

```

1  """
2  Script : choix_seuil_outlier.py
3  Auteur : CISSE Ibrahim
4  Objectif : Tester plusieurs méthodes de détection d'outliers sur une variable
numérique et choisir le seuil le plus pertinent.
5
6  Méthodologie :
7  1. Chargement d'un jeu de données simulé
8  2. Application de trois méthodes de détection :
9      - Écart-type ( $\pm 2\sigma$ )
10     - IQR (boîte à moustaches)
11     - Quantiles extrêmes (1% et 99%)
12  3. Comparaison du nombre d'outliers détectés
13  4. Visualisation des seuils et des points extrêmes
14  5. Choix du seuil optimal selon le contexte métier
15  """
16
17  import numpy as np
18  import pandas as pd
19  import matplotlib.pyplot as plt
20
21  # 1. Génération de données simulées
22  np.random.seed(42)
23  data = np.random.normal(loc=100, scale=15, size=1000)
24  data = np.append(data, [30, 200, 250]) # Ajout d'outliers
25  df = pd.DataFrame({"valeurs": data})
26
27  # 2. Méthode écart-type
28  mean = df["valeurs"].mean()
29  std = df["valeurs"].std()
30  outliers_std = df[(df["valeurs"] < mean - 2*std) | (df["valeurs"] > mean + 2*std)]
31
32  # 3. Méthode IQR
33  q1 = df["valeurs"].quantile(0.25)
34  q3 = df["valeurs"].quantile(0.75)
35  iqr = q3 - q1
36  outliers_iqr = df[(df["valeurs"] < q1 - 1.5*iqr) | (df["valeurs"] > q3 + 1.5*iqr)]
37
38  # 4. Méthode quantiles extrêmes
39  low = df["valeurs"].quantile(0.01)
40  high = df["valeurs"].quantile(0.99)
41  outliers_quant = df[(df["valeurs"] < low) | (df["valeurs"] > high)]
42
43  # 5. Comparaison
44  print("Méthode écart-type : ", len(outliers_std), " outliers")
45  print("Méthode IQR : ", len(outliers_iqr), " outliers")
46  print("Méthode quantiles : ", len(outliers_quant), " outliers")
47
48  # 6. Visualisation
49  plt.figure(figsize=(10, 6))
50  plt.boxplot(df["valeurs"], vert=False)
51  plt.title("Détection des outliers par boîte à moustaches")
52  plt.xlabel("Valeurs")
53  plt.grid(True)
54  plt.show()
55

```