

Credit Scoring : biais d'échantillon ou *réintégration des refusés*

Adrien Ehrhardt

Christophe Biernacki, Vincent Vandewalle,
Philippe Heinrich, Sébastien Beben

Crédit Agricole Consumer Finance
INRIA Lille - Nord-Europe

Modal's days 2017

19 janvier 2017

1 Contexte

- Entreprise
- Produits et demandes de crédit
- Système d'acceptation
- Credit Scoring

2 Réintégration des refusés

- Exemple Fuzzy Augmentation
- Etat de l'art
- Formalisation du problème
- Mécanisme de manque
- Bon/mauvais modèle
- Quelle case en Credit Scoring ?
- Améliorer l'estimateur $\hat{\theta}^f$?
- Résultats
- Augmentation de la dimension

3 Conclusion

Contexte

Contexte : Entreprise



La Redoute



Quelques chiffres :

- Encours gérés (2015) : 71 milliards € (26 milliards € en France)
- Produit Net Bancaire : 2 milliards €
- Résultat : +380 M€
- **Coût du risque** France : 200 M€

Produits :

- Facilités de paiement (exemple : en 3 fois)
- Prêt personnel (mensualité + échéance)
- Carte de crédit (revolving)

Contexte : Produits et demandes de crédit II

Cotisation Fidélité * 1an * 3ans *		Offre Carte Castorama Simplifiée	
Emprunteur:	M <input checked="" type="radio"/> Mme <input type="radio"/> Mlle <input type="radio"/>	Conjoint:	M <input type="radio"/> Mme <input type="radio"/> Mlle <input type="radio"/>
Pièce d'identité*	CARTE NATIONALE D'IDENTITE	Pièce d'identité*	
Nationalité*	FRANCE	Nationalité*	FRANCE
Fin de séjour		Fin de séjour	
Nom*	MODAL	Nom*	
Prénom*	DAYZ	Prénom*	
Nom de jeune fille		Nom de jeune fille	
Date de naissance*	01 01 1980	Date de naissance*	
Ville de naissance*	LILLE	Ville de naissance*	
Département de naissance*	59	Département de naissance*	
Pays de naissance*	FRANCE	Pays de naissance*	FRANCE
Situation familiale*	Union libre <input type="checkbox"/> Enfants* 2	Co-emprunteur:	<input type="radio"/> Oui <input type="radio"/> Non
Adresse			
Voie*	AVENUE DU HALLEY	Complément	
Lieu-dit		Ville*	Villeneuve-d'Ascq
Code Postal*	59650	Pays*	FRANCE
Téléphone domicile		Email *	aehrhadt@ca-cf.fr
Téléphone Portable		Vérification de l'email *	aehrhadt@ca-cf.fr
Code habitat*	LOGE ADMINISTR.	Depuis le*	09 2012
Relevé de compte dématérialisé* ?	<input type="radio"/> Oui <input type="radio"/> Non		

Contexte : Produits et demandes de crédit III

Situation professionnelle emprunteur		Situation professionnelle conjoint	
Profession*	ARTISTE, PRESSE 35	Profession	
Employeur		Employeur	
Depuis*	09 2012	Depuis	
Adresse		Adresse	
Code Postal		Code Postal	
Ville		Ville	
Pays	FRANCE	Pays	FRANCE
Téléphone		Téléphone	
Employeur		Employeur	

Revenus :		Charges :	
Emprunteur :		Charges mensuelles d'habitation*	600
Revenus nets mensuels*	1800	Autres charges mensuelles*	150
Autres revenus nets mensuels*	200	Autres crédits*	0
Conjoint :			
Revenus nets mensuels			
Autres revenus nets mensuels			
Prestations familiales*	50		
Total	2050,00 €	Total	750,00 €

Carte facultative emprunteur : Oui Non

Assurances DIM CACIT1 Offre promo Oui Non

Coordonnées bancaires ou postales

BIC IBAN Titulaire

DAYZ MODAL

N° de dossier 49800146566

ETUDE COMPLEMENTAIRE

Capital attribué : 2000 euros

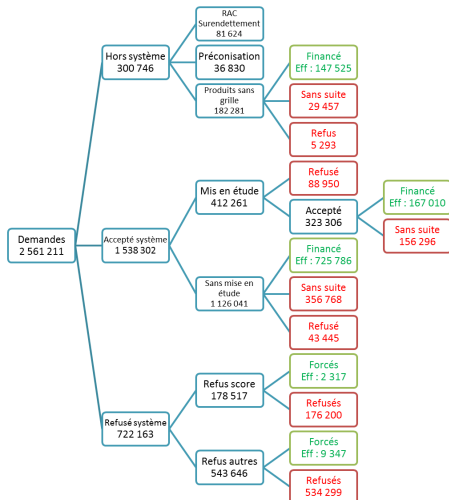


Imprimer le dossier de crédit



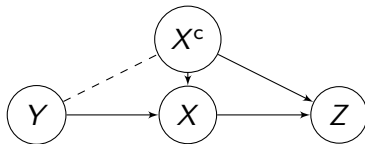
Imprimer via Internet

Contexte : Système d'acceptation



X : vecteur a. des caractéristiques
 Y dans $\{0, 1\}$: remboursement
 A_S dans $\{a, \bar{a}\}$: $f(X)$
 A_C dans $\{a, \bar{a}\}$: $g(X, X^c)$
 Z dans $\{f, nf\}$: v. a. de financement

On oublie A_S et A_C dans la suite.



$$\exists \theta = (\theta_0, \dots, \theta_d) \in \mathbb{R}^{d+1} \text{ s.t. } \forall \mathbf{x}, \ln \left(\frac{P(Y = 1|\mathbf{x})}{P(Y = 0|\mathbf{x})} \right) = \theta_0 + \theta_1 x^1 + \dots + \theta_d x^d$$

n clients financés ($Z = f$)

m clients non financés ($Z = nf$)

\mathbf{x} : *features* observées des clients

\mathbf{y} : remboursement observé

\mathbf{x}^f (\mathbf{x}^{nf}) : *features* observées des clients financés (non financés)

\mathbf{y}^f (\mathbf{y}^{nf}) : remboursement (**non**) observé des clients financés (non financés)

$$\underbrace{\ell(\theta; \mathbf{x}, \mathbf{y})}_{\substack{\text{vraisemblance} \\ \text{complète}}} = \sum_{i=1}^n \ln(p_{\theta}(y_i|x_i)) + \sum_{i=n+1}^{n+m} \ln(p_{\theta}(y_i|x_i)) = \underbrace{\ell(\theta; \mathbf{x}^f, \mathbf{y}^f)}_{\substack{\text{vraisemblance} \\ \text{observée}}} + \ell(\theta; \mathbf{x}^{nf}, \mathbf{y}^{nf})$$

Quel intérêt à utiliser \mathbf{x}^{nf} ?

Quel risque à n'utiliser que $(\mathbf{x}^f, \mathbf{y}^f)$?

Réintégration des refusés

Réintégration des refusés : Exemple Fuzzy Augmentation I

$$\mathbf{y}^f \begin{pmatrix} y_1 \\ \vdots \\ y_n \\ \text{NA} \\ \vdots \\ \text{NA} \end{pmatrix}$$
$$\mathbf{y}^{nf}$$

$$\mathbf{x}^f \begin{pmatrix} x_1^1 & \cdots & x_1^d \\ \vdots & \vdots & \vdots \\ x_n^1 & \cdots & x_n^d \\ x_{n+1}^1 & \cdots & x_{n+1}^d \\ \vdots & \vdots & \vdots \\ x_{n+m}^1 & \cdots & x_{n+m}^d \end{pmatrix}$$
$$\mathbf{x}^{nf}$$

Abandon de \mathbf{x}^{nf} et construction de $\hat{\theta}^{\text{f}}$

$$\begin{array}{c}
 \mathbf{y}^{\text{f}} \\
 \mathbf{y}^{\text{nf}}
 \end{array}
 \begin{pmatrix}
 y_1 \\
 \vdots \\
 y_n \\
 \text{NA} \\
 \vdots \\
 \text{NA}
 \end{pmatrix}
 \qquad
 \begin{array}{c}
 \mathbf{x}^{\text{f}} \\
 \mathbf{x}^{\text{nf}}
 \end{array}
 \begin{pmatrix}
 x_1^1 & \cdots & x_1^d \\
 \vdots & \vdots & \vdots \\
 x_n^1 & \cdots & x_n^d \\
 x_{n+1}^1 & \cdots & x_{n+1}^d \\
 \vdots & \vdots & \vdots \\
 x_{n+m}^1 & \cdots & x_{n+m}^d
 \end{pmatrix}$$

Réintégration des refusés : Exemple Fuzzy Augmentation III

Remplacement de \mathbf{y}^{nf} par les proba données par $\hat{\theta}^{\text{f}}$

$$\begin{array}{c} \mathbf{y}^{\text{f}} \\ \mathbf{y}^{\text{nf}} \end{array} \begin{pmatrix} y_1 \\ \vdots \\ y_n \\ p_{\hat{\theta}^{\text{f}}}(y_{n+1} = 1 | x_{n+1}) \\ \vdots \\ p_{\hat{\theta}^{\text{f}}}(y_{n+m} = 1 | x_{n+m}) \end{pmatrix} \quad \begin{array}{c} \mathbf{x}^{\text{f}} \\ \mathbf{x}^{\text{nf}} \end{array} \begin{pmatrix} x_1^1 & \cdots & x_1^d \\ \vdots & \vdots & \vdots \\ x_n^1 & \cdots & x_n^d \\ x_{n+1}^1 & \cdots & x_{n+1}^d \\ \vdots & \vdots & \vdots \\ x_{n+m}^1 & \cdots & x_{n+m}^d \end{pmatrix}$$

Apprendre $\hat{\theta}^{\text{fuzzy}}$ sur le dataset résultant.

Problème : $\hat{\theta}^{\text{fuzzy}} = \hat{\theta}^{\text{f}}$

[Feelders, 2000] : formalisation et approche par données manquantes.

[Viennet et al., 2006, Guizani et al., 2013, Banasik and Crook, 2007, Nguyen, 2016] : description de méthodes empiriques pour utiliser x^{nf} .

[Viennet et al., 2006, Banasik and Crook, 2007] : la réintégration a un intérêt.

[Nguyen, 2016, Guizani et al., 2013] : la réintégration n'a pas d'intérêt.

[Kiefer and Larson, 2006] : étude de cas du mauvais modèle.

Réintégration des refusés : Formalisation du problème

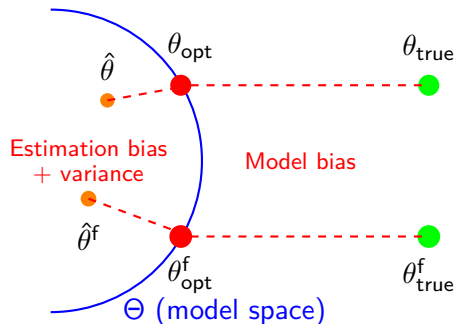
Objet d'intérêt : $p_{\text{true}}(y|x)$

Proposition d'un modèle : $p_{\theta}(y|x)$

Problème (sans surprise...) : estimer θ

Données :

- 1 Cas idéal : $\mathbf{x}^f, \mathbf{x}^{nf}$ et $\mathbf{y}^f, \mathbf{y}^{nf}$
- 2 Cas CACF : \mathbf{x}^f et \mathbf{y}^f



Réintégration des refusés : Formalisation du problème

Estimateurs :

① Cas idéal : $\sqrt{n+m}(\hat{\theta} - \theta_{\text{opt}}) \xrightarrow[n, m \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}})$

② Cas CACF : $\sqrt{n}(\hat{\theta}^f - \theta_{\text{opt}}^f) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}_{d+1}(0, \Sigma_{\theta_{\text{opt}}^f}^f)$

Question 1 : propriétés asymptotiques des deux estimateurs

$$\text{(Q1) } \theta_{\text{opt}} \stackrel{?}{=} \theta_{\text{opt}}^f$$

$$\text{(Q2) } \text{ARE}(\hat{\theta}^f, \hat{\theta}) = \left(\varphi \frac{|\Sigma_{\theta_{\text{opt}}}|}{|\Sigma_{\theta_{\text{opt}}^f}^f|} \right)^{\frac{1}{d+1}}$$

$$\text{(Q2a) } \frac{n}{n+m} \xrightarrow[n, m \rightarrow \infty]{\text{a.s.}} \varphi$$

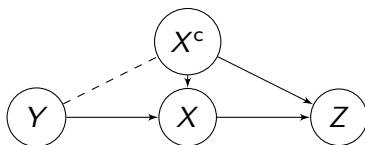
$$\text{(Q2b) } \Sigma_{\theta_{\text{opt}}} \stackrel{?}{=} \Sigma_{\theta_{\text{opt}}^f}^f$$

Réponse 1 (toujours sans surprise...) : ça dépend.

- ① Mécanisme des données manquantes $p_{\text{true}}(z|x)$
- ② La vraie loi prédictive $p_{\text{true}}(y|x)$ appartient-elle à Θ ?

Réintégration des refusés : Mécanisme de manque

- **MCAR** : $\forall x, y, z, p_{\text{true}}(z|x, y) = p_{\text{true}}(z)$
→ Inadapté au Credit Scoring.
- **MAR** : $\forall x, y, z, p_{\text{true}}(z|x, y) = p_{\text{true}}(z|x)$
- **MNAR** : $\exists x, y, z, p_{\text{true}}(z|x, y) \neq p_{\text{true}}(z|x)$
→ Influence du "feeling" des conseillers X^c .



Ignorabilité :

Les paramètres de $p_{\text{true}}(z|x, y)$ ne dépendent pas fonctionnellement de θ_{true} .

Réintégration des refusés : Bon/mauvais modèle

- Bon modèle : $\exists \theta_{\text{true}}, p_{\text{true}}(y|x) = p_{\theta_{\text{true}}}(y|x)$.
→ Données réelles \Rightarrow hypothèse peu probable.
- Mauvais modèle : θ_{opt} minimise l'ignorance sur la vraie loi.
→ Utilisation de la régression logistique pour sa robustesse à la misspecification.

$p_{\theta}(y x, z)$ \ $p_{\text{true}}(z x)$	MCAR	MAR	MNAR
True	$\theta_{\text{opt}}^f = \theta_{\text{opt}}$	$\theta_{\text{opt}}^f = \theta_{\text{opt}}$ $\Sigma_{\theta_{\text{opt}}^f} \neq \Sigma_{\theta_{\text{opt}}}$	$\theta_{\text{opt}}^f \neq \theta_{\text{opt}}$
Misspecified	$\Sigma_{\theta_{\text{opt}}^f} = \Sigma_{\theta_{\text{opt}}}$	$\theta_{\text{opt}}^f \neq \theta_{\text{opt}}$ $\Sigma_{\theta_{\text{opt}}^f} \neq \Sigma_{\theta_{\text{opt}}}$	$\Sigma_{\theta_{\text{opt}}^f} \neq \Sigma_{\theta_{\text{opt}}}$

Mécanisme de manque :

- MCAR : la sélection des clients n'est pas aléatoire !
- MAR : vrai si X détermine Z (cas du score)
- MNAR : vrai si conseillers = oracles

On ne peut pas tester MNAR.

Bon modèle : non.

Mais toutes les variables sont discrétisées : on peut approcher n'importe quelle relation fonctionnelle entre la variable cible et chaque variable continue.

Question 2 : Comment "sauter" dans une meilleure case ?

Leviers :

- ~~Changer le modèle (i.e. l'espace Θ)~~ régression logistique,
- ~~Modéliser la sélection (i.e. $p_{\alpha}(z|x, y)$)~~ relève de la croyance,
- Utiliser x^{nf} .

Case 1 : MAR et bon modèle

Idée : Améliorer la vitesse de convergence

Proposition : Modèle génératif avec données manquantes

Case 2 : MAR et mauvais modèle

Idée : "Dé-biaiser" le modèle en corrigeant $p_{\theta}(y|x)$

Proposition : Méthode "Augmentation" proche de l'Importance Sampling

Case 3 : MNAR

Idée : "Dé-biaiser" le modèle en corrigeant $p_{\theta}(y|x, z)$

Proposition : Méthode "Parcelling" : hypothèses invérifiables

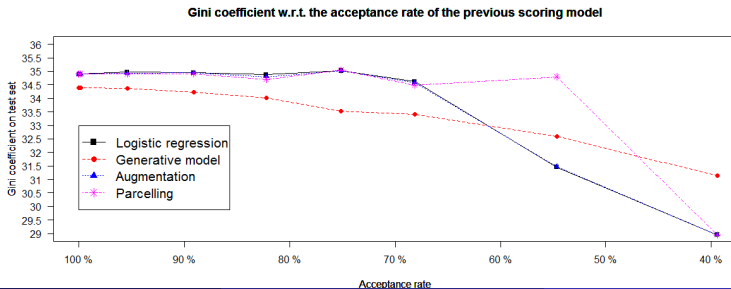
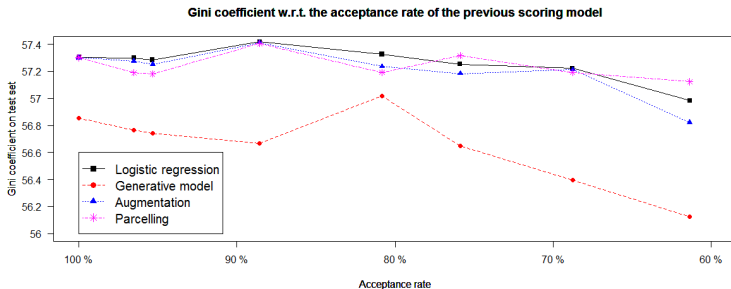
Certaines méthodes sont ré-interprétables comme des hypothèses partielles sur la loi jointe.

6 méthodes testées :

- 1 2 méthodes conduisent à θ_{opt}^f .
- 2 1 méthode utilisant un algorithme sans garantie de convergence.
- 3 Modèle génératif : biais de modèle important.
- 4 Augmentation (**MAR/mauvais modèle**) : hypothèses non satisfaites.
- 5 Parcelling (**MNAR**) : introduction d'*a priori*.

Réponse 2 : arbitrer entre (pas de réintégration) et (méthode Toto de réintégration), c'est comparer deux "prises de risque" inquantifiables.

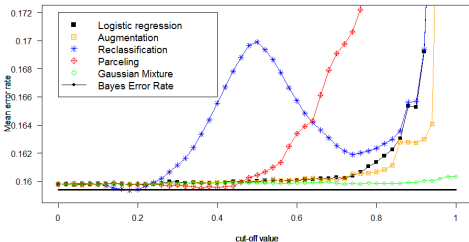
Réintégration des refusés : Résultats



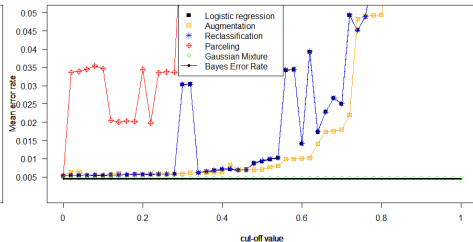
Réintégration des refusés : Augmentation de la dimension

Gains hypothétiques de l'utilisation de x^{nf} faibles

Mean error rate on test set for 100 learning sets w.r.t. the cut-off value



Mean error rate on test set for 100 learning sets w.r.t. the cut-off value



Grands gains potentiels de l'augmentation du nombre de prédicteurs :

- On agit sur le modèle.
- On abaisse le taux d'erreur de Bayes (classes séparables).

Conclusion

- Le choix d'une méthode relève de la croyance.
- La formalisation a permis de clotûrer un débat resté empirique jusqu'alors.
- Article en préparation
- Sujet en cours : discrétisation

Merci pour votre attention !

Questions ?



Banasik, J. and Crook, J. (2007).

Reject inference, augmentation, and sample selection.

[European Journal of Operational Research](#), 183(3):1582–1594.



Feelders, A. (2000).

Credit scoring and reject inference with mixture models.

[International Journal of Intelligent Systems in Accounting, Finance & Management](#), 9(1):1–8.



Guizani, A., Souissi, B., Ammou, S. B., and Saporta, G. (2013).

Une comparaison de quatre techniques d'inférence des refusés dans le processus d'octroi de crédit.

In [45 emes Journ'ees de statistique](#), page pp.



Kiefer, N. M. and Larson, C. E. (2006).

Specification and informational issues in credit scoring.

[Available at SSRN 956628](#).



Nguyen, H. T. (2016).

Reject inference in application scorecards : evidence from France.
Technical report, University of Paris West-Nanterre la Défense,
EconomiX.



Viennet, E., Soulié, F. F., and Rognier, B. (2006).

Evaluation de techniques de traitement des refusés pour l'octroi de
crédit.

[arXiv preprint cs/0607048.](https://arxiv.org/abs/cs/0607048)