

Projet 3 : Préparez les données

pour un organisme de Santé Publique



Introduction

OPENCLASSROOMS

Gael Delescluse

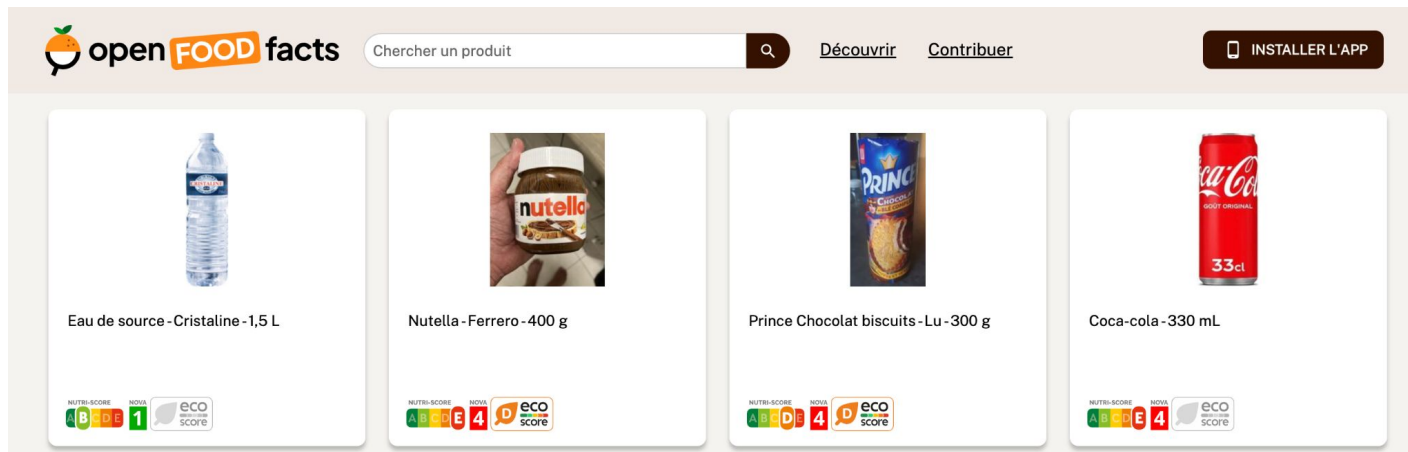
0

Table des matières

Présentation
Analyse Dataset
Données aberrantes
Données manquantes
Analyse univariées
Analyse bivariées
Analyse multivariées

Introduction

Présentation



NUTRI-SCORE



+ 800 millions d'adultes

+ 2,8 millions décès

Dataset

Analyse métier

```
product_name  
nutrition_grade_fr  
energy_100g  
proteins_100g  
sugars_100g  
fat_100g  
saturated-fat_100g  
salt_100g  
sodium_100g  
fiber_100g  
fruits-vegetables-nuts_100g  
dtype: object
```

320.772 individus

162 variables

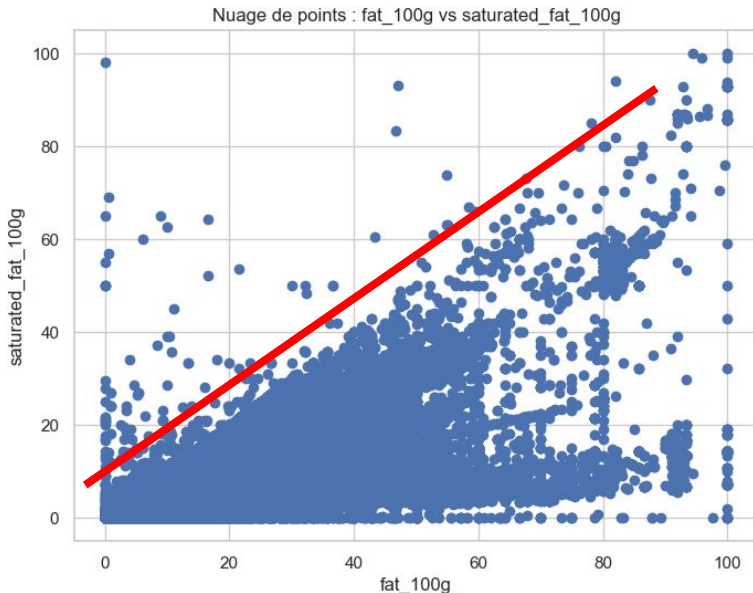
CSV

données manquantes

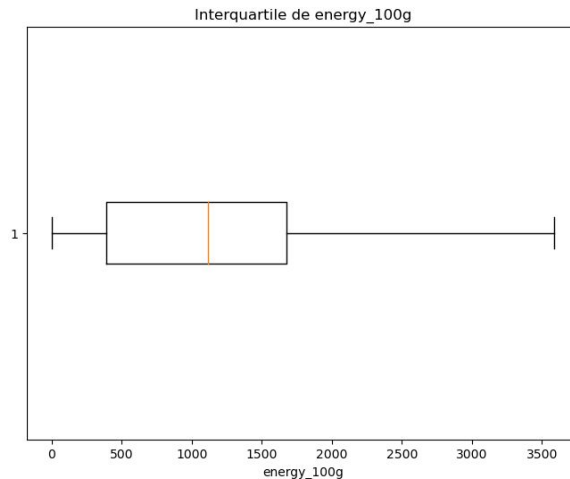
Séparateur \t

Nettoyage du Dataset

Valeurs aberrantes

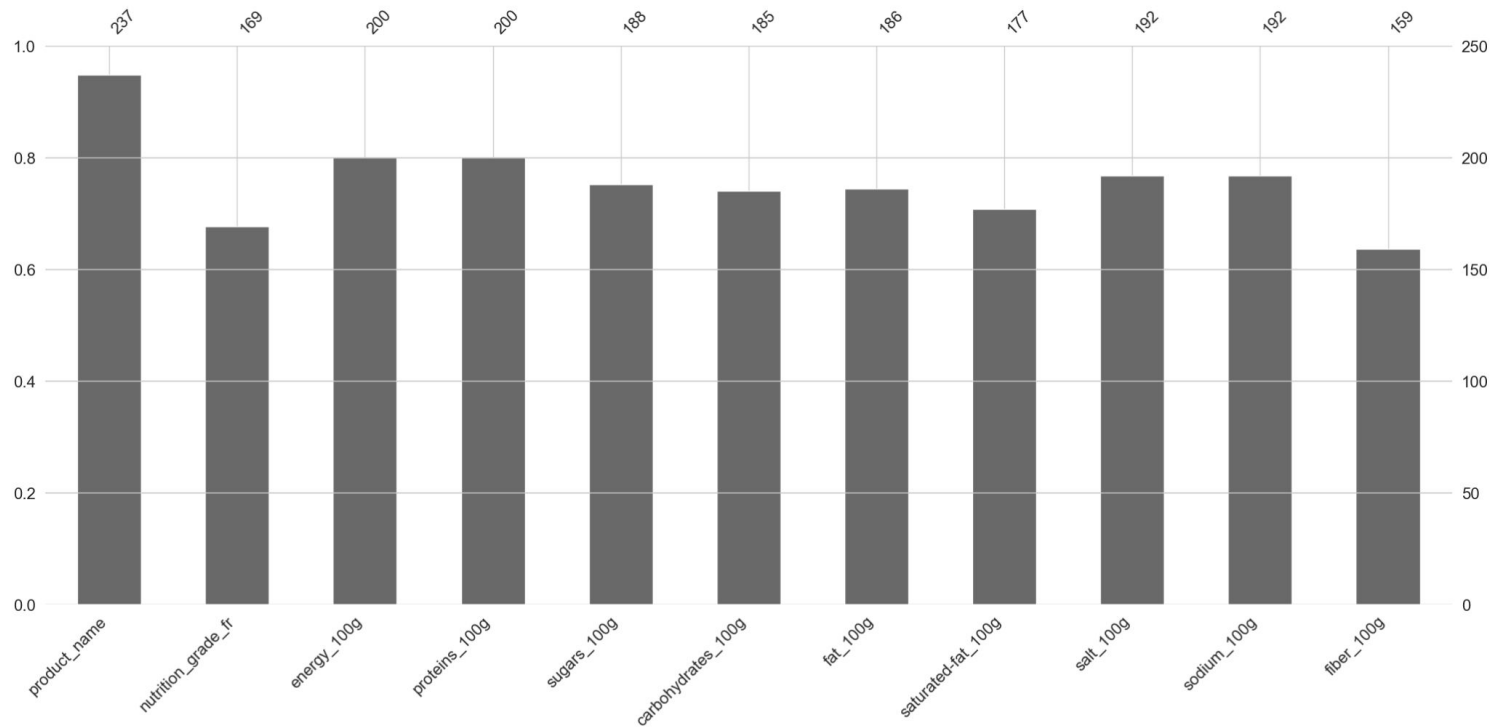


```
index_drop_rogue_value = df[df['proteins_100g'] > 100].index
df.drop(index_drop_rogue_value, inplace=True)
index_drop_rogue_value = df[df['sugars_100g'] > 100].index
df.drop(index_drop_rogue_value, inplace=True)
index_drop_rogue_value = df[df['fat_100g'] > 100].index
df.drop(index_drop_rogue_value, inplace=True)
index_drop_rogue_value = df[df['saturated_fat_100g'] > 100].index
df.drop(index_drop_rogue_value, inplace=True)
index_drop_rogue_value = df[df['salt_100g'] > 100].index
df.drop(index_drop_rogue_value, inplace=True)
index_drop_rogue_value = df[df['fiber_100g'] > 100].index
df.drop(index_drop_rogue_value, inplace=True)
index_drop_rogue_value = df[df['proteins_100g'] < 0].index
df.drop(index_drop_rogue_value, inplace=True)
index_drop_rogue_value = df[df['sugars_100g'] < 0].index
df.drop(index_drop_rogue_value, inplace=True)
index_drop_rogue_value = df[df['fiber_100g'] < 0].index
df.drop(index_drop_rogue_value, inplace=True)
```



Huile d'avocat 3765 kJ

Données manquantes



Corrélation Pearson

fat_100g

salt_100g

sodium_100g

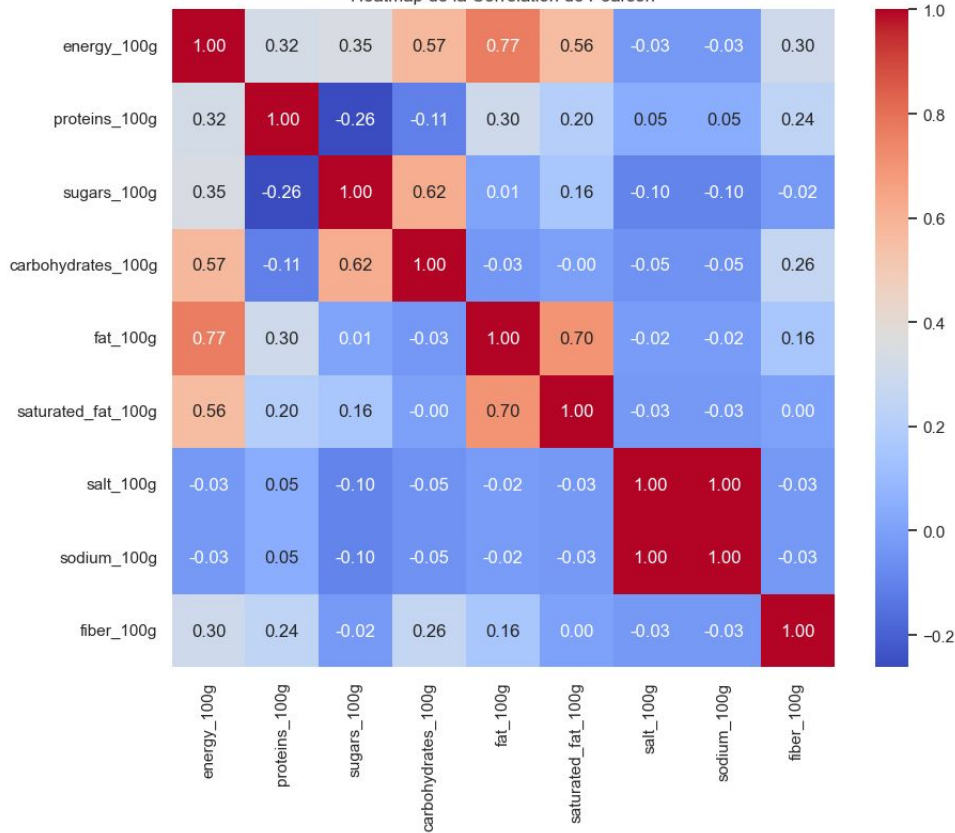
energy_100g

saturated_fat_100g

sugars_100g

carbohydrates_100g

Heatmap de la Corrélation de Pearson



IterativeImputer

proteins_100g -> 1%

salt_100g -> 3%

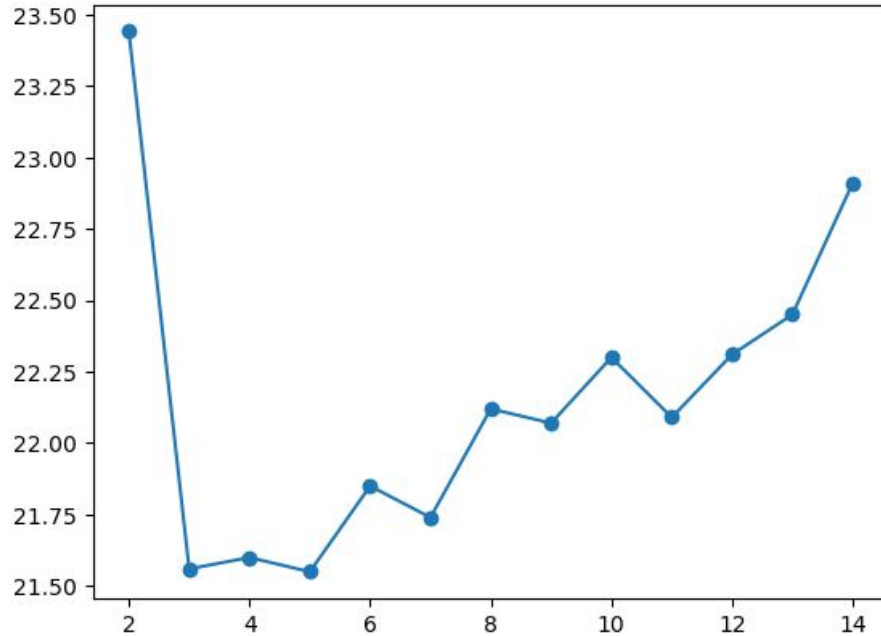
sodium_100g -> 3%

```
product_name      0.00
nutrition_grade_fr 13.04
energy_100g       0.00
proteins_100g     0.94
sugars_100g       0.00
carbohydrates_100g 0.00
fat_100g          0.00
saturated_fat_100g 0.00
salt_100g         2.31
sodium_100g       2.33
fiber_100g        21.21
dtype: float64
```

```
imputer = IterativeImputer()
df['fat_100g'] = imputer.fit_transform(df[['fat_100g', 'saturated_fat_100g', \
                                         'energy_100g']])[:, 0]
df['saturated_fat_100g'] = imputer.fit_transform(df[['saturated_fat_100g', 'fat_100g', \
                                         'energy_100g']])[:, 0]
df['energy_100g'] = imputer.fit_transform(df[['energy_100g', 'fat_100g', 'saturated_fat_100g', \
                                         'carbohydrates_100g']])[:, 0]
df['carbohydrates_100g'] = imputer.fit_transform(df[['carbohydrates_100g', 'sugars_100g', \
                                         'energy_100g']])[:, 0]
df['sugars_100g'] = imputer.fit_transform(df[['sugars_100g', 'carbohydrates_100g']])[:, 0]
```

K nearest neighbor

```
df_sample = df_without_Nan.sample(50000, replace=False)
```

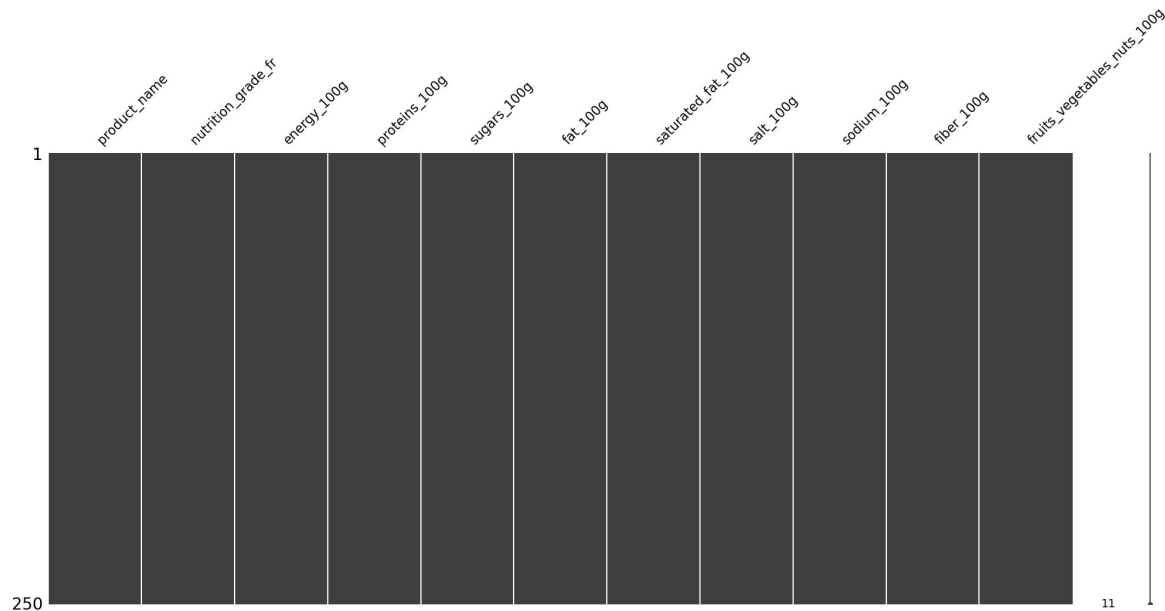


```
n_neighbors=3
```

K nearest neighbor

Accuracy

84%

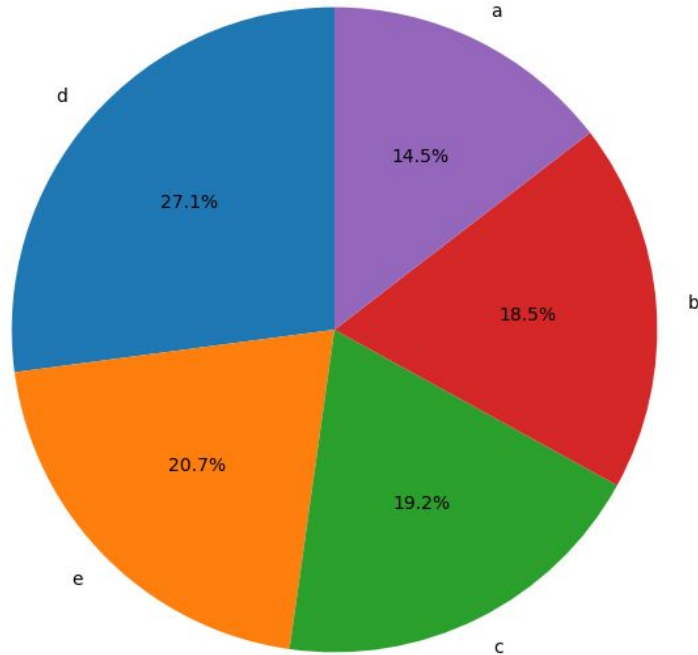


29%

suppression
d'individus

Analyse exploratoire du Dataset

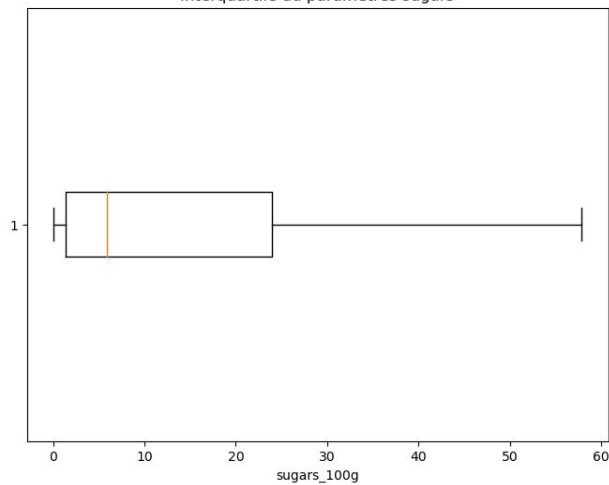
Répartition du nutriscore



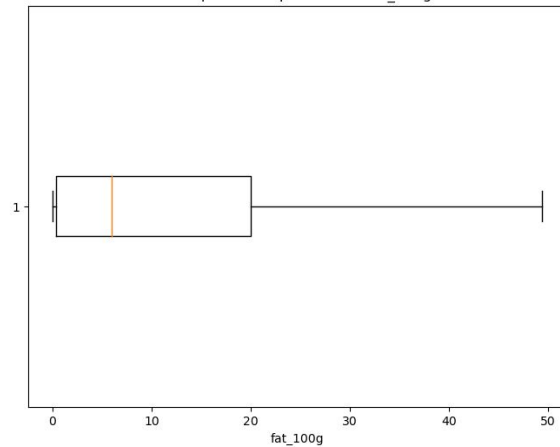
Analyse univariée

Analyse univariée

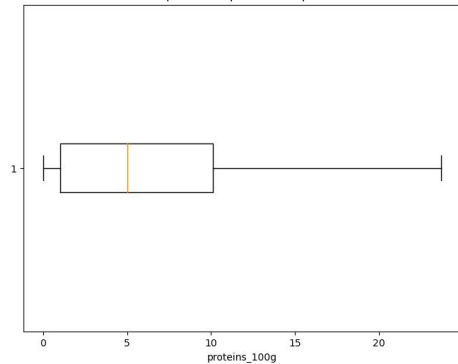
Interquartile du paramètres sugars



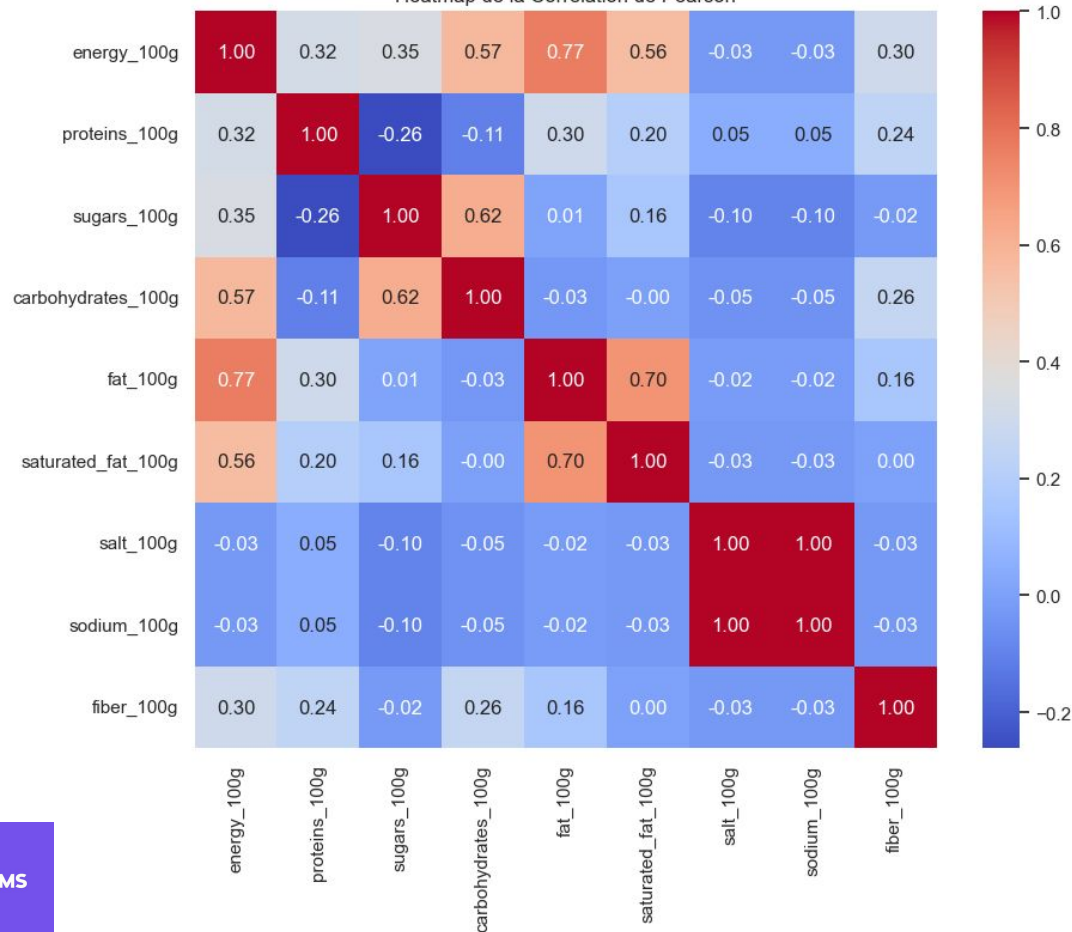
Interquartile du paramètres fat_100g



Interquartile du paramètres proteins

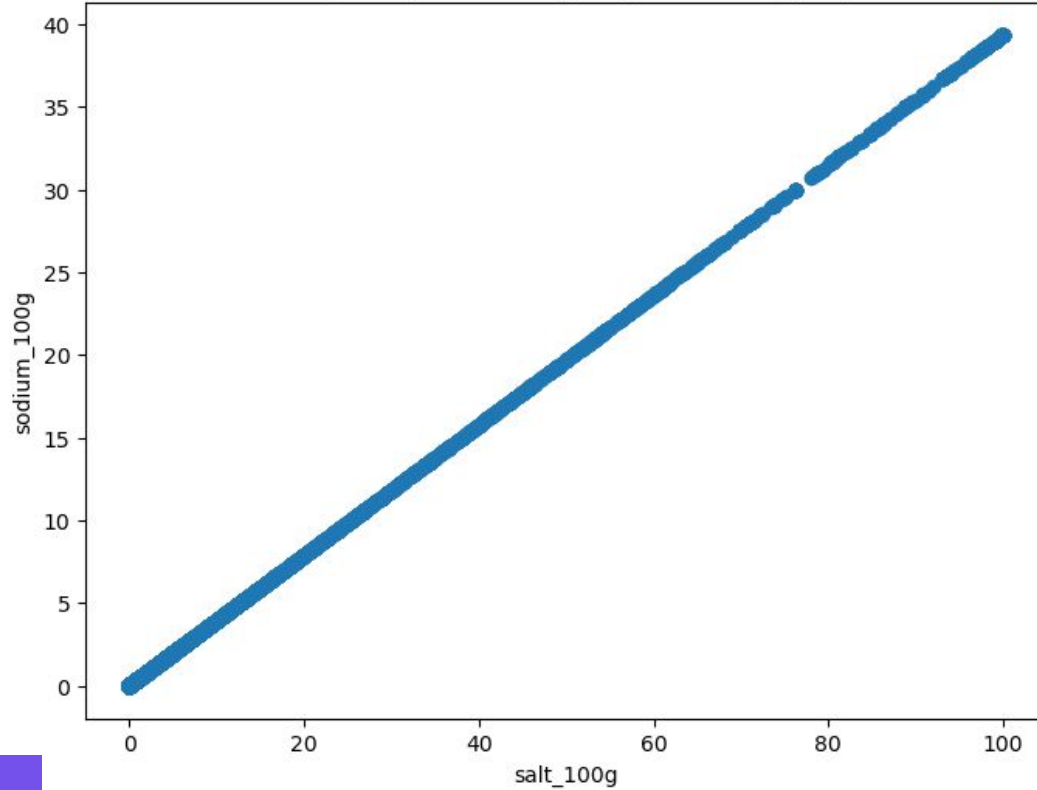


Heatmap de la Corrélation de Pearson

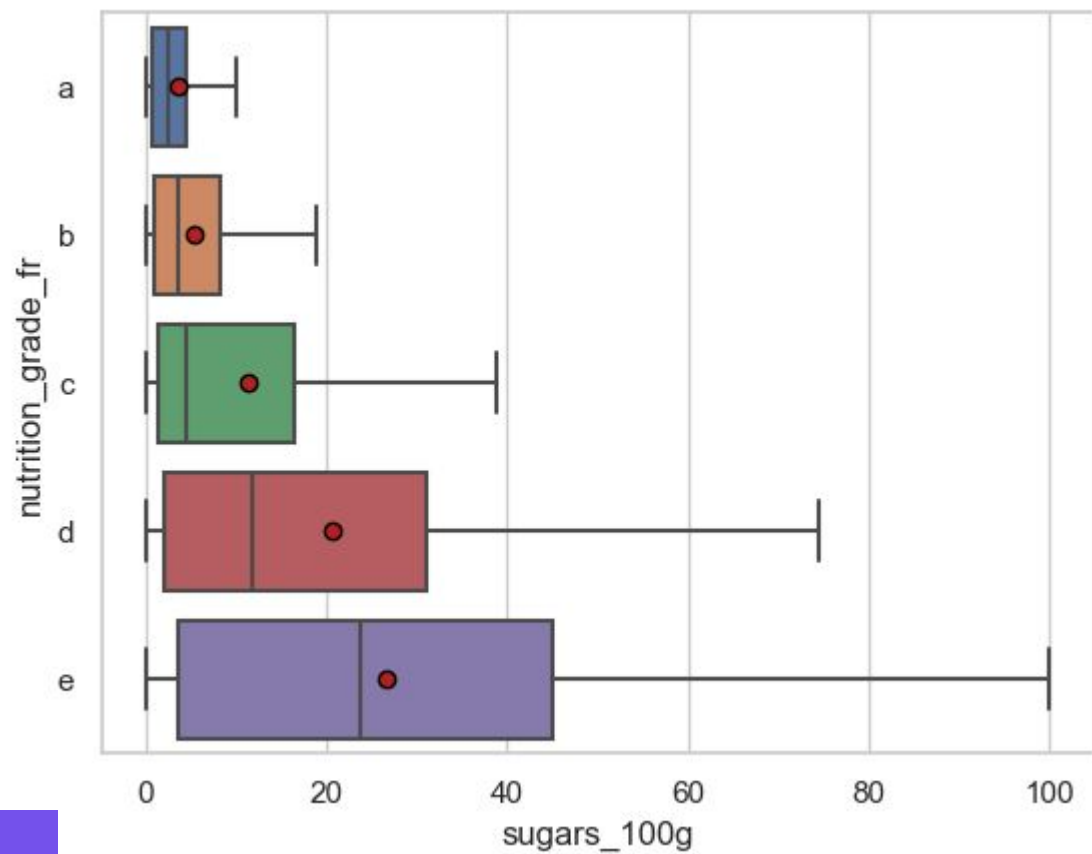


Analyse bvariée

Nuage de points : salt_100g vs sodium_100g

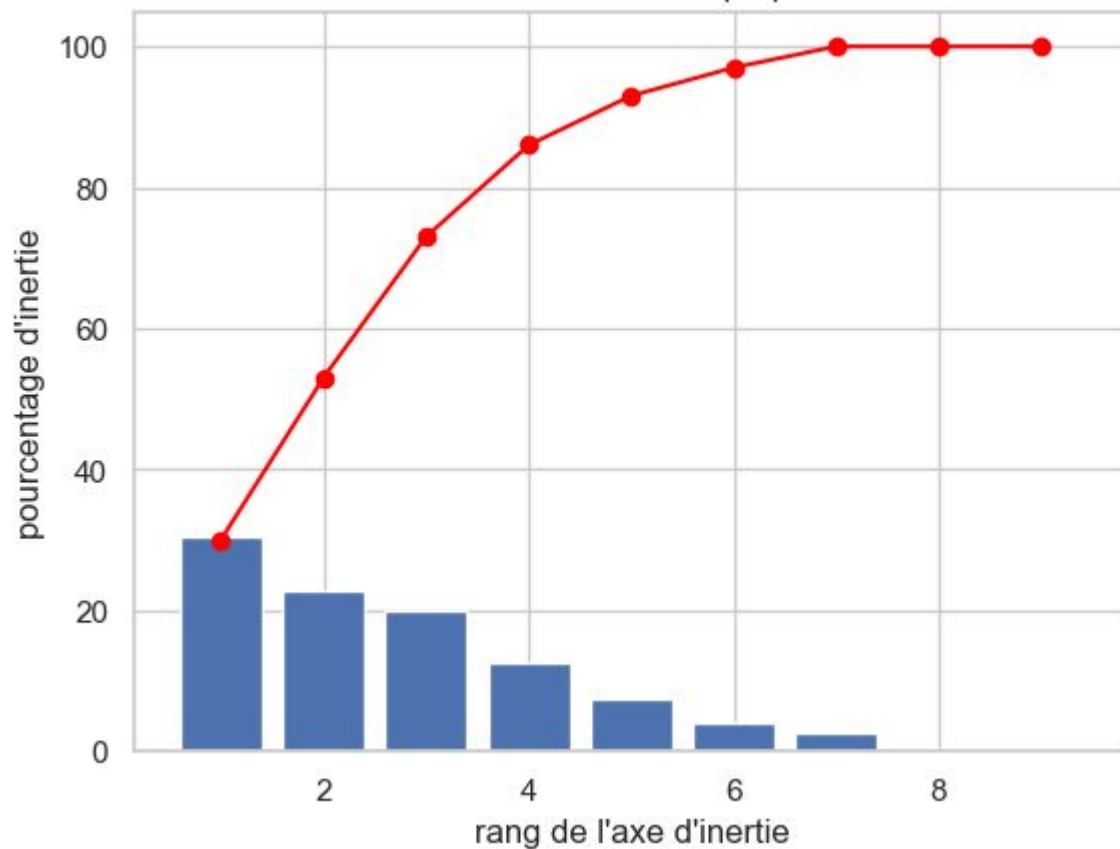


Analyse bivariable

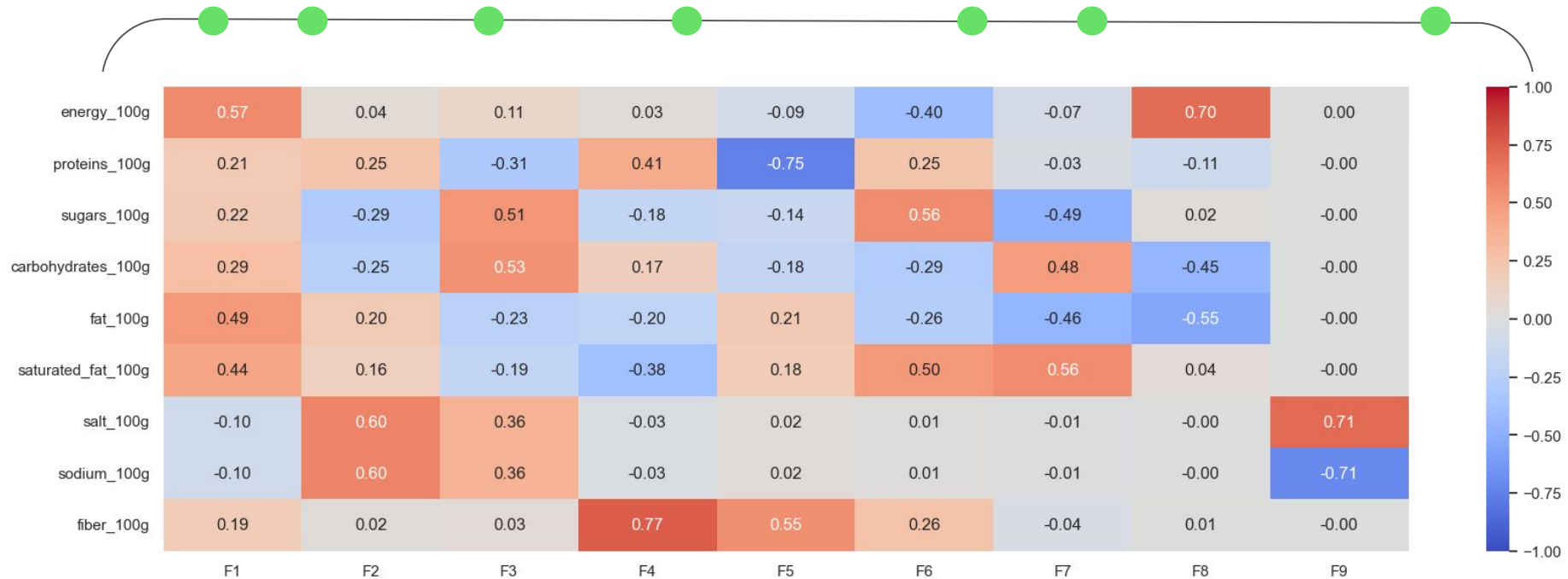


Analyse bivariable

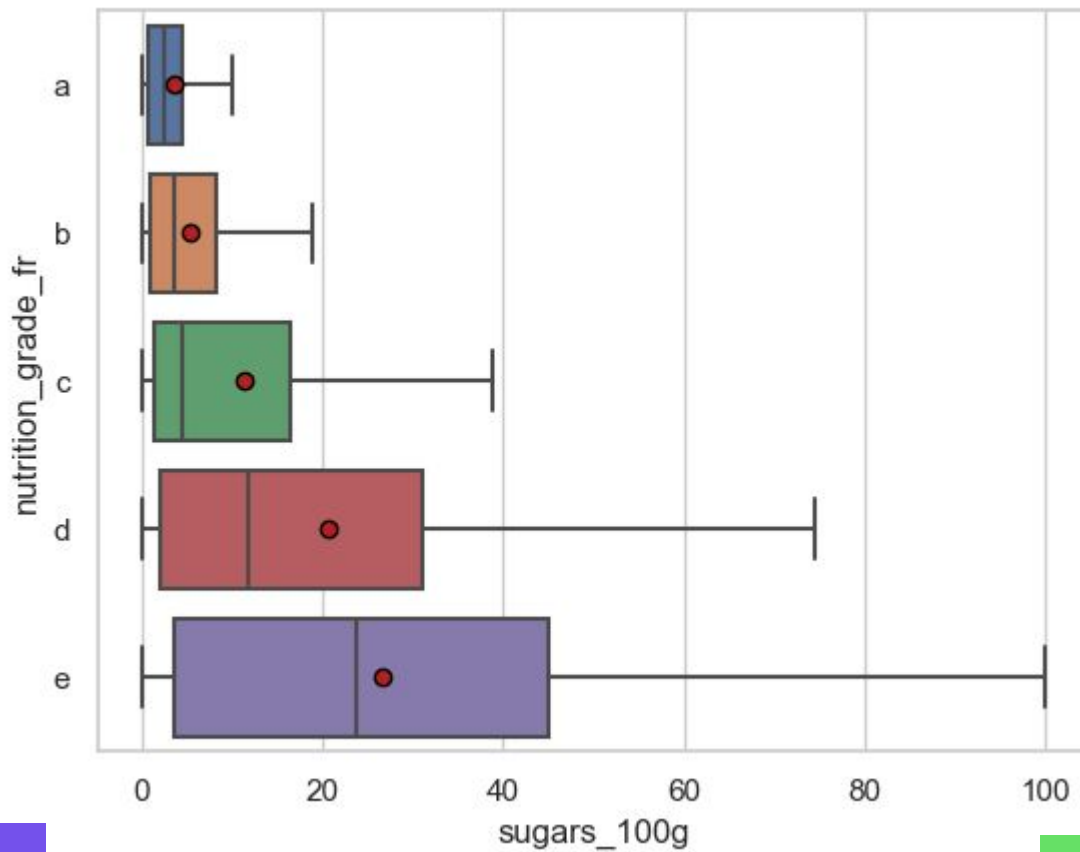
Eboulis des valeurs propres



ACP



ACP Heatmap



```
from scipy.stats import kruskal

groups = [group['sugars_100g'] for name, group in df.groupby('nutrition_grade_fr')]
h_statistic, p_value = kruskal(*groups)
print(f"Statistique H : {h_statistic}")
print(f"p-valeur : {p_value}")

alpha = 0.05 # Niveau de signification
if p_value < alpha:
    print("Les distributions des groupes sont différentes.")
else:
    print("Nous ne pouvons pas conclure que les distributions des groupes sont différentes.")
```



```
Statistique H : 31056.540053574132
p-valeur : 0.0
Les distributions des groupes sont différentes.
```

Protection des données personnelles (RGPD)



Licéité, loyauté et transparence

Limitation des finalités

Minimisation des données

Exactitude des données

Limitation de la conservation

Conclusion de l'étude



Corrélation entre différentes variables

Distribution des valeurs impactent le résultat du nutriscore

Synthétisation de variable

Imputation valeurs manquantes