



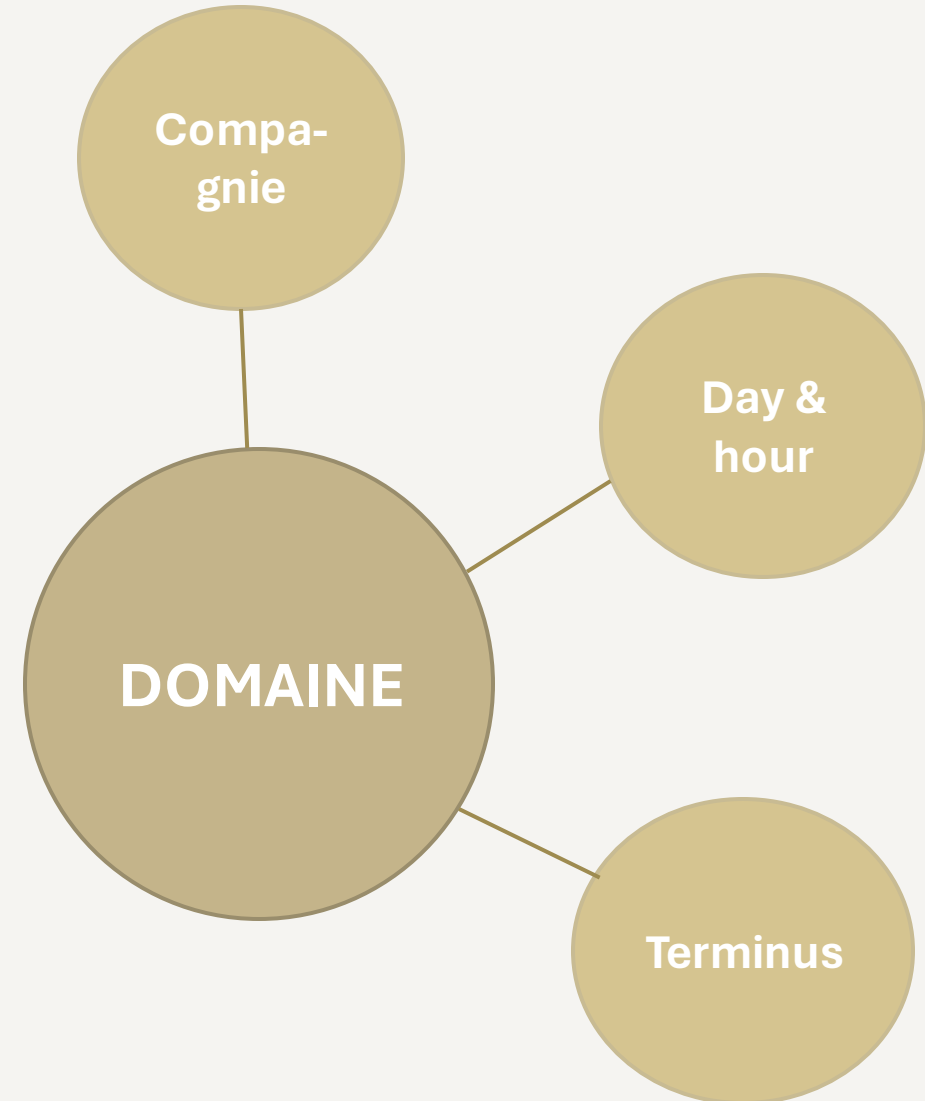
# Présentation sprint : Compagnie de trains polonaise

PAR MAUD TISSOT, ARTHUR SCelles, ALEXANDRE SMADJA,  
AURÉLIEN POUXVIEL

# Peut-on prévoir un jour à l'avance le retard d'un train ?



12/13/2023



# I. Traduction du problème en code

Après traitement des données nous avons donc :

- 241 916 lignes
- 24 variables

## TARGET

- retard OUI/NON si  $>10\text{min}$

## PREPROCESS

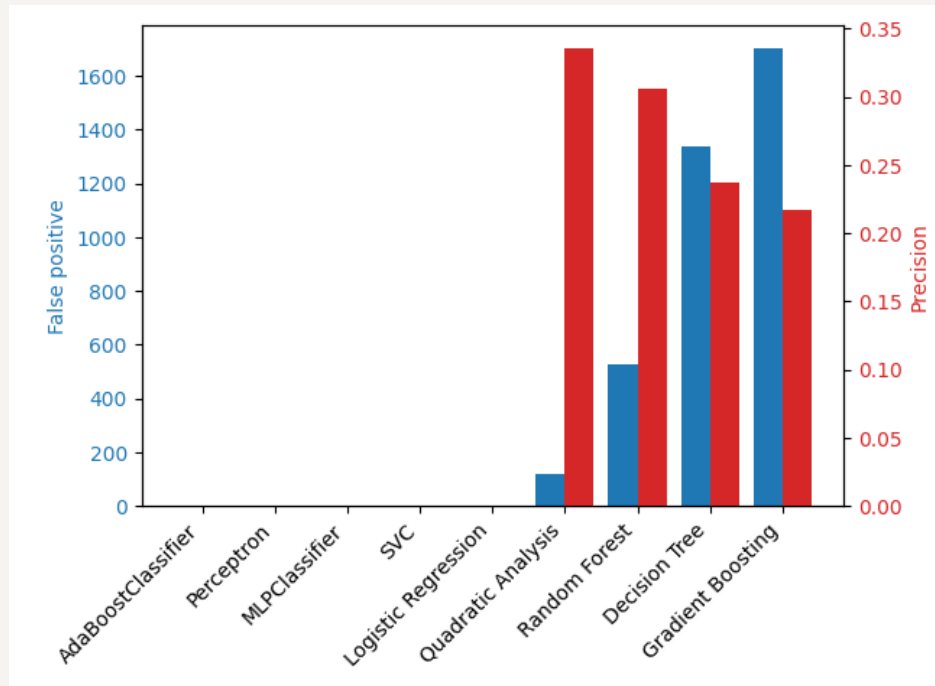
- Filtre 1 trajet = 1 ligne
- New features

## TRAIN / TEST SLIT

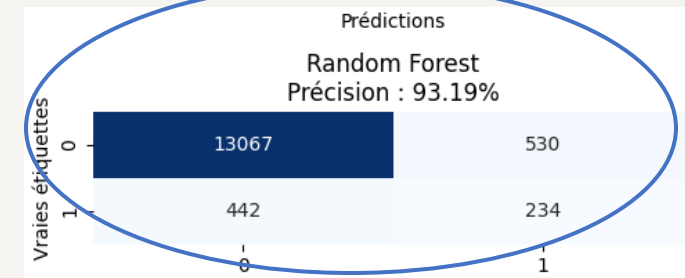
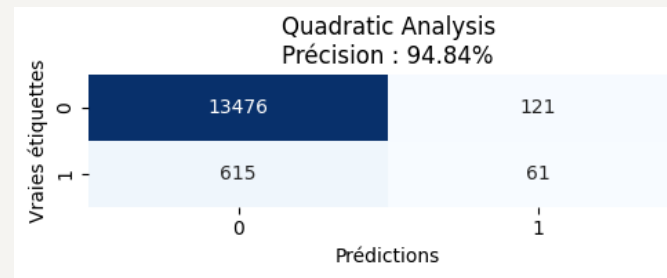
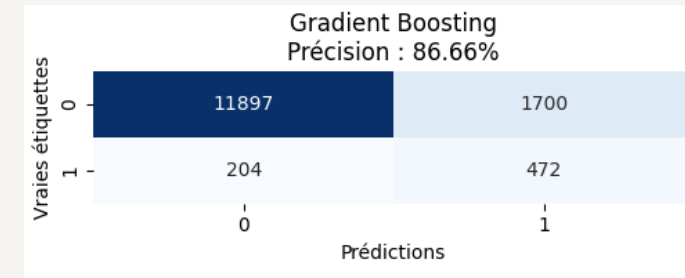
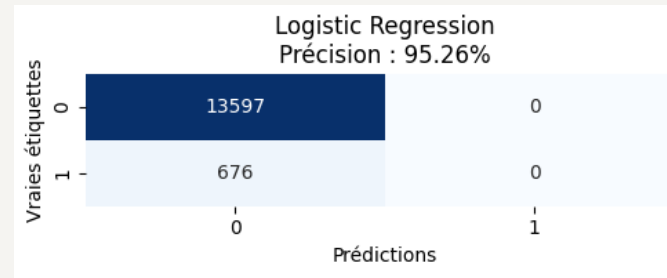
- Test = 1 jour
- Train = dataset avant j-2

## 2. Prédiction du retard un jour à l'avance

- Nous avons essayé différents modèles de différents types
- Choix du modèle



$$Precision = \frac{TP}{TP + FP}$$

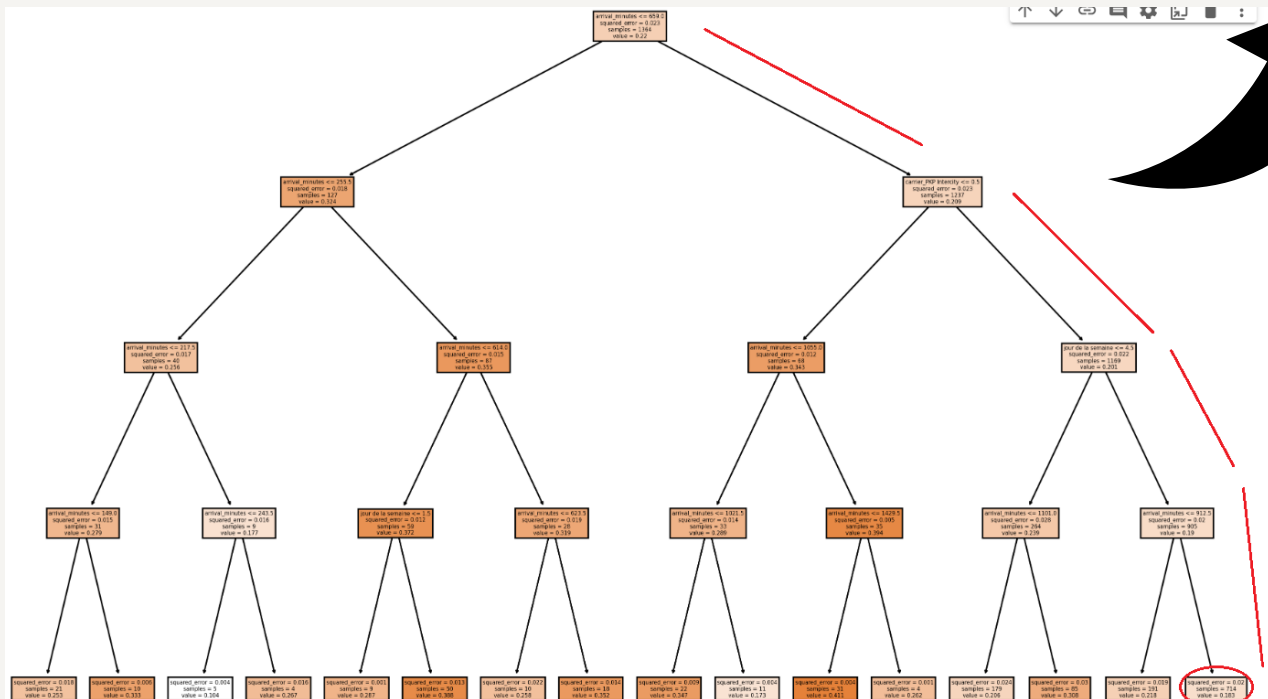


Justement, on ne se fie pas **aveuglément** à la precision et on ne prend pas de modèle 'naif'



# 3. Arbre 1 - Identifier les domaines pour lesquels le retard est prévisible

- Calcul de l'erreur
- 1ere arbre sur l'horaire et les jours = Périodicité des trains



= Trains à partir de 15 h de PKP, pour la fin de semaine (vendredi, samedi, dimanche). 70% de précision (714 / 1100 retards)

arrival\_minutes <= 659.0  
squared\_error = 0.023  
samples = 1364  
value = 0.22

FALSE

carrier\_PKP Intercity <= 0.5  
squared\_error = 0.023  
samples = 1237  
value = 0.209

FALSE

jour de la semaine <= 4.5  
squared\_error = 0.022  
samples = 1169  
value = 0.201

FALSE

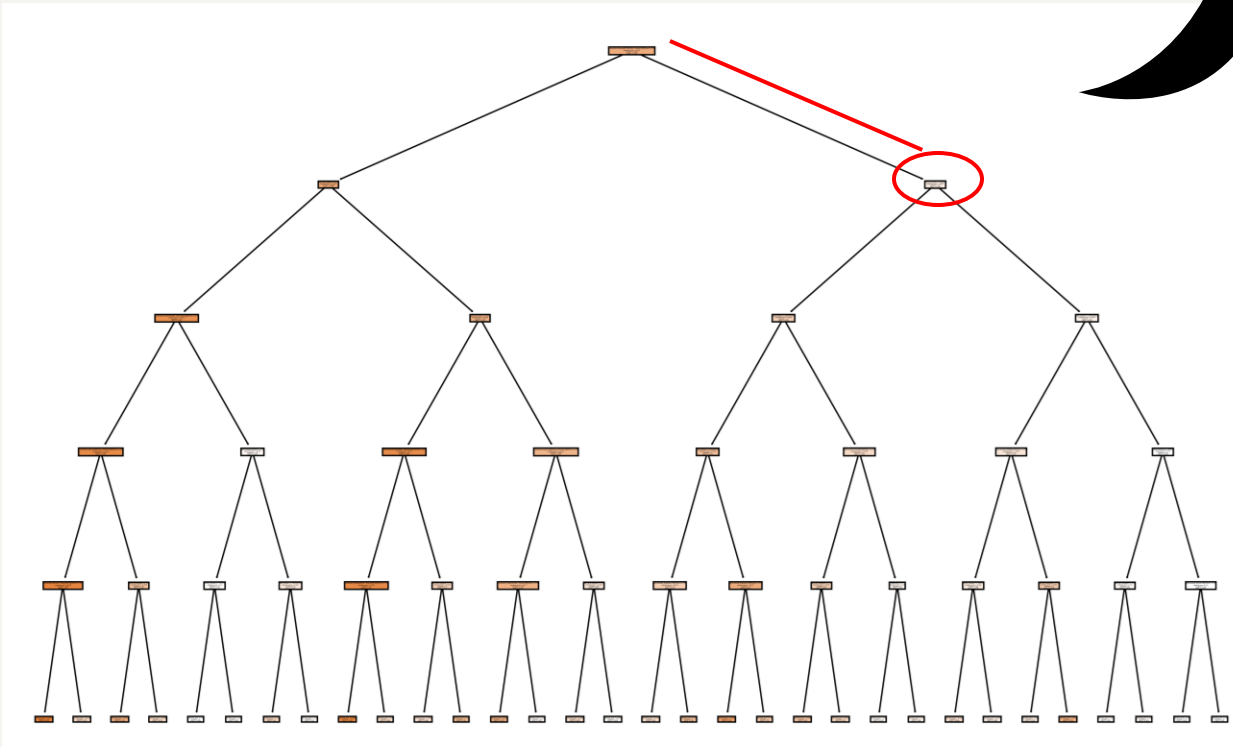
arrival\_minutes <= 912.5  
squared\_error = 0.02  
samples = 905  
value = 0.19

FALSE

squared\_error = 0.02  
samples = 714  
value = 0.183

### 3. Arbre 2 - Géographie (gares, connections)

- Trop de profondeurs = pas toujours pertinent



```
connection_Przemyśl Główny - Gdynia Główna <= 0.5  
squared_error = 0.023  
samples = 1364  
value = 0.224
```

FALSE

```
arrival_minutes <= 1018.5  
squared_error = 0.006  
samples = 84  
value = 0.055
```

= Trains de la connexion **Przemyśl Główny - Gdynia Główna**

**202 trains.**

**97 retards**

**84 prédits sur 97 = 86% de précision**

# Conclusion

## Prédiction large

Nous sommes capables par exemple de proposer à une entreprise comme PKP intercity de prédire leurs retards en fin de journée 24h à l'avance avec une précision de 70%

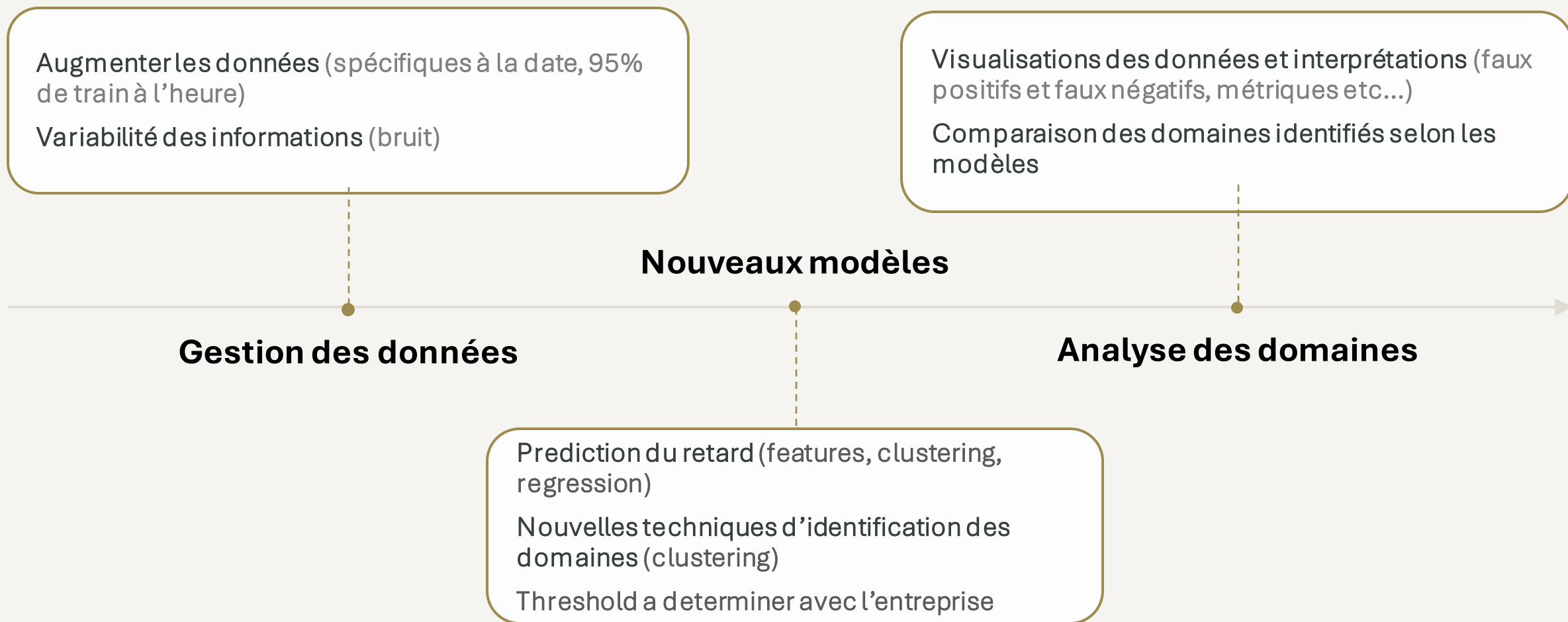
## Proposition d'un domaine bien plus précis

La ligne **Przemyśl Główny - Gdynia Główna** :

- 202 trajets en une semaine
- 1 train sur 2 est en retard

Nous proposons de prédire 24h à l'avance les retards de cette ligne avec une précision de 86%

# Critiques et perspectives





**Merci de votre écoute**

# ANNEXES

