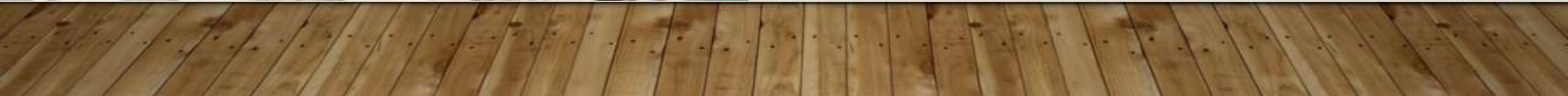




# DÉTECTEZ DES FAUX BILLETS AVEC R OU PYTHON

Par Claude Olukoya



# MA MISSION

---

Mon entreprise vient de décrocher un contrat avec l'Organisation nationale de lutte contre le faux-monnayage (ONCFM) et elle souhaite m'y envoyer en tant que senior data analyst pour cette mission.

Ils mettent en place des méthodes d'identification des faux billets en euros pour lutter contre la contrefaçon et construire une application de machine learning. Celle-ci leur permettra, après avoir scanné des billets (longueur, hauteur, largeur, etc.), de faire une prédiction sur la nature des billets (vrai billet ou faux billet). 4 algorithmes à tester :

1. K-Means

2. Régression logistique

3. KNN

4. Random Forest

# LE FICHIER (CSV) DE DÉPART

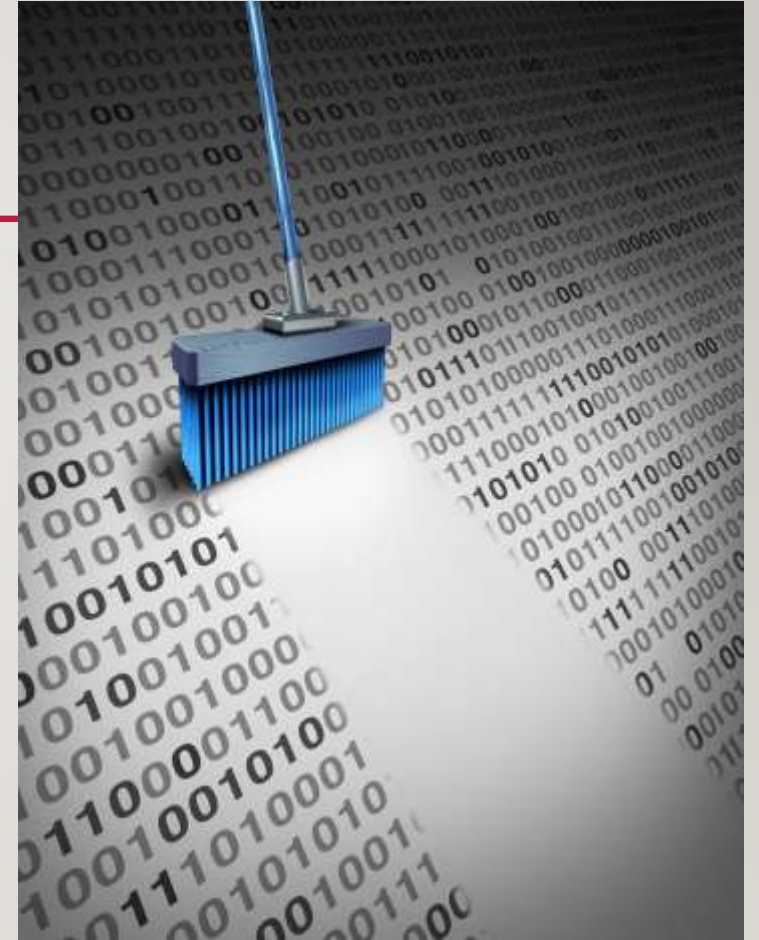
---

Billets.csv



---

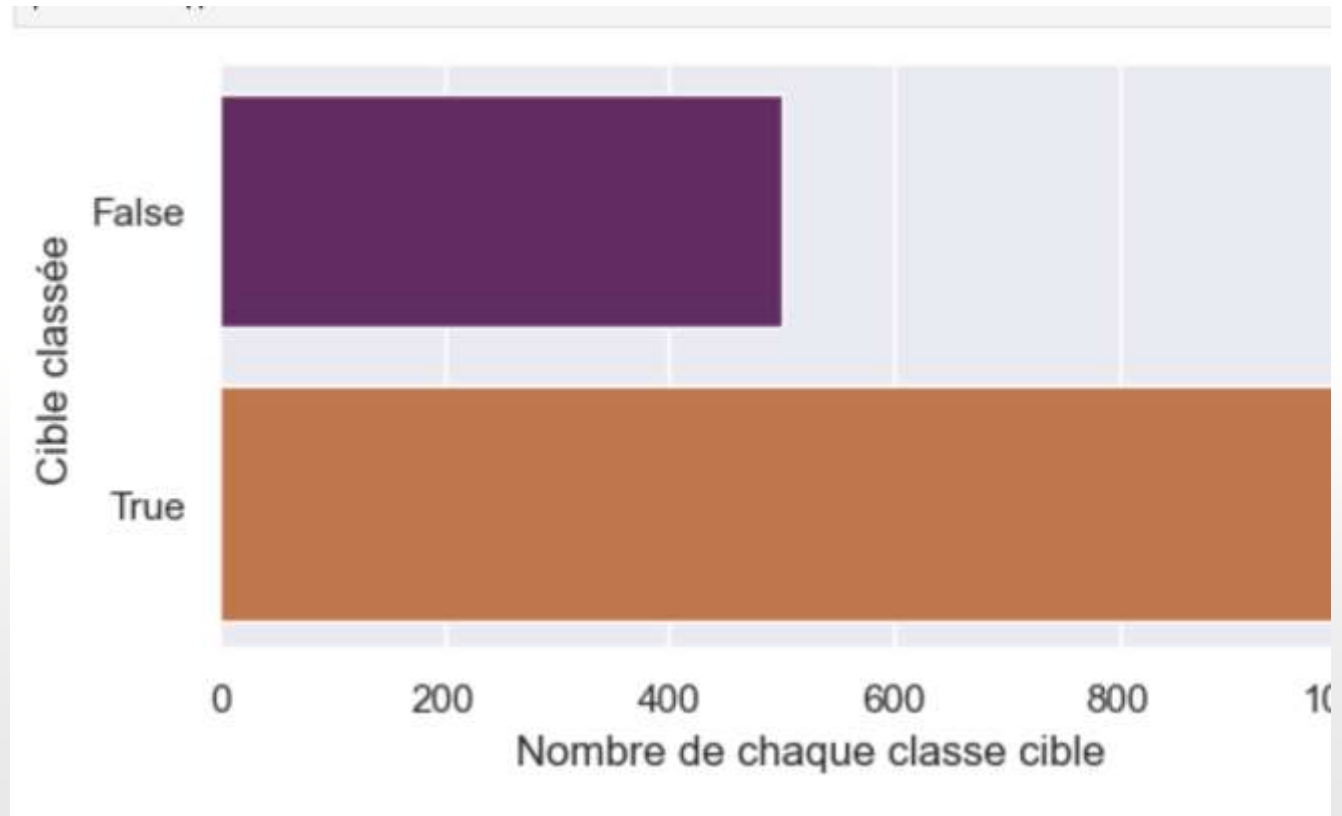
# **PARTIE I : NETTOYAGE & EXPLORATION**



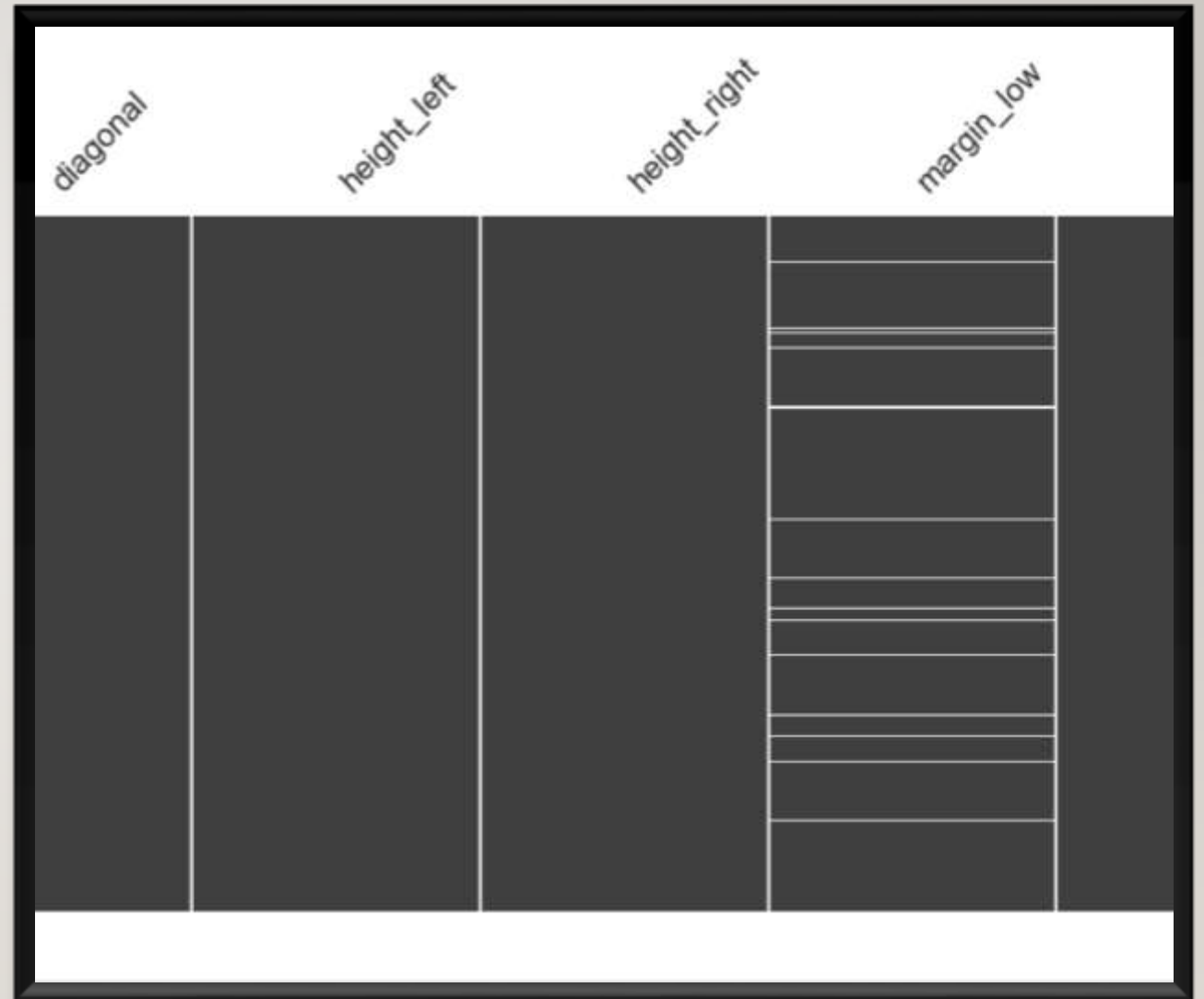


Nous constatons qu'il y a 1,000 vrais billets et 500 faux billets

---



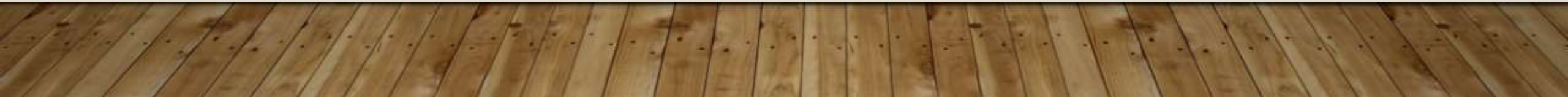
- 
- En utilisant le package *missingno* on dessine une matrice de nullité. Nous constatons qu'il y a 37 valeurs manquantes dans la colonne `margin_low`.



# IMPUTATION DES VALEURS MANQUANTES AVEC 'OLS'

**Étape 1:** *Backward elimination (Elimination inverse)* : utilisée en apprentissage automatique pour identifier le meilleur sous-ensemble de caractéristiques parmi un ensemble donné. Elle consiste à supprimer de manière itérative les caractéristiques qui ne sont pas prédictives de la variable cible ou qui ont le moins de pouvoir prédictif.

**Étape 2:** Une fois qu'on a choisi le nombre de features : *is\_genuine* & *margin\_up*, on fait la prédiction à l'aide de OLS (Ordinary Least Squares) nous aide à trouver la ligne la plus adaptée qui prédit le résultat en fonction des données dont nous disposons.



# IMPUTATION DES VALEURS MANQUANTES

## La méthode manuelle

*Étape 1* : Séparer les valeurs nulles du DataFrame et les considérer comme test\_df

*Étape 2* : Retirer les valeurs nulles du DataFrame et les considérer comme train\_data

*Étape 3* : Créer X\_train et y\_train à partir du train\_data

*Étape 4* : Construire le modèle de Régression Linéaire

*Étape 5* : Créer le X\_test à partir test\_df

*Étape 6* : Appliquer le modèle sur X\_test du test\_df et formuler des prédictions

*Étape 7* : Remplacer les valeurs manquantes par des valeurs prédites

[ ]:



# DIFFERENCE ENTRE LES VRAIS ET FAUX BILLETS

	diagonale	height_left	height_right	margin_low	margin_up	length
<b>Les différences entre vrais &amp; faux billets</b>	0.08271	-0.21727	-0.30512	-1.006586	-0.26659	1.42647

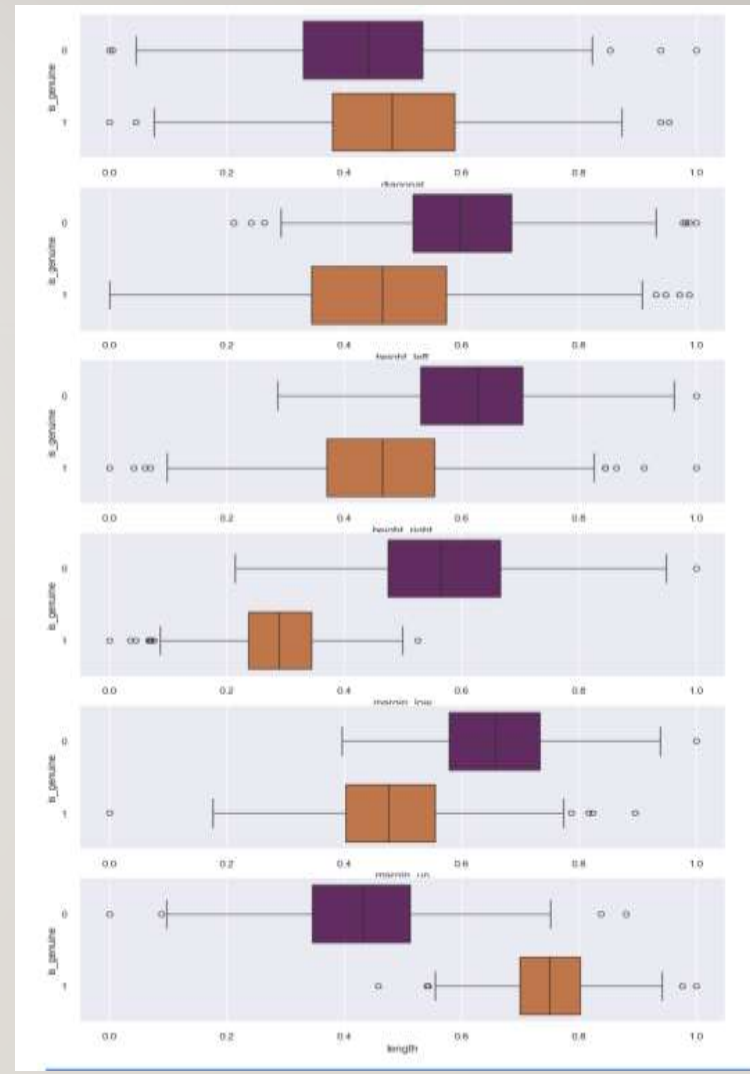
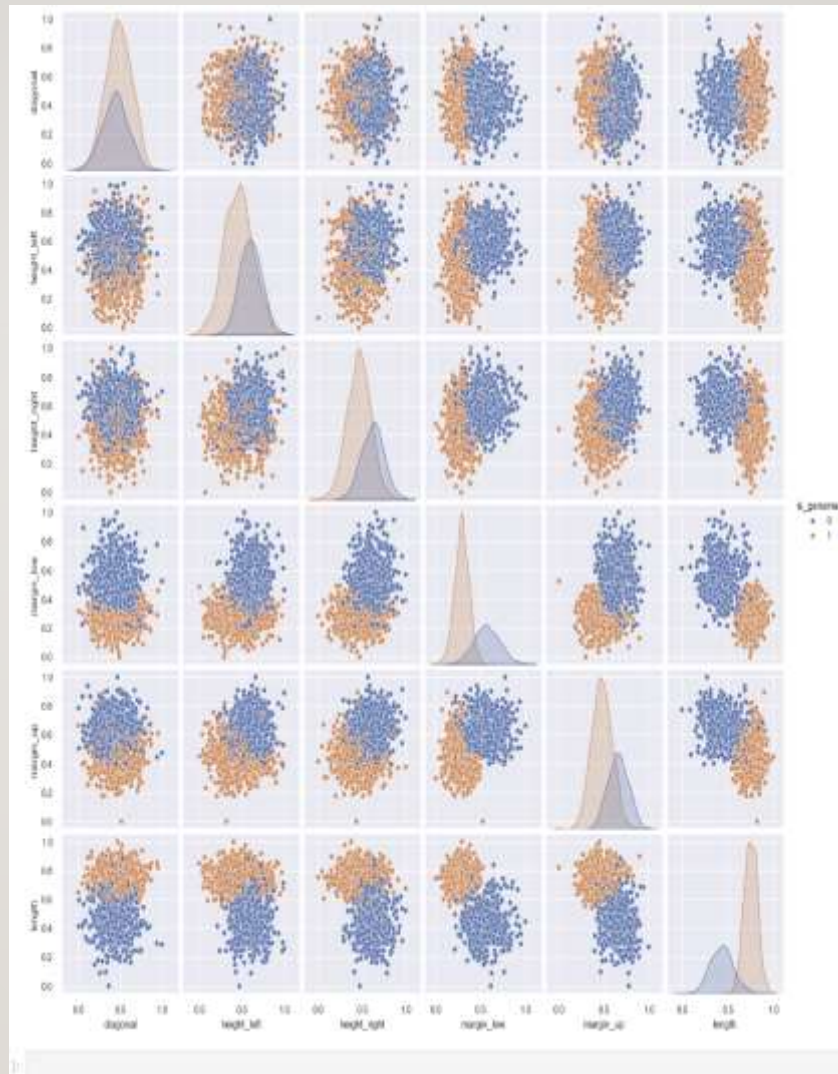
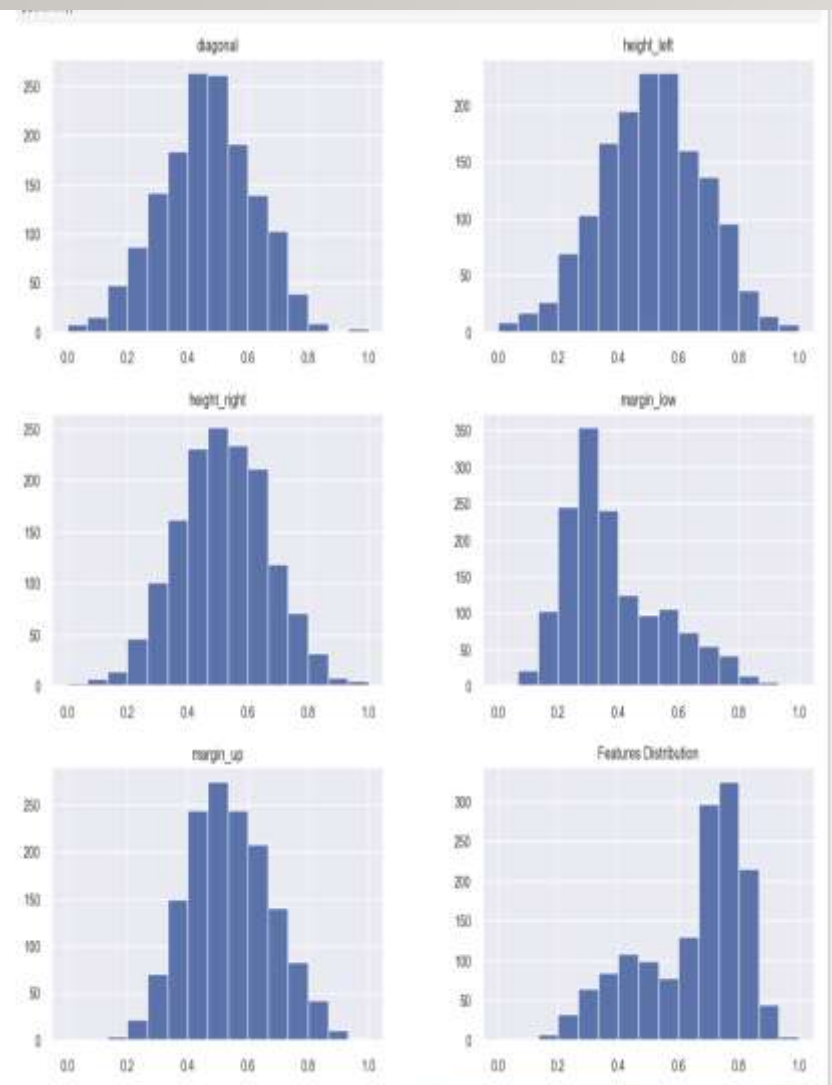
On peut en déduire que :

Les vrais billets ont une diagonale légèrement plus importante

La hauteur et la marge des faux billets sont légèrement plus importantes

Les vrais billets sont plus grands en largeur

# LES DISTRIBUTIONS

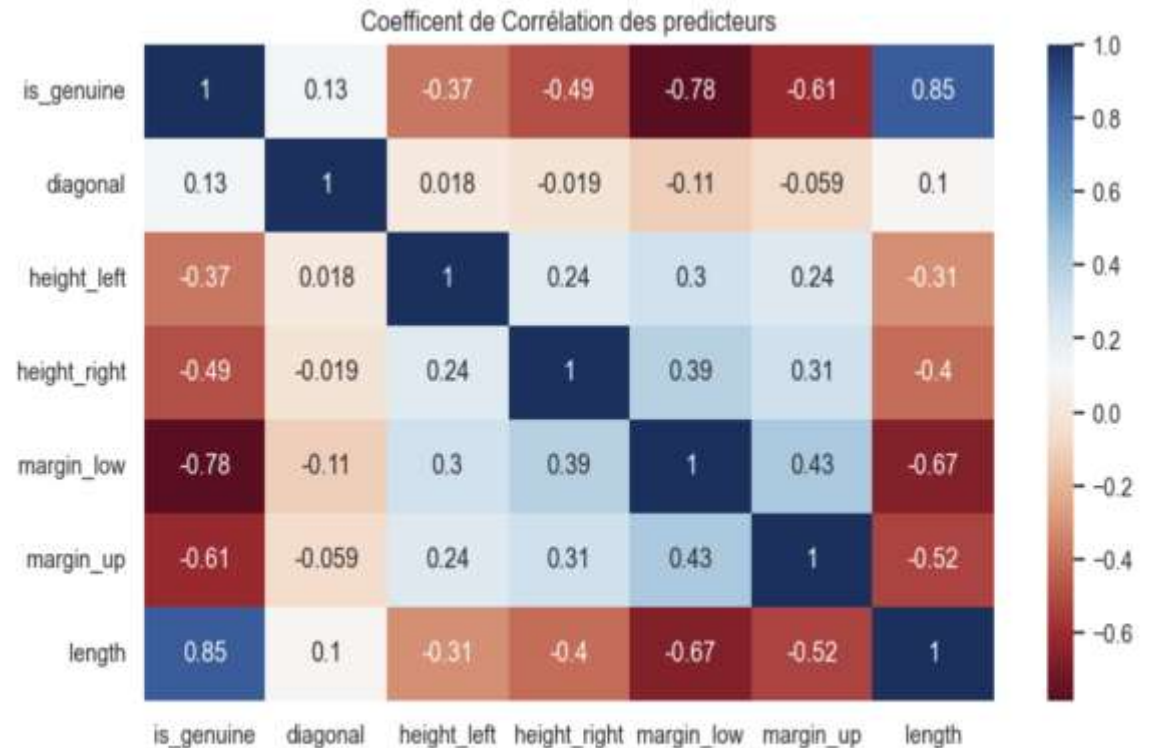


# LES HYPOTHÈSES DE VALIDITÉ À VÉRIFIER

## I. COLINÉARITE DES VARIABLES :

*signifie qu'une forte corrélation existe entre elles, ce qui rend difficile, voire impossible, l'estimation fiable de leurs coefficients de régression individuels*

Nous constatons que les variables *length* et *is\_genuine* ont une très forte corrélation.



```
86]: # On note que la variable "length" a une très forte corrélation avec la variable "is_genuine" avec 0.85
```



## 2. HOMOSCÉDASTICITÉ :

*fait référence à une condition dans laquelle la variance du résidu, ou terme d'erreur, dans un modèle de régression est constante.*

Avec une statistique de test de 29.98 et une valeur p de  $9.01 \times 10^{-24}$ , l'assumption n'est pas validée.

```
] import numpy as np
from scipy.stats import levene

# Récupération des résidus
residuals = ols_model.resid

# La valeur prédite
predicted = ols_model.fittedvalues

# Division des résidus en groupes en fonction de la variable prédite
groups = pd.qcut(predicted, q=5, labels=False)

# Le test de Levene sur les résidus groupés
levene_test = levene(*[residuals[groups == i] for i in np.unique(groups)], center='mean')

# La statistique de test et la valeur p
test_statistic = levene_test.statistic
p_value = levene_test.pvalue

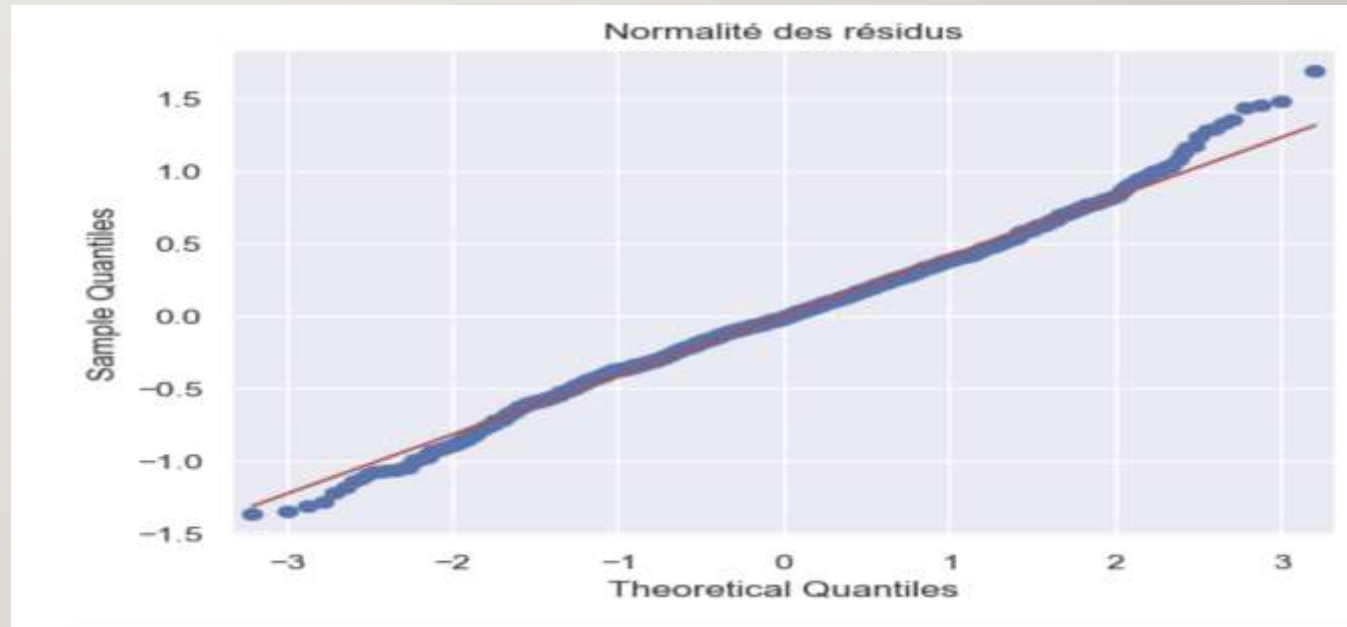
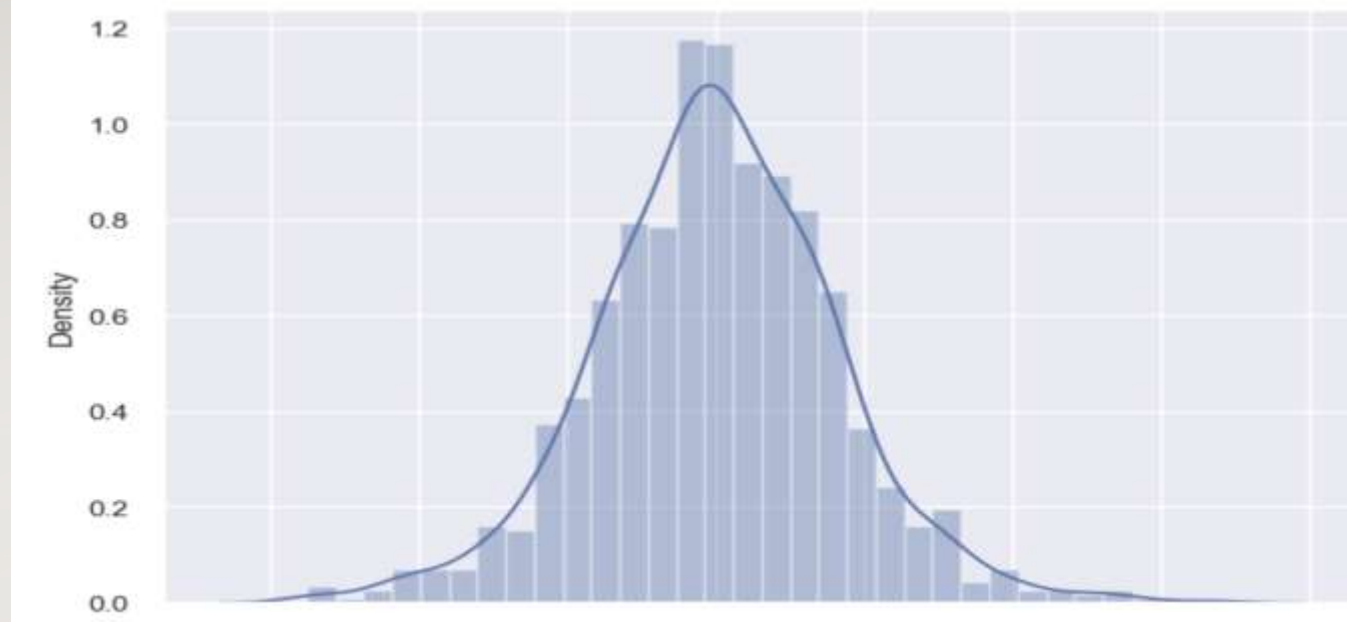
# On affiche les résultats
print("Test de Levene:")
print("Statistique de test:", test_statistic)
print("Valeur p:", p_value)

Test de Levene:
Statistique de test: 29.983642592506712
Valeur p: 9.014373222445792e-24
```

### 3. NORMALITÉ DES RÉSIDUS

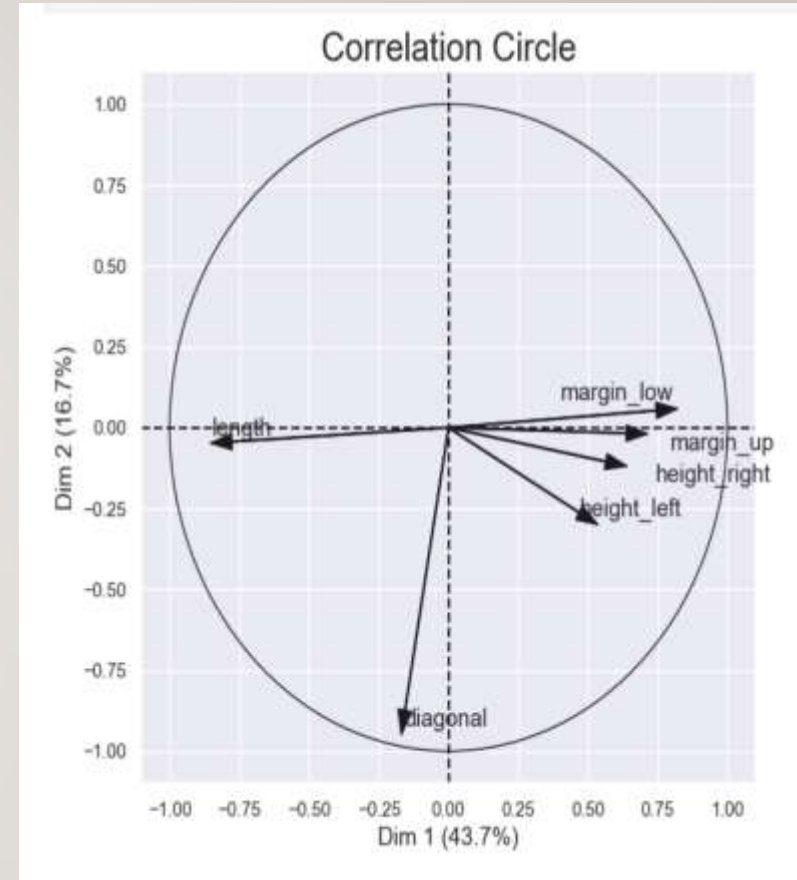
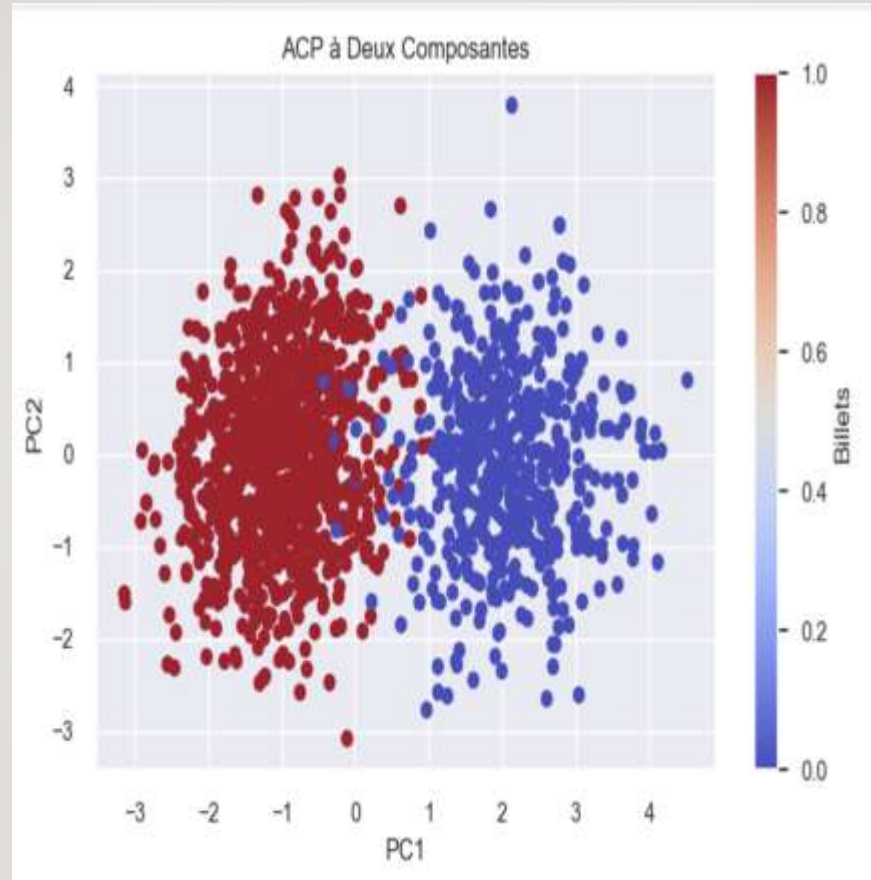
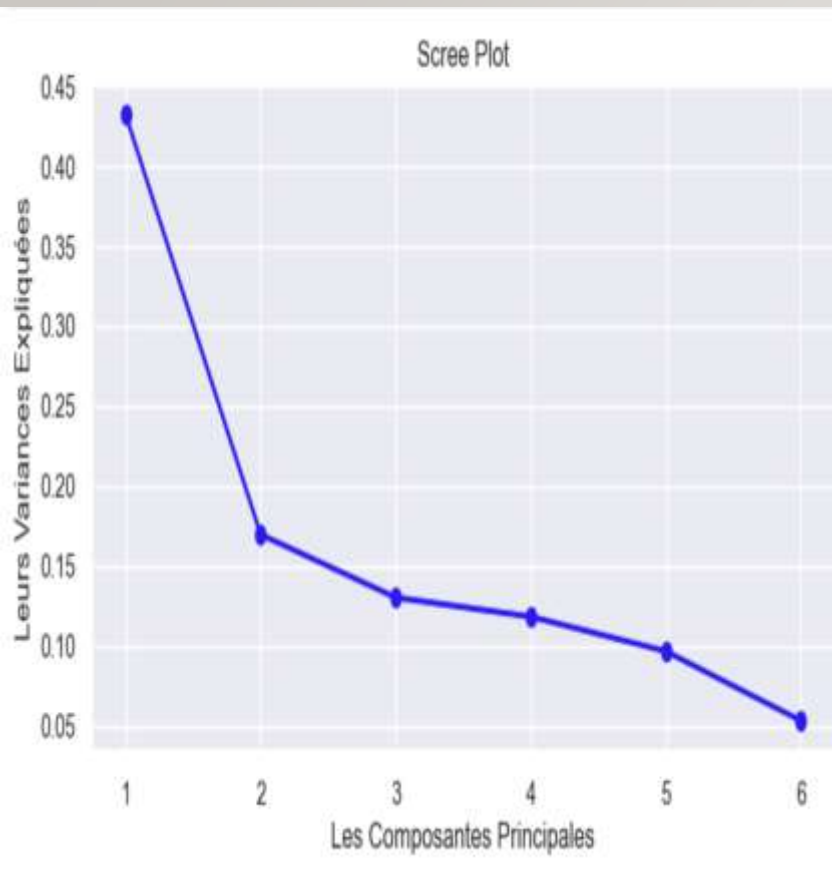
*La normalité est l'hypothèse selon laquelle les résidus sous-jacents sont normalement distribués, ou approximativement ainsi.*

Le distplot montre une distribution normale mais le test Shapiro ne valide pas cette assumption.





# ACP

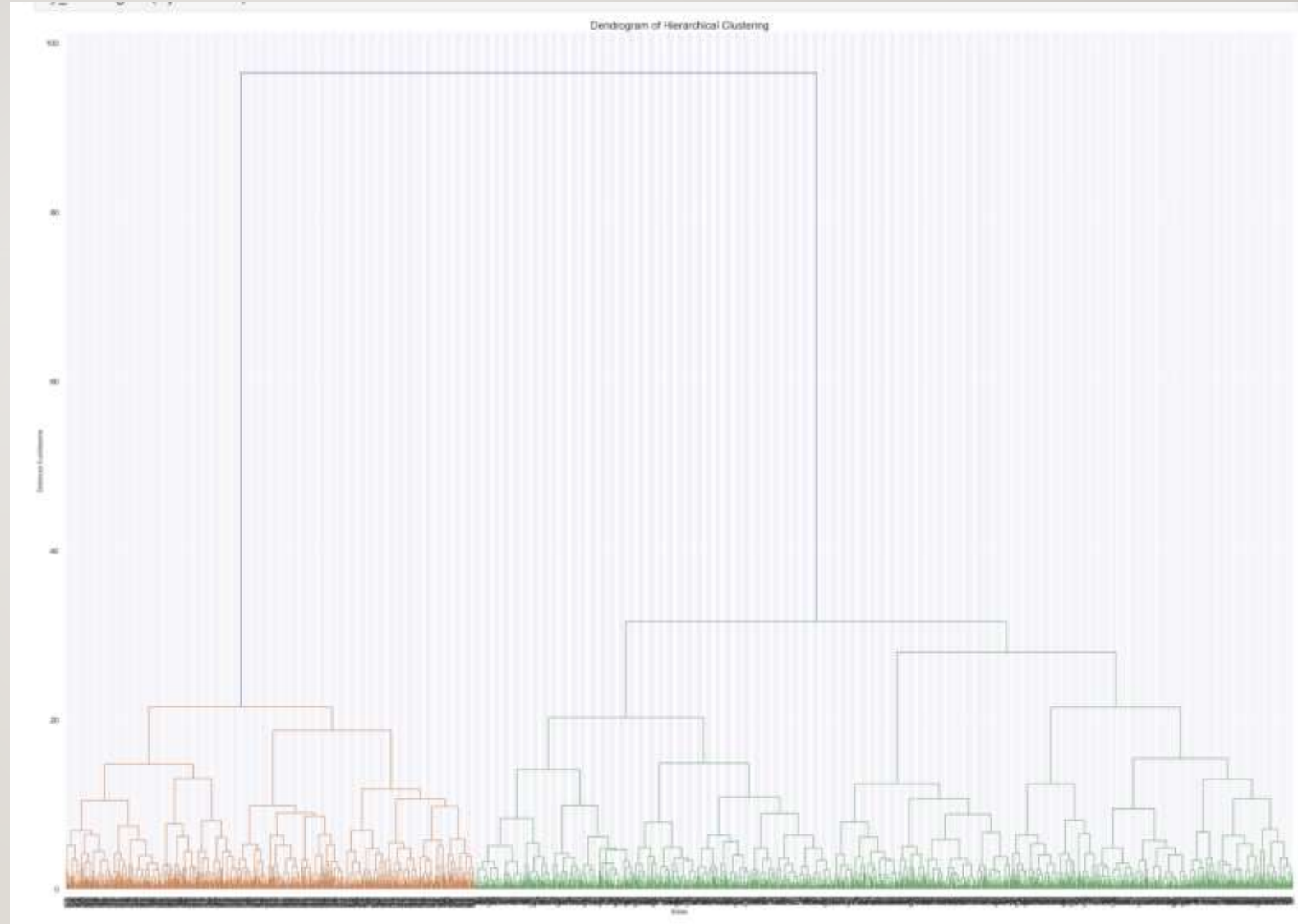


**Du scree Plot, on prend les 2 premières composantes principales (60% de la variance totale)**

**Sur la projection sur les individus PC1 et PC2, on voit clairement deux clusters qui expliquent les vrais et faux billets.**

# Classification Ascendante Hierarchique (CAH)

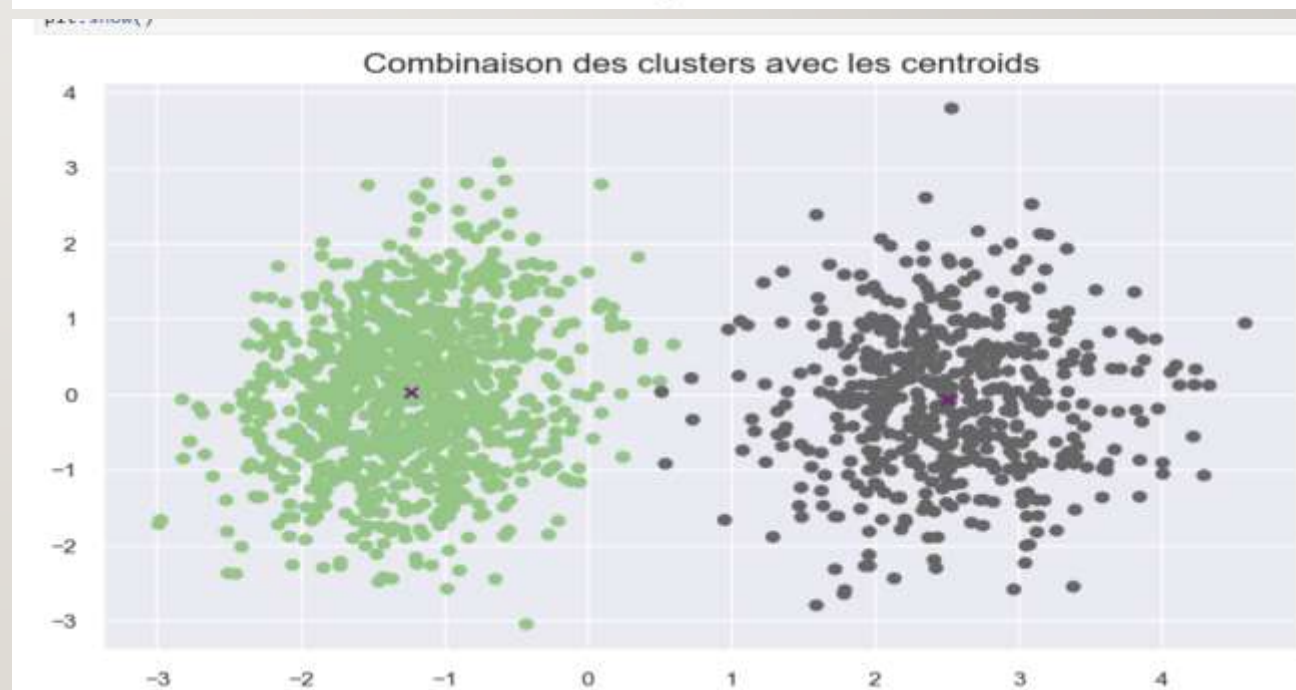
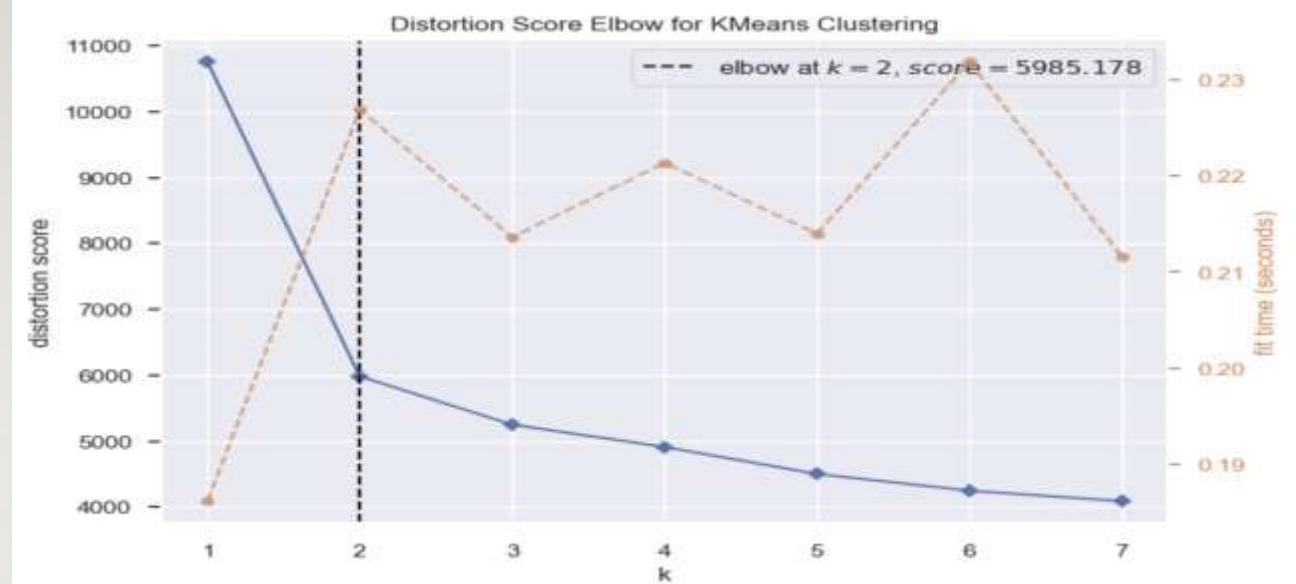
J'ai essayé une autre méthode de clustering, en l'occurrence CAH pour voir les clusters. On voit clairement 2 clusters dans ce dendrogramme



# K-Means

De la 2<sup>ème</sup> méthode de clustering K-Means, on en déduit que le nombre optimal de  $k = 2$  avec la méthode du Coude.

Après avoir entraîné le modèle sur les données, on voit clairement les deux clusters avec les centroides qui distinguent les vrais et faux billets.



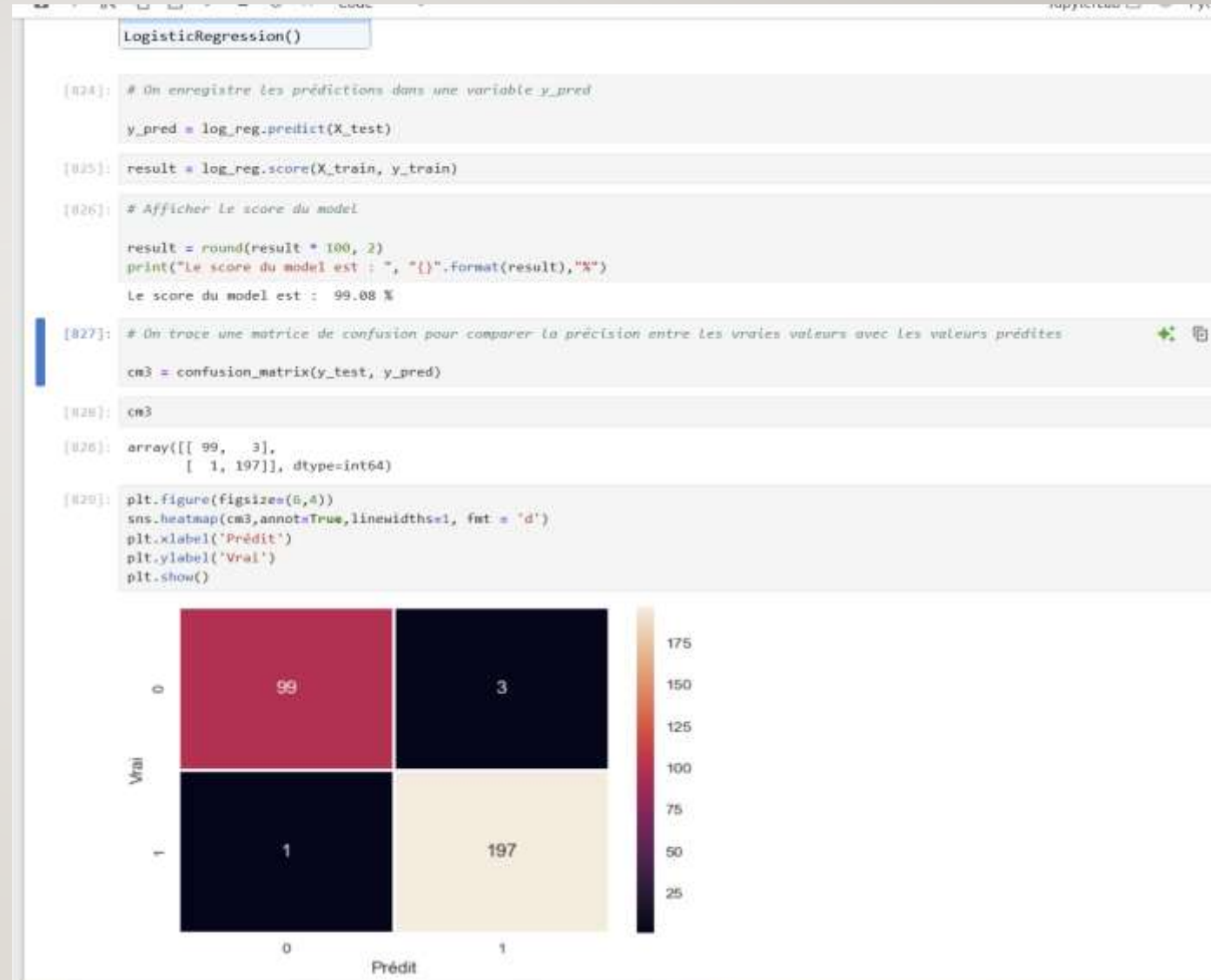


# PRÉDICTION DES BILLETS : 4 ALGORITHMES À TESTER

## I. LA REGRESSION LOGISTIQUE

*La régression logistique est un algorithme d'apprentissage automatique supervisé qui accomplit des tâches de classification binaire en prédisant la probabilité d'un résultat, d'un événement ou d'une observation.*

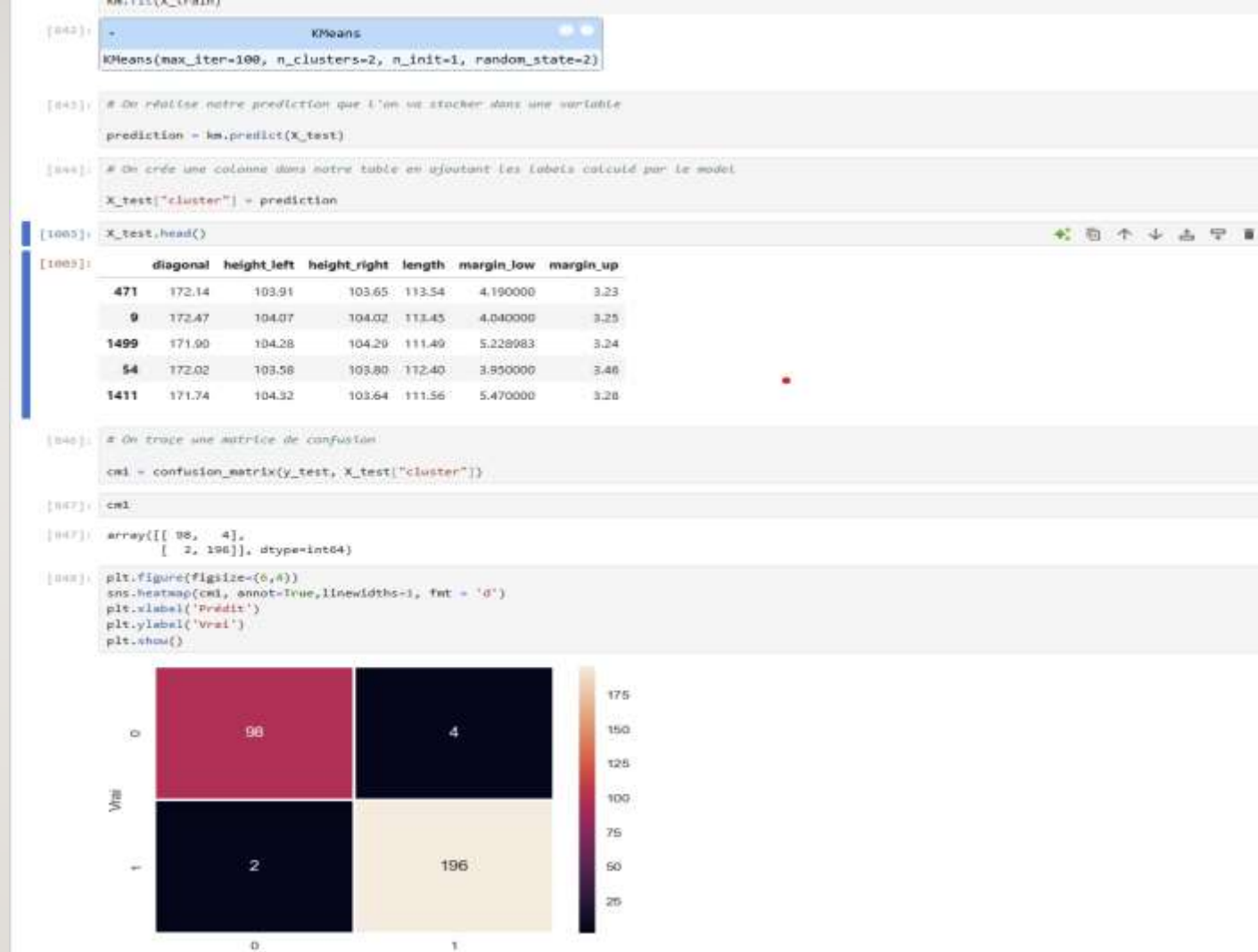
- Vérifier les 6 assomptions avant d'entraîner le modèle
- Utilisation de DecisionTreeClassifier pour sélectionner les features les plus pertinents
- L'évaluation de la performance donne un F1 score de 98.99%



## 2. K-MEANS

*Bien que K-means lui-même ne prédise pas directement, il peut s'agir d'une étape puissante de prétraitement ou d'ingénierie des fonctionnalités qui améliore les modèles prédictifs en : Identifiant des clusters ou des segments de clientèle significatifs.*

- On choisit le nombre optimal de K clusters (2) avec la méthode Elbow
- On crée une colonne dans notre table en ajoutant les labels calculés par le modèle
- Le f1 score donne 98.49%

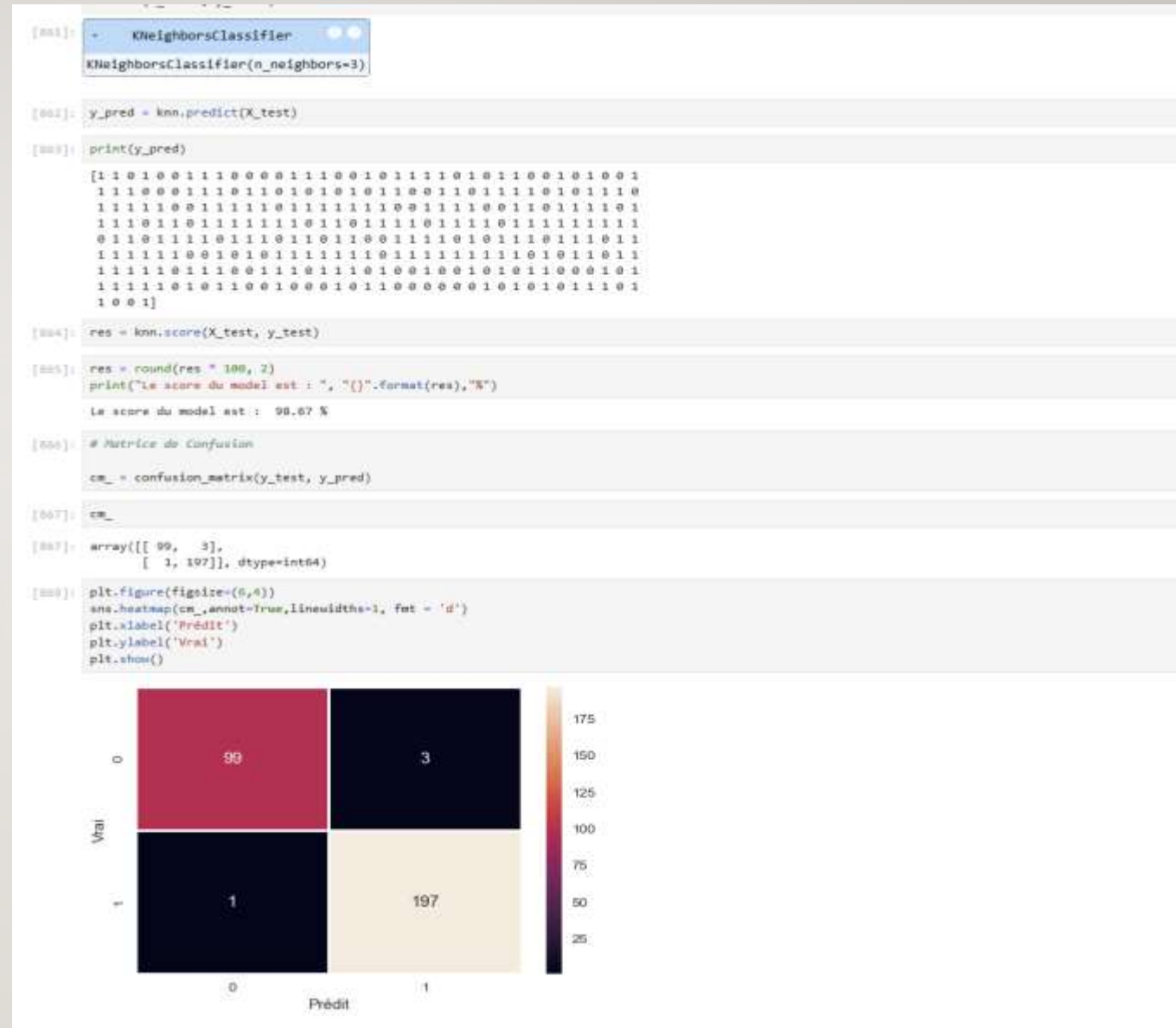




### 3. K-Nearest Kneighbors Classifier (KNN)

*L'algorithme des k-plus proches voisins (KNN) est un classificateur d'apprentissage supervisé non paramétrique, qui utilise la proximité pour effectuer des classifications ou des prédictions sur le regroupement d'un point de données individuel.*

- On partitionne notre dataset en X-train, X-test
- On entraîne notre modèle
- On prédit X\_test
- Notre modèle nous donne un score de 98.99%



## 4. RANDOM FOREST

*Un algorithme qui génère une 'forêt d'arbres'. Il les combine ensuite pour éviter le surapprentissage et produire des prédictions plus précises.*

- On partitionne notre dataset en X-train, X-test
- On entraîne notre modèle
- On prédit X\_test
- Notre modèle nous donne un score de 99.24%

```
[879]: rf.fit(X_train.values, y_train)
[879]: RandomForestClassifier
RandomForestClassifier()

[880]: y_pred = rf.predict(X_test)

Evaluer la performance du modèle

[881]: new_res = rf.score(X_test, y_test)

[882]: new_res = round(new_res * 100, 2)
print("Le score du modèle est : ", "{}".format(new_res), "%")
Le score du modèle est : 99.24 %

[883]: from sklearn.metrics import classification_report

[884]: print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.99	0.98	0.99	102
1	0.99	0.99	0.99	198
accuracy			0.99	300
macro avg	0.99	0.99	0.99	300
weighted avg	0.99	0.99	0.99	300

```
[885]: # On visualise l'importance de chaque features par dans le modèle
features = pd.DataFrame(rf.feature_importances_, index = X.columns)

[886]: features
```

	0
diagonal	0.007617
height_left	0.027538
height_right	0.039911
length	0.443080
margin_low	0.362590
margin_up	0.119284

```
[887]: cm2 = confusion_matrix(y_test, y_pred)
[888]: cm2
```

array([[100, 2],
[ 1, 197]], dtype=int64)

```
[889]: plt.figure(figsize=(6,4))
sns.heatmap(cm2, annot=True, linewidths=1, fmt='d')
plt.xlabel('Prédit')
plt.ylabel('Vrai')
plt.show()
```

	0	1
0	100	2
1	1	197

```
[890]: # Sauvegarder le modèle pour une utilisation future
```