

Détection de Fraude Bancaire

Guide Complet : de l'Exploration des Données à la Modélisation

Fiche Projet – Data Science

2026

1. Introduction et Contexte

⌚ Objectif du Projet

Ce projet vise à développer un système de détection automatique de fraudes bancaires en utilisant des techniques de Machine Learning. L'objectif est de :

- Analyser les caractéristiques des transactions frauduleuses,
- Construire un modèle prédictif performant,
- Évaluer rigoureusement sa capacité de détection.

La fraude bancaire représente un enjeu financier majeur pour les institutions financières. Détecter ces transactions frauduleuses de manière automatisée permet de :

- Réduire les pertes financières,
- Protéger les clients,
- Améliorer la confiance dans le système bancaire.

2. Phase 1 : Exploration des Données (EDA)

L'analyse exploratoire des données est une étape fondamentale qui précède toute modélisation.

2.1 Chargement et Inspection Initiale

```

1 import pandas as pd
2 import numpy as np
3
4 # Chargement du dataset
5 data_credit = pd.read_csv('creditcard_2023.csv')
6
7 # Aperçu des données
8 print(data_credit.head())
9 print(data_credit.info())
10 print(data_credit.describe())

```

💡 Points Clés

Cette première inspection permet d'identifier :

- **La variable cible** : `Class` (0 = Normal, 1 = Fraude)
- **Les features** : Variables anonymisées V_1, V_2, \dots, V_{28} et `Amount`
- **Les types de données** et la présence éventuelle de valeurs manquantes

2.2 Analyse de la Distribution des Classes

```

1 # Vérification de l'équilibre des classes
2 counts = data_credit['Class'].value_counts()
3 print(f"Transactions normales (0) : {counts[0]}")
4 print(f"Transactions frauduleuses (1) : {counts[1]}")
5
6 # Visualisation
7 import matplotlib.pyplot as plt
8 data_credit['Class'].value_counts().plot(kind='bar')
9 plt.title('Distribution des Classes')
10 plt.xlabel('Class')
11 plt.ylabel('Nombre de transactions')
12 plt.show()

```

⚠️ Attention - Contexte Pédagogique

Particularité de ce dataset pédagogique :

Les classes sont équilibrées (50% normales / 50% frauduleuses). Cette configuration simplifie l'apprentissage mais ne reflète **pas la réalité**.

2.3 Analyse des Corrélations

La corrélation de Pearson mesure la relation linéaire entre chaque variable et la fraude :

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} \quad (1)$$

```

1 # Calcul des corrélations avec la cible
2 correlation = data_credit.corr()['Class'].sort_values(ascending
   =False)
3 print(correlation)
4
5 # Visualisation heatmap (optionnel)
6 import seaborn as sb
7 plt.figure(figsize=(12, 8))
8 sb.heatmap(data_credit.corr(), cmap='coolwarm', center=0)
9 plt.title('Matrice de Corrélation')
10 plt.show()

```

💡 Points Clés

Les variables avec $|\rho|$ élevé (ex : V14, V4, V10) sont des candidats prioritaires pour la modélisation car elles contiennent une forte information prédictive.

2.4 Visualisation des Variables Clés

```

1 # Boxplots pour comparer les distributions
2 import seaborn as sb
3
4 fig, axes = plt.subplots(2, 2, figsize=(14, 10))
5
6 sb.boxplot(x='Class', y='V14', data=data_credit, ax=axes[0,0])
7 sb.boxplot(x='Class', y='V4', data=data_credit, ax=axes[0,1])
8 sb.boxplot(x='Class', y='V10', data=data_credit, ax=axes[1,0])
9 sb.boxplot(x='Class', y='Amount', data=data_credit, ax=axes
   [1,1])
10
11 plt.tight_layout()
12 plt.show()

```

💡 À Retenir

Une bonne séparation visuelle entre les deux classes indique un fort pouvoir discriminant de la variable.

3. Phase 2 : Prétraitement des Données

Le prétraitement transforme les données brutes en format exploitable par les algorithmes de ML.

3.1 Les Quatre Piliers du Preprocessing

1. Sélection des Features (Feature Selection)

Conservation uniquement des variables à fort pouvoir prédictif pour réduire le bruit.

2. Mise à l'Échelle (Scaling)

Standardisation des variables pour éviter qu'une variable à grande échelle domine le modèle.

3. Nettoyage

Vérification et traitement des valeurs manquantes et des doublons.

4. Séparation Train/Test

Division du dataset pour entraîner et évaluer objectivement le modèle.

3.2 Standardisation (Z-score)

La standardisation ramène chaque variable à une distribution centrée réduite :

$$z = \frac{x - \mu}{\sigma} \quad (2)$$

où x = valeur d'origine, μ = moyenne, σ = écart-type.

3.3 Implémentation Complète

```

1  from sklearn.preprocessing import StandardScaler
2  from sklearn.model_selection import train_test_split
3
4  # 1. Selection des features importantes
5  important_features = ['V14', 'V4', 'V10', 'V12', 'V17',
6                      'V2', 'V3', 'V9', 'V7', 'Amount']
7  X = data_credit[important_features].copy()
8  y = data_credit['Class']
9
10 # 2. Verification des donnees manquantes
11 print(f"Valeurs manquantes : {X.isnull().sum().sum()}")
12
13 # 3. Standardisation de la variable Amount
14 scaler = StandardScaler()
15 X['Amount'] = scaler.fit_transform(X[['Amount']])
16
17 # 4. Separation Train/Test (80/20)
18 X_train, X_test, y_train, y_test = train_test_split(
19     X, y, test_size=0.2, stratify=y, random_state=42
20 )
21
22 print(f"Taille train : {X_train.shape}")
23 print(f"Taille test : {X_test.shape}")

```

À Retenir**Pourquoi stratify=y ?**

Cette option garantit que la proportion de fraudes est identique dans les ensembles d'entraînement et de test (ici 50/50), assurant une évaluation fiable.

4. Phase 3 : Modélisation par Régression Logistique

4.1 Principe de la Régression Logistique

La régression logistique modélise la probabilité qu'une transaction soit frauduleuse via la fonction sigmoïde :

$$P(Y = 1 | X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)}} \quad (3)$$

Le modèle estime les coefficients β qui maximisent la vraisemblance des données observées.

4.2 Pourquoi la Régression Logistique ?

- **Interprétabilité** : Les coefficients indiquent l'influence de chaque variable
- **Probabilités** : Permet d'ajuster le seuil de décision selon les coûts métier
- **Efficacité** : Rapide à entraîner, même sur de gros volumes
- **Baseline** : Modèle de référence en finance

4.3 Entraînement du Modèle

```

1  from sklearn.linear_model import LogisticRegression
2
3  # Initialisation et entraînement
4  model = LogisticRegression(
5      solver='lbfgs',
6      max_iter=1000,
7      random_state=42
8  )
9
10 model.fit(X_train, y_train)
11
12 # Predictions
13 y_pred = model.predict(X_test)
14 y_proba = model.predict_proba(X_test)[:, 1]

```

5. Phase 4 : Évaluation des Performances

5.1 Matrice de Confusion

	Prédit Normal	Prédit Fraude
Réel Normal	TN (Vrai Négatif)	FP (Faux Positif)
Réel Fraude	FN (Faux Négatif)	TP (Vrai Positif)

```

1 from sklearn.metrics import confusion_matrix,
2     classification_report
3
4 # Matrice de confusion
5 cm = confusion_matrix(y_test, y_pred)
6 sb.heatmap(cm, annot=True, fmt='d', cmap='Blues')
7 plt.title('Matrice de Confusion')
8 plt.ylabel('Classe Réelle')
9 plt.xlabel('Classe Prédite')
10 plt.show()

```

5.2 Métriques de Performance

Points Clés

Les 4 Métriques Essentielles :

- **Accuracy** : $\frac{TP+TN}{TP+TN+FP+FN}$ – Taux de bonnes prédictions global
- **Precision** : $\frac{TP}{TP+FP}$ – Fiabilité des alertes fraude
- **Recall (Sensibilité)** : $\frac{TP}{TP+FN}$ – Capacité à détecter les fraudes réelles
- **F1-Score** : $2 \times \frac{Precision \times Recall}{Precision + Recall}$ – Compromis optimal

```

1 from sklearn.metrics import accuracy_score, precision_score,
2     recall_score, f1_score
3
4 print(f"Accuracy : {accuracy_score(y_test, y_pred):.4f}")
5 print(f"Precision : {precision_score(y_test, y_pred):.4f}")
6 print(f"Recall : {recall_score(y_test, y_pred):.4f}")
7 print(f"F1-Score : {f1_score(y_test, y_pred):.4f}")
8
9 # Rapport détaillé
10 print("\n", classification_report(y_test, y_pred))

```

5.3 Courbe ROC et AUC

```

1  from sklearn.metrics import roc_curve, roc_auc_score
2
3  # Calcul de l'AUC
4  auc = roc_auc_score(y_test, y_proba)
5  print(f"AUC-ROC : {auc:.4f}")
6
7  # Trace de la courbe ROC
8  fpr, tpr, thresholds = roc_curve(y_test, y_proba)
9  plt.figure(figsize=(8, 6))
10 plt.plot(fpr, tpr, label=f'ROC Curve (AUC = {auc:.4f})')
11 plt.plot([0, 1], [0, 1], 'k--', label='Random Classifier')
12 plt.xlabel('Taux de Faux Positifs')
13 plt.ylabel('Taux de Vrais Positifs')
14 plt.title('Courbe ROC')
15 plt.legend()
16 plt.grid(alpha=0.3)
17 plt.show()

```

$$0.5 \leq AUC \leq 1 \quad (4)$$

Plus l'AUC est proche de 1, meilleure est la capacité discriminante du modèle.

6. Résultats du Projet

6.1 Performances Obtenues

Métrique	Valeur
Accuracy	0.9580
Precision	0.9791
Recall	0.9360
F1-score	0.9571
AUC-ROC	0.9907

6.2 Validation Croisée

```

1  from sklearn.model_selection import cross_val_score
2
3  # Validation croisée 5-fold
4  cv_scores = cross_val_score(model, X, y, cv=5, scoring='f1')
5  print(f"F1-Score moyen : {cv_scores.mean():.4f} (+/- {cv_scores.std():.4f})")

```

Métrique	Moyenne ± Écart-type
F1-score (CV)	0.9573 ± 0.0010

💡 Points Clés**Interprétation des Résultats :**

- Recall de 93,6% : le modèle détecte correctement plus de 9 fraudes sur 10
- Precision de 97,9% : très peu de fausses alertes
- AUC de 0.99 : excellente capacité de discrimination
- Faible variance : modèle stable et robuste

7. Limites et Perspectives

7.1 Contexte Pédagogique vs Réalité

⚠️ Attention - Contexte Pédagogique**Différences Majeures avec la Réalité Bancaire :**

1. **Déséquilibre extrême** : En pratique, les fraudes représentent moins de 0,1% à 1% des transactions (et non 50%). Cela nécessite :
 - Des techniques de rééquilibrage (SMOTE, sous-échantillonnage)
 - Des métriques adaptées (Precision-Recall plutôt qu'Accuracy)
 - Un seuil de décision ajusté selon les coûts métier
2. **Coûts asymétriques** : Une fraude non détectée (FN) coûte bien plus cher qu'une fausse alerte (FP). Le modèle doit privilégier le Recall.
3. **Évolution temporelle** : Les fraudeurs adaptent leurs techniques. Le modèle nécessite un réentraînement régulier.

7.2 Méthodes Plus Avancées

En pratique, les institutions financières utilisent des approches plus sophistiquées :

- **Random Forest** : Capture les interactions non-linéaires entre variables
- **XGBoost / LightGBM** : Performance supérieure grâce au boosting
- **Réseaux de Neurones** : Deep Learning pour patterns complexes
- **Isolation Forest** : Détection d'anomalies non supervisée
- **Approches hybrides** : Combinaison de plusieurs modèles (ensemble)

7.3 Améliorations Possibles

1. **Feature Engineering** : Créer de nouvelles variables (fréquence, montants moyens, géolocalisation)
2. **Ajustement du seuil** : Optimiser le seuil de décision selon les coûts métier

3. **Analyse temporelle** : Intégrer l'historique des transactions
4. **Explainability** : Utiliser SHAP ou LIME pour interpréter les prédictions

8. Conclusion

✓ Synthèse du Projet

Ce projet a permis de :

- Maîtriser le pipeline complet d'un projet ML : EDA → Preprocessing → Modélisation → Évaluation
- Construire un modèle de régression logistique performant ($AUC = 0.99$)
- Comprendre les métriques adaptées à la détection de fraude
- Identifier les différences entre contexte pédagogique et réalité opérationnelle

❶ À Retenir

Points à retenir :

La régression logistique constitue une excellente **baseline** pour la détection de fraude grâce à son interprétabilité et sa rapidité. Toutefois, dans un contexte réel avec déséquilibre extrême (99% vs 1%), des modèles plus sophistiqués (XGBoost, Neural Networks) et des techniques de rééquilibrage sont nécessaires pour maximiser le Recall tout en contrôlant les faux positifs.