

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.preprocessing import LabelEncoder
4 from sklearn.experimental import enable_iterative_imputer
5 from sklearn.impute import IterativeImputer
6 from sklearn.preprocessing import StandardScaler
7
8 class Data_Cleaning:
9     def __init__(self,csv_file, sep = ','):
10         self.csv_file = csv_file
11         self.df = None
12         self.sep = sep
13         self.cols = None
14         self.x_num = None
15         self.x_str = None
16         self.x = None
17         self.y = None
18         self.target = None
19         self.reply = None
20
21     # Est ce que le df contient de header ou non
22     def isHeader(self):
23         if self.reply == 'yes':
24             self.df = pd.read_csv(self.csv_file, sep = self.sep)
25         else:
26             self.df = pd.read_csv(self.csv_file, header = None)
27             self.cols = [c.strip() for c in self.cols.split(',')]
28             self.df.columns = self.cols
29
30     # On fait la separation les colonnes numeriques et les colonnes catgoriales
31     def separation_xnum_xstr(self):
32         self.x_num = self.df.select_dtypes(include = [np.number])
33         self.x_str = self.df.select_dtypes(exclude = [np.number])
34
35         for i in self.x_str.columns:
36             self.x_str[i] = self.x_str[i].fillna(self.x_str[i].mode()[0])
37
38     # On fait l'encodeage des colonnes categoriale par LabelEncoder
39     def encodage(self):
40         encoder = LabelEncoder()
41
42         for i in self.x_str.columns:
43             self.x_str[i] = encoder.fit_transform(self.x_str[i])
44         # On fait la concatination des deux colonnes encodées (initiallement categoriale) avec les colonnes numerique
45         self.x = pd.concat([self.x_num, self.x_str], axis = 1)
46
47     # Choisir la methode de remplissage des valeurs manquantes
48     def val_mang(self):
49         if self.reply == 'intelligente':
50             impute = IterativeImputer(random_state = 42)
51             self.x = pd.DataFrame(impute.fit_transform(self.x), columns = self.x.columns)
52         elif self.reply == 'moyenne':
53             self.x = self.x.fillna(self.x.mean())
54         else:
55             self.x = self.x.fillna(self.x.median())
56
57     # La suppression des observation doublons
58     def duplication(self):
59         self.x = self.x.drop_duplicates()
60
61     # Pour traiter les valeurs aburantes, j'ai utilisé le "capping"
62     def outlier(self,col):
63         Q1 = col.quantile(0.25)
64         Q3 = col.quantile(0.75)
65
66         IQR = Q3 - Q1
67
68         Fb = Q1 - 1.5 * IQR
69         Fh = Q3 + 1.5 * IQR
70
71         return col.clip(lower = Fb, upper = Fh)
72
73     def remp_outlier(self):
74         for col in self.x.columns :
75             self.x[col] = self.outlier(self.x[col])
76

```

```
77 # Séparation des caractéristiques (x) et la variable cible (y)
78 def separation_x_y(self):
79     self.y = self.x[self.target]
80     self.x = self.x.drop(columns = self.target)
81
82 # Est ce que l'utilisateur souhaite de faire la standardisation ou non
83 def standarisation(self):
84     scale = StandardScaler()
85
86     if self.reply == 'yes':
87         self.x = pd.DataFrame(scale.fit_transform(self.x),columns = self.x.columns)
88
89 # Retourner la nouvelle dataframe nettoyée
90 def df_final(self):
91     self.df = pd.concat([self.x,self.y], axis = 1)
92
93     return self.df
```