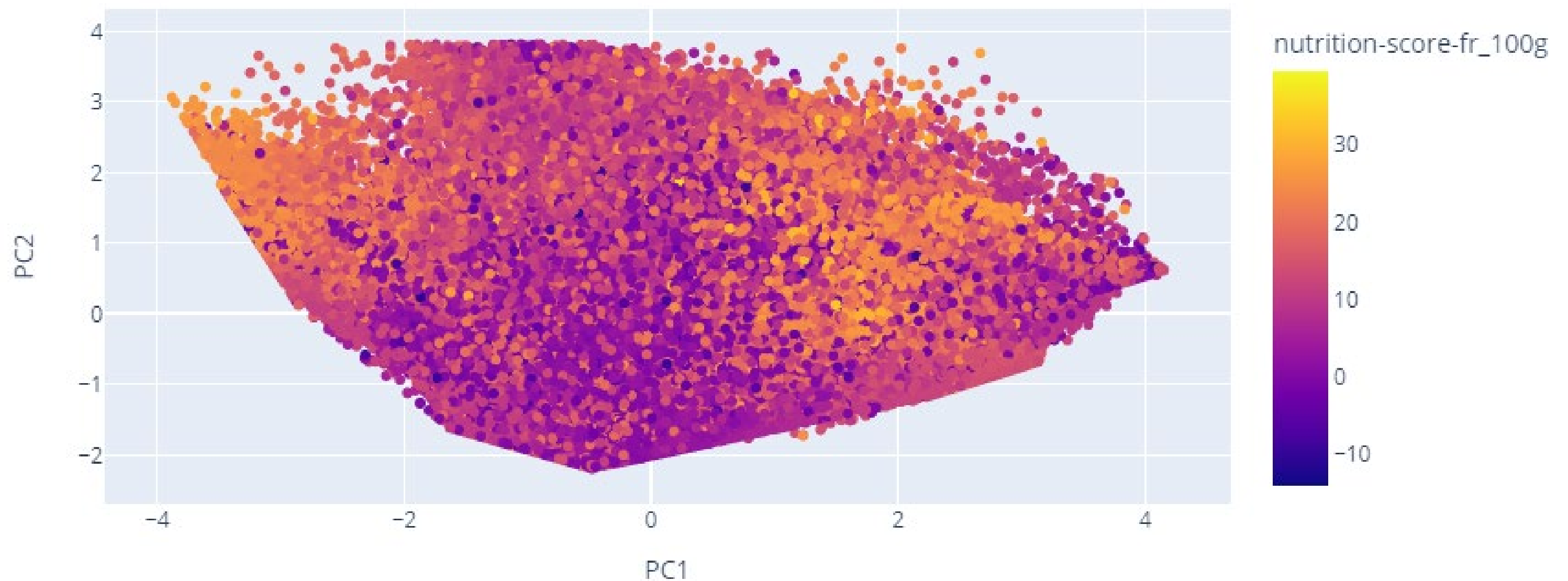


# Machine Learning Engineer

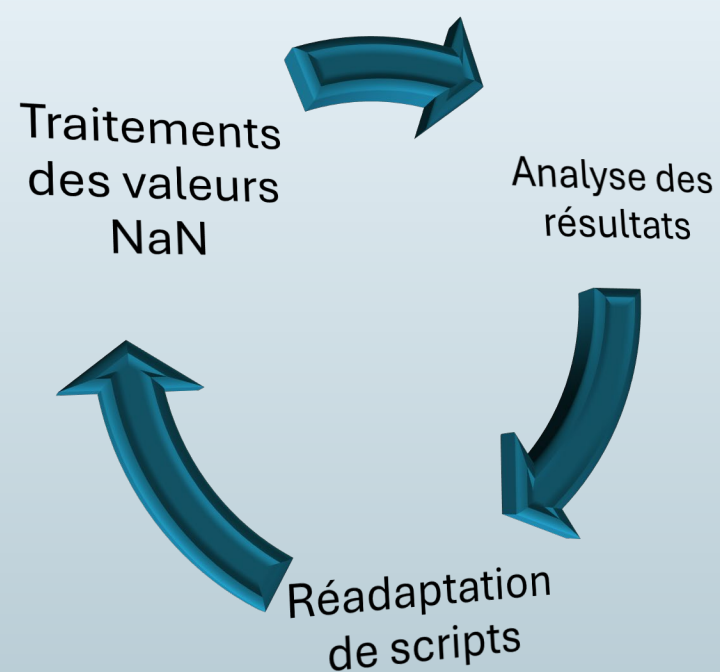
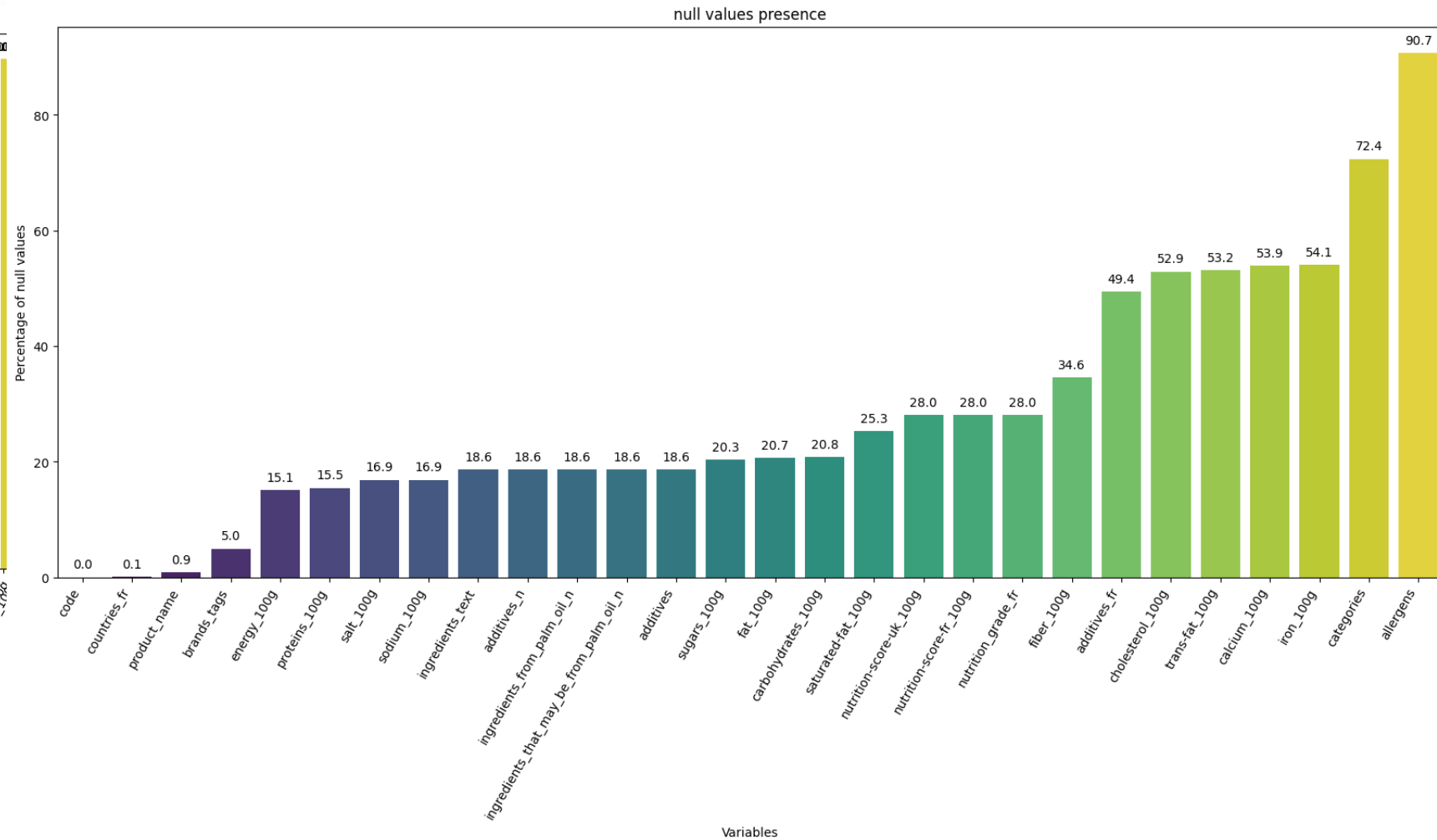
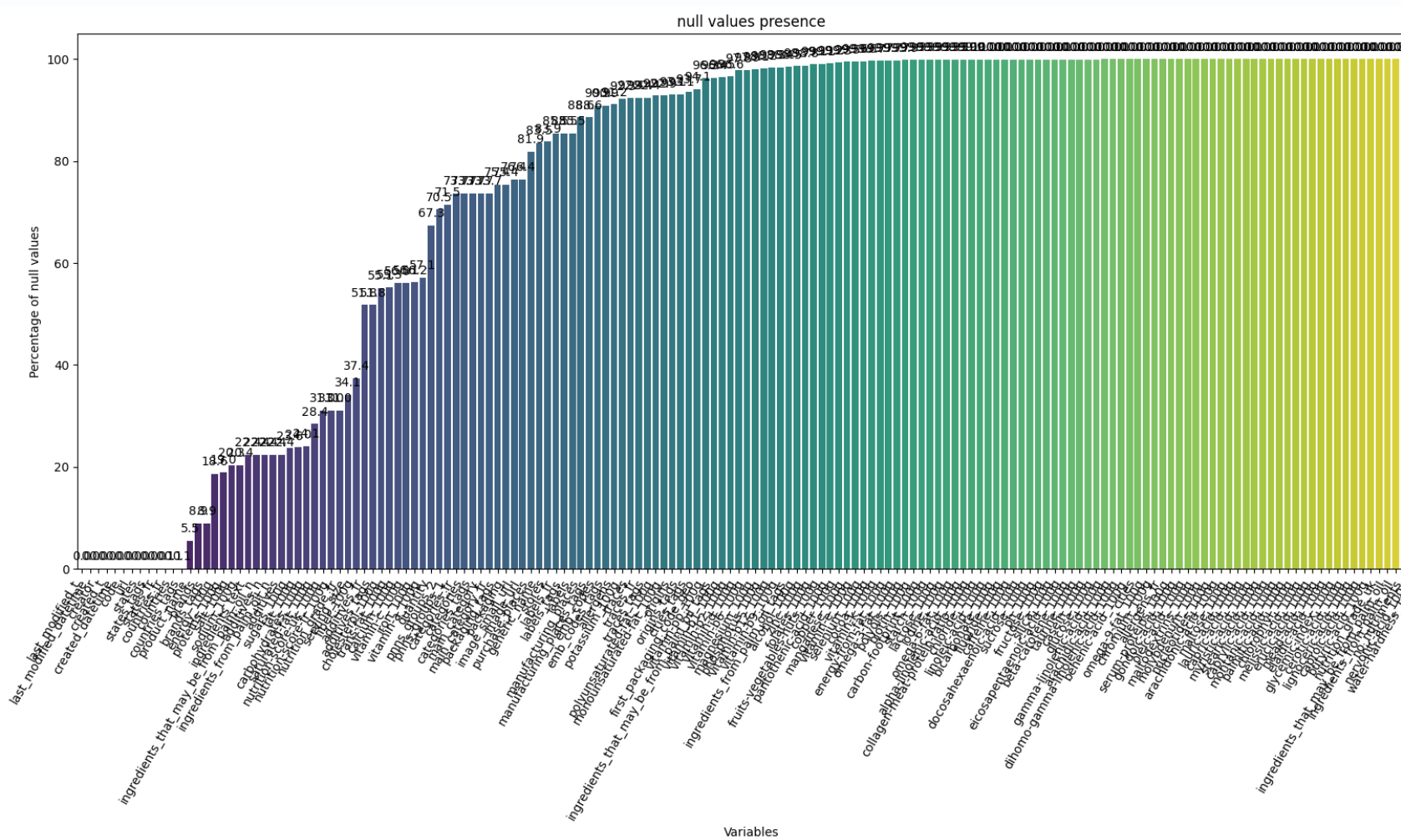
PCA Projection Colored by Nutrition Score



**Préparez des données pour un organisme de santé publique**

Faisabilité d'une application pour la gestion d'ajout des données openFood

# Machine Learning Engineer



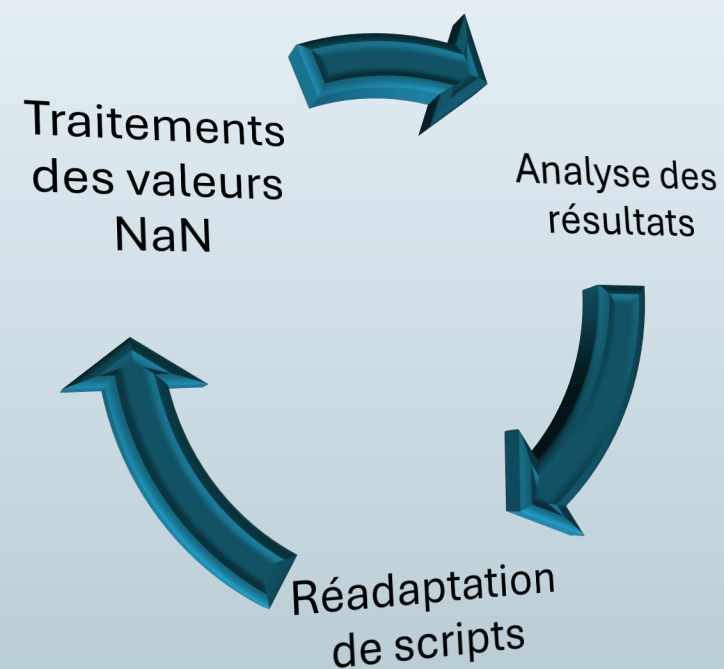
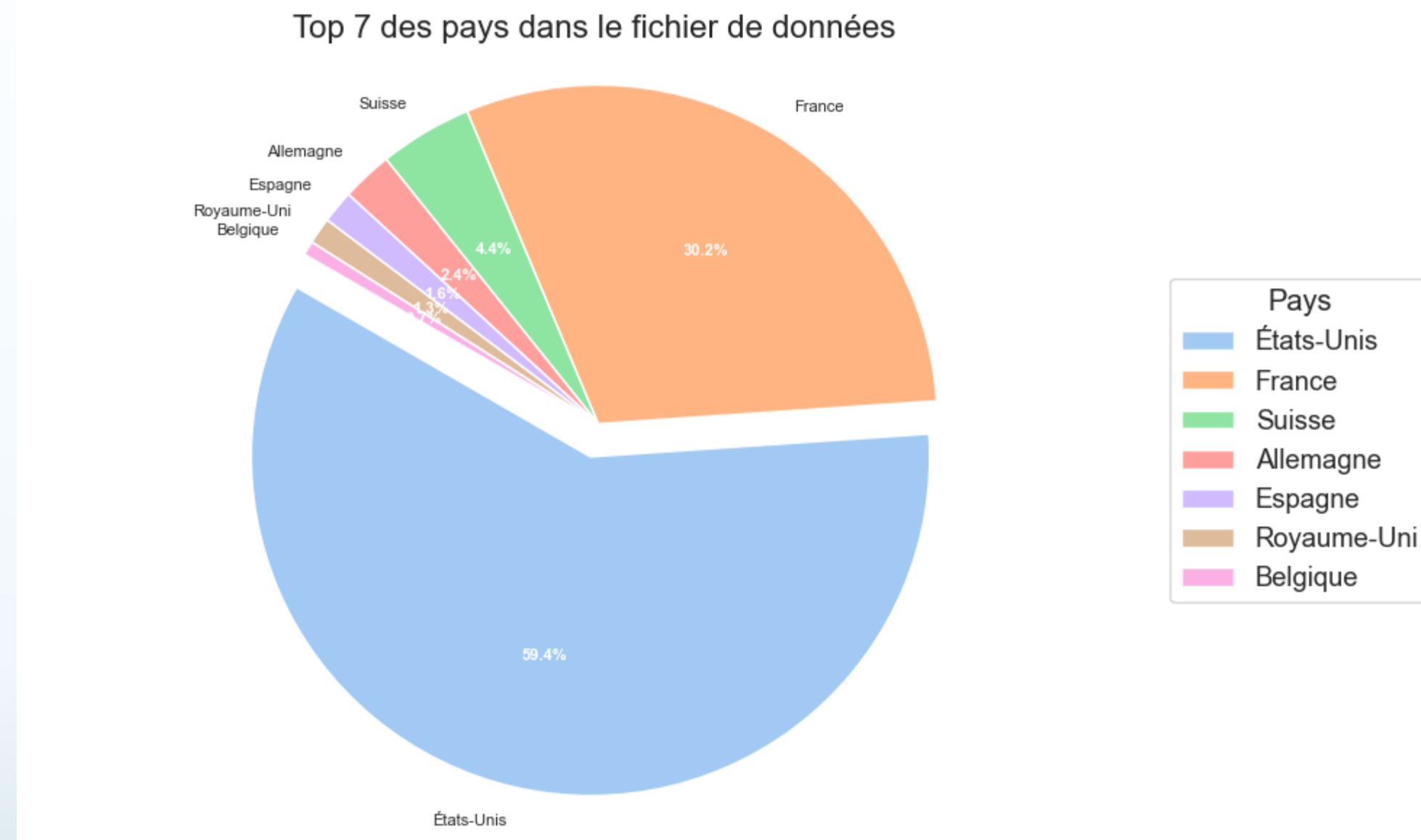
- Traitement large des valeurs null sur les données
- Respect de la RGPD dans la conservation de certaines données

# La préparation des données

## Suppression et Imputation de données

# Machine Learning Engineer

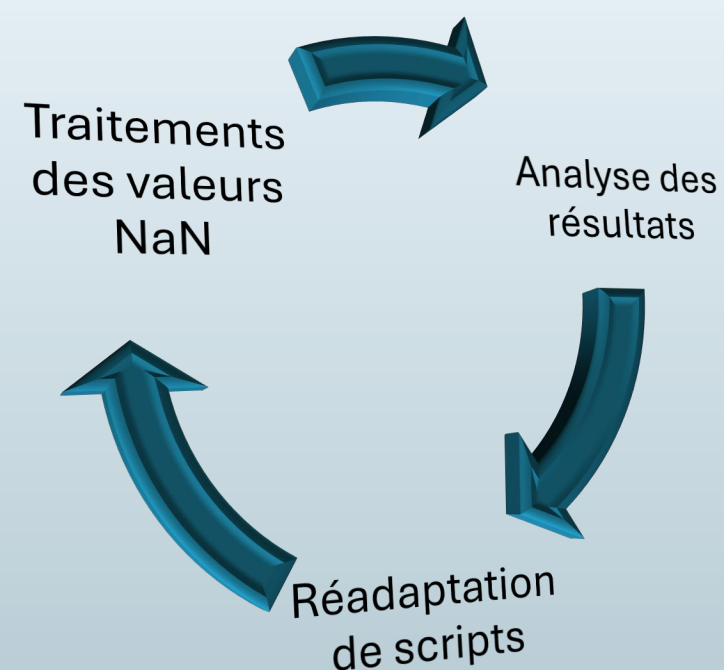
- Conservation des 7 pays les plus alimentés en données



**La préparation des données**  
**Suppression et Imputation de données**

# Machine Learning Engineer

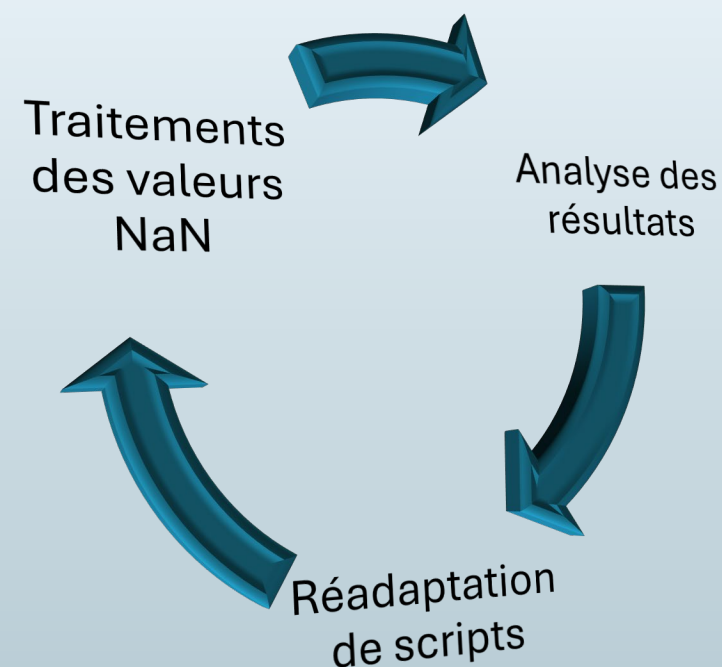
- Traitement des doublons : chaque produit est unique au sein de sa marque
- .Nettoyage de la colonne Ingredient\_text dans une démarche d'anticipation



**La préparation des données**  
**Suppression et Imputation de données**

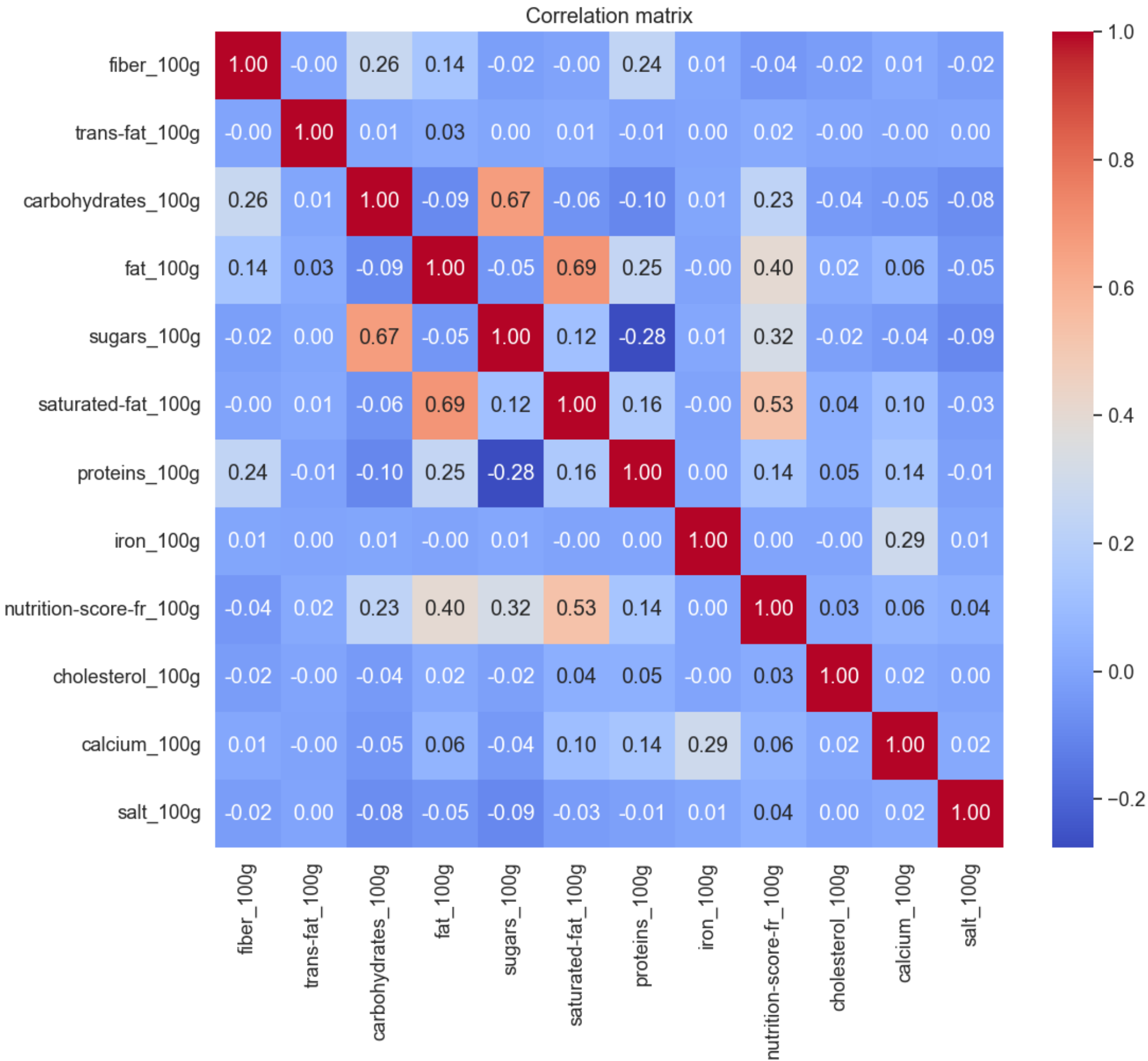
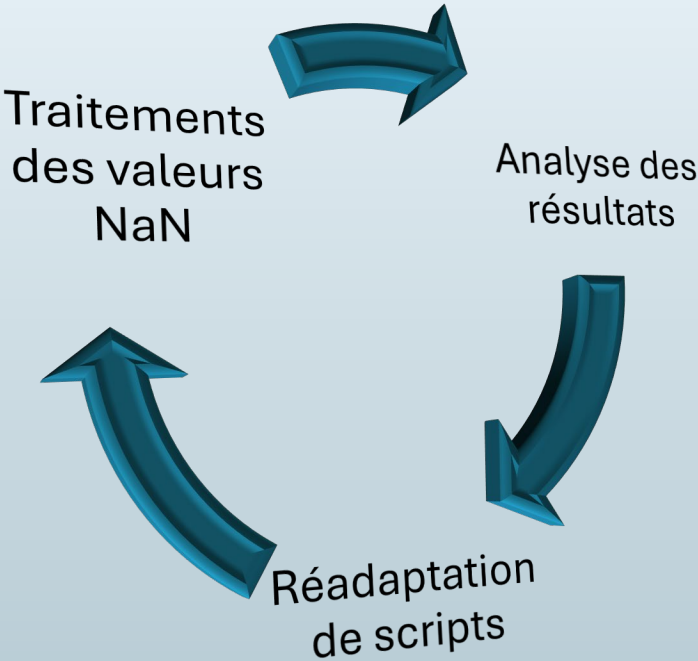
# Machine Learning Engineer

- Traitement général de valeur null
- Limitation générale des valeurs aberrantes
  - Les calories joules ne sont pas supérieurs à 3900 unités pour 100g
  - Les variables numériques pour 100g ne peuvent pas être inférieur à 0 ni supérieur à 100



**La préparation des données**  
**Suppression et Imputation de données**

Matrix de Correlation



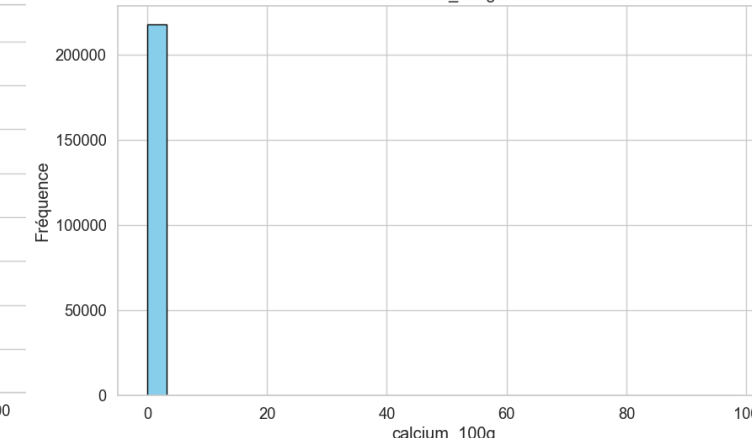
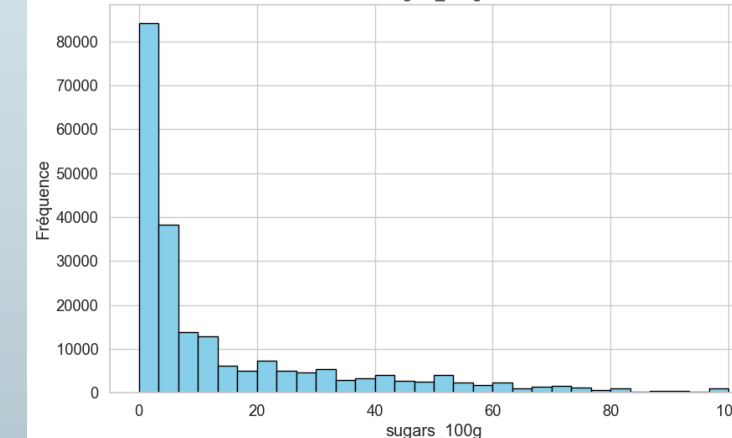
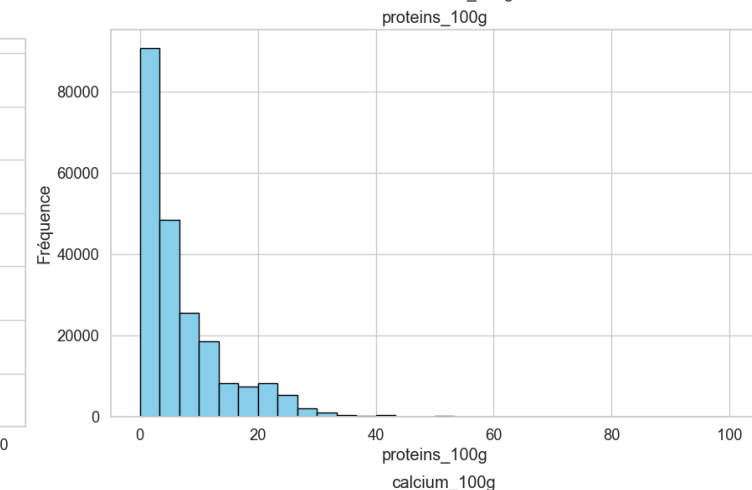
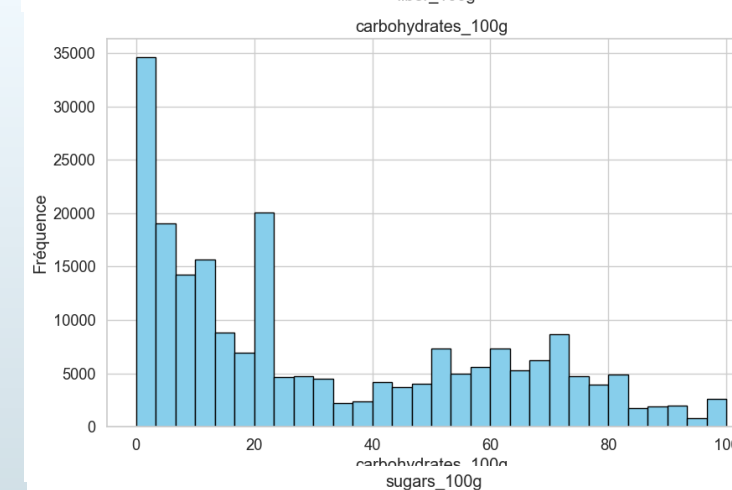
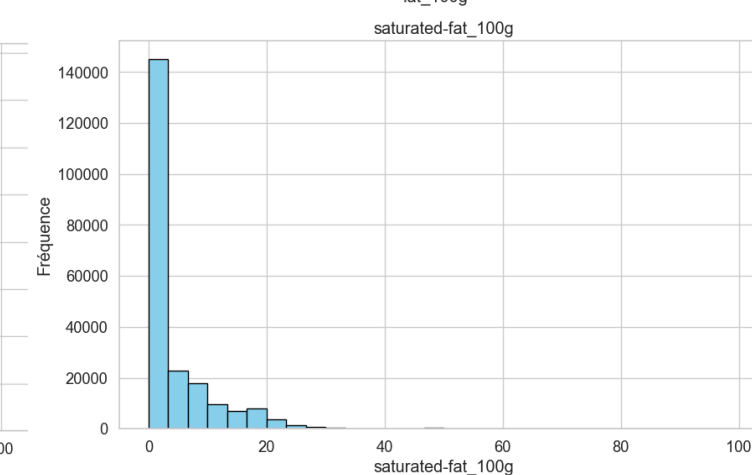
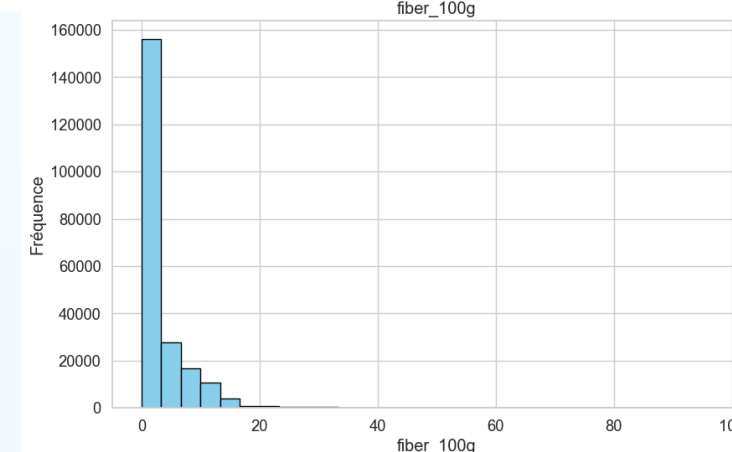
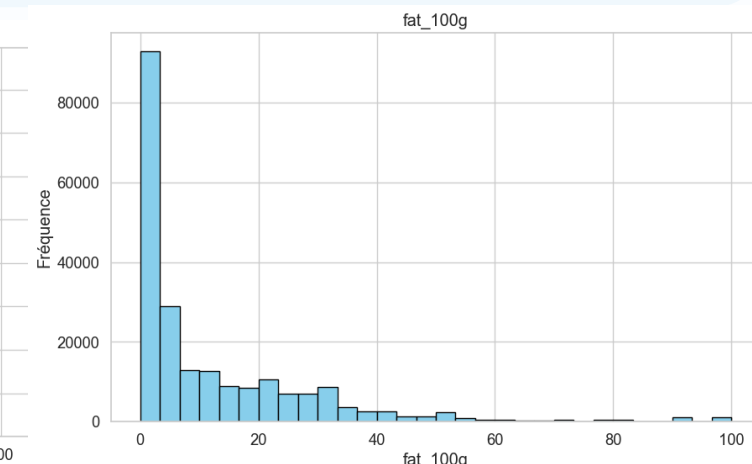
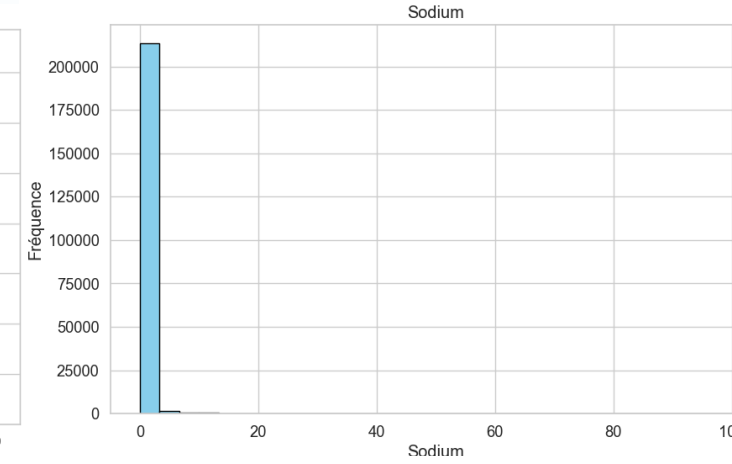
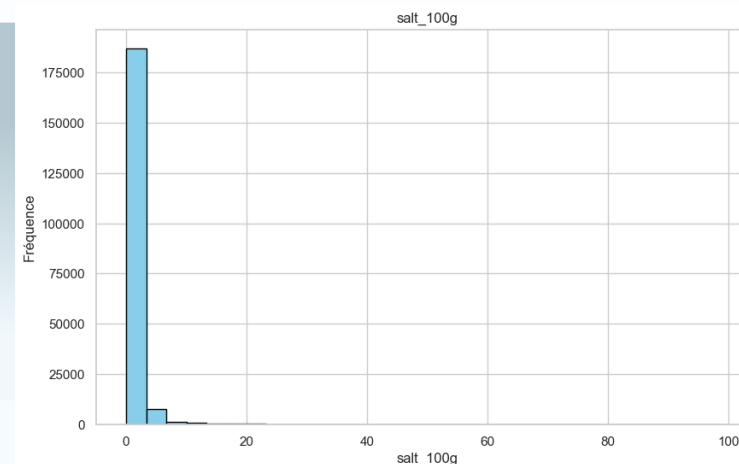
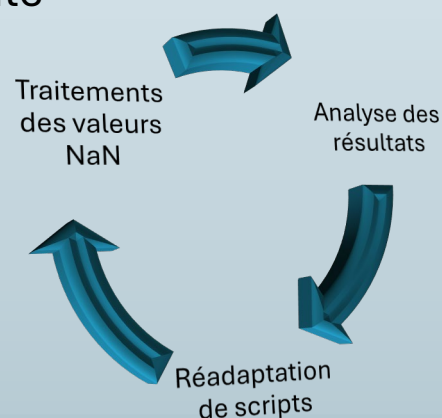
Vérification de la relation entre les données  
Analyse et Imputation de données



# Machine Learning Engineer

## Imputation des valeurs Nan

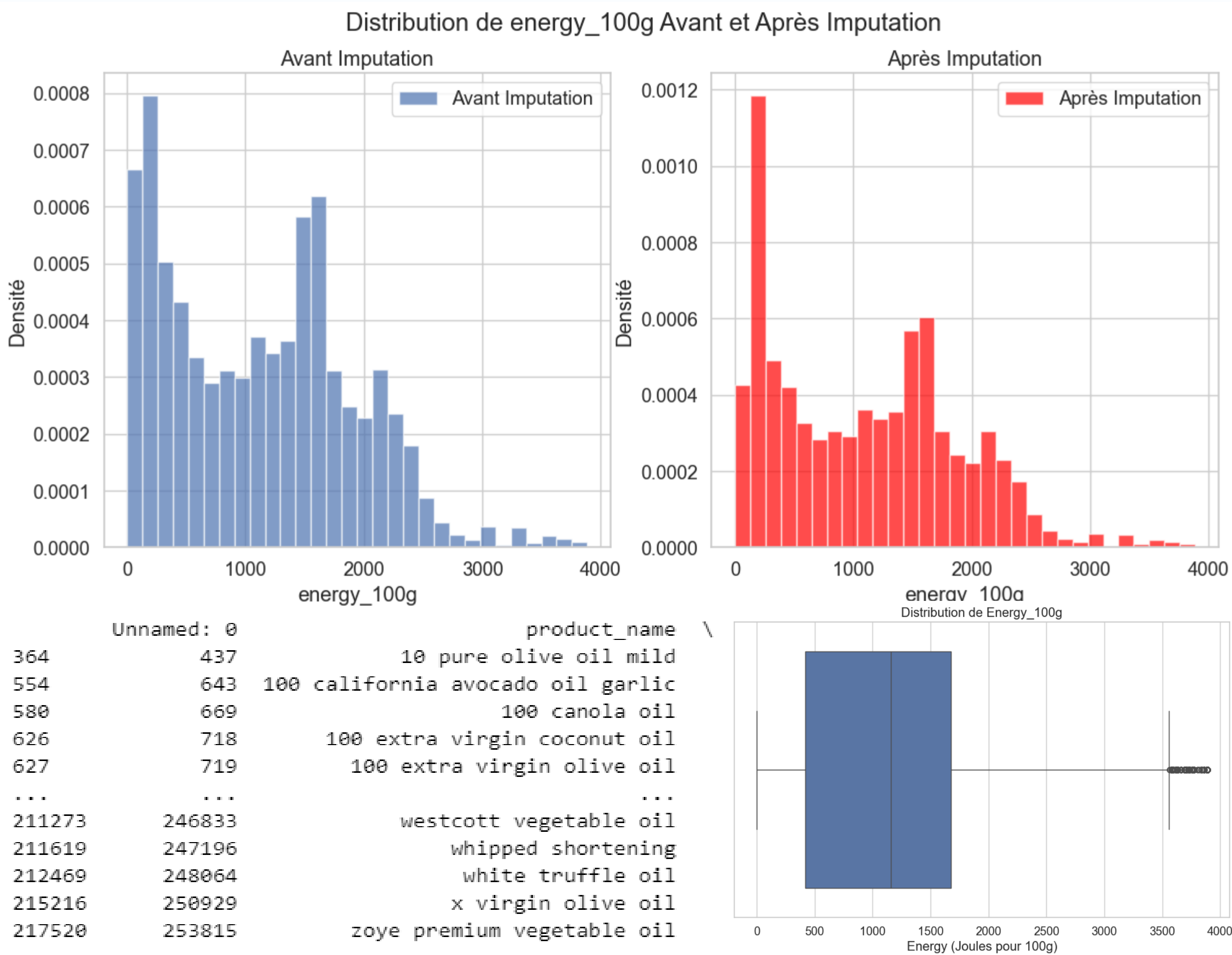
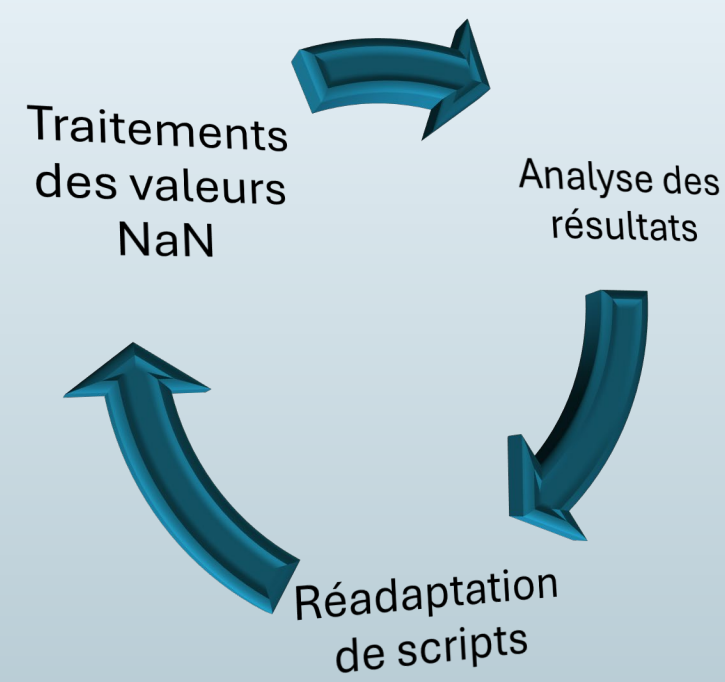
- La variable salt\_100g : la médiane dont la moitié est diminué si le sucre est conséquent pour ses valeurs NaN
- La variable sodium\_100g : un ratio de 40% du poids du sel
- Les variables sugar\_100g, carbohydrates\_100g, fiber\_100g, leur valeur NaN sont imputé en relation les un avec les autres.
- Les variables fat\_100g et saturated-fat\_100g sont corrélées, leur imputation se fait par ratio selon les valeurs présentes.
- La variable protein\_100g et celle du calcium\_100g, leur Nan sont imputé et influencé par la présence du sucre
- La variable Cholesterol\_100g : nous estimons que si elle est null, elle est absente



**Vérification de la relation entre les données**  
**Analyse et Imputation de données**

# Distribution des calories

- Constat de la presence de 2 unites de mesure : Les joules et calories
- Imputation avec la méthode KNN : influence des plus proches voisins
- Certaines valeurs aberrantes concernent des produits en rapport avec le domaine des huiles, donc légitime



## La préparation des données

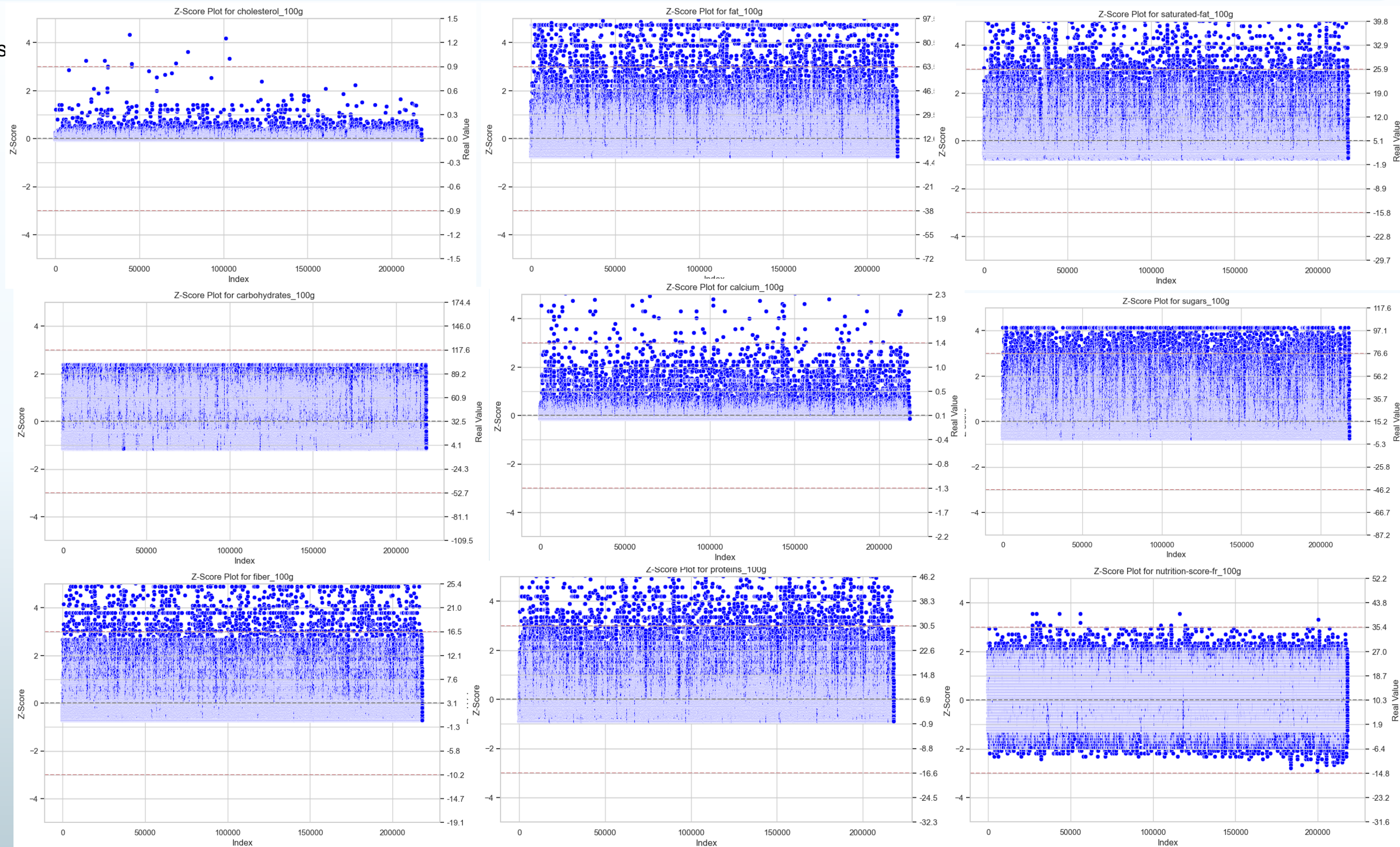
### Analyse et Imputation de données



# Z-Score

# Machine Learning Engineer

- Travaile sur les valeurs aberrantes :



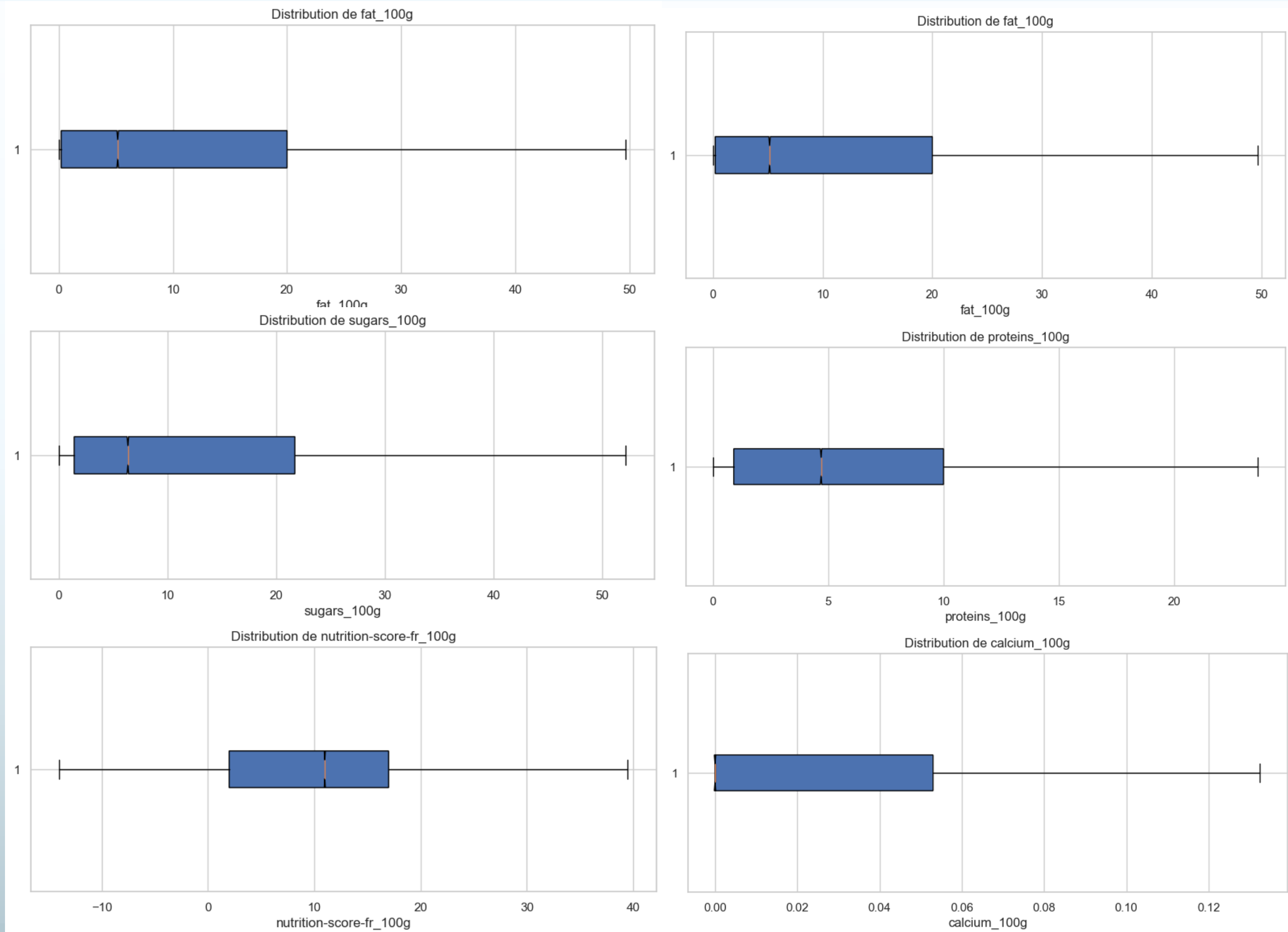
**La préparation des données**  
**Analyse et Traitement des valeurs aberrantes**

# Machine Learning Engineer

## Z-Score

- Travail sur les valeurs aberrantes
- Nous repérons la valeur non aberrante la plus forte et nous adaptons l'extrême à cette valeur, cela permet de respecter la variance entre les données

Après  
traitement :

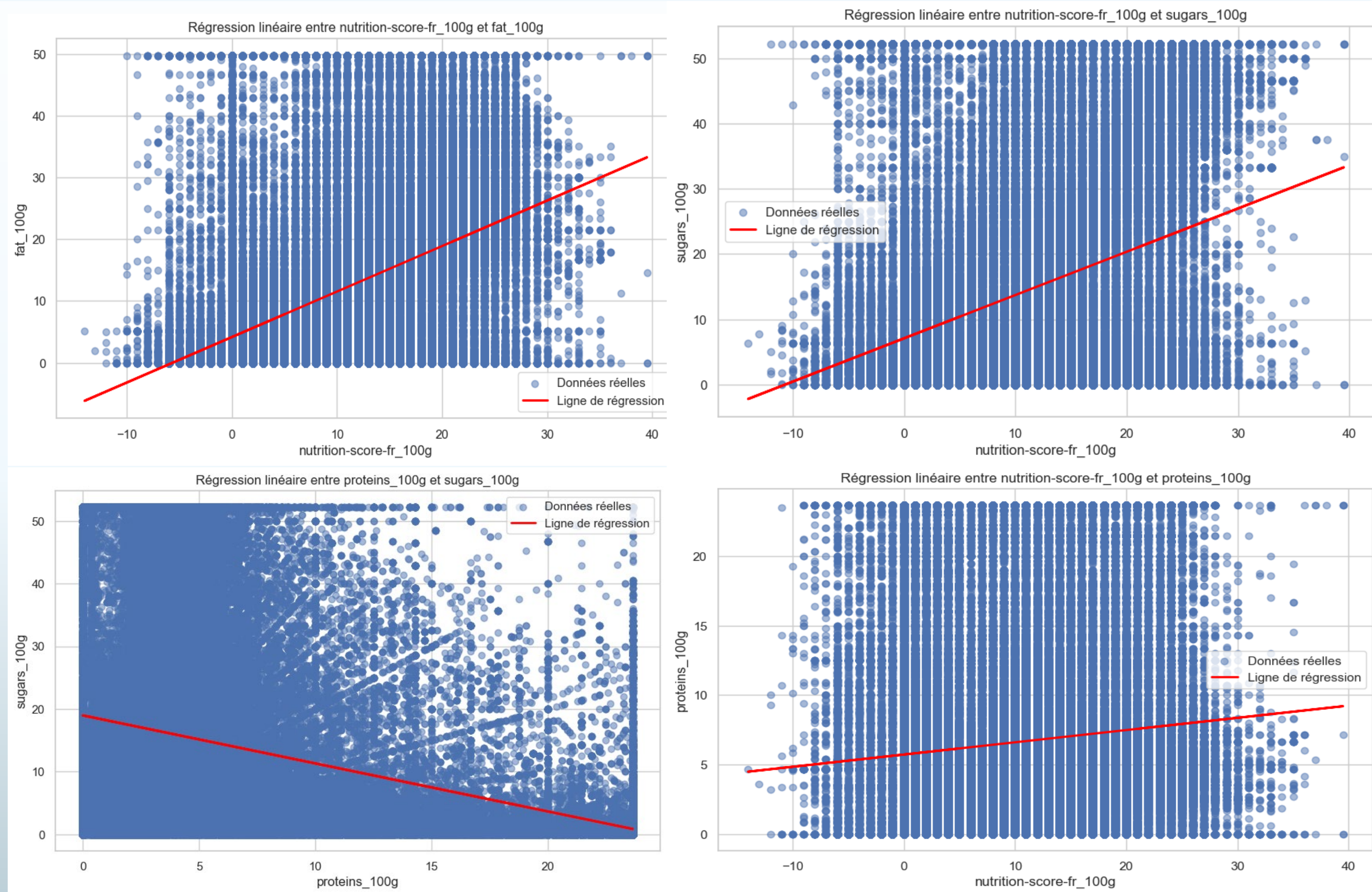


**La préparation des données**  
**Analyse et Traitement des valeurs aberrantes**



# Machine Learning Engineer

- Le gras et le sucre influencent légèrement le score de nutrition
- La proteine ne semble pas impacter le score de nutrition



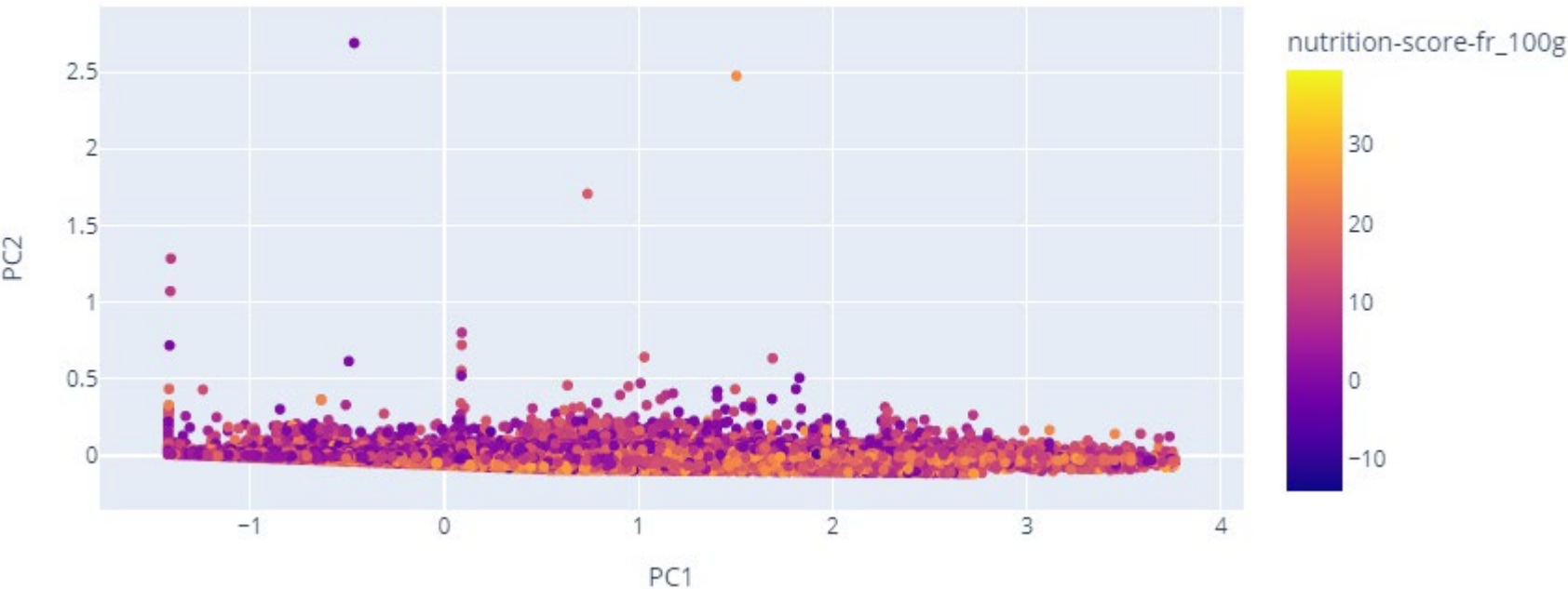
**Regression Lineaire**

**Etude des tendances  
Analyse Bivariée**



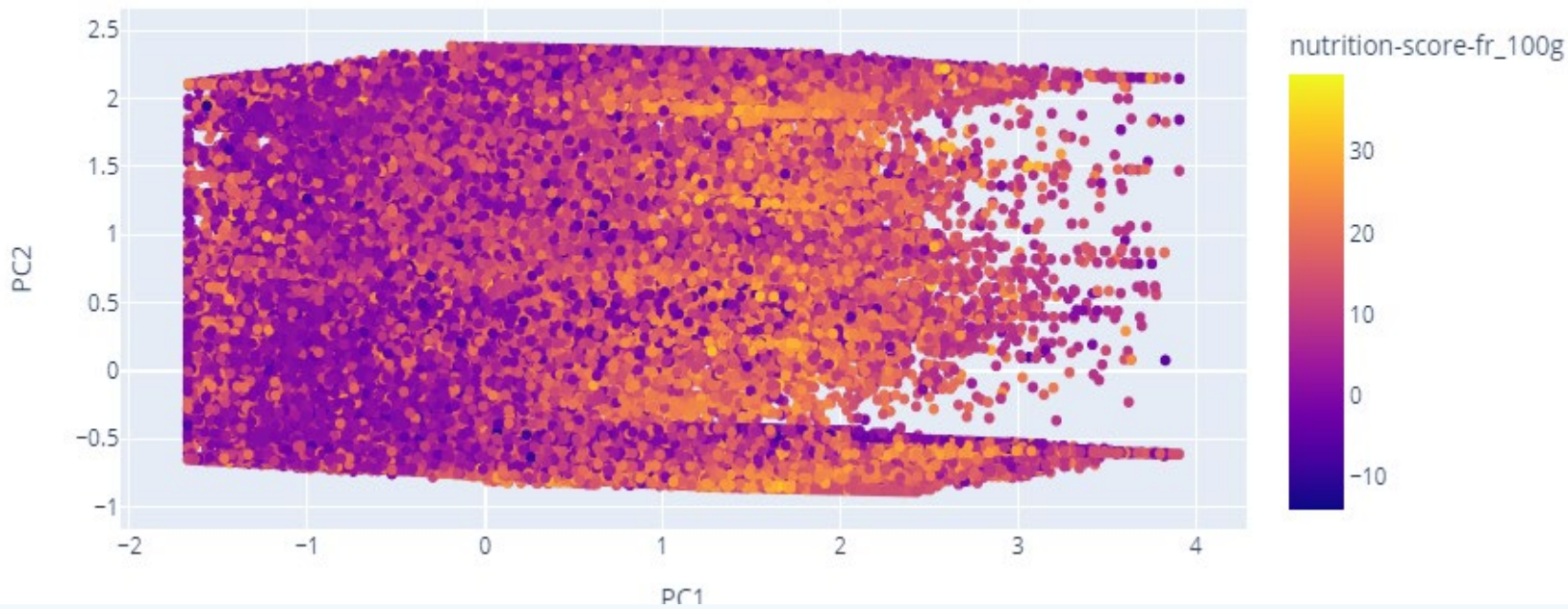
# Machine Learning Engineer

PCA Projection Colored by Nutrition Score



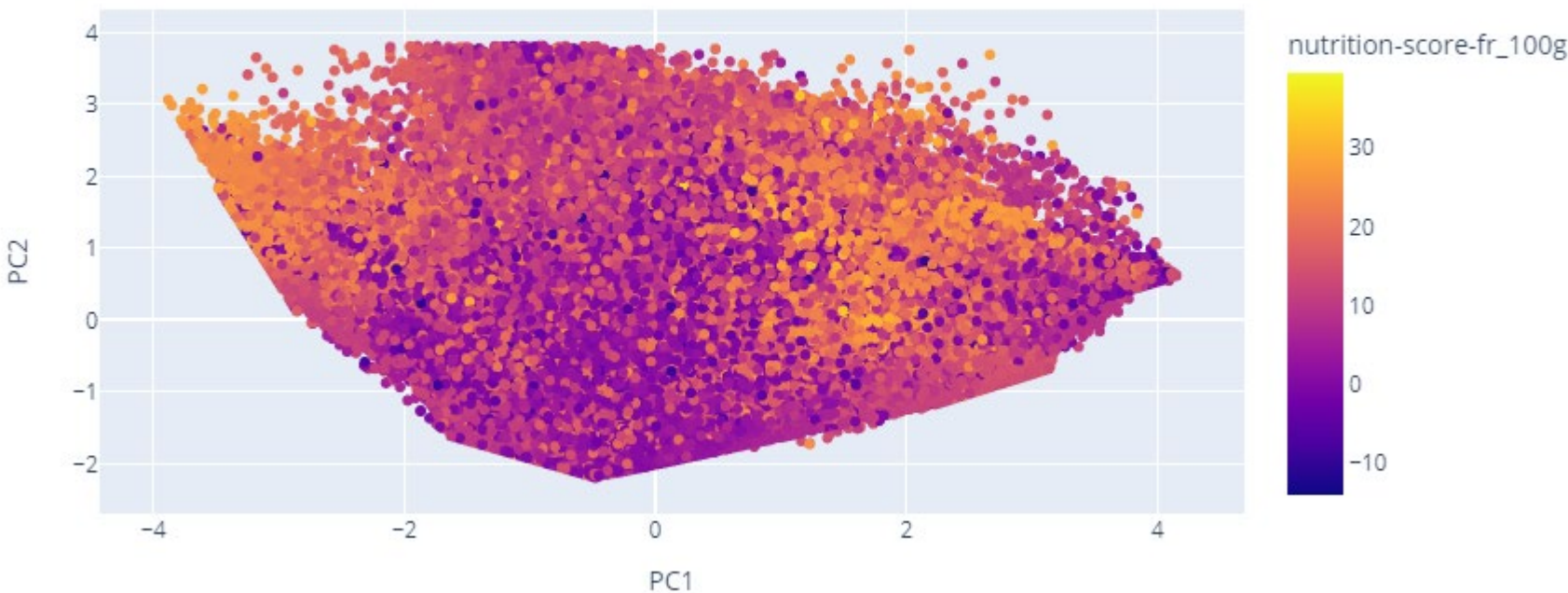
```
['fat_100g', 'proteins_100g', 'iron_100g',  
                                     'cholesterol_100g',  
'calcium_100g'], 'nutrition-score-fr_100g', 'product_name')
```

PCA Projection Colored by Nutrition Score



```
['fiber_100g', 'carbohydrates_100g', 'sugars_100g',  
'calcium_100g'], 'nutrition-score-fr_100g'
```

PCA Projection Colored by Nutrition Score

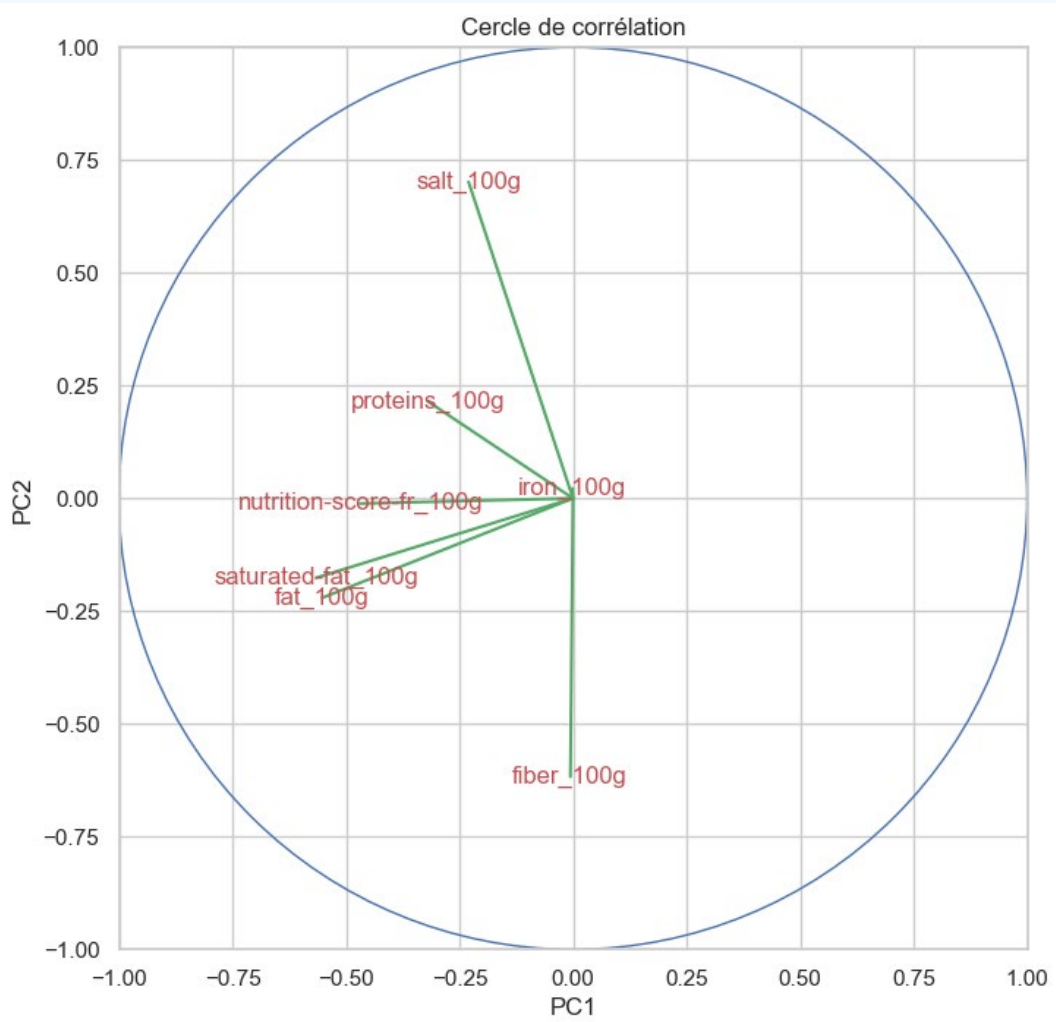
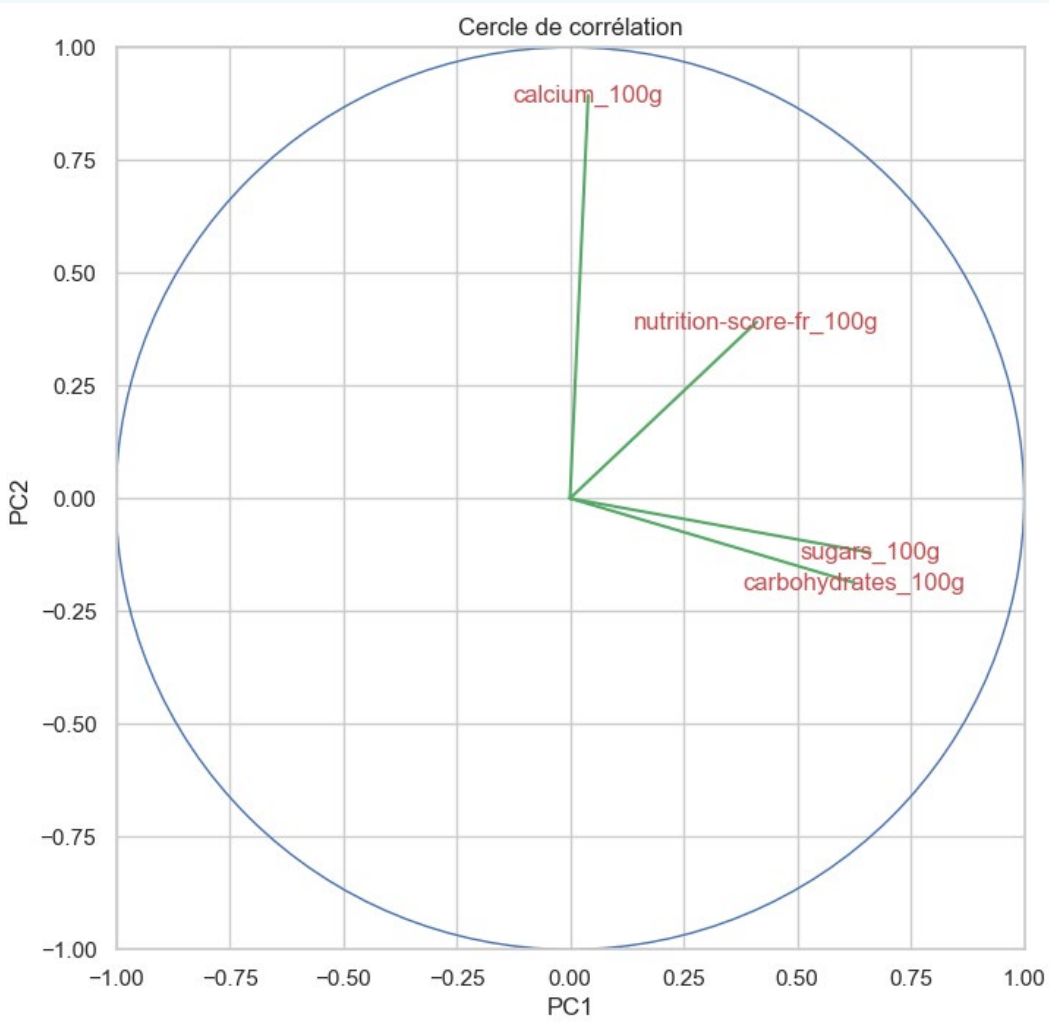
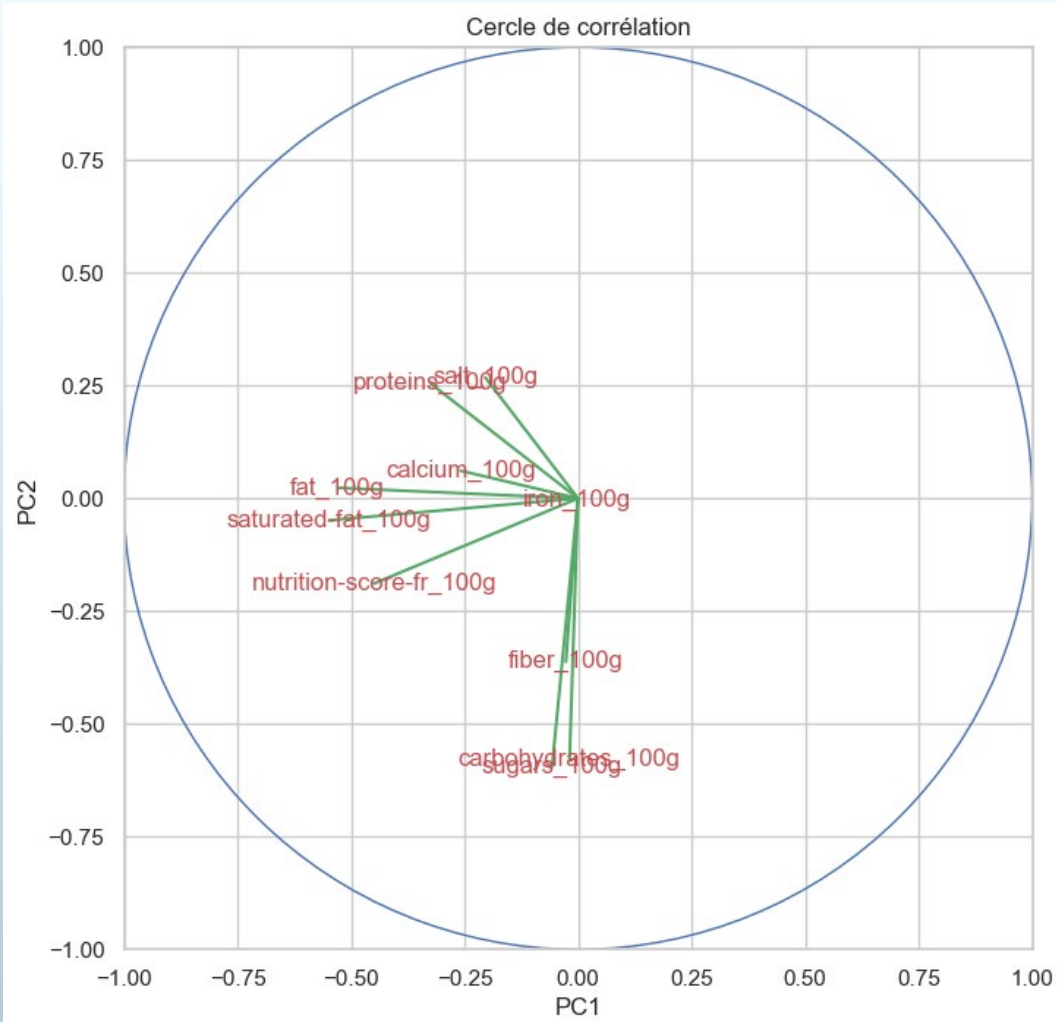


```
['fiber_100g', 'carbohydrates_100g', 'sugars_100g', 'calcium_100g', 'fat_100g',  
'proteins_100g', 'iron_100g',  
                                     'cholesterol_100g', 'salt_100g'], 'nutrition-score-fr_100g'
```

ACP

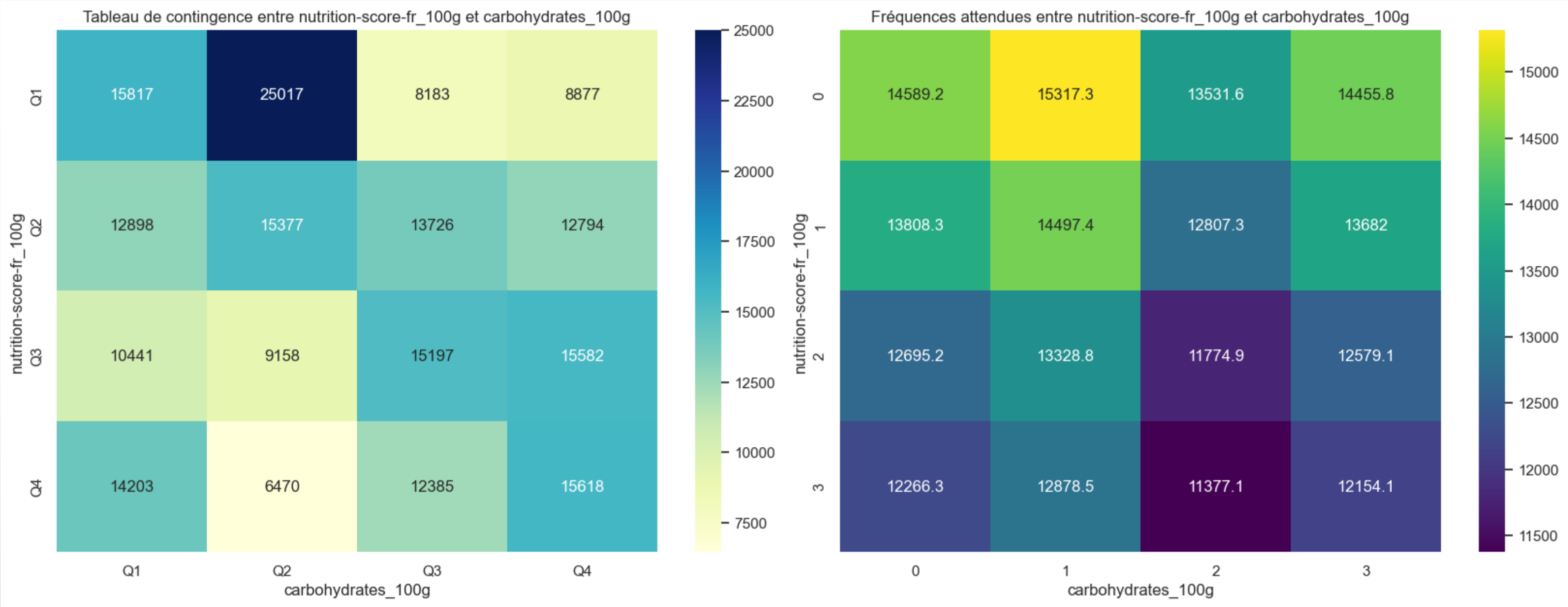
Analyse des Composants Principaux

# Machine Learning Engineer



## Cercle de Corrélation

# Machine Learning Engineer

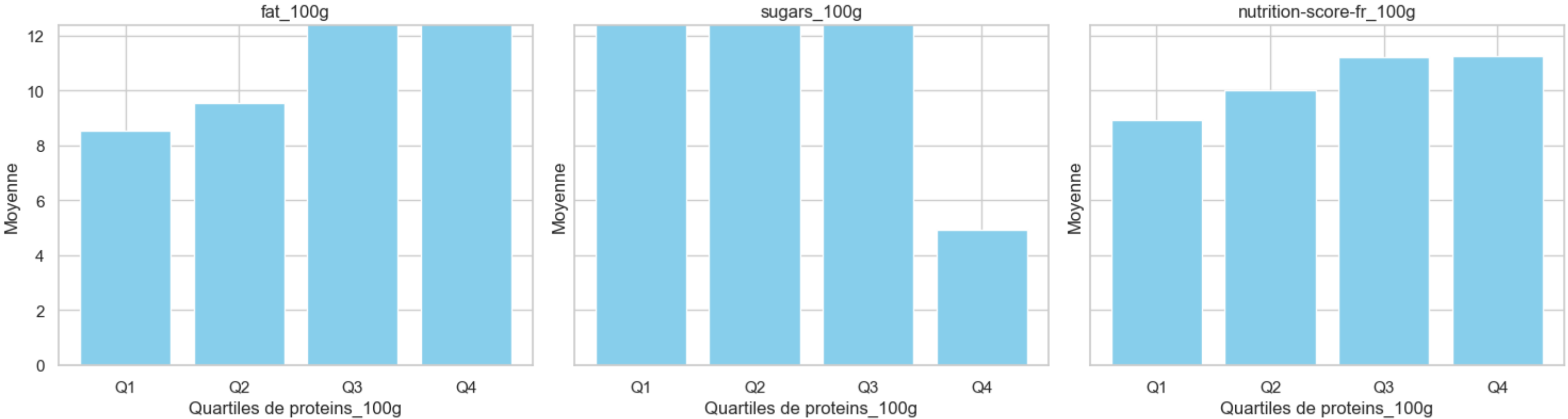


## Le test Chi-carre

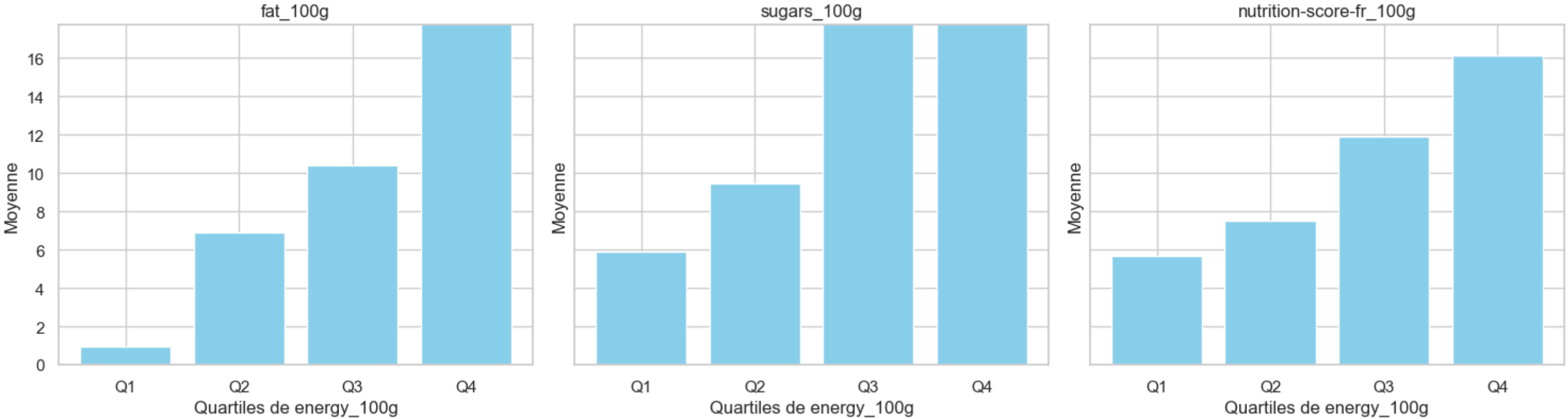


# Machine Learning Engineer

Moyennes par quartiles de protéine pour chaque variable de réponse



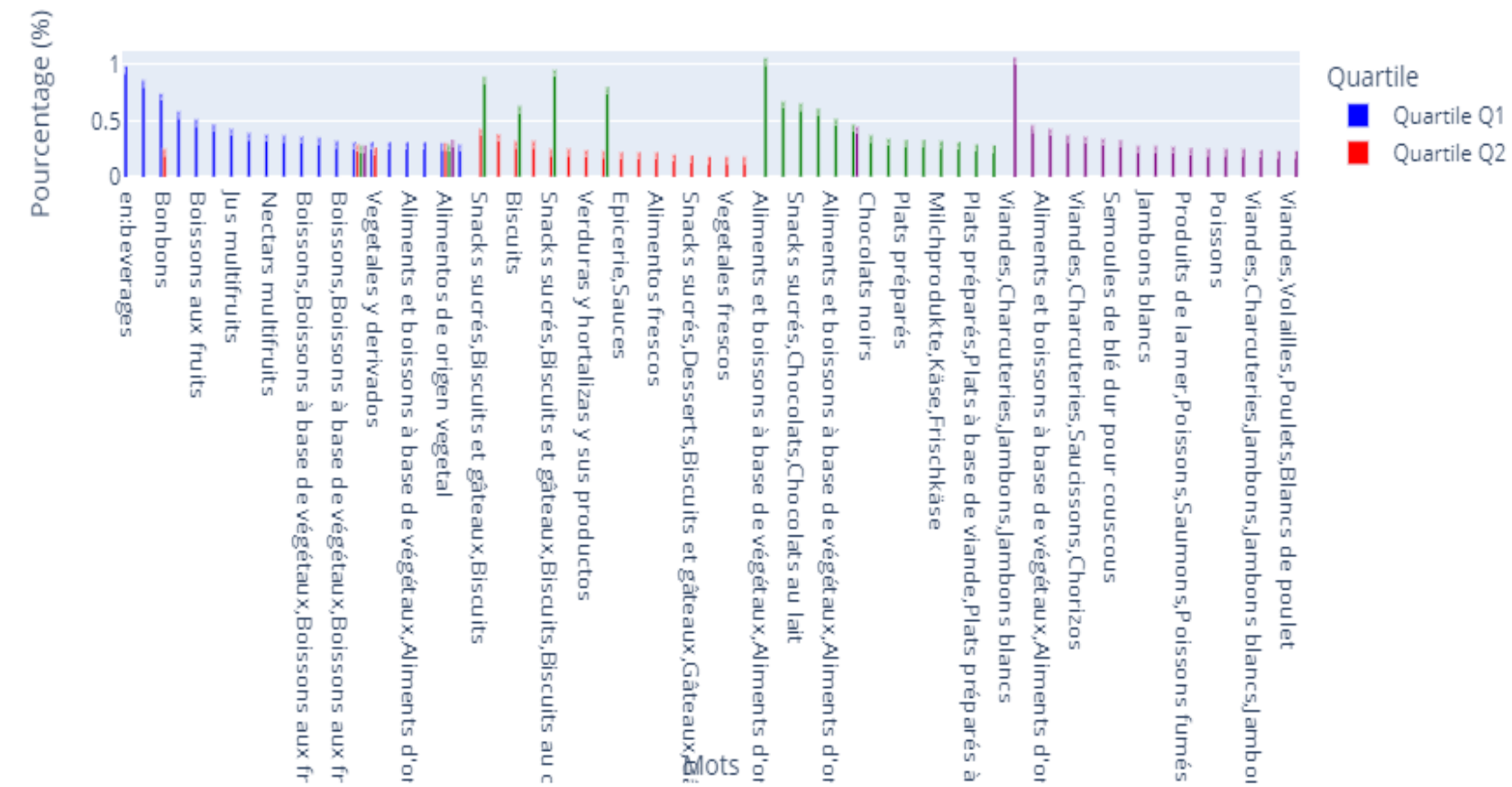
Moyennes par quartiles de protéine pour chaque variable de réponse



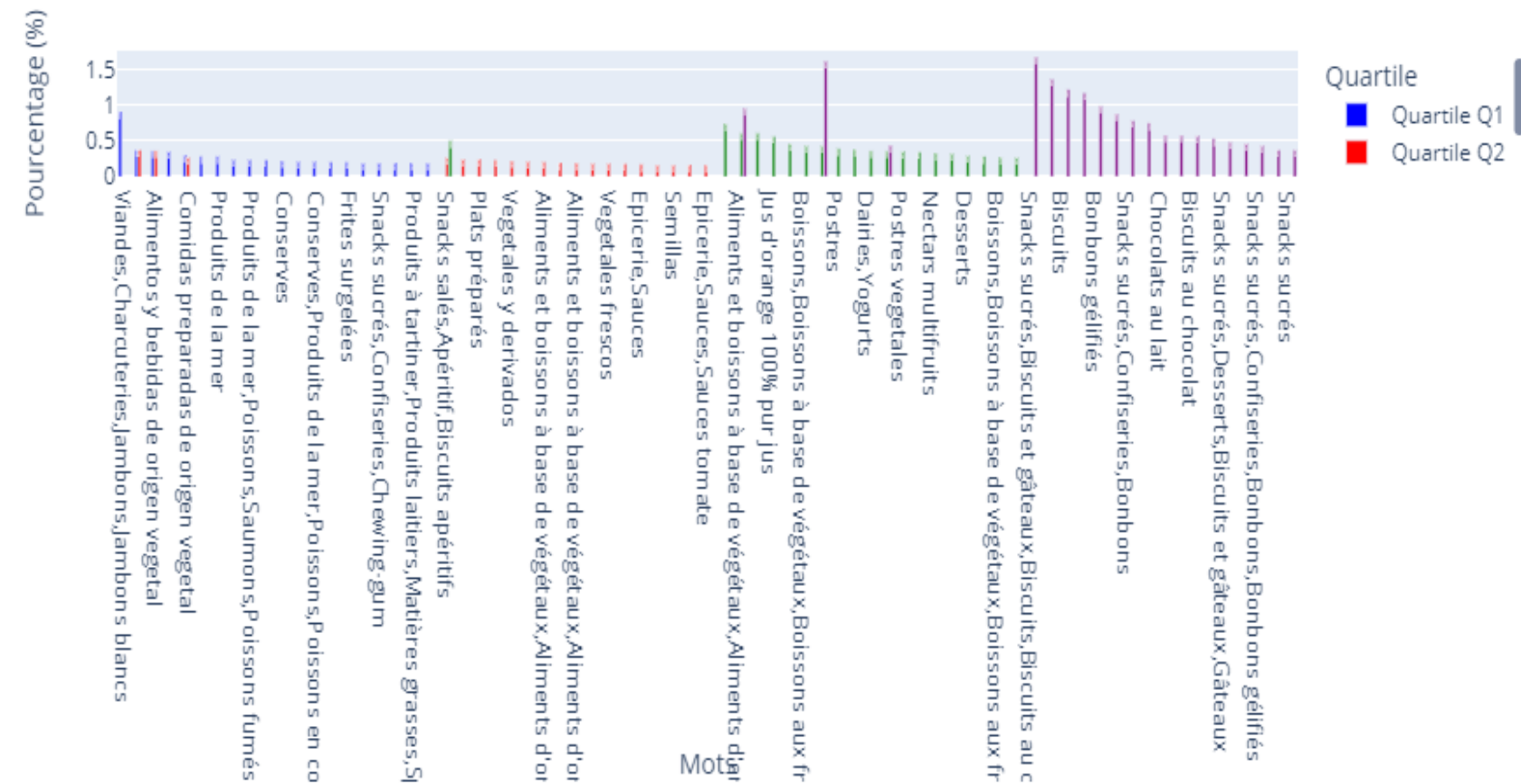
ANOVA

# Machine Learning Engineer

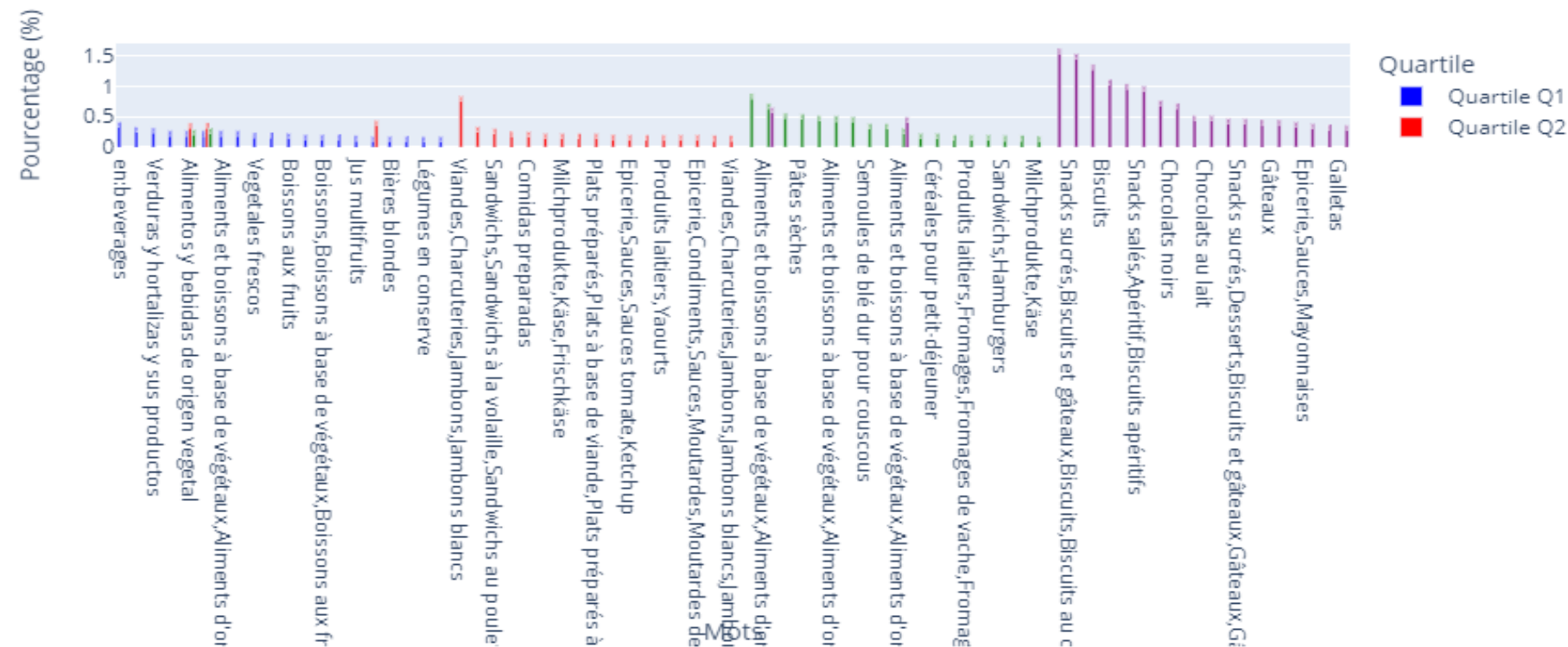
### Pourcentage des mots les plus fréquents par quartile de protéine



### Pourcentage des mots les plus fréquents par quartile de sucre



### Pourcentage des mots les plus fréquents par quartile de calorie



# Q1

## Q2

Q3

Q4

# ANOVA

## Faisabilite de l'application