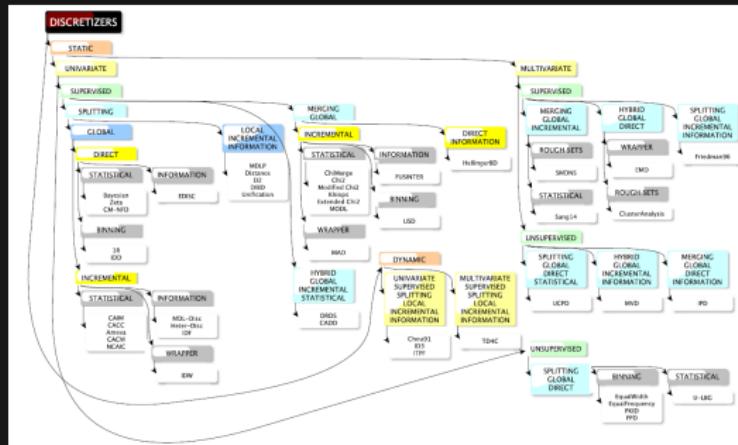
Some more notations III

Grouping

$$\bigsqcup_{h=1}^{m_j} C_{j,h} = \{1, \dots, l_j\}.$$



Feature quantization: Existing approaches



You maximize an *ad hoc* criterion:

$$\hat{\mathbf{q}} = \underset{\mathbf{q}}{\operatorname{argmax}} \text{CRIT}(\mathcal{T}_f),$$

and hope that it's aligned with your original goal:

$$\hat{\theta}_{\hat{\mathbf{q}}} = \underset{\theta_{\hat{\mathbf{q}}}}{\operatorname{argmax}} \ell(\theta_{\hat{\mathbf{q}}}; \mathcal{T}_f).$$

Feature quantization: Approximation

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

Feature quantization: Approximation

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

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

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

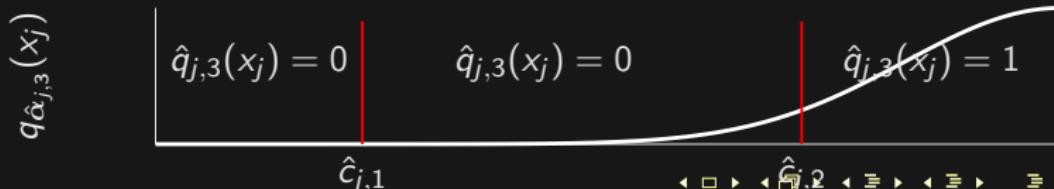
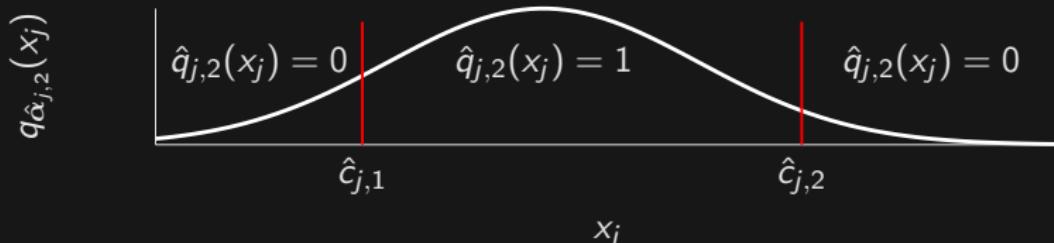
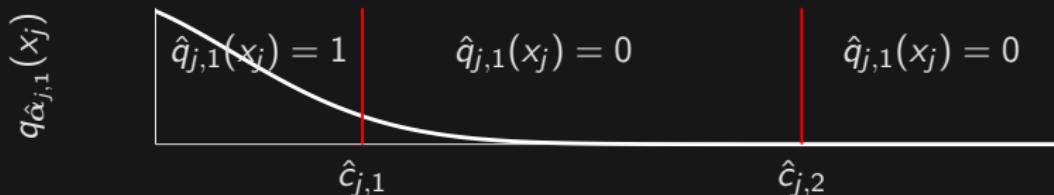
For categorical features, we set for

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

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

Feature quantization: Estimation MAP

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



Feature quantization: Neural networks

We wish to maximize the following likelihood:

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

If there is a true quantization \mathbf{q}^* , then $\alpha^* = \lim_{n \rightarrow \infty} \hat{\alpha}$ is such that $\mathbf{q}_{\alpha^*} = \mathbf{q}^*$.

Feature quantization: Neural networks

We wish to maximize the following likelihood:

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

If there is a true quantization \mathbf{q}^* , then $\alpha^* = \lim_{n \rightarrow \infty} \hat{\alpha}$ is such that $\mathbf{q}_{\alpha^*} = \mathbf{q}^*$.

If not, \mathbf{q}^{MAP} is “guaranteed” to be a good candidate quantization.

Feature quantization: Neural networks

We wish to maximize the following likelihood:

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

If there is a true quantization \mathbf{q}^* , then $\alpha^* = \lim_{n \rightarrow \infty} \hat{\alpha}$ is such that $\mathbf{q}_{\alpha^*} = \mathbf{q}^*$.

If not, \mathbf{q}^{MAP} is “guaranteed” to be a good candidate quantization.

Problem: $\ell(\theta, \alpha; \mathbf{x}, \mathbf{y})$ cannot be directly maximized (it's not even convex).

Feature quantization: Neural networks

We wish to maximize the following likelihood:

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

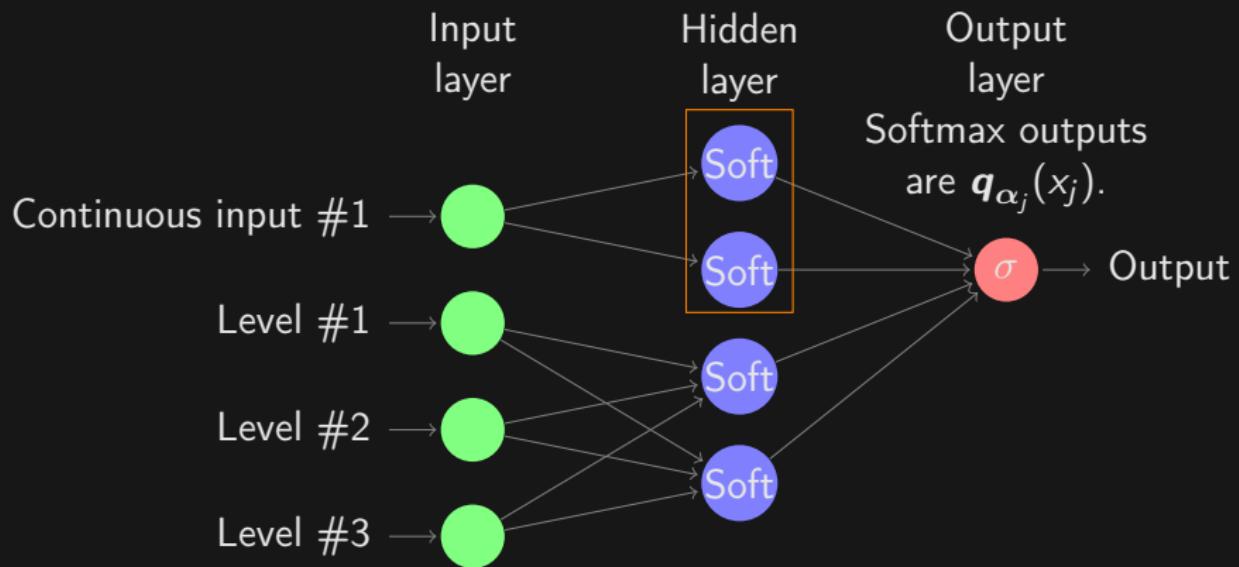
If there is a true quantization \mathbf{q}^* , then $\alpha^* = \lim_{n \rightarrow \infty} \hat{\alpha}$ is such that $\mathbf{q}_{\alpha^*} = \mathbf{q}^*$.

If not, \mathbf{q}^{MAP} is “guaranteed” to be a good candidate quantization.

Problem: $\ell(\theta, \alpha; \mathbf{x}, \mathbf{y})$ cannot be directly maximized (it's not even convex).

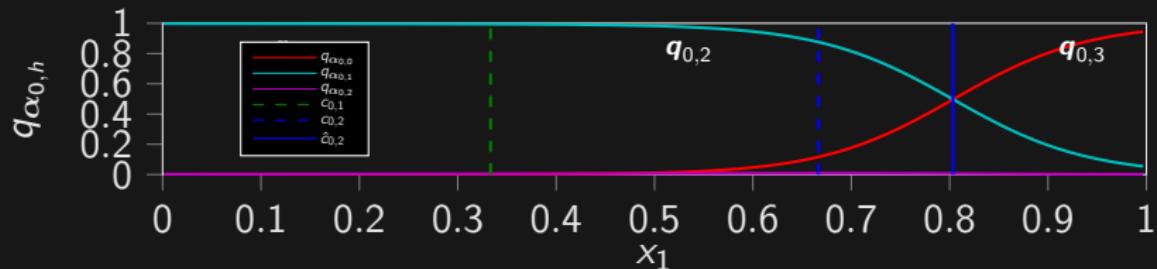
Solution: Resort to gradient descent (not guaranteed to converge to a global maximum!).

Feature quantization: Neural networks



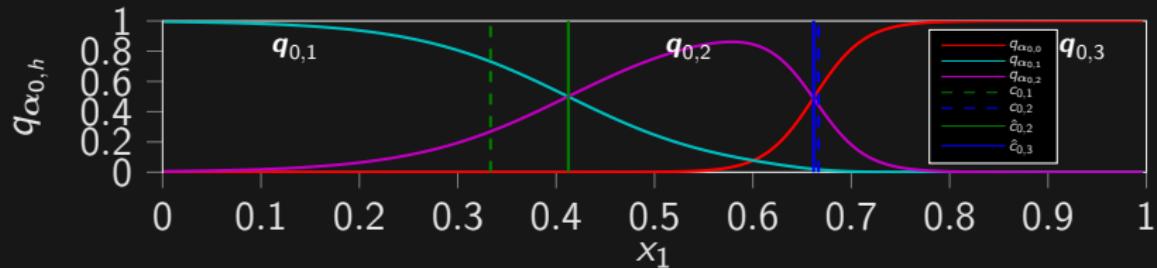
Estimation via neural networks

Continuous feature 0 at iteration 5



(a) Quantization $\hat{q}_1^{(s)}(x_1)$ resulting from the MAP at iter $t = 5$ and $m_{\max} = 3$.

Continuous feature 0 at iteration 300



(b) Quantizations $\hat{q}_1^{(s)}(x_1)$ resulting from the MAP at iter $t = 300$ and $m_{\max} = 3$.

Feature quantization: Model = quantization selection

New model selection criterion

We have drastically restricted the search space to clever candidates $\mathbf{q}^{\text{MAP}(1)}, \dots, \mathbf{q}^{\text{MAP}(\text{iter})}$ resulting from the gradient descent steps.

$$(\mathbf{q}^*, \boldsymbol{\theta}^*) = \operatorname{argmin}_{\hat{\mathbf{q}} \in \{\mathbf{q}^{\text{MAP}(1)}, \dots, \mathbf{q}^{\text{MAP}(\text{iter})}\}, \boldsymbol{\theta} \in \Theta_m} \text{BIC}(\hat{\boldsymbol{\theta}}_{\hat{\mathbf{q}}})$$

Feature quantization: Model = quantization selection

New model selection criterion

We have drastically restricted the search space to clever candidates $\mathbf{q}^{\text{MAP}(1)}, \dots, \mathbf{q}^{\text{MAP}(\text{iter})}$ resulting from the gradient descent steps.

$$(\mathbf{q}^*, \boldsymbol{\theta}^*) = \operatorname{argmin}_{\hat{\mathbf{q}} \in \{\mathbf{q}^{\text{MAP}(1)}, \dots, \mathbf{q}^{\text{MAP}(\text{iter})}\}, \boldsymbol{\theta} \in \Theta_m} \text{BIC}(\hat{\boldsymbol{\theta}}_{\hat{\mathbf{q}}})$$

We would still need to loop over candidates m !

Feature quantization: Model = quantization selection

New model selection criterion

We have drastically restricted the search space to clever candidates $\mathbf{q}^{\text{MAP}(1)}, \dots, \mathbf{q}^{\text{MAP}(\text{iter})}$ resulting from the gradient descent steps.

$$(\mathbf{q}^*, \boldsymbol{\theta}^*) = \operatorname{argmin}_{\hat{\mathbf{q}} \in \{\mathbf{q}^{\text{MAP}(1)}, \dots, \mathbf{q}^{\text{MAP}(\text{iter})}\}, \boldsymbol{\theta} \in \Theta_m} \text{BIC}(\hat{\boldsymbol{\theta}}_{\hat{\mathbf{q}}})$$

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.

Feature quantization: Model = quantization selection

New model selection criterion

We have drastically restricted the search space to clever candidates $\mathbf{q}^{\text{MAP}(1)}, \dots, \mathbf{q}^{\text{MAP}(\text{iter})}$ resulting from the gradient descent steps.

$$(\mathbf{q}^*, \boldsymbol{\theta}^*) = \operatorname{argmin}_{\hat{\mathbf{q}} \in \{\mathbf{q}^{\text{MAP}(1)}, \dots, \mathbf{q}^{\text{MAP}(\text{iter})}\}, \boldsymbol{\theta} \in \Theta_m} \text{BIC}(\hat{\boldsymbol{\theta}}_{\hat{\mathbf{q}}})$$

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

Feature quantization: SEM-Gibbs

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

Feature quantization: SEM-Gibbs

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

- ▶ q is considered a latent (unobserved) feature;

Feature quantization: SEM-Gibbs

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

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

Feature quantization: SEM-Gibbs

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

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

SEM-Gibbs: estimation

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

SEM-Gibbs: 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: 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: 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: 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: algorithm

Initialization

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

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^{MAP}

$$\begin{pmatrix} q^{\text{MAP}, 1,j} \\ \vdots \\ q^{\text{MAP}, n,j} \end{pmatrix} \xrightarrow{\text{MAP estimate}} \begin{pmatrix} \underset{q_j}{\text{argmax}} p_{\alpha_j}(q_j | x_{1,j}) \\ \vdots \\ \underset{q_j}{\text{argmax}} p_{\alpha_j}(q_j | x_{n,j}) \end{pmatrix}$$

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

UCI data

Table: Gini indices (the greater the value, the better the performance) of our proposed quantization algorithm *glmdisc* and two baselines: ALLR and MDLP / χ^2 tests obtained on several benchmark datasets from the UCI library.

Dataset	ALLR	MDLP/ χ^2	<i>glmdisc</i>
Adult	81.4 (1.0)	85.3 (0.9)	80.4 (1.0)
Australian	72.1 (10.4)	84.1 (7.5)	92.5 (4.5)
Bands	48.3 (17.8)	47.3 (17.6)	58.5 (12.0)
Credit	81.3 (9.6)	88.7 (6.4)	92.0 (4.7)
German	52.0 (11.3)	54.6 (11.2)	69.2 (9.1)
Heart	80.3 (12.1)	78.7 (13.1)	86.3 (10.6)

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 of Table 4 and the current scorecard (manual / expert representation) obtained on several portfolios of Crédit Agricole Consumer Finance.

Portfolio	ALLR	Current	MDLP/ χ^2	<i>glmdisc</i>
Automobile	59.3 (3.1)	55.6 (3.4)	59.3 (3.0)	58.9 (2.6)
Renovation	52.3 (5.5)	50.9 (5.6)	54.0 (5.1)	56.7 (4.8)
Standard	39.7 (3.3)	37.1 (3.8)	45.3 (3.1)	44.0 (3.1)
Revolving	62.7 (2.8)	58.5 (3.2)	63.2 (2.8)	62.3 (2.8)
Mass retail	52.8 (5.3)	48.7 (6.0)	61.4 (4.7)	61.8 (4.6)
Electronics	52.9 (11.9)	55.8 (10.8)	56.3 (10.2)	72.6 (7.4)

See [this gist](#) for χ^2 automated grouping tests.

Bivariate interactions

Bivariate interactions: Notations

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

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

Bivariate interactions: Notations

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

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

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

$$(\boldsymbol{\theta}^*, \boldsymbol{\delta}^*) = \underset{\boldsymbol{\theta}, \boldsymbol{\delta} \in \{0,1\}^{\frac{d(d-1)}{2}}}{\operatorname{argmax}} \sum_{i=1}^n \ln p_{\boldsymbol{\theta}}(y_i | \mathbf{q}(\mathbf{x}_i), \boldsymbol{\delta}) - \text{penalty}(n; \boldsymbol{\theta})$$

Bivariate interactions: Notations

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

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

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

$$(\boldsymbol{\theta}^*, \boldsymbol{\delta}^*) = \underset{\boldsymbol{\theta}, \boldsymbol{\delta} \in \{0,1\}^{\frac{d(d-1)}{2}}}{\operatorname{argmax}} \sum_{i=1}^n \ln p_{\boldsymbol{\theta}}(y_i | \mathbf{q}(\mathbf{x}_i), \boldsymbol{\delta}) - \text{penalty}(n; \boldsymbol{\theta})$$

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

Bivariate interactions: Model proposal

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

Bivariate interactions: Model proposal

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

Strategy used here: Metropolis-Hastings sampling algorithm.

Idea: Propose “clever” interactions and accept / reject them based on the BIC criterion of the resulting logistic regression.

Bivariate interactions: Model proposal

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

Strategy used here: Metropolis-Hastings sampling algorithm.

Idea: Propose “clever” 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

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

Strategy used here: Metropolis-Hastings sampling algorithm.

Idea: Propose “clever” interactions and accept / reject them based on the BIC criterion of the resulting logistic regression.

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

Bivariate interactions: Model proposal

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

Strategy used here: Metropolis-Hastings sampling algorithm.

Idea: Propose “clever” interactions and accept / reject them based on the BIC criterion of the resulting logistic regression.

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

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

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 (provably) 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 (provably) good chance that it will be in the full multivariate model.

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

Bivariate interactions: Results

Données UCI

Table: Gini indices (the greater the value, the better the performance) of our proposed quantization algorithm *glmdisc* and two baselines: ALLR and MDLP / χ^2 tests obtained on several benchmark datasets from the UCI library.

Dataset	ALLR	<i>ad hoc</i> methods	Our proposal: <i>glmdisc</i> -NN	Our proposal: <i>glmdisc</i> -SEM	<i>glmdisc</i> -SEM w. interactions
Adult	81.4 (1.0)	85.3 (0.9)	80.4 (1.0)	81.5 (1.0)	81.5 (1.0 - no interaction)
Australian	72.1 (10.4)	84.1 (7.5)	92.5 (4.5)	100 (0)	100 (0 - no interaction)
Bands	48.3 (17.8)	47.3 (17.6)	58.5 (12.0)	58.7 (12.0)	58.8 (13.0)
Credit	81.3 (9.6)	88.7 (6.4)	92.0 (4.7)	87.7 (6.4)	87.7 (6.4 - no interaction)
German	52.0 (11.3)	54.6 (11.2)	69.2 (9.1)	54.5 (10)	
Heart	80.3 (12.1)	78.7 (13.1)	86.3 (10.6)	82.2 (11.2)	84.5 (10.8)

Bivariate interactions: Results

Medicine data

Table: Gini indices of our proposed quantization algorithm *glmdisc*-SEM and two baselines: ALLR and ALLR with all pairwise interactions on several medicine-related benchmark datasets.

	Pima	Breast	Birthwt
ALLR	73.0	94.0	34.0
ALLR LR w. interactions	60.0	51.0	15.0
glmdisc	57.0	93.0	18.0
glmdisc w. interactions	62.0	95.0	54.0

Bivariate interactions: Results

CACF data

Table: Gini indices (the greater the value, the better the performance) of our proposed quantization algorithm *glmdisc*, the two baselines of Table 4 and the current scorecard (manual / expert representation) obtained on several portfolios of Crédit Agricole Consumer Finance.

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)	64.8 (2.0)
Renovation	52.3 (5.5)	50.9 (5.6)	54.0 (5.1)	56.7 (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)	47.2 (2.8)
Revolving	62.7 (2.8)	58.5 (3.2)	63.2 (2.8)	62.3 (2.8)	60.7 (2.8)	67.2 (2.5)
Mass retail	52.8 (5.3)	48.7 (6.0)	61.4 (4.7)	61.8 (4.6)	61.0 (4.7)	60.3 (4.8)
Electronics	52.9 (11.9)	55.8 (10.8)	56.3 (10.2)	72.6 (7.4)	62.0 (9.5)	63.7 (9.0)

Bivariate interactions: Results

Older results

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)	≈ 15	30	12.2
Basel II (n=70,000 ; d=13)	70	71.3	19

Segmentation: logistic regression trees

Segmentation: logistic regression trees

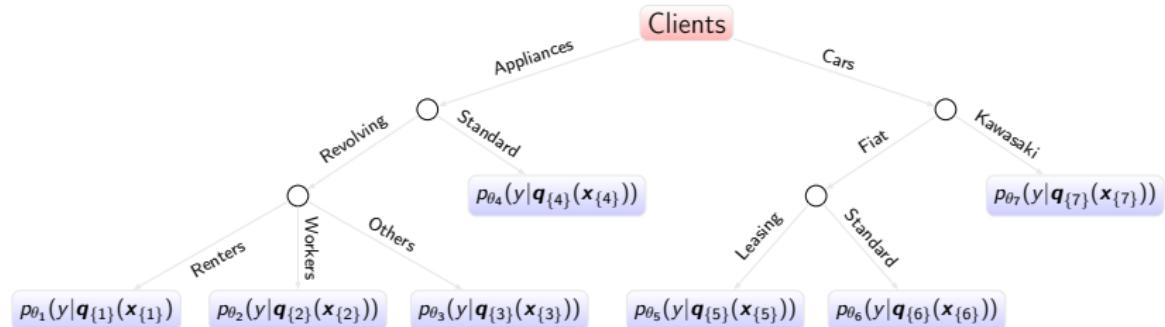


Figure: Scorecards tree structure in acceptance system.

Segmentation: logistic regression trees: Notations

K segments.

Segmentation: logistic regression trees: Notations

K segments.

$c \in \{1, \dots, K\}$: latent feature of the client's segment.

Segmentation: logistic regression trees: Notations

K segments.

$c \in \{1, \dots, K\}$: latent feature of the client's segment.

We suppose there is a true segmentation c^* , K^* and logistic regressions θ^{*,c^*} at its leaves.

Segmentation: logistic regression trees: Notations

K segments.

$c \in \{1, \dots, K\}$: latent feature of the client's segment.

We suppose there is a true segmentation c^* , K^* and logistic regressions θ^{*,c^*} at its leaves.

If we could evaluate all segmentations, the true one would be selected by

$$\operatorname{argmax}_{c,K} \sum_{c=1}^K \text{BIC}(\hat{\theta}^c),$$

where $\hat{\theta}^c$ is the MLE of the logistic regression on .

Segmentation: logistic regression trees: Model proposal

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

Segmentation: logistic regression trees: Model proposal

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

$$p(y|\mathbf{x}) = \sum_{c=1}^K p_\theta(y|\mathbf{x}; c) p_\beta(c|\mathbf{x}).$$

Segmentation: logistic regression trees: Model proposal

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

$$p(y|\mathbf{x}) = \sum_{c=1}^K p_{\theta}(y|\mathbf{x}; c) p_{\beta}(c|\mathbf{x}).$$

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

Segmentation: logistic regression trees: Model proposal

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

$$p(y|\mathbf{x}) = \sum_{c=1}^K p_{\theta}(y|\mathbf{x}; c) p_{\beta}(c|\mathbf{x}).$$

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

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

Segmentation: logistic regression trees: Model proposal

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

$$p(y|\mathbf{x}) = \sum_{c=1}^K p_{\theta}(y|\mathbf{x}; c) p_{\beta}(c|\mathbf{x}).$$

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

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

$$\beta^{(s+1)} = C4.5(\mathbf{c}^{(s+1)}, \mathbf{x}).$$

Segmentation: logistic regression trees: Some results

Oracle = ALLR		<i>glmtree</i> -SEM	FAMD	PLS	LMT	MOB
Gini	69.7	69.7	65.3	47.0	69.7	64.8

Oracle		ALLR	<i>glmtree</i> -SEM	FAMD	PLS	LMT	MOB
Gini	69.7	25.8	69.7	17.7	48.4	65.8	69.7

Bonus

Big “unstructured” data

Some theoretical results about an ever bigger d (not the one you think about though).

Online logistic regression

What if we dynamically adjusted logistic regression coefficients of a given scorecard (still learnt on a cold database) on new data as they come in?

Profitability

Good / bad label is merely a proxy of the true performance measure: profitability.

Already done: weighting observations by the amount of the loan gives rise to roughly the same logistic regression coefficients.

Predicting IR3 in 2 months based on the month's applications

Current process: finance people, wait 3 months, if risk \neq budget then adjust acceptance policy, wait 3 months again and repeat.
Couldn't we anticipate by looking at the quality (e.g. through the score) of the applications?

Thanks!

References |

-  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).
-  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 (visited on 08/25/2016).
-  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 (visited on 08/25/2016).

References II

-  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.
-  Bianca Zadrozny. "Learning and evaluating classifiers under sample selection bias". In: Proceedings of the twenty-first international conference on Machine learning. ACM. 2004, p. 114.