

## Exercice – Préparation et nettoyage des données (en Python)

### 1. Découpage des données

Séparer le jeu de données en **Train / Validation / Test** (par exemple : 60% / 20% / 20%).

### 2. Typage

Convertir les colonnes numériques importées en texte en types numériques appropriés :

Montant\_raw, RevenuMensuel\_raw, Dette\_raw, etc.

### 3. Unités / devises

- Nettoyer la colonne Montant\_raw.
- Convertir tous les montants en **HTG** (taux : **1 USD = 130 HTG**).
- Gérer différents formats, par exemple :
  - '1,200.50'
  - '20k'
  - '500 USD'
  - '12.50USD'

### 4. Dates

- Parser la colonne DateTransaction\_raw (formats multiples + valeurs invalides).
- Créer des variables temporelles :
  - jour\_semaine
  - mois
  - heure (si l'heure peut être reconstruite)

### 5. Doublons

DéTECTER et supprimer les lignes dupliquées (certaines sont **exactement identiques**).

### 6. Valeurs manquantes

- Traiter les valeurs manquantes (NaN, 'N/A', 'Unknown').
- Ajouter des indicateurs binaires si pertinent (ex. : revenu manquant).

### 7. Variables catégorielles

- Normaliser Ville\_raw (variantes PAP, P-au-P, etc.).
- Appliquer un **One-Hot Encoding** pour :

- Canal
  - Device
  - Ville\_norm
- Appliquer un **Ordinal Encoding** pour NiveauEtude.

## 8. Valeurs aberrantes (Outliers)

Identifier et traiter les valeurs aberrantes, par exemple :

- Age  $\in [-4, 150]$
- NbTrans\_24h  $\in [120, 250]$
- Montants très élevés

Tester différentes approches :

- Clipping
- log1p
- RobustScaler

## 9. Mise à l'échelle

Standardiser ou normaliser les variables numériques pour les modèles sensibles aux distances ou au gradient :

- Régression logistique
- KNN
- SVM

## 10. Data leakage

Respecter la règle suivante :

- **Fit** des transformations uniquement sur X\_train
- **Transform** ensuite sur X\_validation et X\_test