# Apprentissage de classes déséquilibrées HAX907X - Apprentissage statistique

SAWADOGO Kader GERMAIN Marine LABOURAIL Célia MARIAC Damien

Université Montpellier Département de Mathématique

17 octobre 2025

- Contexte
- 2 Problématique
- Méthodes
- 4 Limites des méthodes
- 6 Application
- 6 Conclusion

# Problème du déséquilibre et motivation du projet

#### Un jeu de données très déséquilibré

- $\bullet$  CréditCard : 284 807 transactions, dont seulement 492 fraudes (0,17 %).
- ⇒ Les modèles ont tendance à ignorer la classe minoritaire.

#### Illustration du problème

Modèle	Recall	Precision	F1-score
Régression Logistique	0.59	0.89	0.70
Random Forest (200 arbres)	0.60	0.95	0.85

- Contexte
- 2 Problématique
- Méthodes
- 4 Limites des méthodes
- 6 Application
- 6 Conclusion

# Les classes déséquilibrées

Problème : difficulté à prédire la classe minoritaire

 $\Rightarrow$  Le modèle a tendance à ignorer cette classe.

#### Exemple général

- 99% vs 1%.
- ullet Un modèle naı̈f prédit la classe majoritaire à une précision de 99 %.
- Mauvais modèle.

# Problématique scientifique

Comment atténuer le déséquilibre des classes pour améliorer la performance réelle du modèle?

- Contexte
- 2 Problématique
- Méthodes
- 4 Limites des méthodes
- 6 Application
- 6 Conclusion

# Random Over-sampling

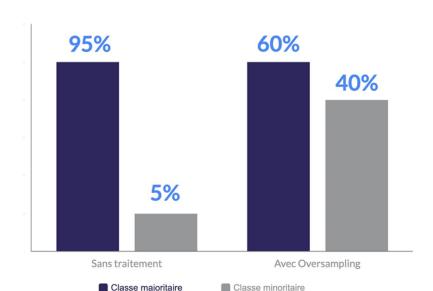


Table – Jeu de données après sur-échantillonnage (ROS )

х	label	source
1	0	original
2	0	original
3	0	original
4	0	original
5	0	original
6	0	original
7	0	original
8	1	original
9	1	original
10	1	original
8	1	duplicated
9	1	duplicated
10	1	duplicated
8	1	duplicated

# Random Under-sampling

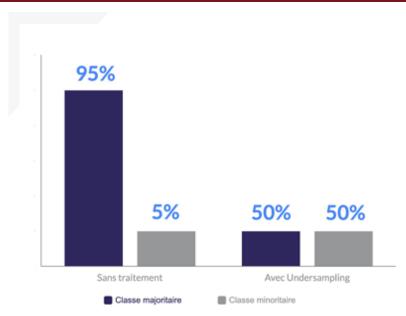


Table – Jeu de données après sous-échantillonnage (RUS )

X	label	source
1	0	supprimé
2	0	suprimé
3	0	suprimé
4	0	original
5	0	original
6	0	original
7	0	original
8	1	original
9	1	original
10	1	original

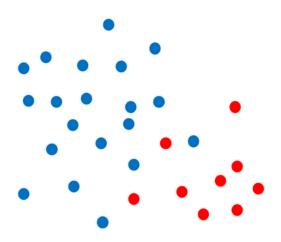


Figure – Schéma de SMOTE

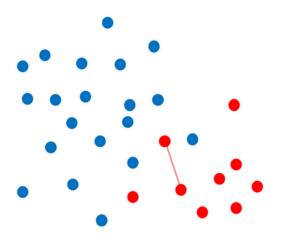


Figure – Schéma de SMOTE

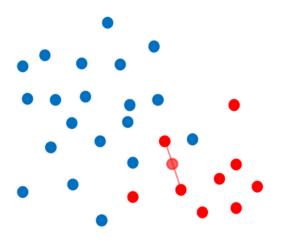


Figure – Schéma de SMOTE

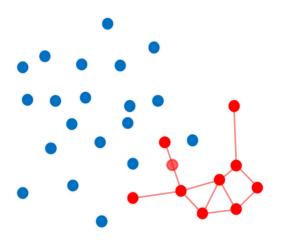


Figure – Schéma de SMOTE

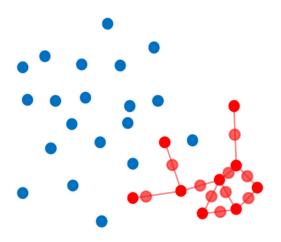


Figure – Schéma de SMOTE

#### **Notations**

- n : nb total d'observations
- n<sub>min</sub>: nb d'observations minoritaires
- d : dimension (nb de variables)

- k : nb de plus proches voisins
- M : nb de points synthétiques générés

# Étapes dominantes & complexité (naïf)

- **Q** Recherche des k-PPV (vers tous les points) : coût d'une distance  $\mathcal{O}(d) \Rightarrow$  comparaison à n points  $\mathcal{O}(n\,d) \Rightarrow$  pour  $n_{\min}$  points minoritaires  $\left[\mathcal{O}(n_{\min}\,n\,d)\right]$ .
- **2 Génération** :  $x_{\text{new}} = x_i + \lambda(x_j x_i)$ ,  $\lambda \sim \mathcal{U}(0, 1)$

 $\mathcal{O}(Md)$ 

#### Synthèse

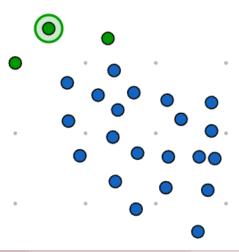
 $T_{\text{SMOTE}} = \mathcal{O}(n_{\min} n d) + \mathcal{O}(M d)$ 

(recherche kNN dominante).

- Contexte
- 2 Problématique
- Méthodes
- 4 Limites des méthodes
- 6 Application
- 6 Conclusion

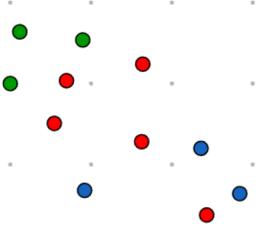
# ROS

overfiting



# **RUS**

• perte d'information

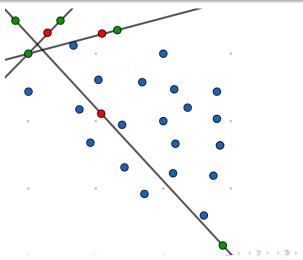


#### **SMOTE**

- temps de calculs
- création de points abérrants
- hyperparamètre k
- variable qualitative

# SMOTE

Points abérrants



- Contexte
- 2 Problématique
- Méthodes
- 4 Limites des méthodes
- 6 Application
- 6 Conclusion

Méthode	Transactions normales	Fraudes
ROS (Over-Sampling)	284 315	284 315
RUS (Under-Sampling)	492	492
SMOTE	284 315	284 315

# Comparaison des méthodes de rééchantillonnage (Régression Logistique)

Méthode	Accuracy	Recall	Balanced Accuracy
RUS	0.956	0.959	0.956
ROS	0.949	0.977	0.949
SMOTE	0.981	0.992	0.970

- Contexte
- 2 Problématique
- Méthodes
- 4 Limites des méthodes
- 6 Application
- 6 Conclusion

# Bilan général des méthodes

Méthode	Points forts	Limites
ROS	Simplicité, conserve	Overfitting, grand
	toutes les données	volume de données
RUS	Rapide et réduit le	perte d'information
	biais	et représentativité
SMOTE	Données	Coût élevé, sensible
	synthétiques variées	aux outliers

#### Aucune méthode n'est universelle :

le choix dépend du jeu de données et du modèle.

# Conclusion et perspectives

- Pour notre jeu de données, la méthode la plus efficace est SMOTE.
- Pour aller plus loin : il serait pertinent de combiner des méthodes existantes ou de pondérer les modèles.

# Merci pour votre attention!