Financial Fraud Detection

Prénom : Denisa Année: M2 Semestre: 9 Nom: **BANULESCU-RADU** Nature: CM Volume horaire: 12 H ECTS / Coef: 2 Notions d'économétrie linéaire et des variables qualitatives Prérequis Bases en statistiques, probabilités et optimisation Connaissance des méthodes de machine learning (supervisé & non supervisé) Ce cours a pour objectif de former les participants aux méthodes économétriques et aux techniques d'apprentissage automatique, supervisées et non supervisées, appliquées à la détection de la fraude financière. Après une introduction aux principales typologies de fraude, l'accent sera mis sur l'analyse et le traitement des bases de données contenant des observations frauduleuses. Deux grandes catégories de modèles seront étudiées, ainsi que leurs mesures de performance : Modèles de prévention (non supervisés) : conçus pour identifier, en amont, les transactions ou clients présentant des comportements atypiques. Résumé Modèles de détection (supervisés) : construits à partir de données historiques afin de classer les nouvelles observations comme frauduleuses ou non. Un enjeu central est le caractère rare de la fraude : les bases de données utilisées sont très volumineuses mais fortement déséquilibrées (par exemple, les fraudes à la carte de crédit représentent souvent moins de 0,5 % des transactions). La dernière partie du cours portera ainsi sur les méthodes permettant de corriger ce déséquilibre et d'améliorer la robustesse des modèles. Le cours sera complété par la réalisation d'un projet pratique mettant en œuvre les techniques étudiées sur une base de données spécifique. Comprendre la définition et les principales typologies de fraude financière Identifier et appliquer les techniques analytiques utilisées pour la détection de la fraude Analyser et traiter des bases de données contenant des cas de fraude Sélectionner les variables pertinentes pour la prévention et la détection de la fraude Mettre en œuvre des méthodes d'apprentissage automatique supervisées et non supervisées adaptées à la détection de la fraude **Objectifs** Évaluer des modèles économétriques et de machine learning à l'aide de mesures de performance adaptées (classification et régression) Utiliser les principales méthodes de rééchantillonnage pour traiter les bases de données déséquilibrées (oversampling, undersampling, SMOTE, etc.) Réaliser une étude de cas pratique appliquée à la détection de la fraude Baesens, B., Van Vlasselaer, V., and Verbeke, W. (2015). Fraud analytics using descriptive, predictive, and social network techniques: a guide to data science for fraud detection. John Wiley & Sons. Hastie, T., Tibshirani, R., and Friedman, J. (2009). The elements of statistical learning: data mining, inference, and prediction. Springer Science & Business **Bibliographie** Fernández, A., García, S., Galar, M., Prati, R. C., Krawczyk, B., and Herrera, F. (2018).

Learning from imbalanced data sets. Springer.

applications. John Wiley & Sons.

He, H. and Ma, Y. (2013). Imbalanced learning: foundations, algorithms, and

PLAN

- **Chapter 1 Introduction**
- Chapter 2 Data
 - o 2.1. Data typologies and sources
 - o 2.2. Operations on data
- Chapter 3 Descriptive analytics for fraud detection Unsupervised learning
 - o 3.1. Outlier detection
 - o 3.2. Clustering
- Chapter 4 Predictive analytics for fraud detection Supervised learning
 - 4.1. Linear regression
 - o 4.2. Logistic regression
 - 4.3. Decision trees
 - 4.4. Ensemble methods: bagging, boosting, random forest
- Chapter 5 Predictive models for skewed datasets
 - o 5.1. Undersampling
 - 5.2. Oversampling
 - 5.3. Adjusting posterior probabilities
 - o 5.4. Cost-sensitive learning
- Chapter 6 Evaluation of predictive models for fraud detection
 - o 6.1. Data splitting
 - o 6.2. Performance measures for classification models
 - o 6.3. Illustration
- **Chapter 7 Cost-sensitive learning for fraud detection**
 - o 7.1. Cost matrix
 - o 7.2. Cost-sensitive Logistic regression
 - o 7.3. Cost-Sensitive evaluation metrics
 - o 7.4. Illustration