



Fouille de Données Massives

Projet : Sujet 1

M2 Informatique - SISE

Guillaume Metzler

Institut de Communication (ICOM)
Université de Lyon, Université Lumière Lyon 2
Laboratoire ERIC UR 3083, Lyon, France

guillaume.metzler@univ-lyon2.fr

Ce projet sera effectué par groupe de deux ou trois étudiants maximum et porte sur l'apprentissage dans un contexte de données déséquilibrées avec une application plus précise sur la détection de fraudes.

Le travail sera à rendre par mail à guillaume.metzler@univ-lyon2.fr avant le **9 Février 2025 au plus tard**. Ce travail se composera d'un dossier retraçant vos démarches et résultats entrepris pour traiter le sujet, de plus amples explications vous sont données ci-dessous. Vous êtes libres d'utiliser le langage de votre choix pour effectuer ce travail, R ou Python. Choisissez celui avec lequel vous êtes le plus à l'aise.

Vous pouvez télécharger les données à l'adresse ci-dessous

[Lien téléchargement](#)

1 A propos des données

Les données sur lesquelles vous allez travailler sont des données réelles. Elles sont issues d'une enseigne de la grande distribution ainsi que de certains organismes bancaires (*FNCI* et *Banque de France*). Chaque ligne représente une transaction effectuée par chèque dans un magasin de l'enseigne quelque part en France, elles ne sont pas brutes et plusieurs variables sont déjà des variables créées, *i.e.* sont issues du *feature engineering*, nous avons un ensemble de 23 variables qui ont la signification suivante :

- **ZIBZIN** : identifiant relatif à la personne, *i.e.* il s'agit de son identifiant bancaire (relatif au chéquier en cours d'utilisation)
- **IDAvisAutorisationCheque** : identifiant de la transaction en cours
- **Montant** : montant de la transaction
- **DateTransaction** : date de la transaction
- **CodeDecision** : il s'agit d'une variable qui peut prendre ici 4 valeurs
 - 0 : la transaction a été acceptée par le magasin
 - 1 : la transaction et donc le client fait partie d'une **liste blanche** (bons payeurs). Vous n'en rencontrerez pas dans cette base de données
 - 2 : le client fait d'une partie d'une **liste noire**, son historique indique cet un mauvais payer (des impayés en cours ou des incidents bancaires en cours), sa transaction est alors automatiquement refusée
 - 3 : client ayant été arrêté par le système par le passé pour une raison plus ou moins fondée
- **VérifianceCPT1** : nombre de transactions effectuées par le même identifiant bancaire au cours du même jour
- **VérifianceCPT2** : nombre de transactions effectuées par le même identifiant bancaire au cours des trois derniers jours
- **VérifianceCPT3** : nombre de transactions effectuées par le même identifiant bancaire au cours des sept derniers jours
- **D2CB** : durée de connaissance du client (par son identifiant bancaire), en jours. Pour des contraintes légales, cette durée de connaissance ne peut excéder deux ans
- **ScoringFP1** : score d'anormalité du panier relatif à une première famille de produits (ex : denrées alimentaires)
- **ScoringFP2** : score d'anormalité du panier relatif à une deuxième famille de produits (ex : électroniques)
- **ScoringFP3** : score d'anormalité du panier relatif à une troisième famille de produits (ex : autres)
- **TauxImpNb_RB** : taux impayés enregistrés selon la région où a lieu la transaction
- **TauxImpNB_CPM** : taux d'impayés relatif au magasin où a lieu la transaction
- **EcartNumCheq** : différence entre les numéros de chèques
- **NbrMagasin3J** : nombre de magasins différents fréquentés les 3 derniers jours
- **DiffDateTr1** : écart (en jours) à la précédente transaction
- **DiffDateTr2** : écart (en jours) à l'avant dernière transaction
- **DiffDateTr3** : écart (en jours) à l'antépénultième transaction
- **CA3TRetMtt** : montant des dernières transactions + montant de la transaction en cours
- **CA3TR** : montant des trois dernières transactions
- **Heure** : heure de la transaction
- **FlagImpaye** : acception (0) ou refus de la transaction (1)

Remarque La variable *CodeDecision* n'est pas une variable à utiliser pour faire de la prédiction car cette information est acquise post-transaction. On peut en revanche s'en servir lors de la phase d'apprentissage pour analyser les données par exemple.

Vous disposez donc d'un jeu de données comprenant 10 mois de transactions du "2017-02-01" au "2017-11-30". A vous de voir si toutes ces informations sont nécessaires ou non pour établir le modèle.

On définira les ensembles de la façon suivante :

- **Apprentissage** : transactions ayant eu lieu entre le "2017-02-01" et le "2017-08-31".
- **Test** : transactions ayant eu lieu entre le "2017-09-01" et le "2017-11-30"

2 Travail à effectuer : première partie

La variable à prédire est la variable *FlagImpaye*, il s'agit d'une variable qui ne peut prendre que deux valeurs possibles : 0 la transaction est acceptée et considérée comme "normale", 1 la transaction est refusée car considérée comme "frauduleuse".

Nous avons vu que plusieurs critères peuvent être utilisés pour évaluer la performance d'un modèle comme l'Accuracy, la précision, le rappel, la F-mesure ou encore l'aire sous la courbe ROC (AUC ROC). Dans le cas présent, vous allez chercher à établir le modèle vous permettant d'obtenir les meilleurs résultats en classification en terme de F-mesure F dont la définition est rappelée ci-dessous :

$$F = \frac{2TP}{2TP + FN + FP}$$

où un TP est une fraude prédite fraude par votre modèle, un FN est une fraude non identifiée comme tel par votre modèle, un FP est une transaction non frauduleuse mais identifiée comme frauduleuse par votre modèle et enfin un TN est une transaction non frauduleuse.

Pour présenter votre travail, vous pourrez vous appuyer sur le modèle présenté en Annexe du présent document et suivre les indications sur le contenu des différentes sections.

3 Quelques suggestions

Idéalement, le travail effectué devrait comprendre au moins 5 procédures ou méthodes différentes vues en cours : une méthode peut être un algorithme de classification seul ou encore couplé ou non à une méthode d'échantillonnage par exemple. Je vous donne ci-dessous une liste non exhaustive des méthodes que vous pourriez utiliser pour votre travail.

- **Pre-Process sur les données** : utilisation d'algorithmes d'over-sampling (random - SMOTE - Adasyn - ...) ou encore des approches d'under-sampling (random - Tomek Link - Edited Nearest Neighbour - NearMiss - ...) vous pouvez aussi utiliser des méthodes dites *cost-sensitive* qui vont accorder plus de poids aux exemples d'une classe donnée, voire même des poids spécifiques à chaque exemple.
- **Algorithmes** : vous pourrez utiliser des algorithmes non supervisés comme des méthodes de clustering (k-means, clustering hiérarchique ou encore les auto-encodeurs) pour détecter les fraudes. Vous disposez également d'un large éventail d'algorithmes de classification supervisés que vous pouvez utiliser : decision trees, random forests, gradient boosting, nearest-neighbors, réseaux de neurones profonds, SVM (linéaire ou non ...), analyse discriminante, boosting, ...
- **Post-traitement** : combinez les résultats issus de différents modèles (bagging) afin de créer un modèle potentiellement plus puissant.

N'hésitez pas à regarder sur internet quelques exemples d'utilisations des algorithmes sus-mentionnés et votre objectif sera de les adapter au contexte des données (consulter des sites comme Kaggle - MachineLearningMastery ou encore Medium qui seront pour vous une bonne source d'inspiration). Vous verrez que toutes les méthodes ne sont pas forcément applicables à ce type de données : si tel est le cas, n'hésitez pas à préciser dans votre rapport pourquoi une méthode n'a pas fonctionné selon vous.

4 Travail à effectuer : deuxième partie

Dans cette dernière partie, on souhaite maximiser la marge de l'enseigne. Cette marge dépend directement du montant de la transaction effectuée, de sa nature (frauduleuse ou non) et de la décision prise par le système d'aide à la décision.

On va noter m le montant de la transaction et on va considérer la matrice des coûts suivante :

- si on accepte une bonne transaction (TN) : la marge **générée** par l'enseigne est égale à 5% du montant de la transaction, *i.e.*, on réalise un gain de $r \times m$ où $r = 0.05$

- si on accepte une mauvaise transaction (FN) : on considère que le perte peut être décrite de la façon suivante :

$$\text{pertes} = \begin{cases} 0 & \text{si } m \leq 20, \\ 0.2 \times m & \text{si } 20 < m \leq 50, \\ 0.3 \times m & \text{si } 50 < m \leq 100, \\ 0.5 \times m & \text{si } 100 < m \leq 200, \\ 0.8 \times m & \text{si } m > 200. \end{cases}$$

- lorsque vous refusez une bonne transaction (FP) : vous **générez** une marge égale à 70% du montant de la transaction multiplié par le taux de marge r , *i.e.* égale à $0.7 \times m \times r$
- lorsque vous refusez une transaction frauduleuse, vous ne perdez ni ne gagnez d'argent.

Résumé

Il s'agit du résumé de votre travail, vous devez donc, en une dizaine de lignes environ, présenter le problème traité ainsi que l'approche proposée dans le rapport. Enfin vous finirez par donnée une idée des conclusions obtenues.

Introduction (et éventuellement état de l'art)

Dans cette section, vous devrez présenter le contexte de l'étude (on pourra par exemple reprendre la description donnée au début du document) ainsi que les difficultés liées à cette étude. On pourra également citer quelques exemples d'applications liées à la problématique traitée.

Si vous avez consulté quelques références, liens sur internet ou autre, c'est le moment de les citer en évoquant en deux ou trois lignes la méthode employée.

Cette section se finira par une présentation de la structure de votre rapport, par exemple.

La Section 2 sera consacrée à l'analyse des données. La Section 3 sera consacrée à la présentation de la méthode proposée pour résoudre ce problème ...

Analyse des données

On évoque rapidement le jeu de données. Vous commencerez par une analyse synthétique des données à l'aide d'outils statistiques élémentaires : vous présenterez quelques graphes pertinents et intéressants s'il y a lieu ainsi que les grandes caractéristiques du jeu de données, sélectionnez les informations intéressantes.

Méthodologie

Vous commencerez par présenter les notations que vous allez employer tout au long de la rédaction de votre rapport.

Vous présenterez ensuite les outils que vous allez utiliser dans la partie expérimentale. On commencera par parler de ce que l'on souhaite maximiser (par exemple AUCROC ou F-mesure) en les définissant avant de s'attaquer à la présentation des algorithmes.

Par exemple, si vous faites une contribution basée sur du boosting et que vous combinez avec des méthodes à noyaux, il faudra rappeler ce qu'est le boosting, ce que sont les méthodes à noyaux (ce sont les approches de bases) et ensuite vous expliquerez comment vous combinez les deux pour résoudre le problème confié.

Il ne faut pas hésiter à présenter le processus de façon *abstraite*, c'est-à-dire avec des notations mathématiques et ne pas être uniquement verbeux.

Un pseudo-code est également appréciable pour synthétiser l'approche proposée.

Dans le cas où vous proposez plusieurs approches à des fins de comparaisons, il faudra prendre soin de présenter les différentes approches et de justifier pourquoi vous intéressez à ces approches là.

Remarque : il n'est pas nécessaire de présenter tous les algorithmes employés, mais uniquement ceux qui servent à l'élaboration d'une version "exotique".

Expériences

Vous dresserez ensuite votre protocole expérimentale qui présentera la ou les méthodes sélectionnées pour répondre à la tâche demandée. Celui-ci comprend en général trois parties

Présentation de données

Vous représentez ici rapidement le jeu de données employées ainsi que ses caractéristiques : nombre d'exemples, dimensions, taux de déséquilibre,...

Protocole expérimental

Vous présentez rapidement les expériences que vous allez faire, *i.e.* les différents algorithmes testés, le range des hyper-paramètres employés ainsi que la façon dont sont optimisées ces hyper-paramètres (cross-validation en k -folds, simple validation ou est-ce que vous faites le choix de les fixer). Quels sont vos ensembles d'entraînement/validation/test ?

Les informations que vous fournissez dans cette section doivent permettre au lecteur de pouvoir reproduire les résultats que vous allez présenter dans la sous-section suivante.

Résultats

Ici vous allez présenter et analyser les résultats obtenus à l'aide de graphiques et/ou tableaux. Outre les performances, on pourra aussi s'intéresser au critère de rapidité d'un algorithme.

L'analyse doit aussi permettre de mettre en exergue les avantages/inconvénients des méthodes proposées. Cela peut passer par l'utilisation d'autres mesures de performances/critères pour évaluer/comparer vos algorithmes.

Conclusions

Il s'agit de conclure quant à votre étude. Reprendre le travail proposé et les principales conclusions. Il est également important de proposer des perspectives à votre travail en fonction des résultats obtenus et de l'approche proposée.