

# Application en lien avec l'alimentation

---

Appel à projet par l'agence "Santé publique France".



Le jeu de données Open Food Fact

Informations générales du produit : nom, date de modification, etc.

Tags : catégorie du produit, localisation, marque, etc.

Valeurs nutritionnelles pour 100 grammes du produit.



# Sommaire

---

- ❖ Idée application

- ❖ Nettoyage

- Variables « intéressantes »
  - Taux de remplissage
  - Colonnes redondantes
  - Localisation

- ❖ Analyse univariée

- ❖ Analyse multivariée

- corrélation
  - réduction dimensionnelle
  - ANOVA

- ❖ Synthèse

# Idée d'application

# Idée d'application

## Informations du data set

- 1 750 000 produits, 184 "informations"
- Code
- Nutriscore
- Categories
- Données nutritionnelles
- Emballages/empreinte Co2
- Photo
- Marque

## Idée d'application

- Scan le produit/photo et envoie sur la fiche produit avec ses information ou bien la créer si le produit n'est pas dans la DataBase
- Evaluation du Nutriscore en fonction des valeurs nutritionnelles si celui n'est pas présent, propositions produits de meme categorie qui a un meilleur score
- Impact Co2

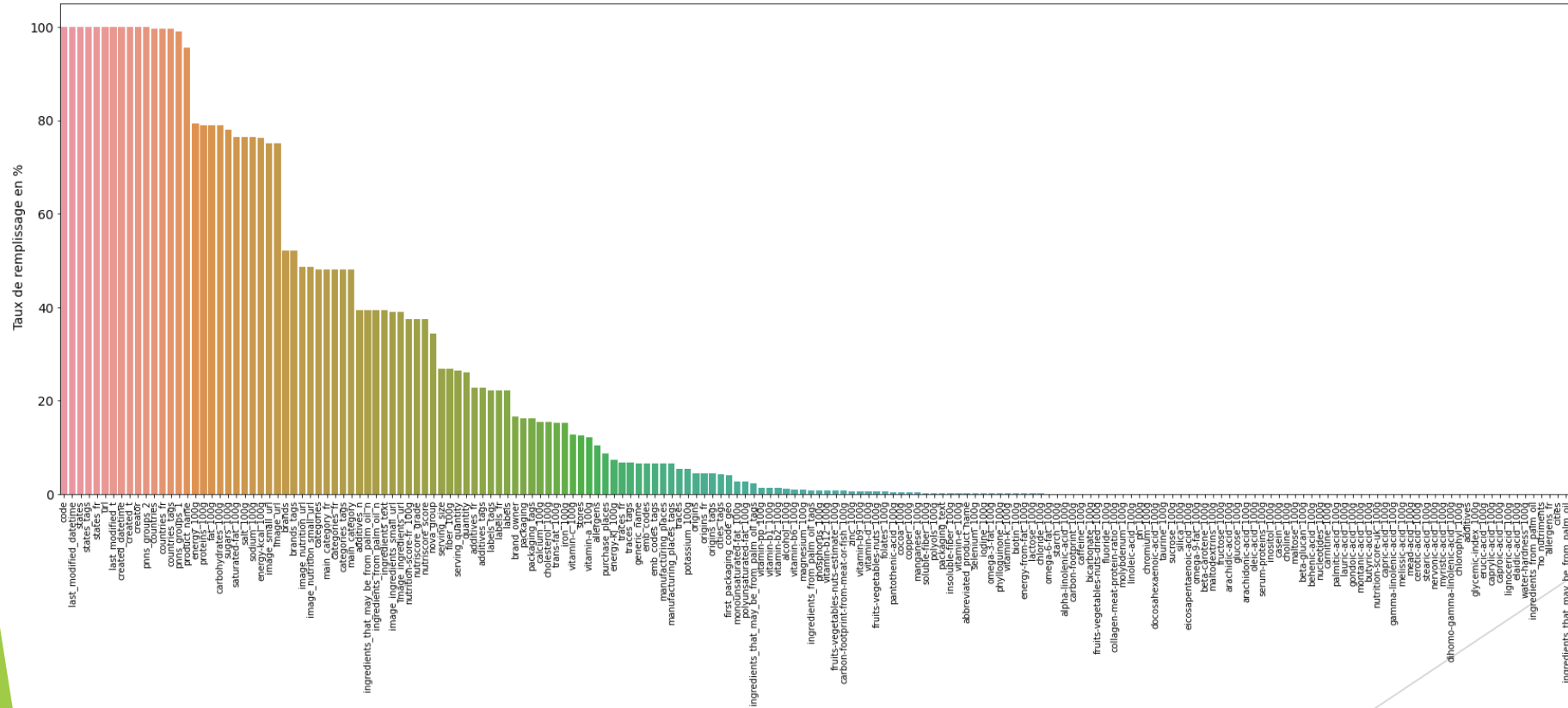
The background features abstract, overlapping green geometric shapes, primarily triangles and polygons, in various shades of green, creating a modern and dynamic visual effect. The word "Nettoyage" is centered in a green serif font.

Nettoyage

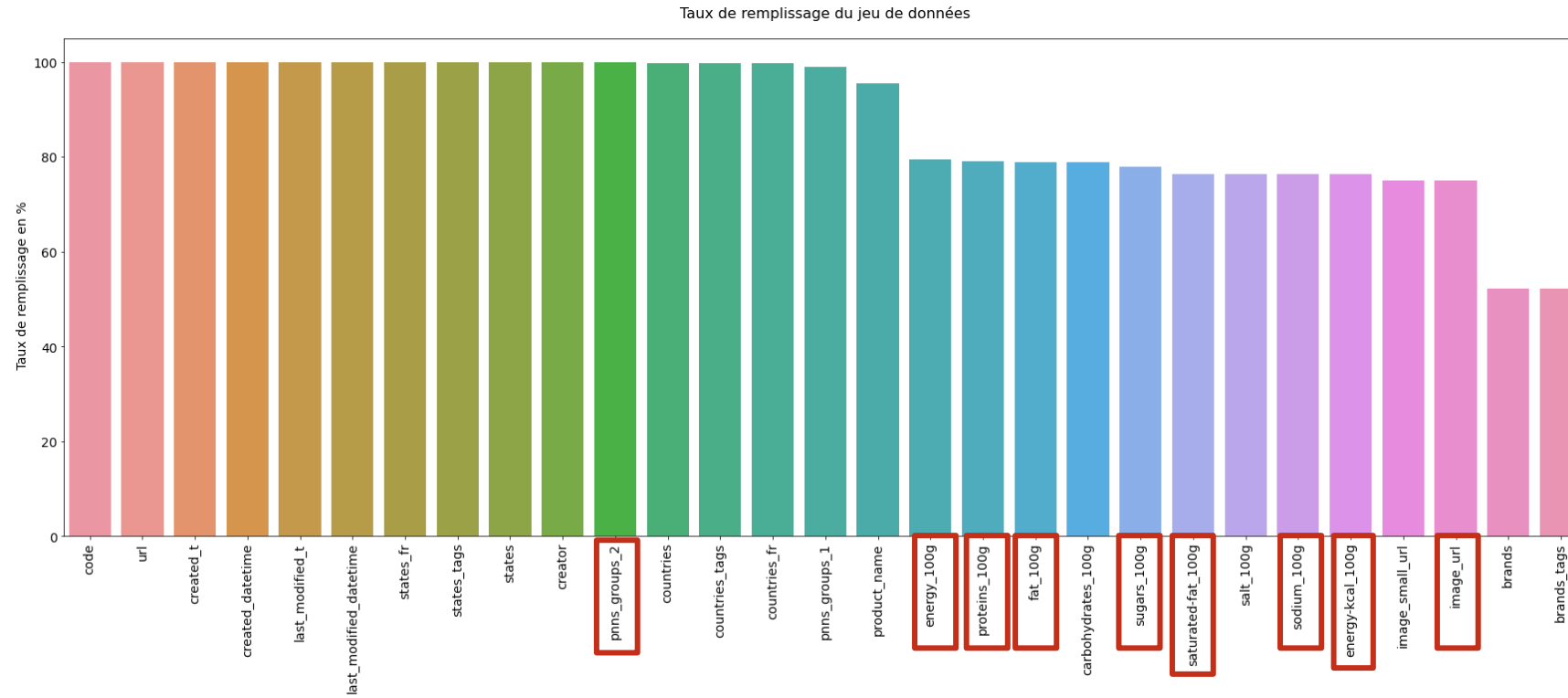
## Identifications de doublons (global/code)

## Identification des variables vides ou interessantes

### Taux de remplissage du jeu de données

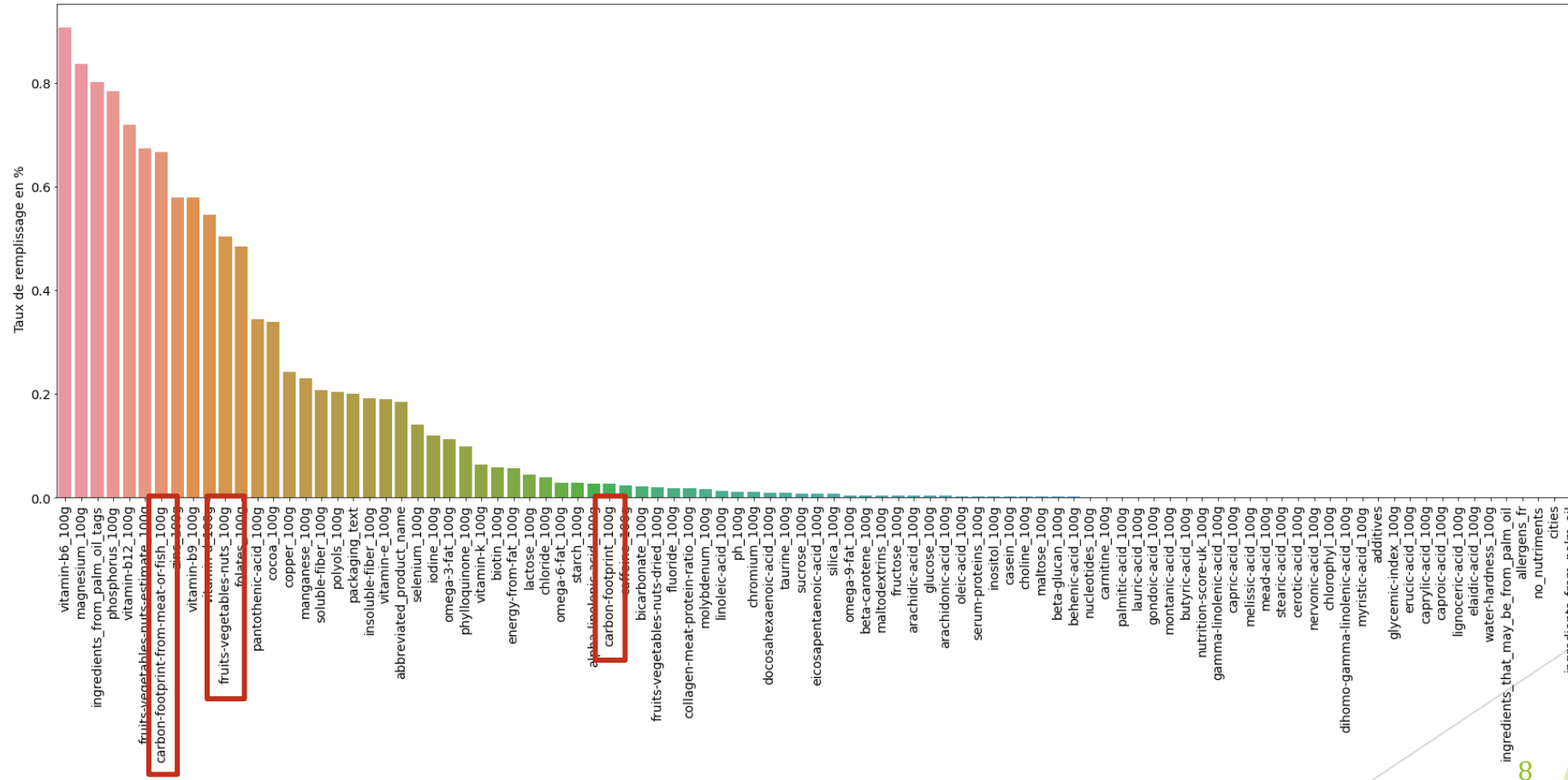


# Nettoyage



# Nettoyage

Taux de remplissage du jeu de données





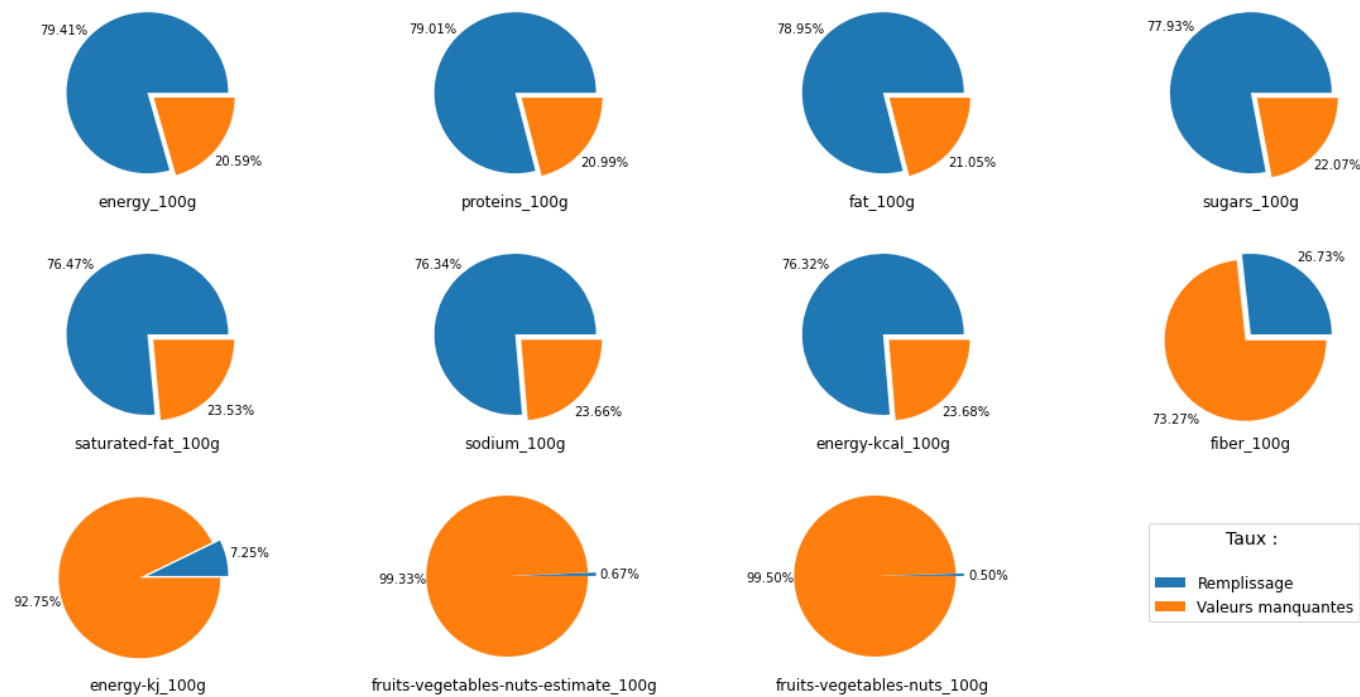
# Nettoyage

Qu'est-ce qui va nous intéresser ?



Valeurs nutritionnelles et nutriscore

Taux de remplissage et de valeurs manquantes

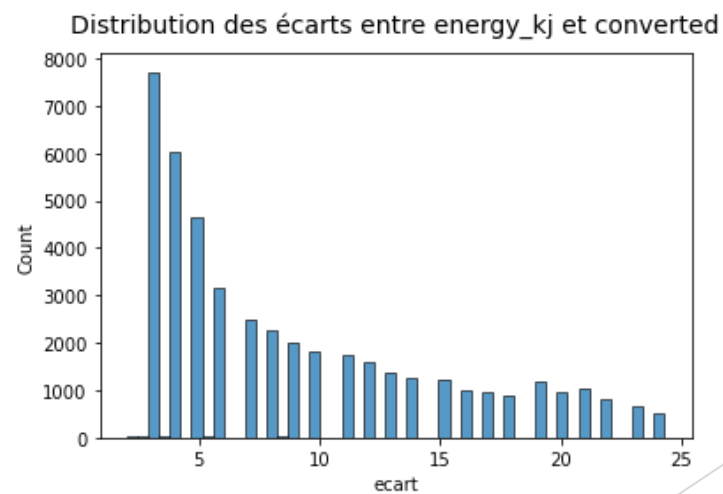
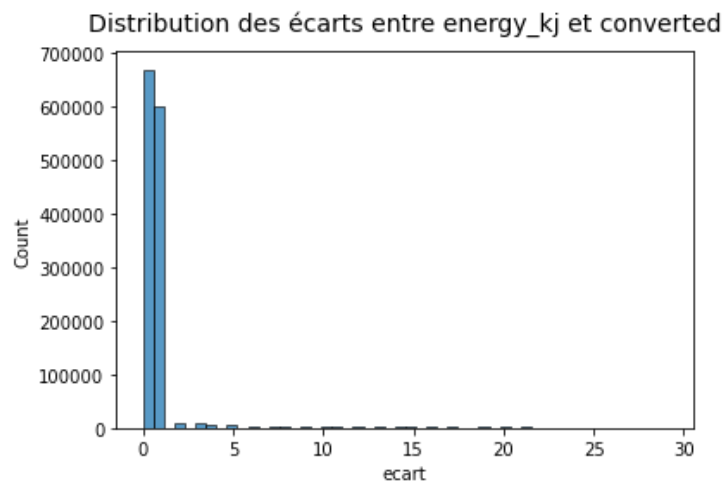


# Nettoyage

Colonnes redondantes

Ernergy

Kj/kcal/sans unité



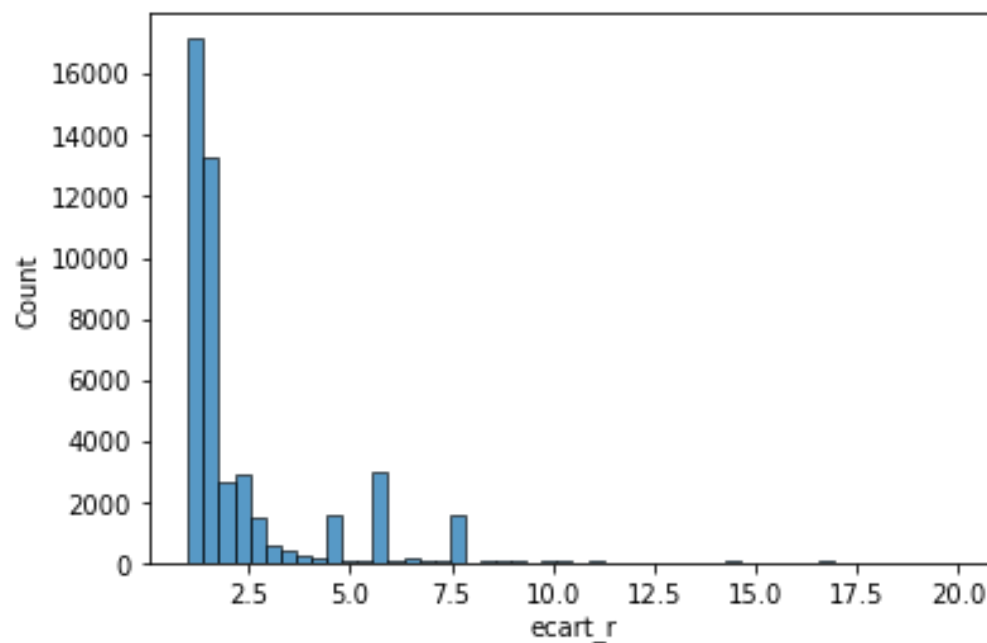
# Nettoyage

Colonnes redondantes

Ernergy

Kj/kcal/sans unité

Distribution des écarts relatifs entre energy\_kj et converted



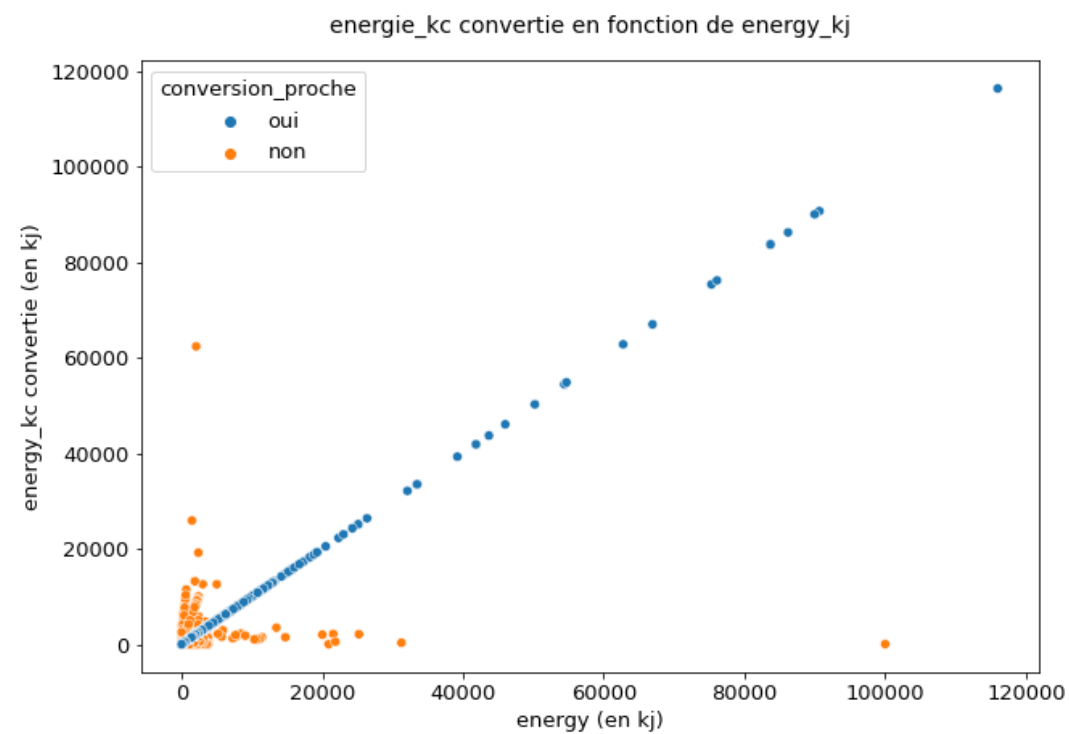
# Nettoyage

Colonnes redondantes

Ernergy

Kj/kcal/sans unité

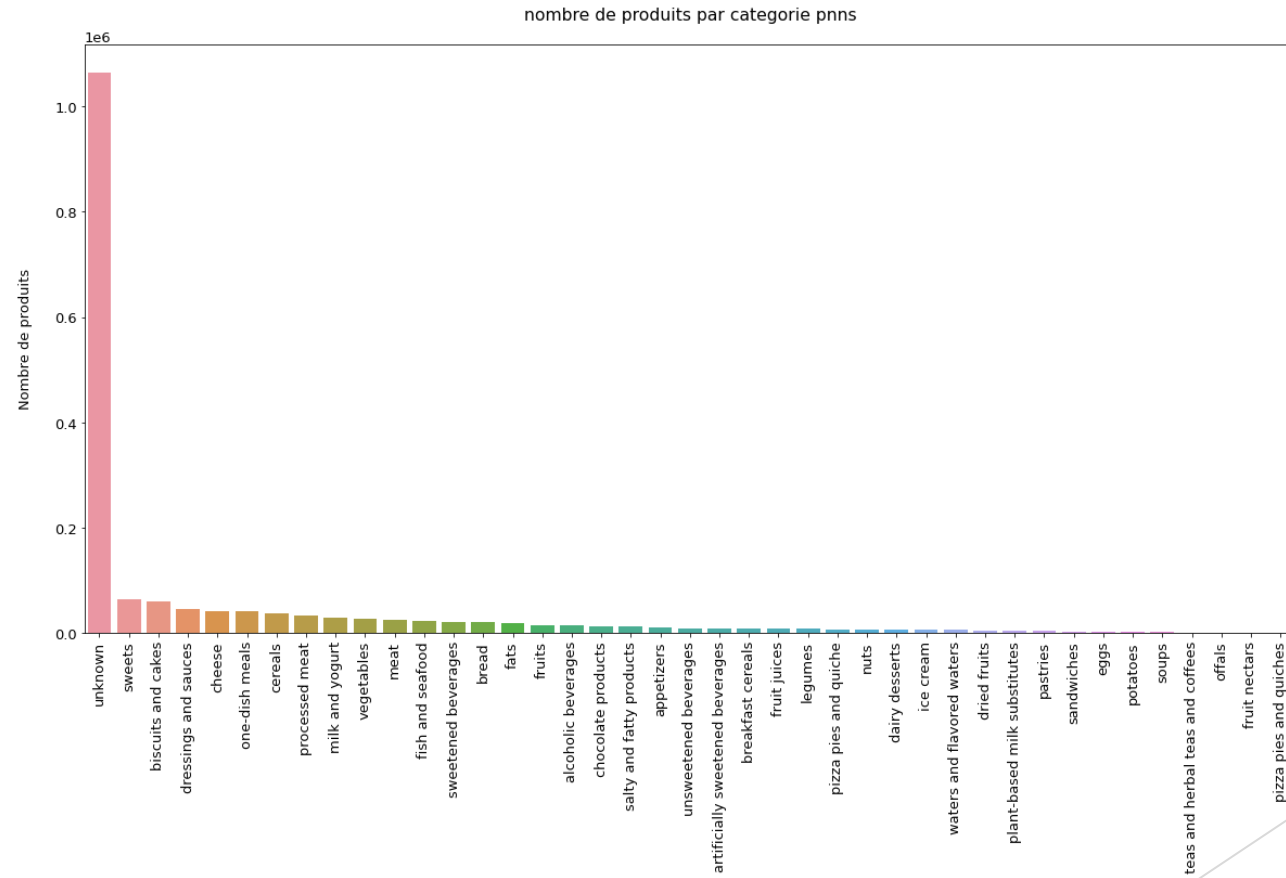
Elimination 3500 produits





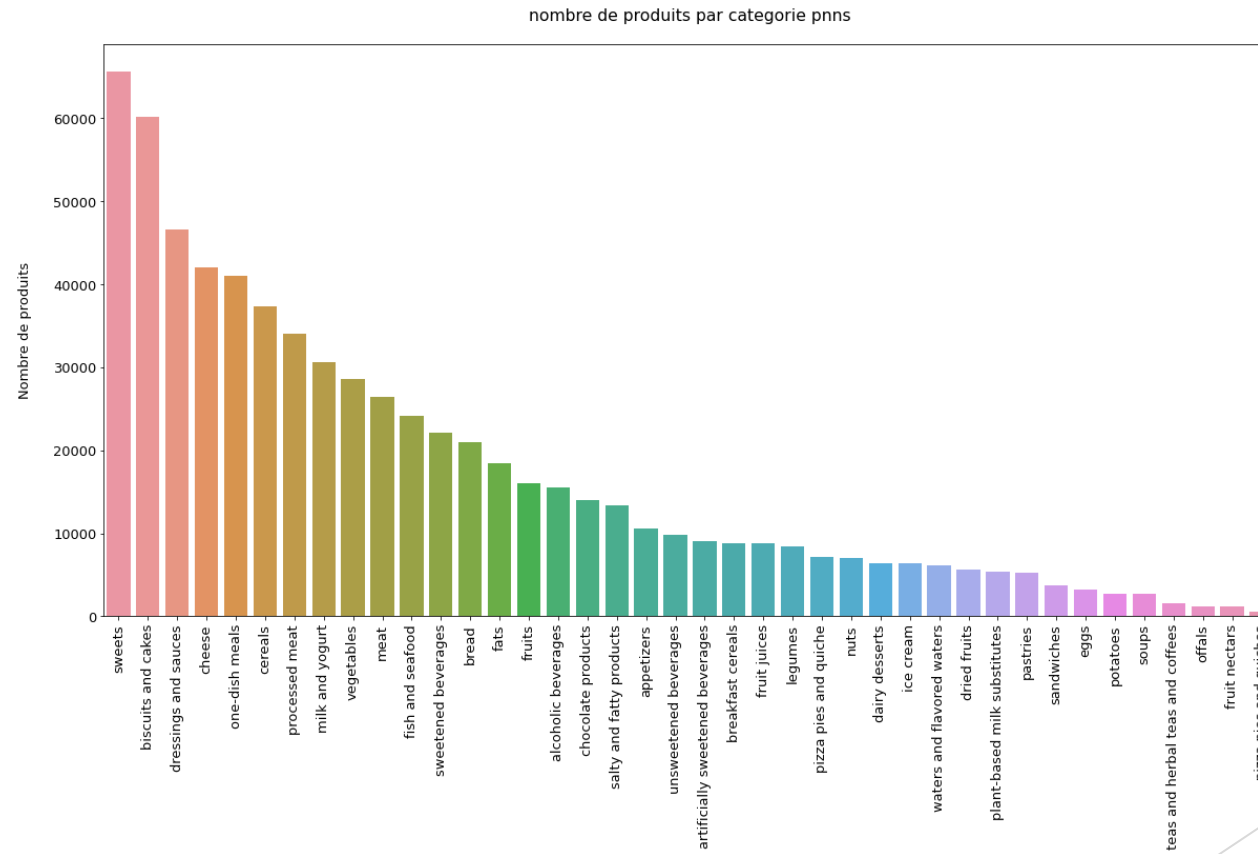
# Nettoyage

Colonnes redondantes  
Pnns Groups



# Nettoyage

Colonnes redondantes  
Pnns Groups



# Nettoyage

## Colonnes redondantes

Nutriscore (score/grade)

Pays ( Selections produits vendus en France)

France	792372	45.45%
Non France	951160	54.55%



# Nettoyage

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 789472 entries, 0 to 1743530
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   code                   789472 non-null object
1   product_name           767443 non-null object
2   brand                   418293 non-null object
3   category                353482 non-null object
4   info_nutri_complete    789472 non-null int64
5   pnns_group              789390 non-null object
6   nutriscore              275387 non-null float64
7   nutriscore_grade       275383 non-null object
8   energy_kj              617648 non-null float64
9   sugars                 613736 non-null float64
10  fat                    611579 non-null float64
11  saturated_fat          614549 non-null float64
12  sodium                 598321 non-null float64
13  fiber                  137888 non-null float64
14  proteins                613557 non-null float64
15  fruits_vegetables      16454 non-null  float64
dtypes: float64(9), int64(1), object(6)
memory usage: 102.4+ MB
```

DataSet à la fin du nettoyage

De 1 750 000 à 790 000 produits

De 184 à 16 colonnes

Conservations des informations  
relatives au nutriscore et  
catégories de produits

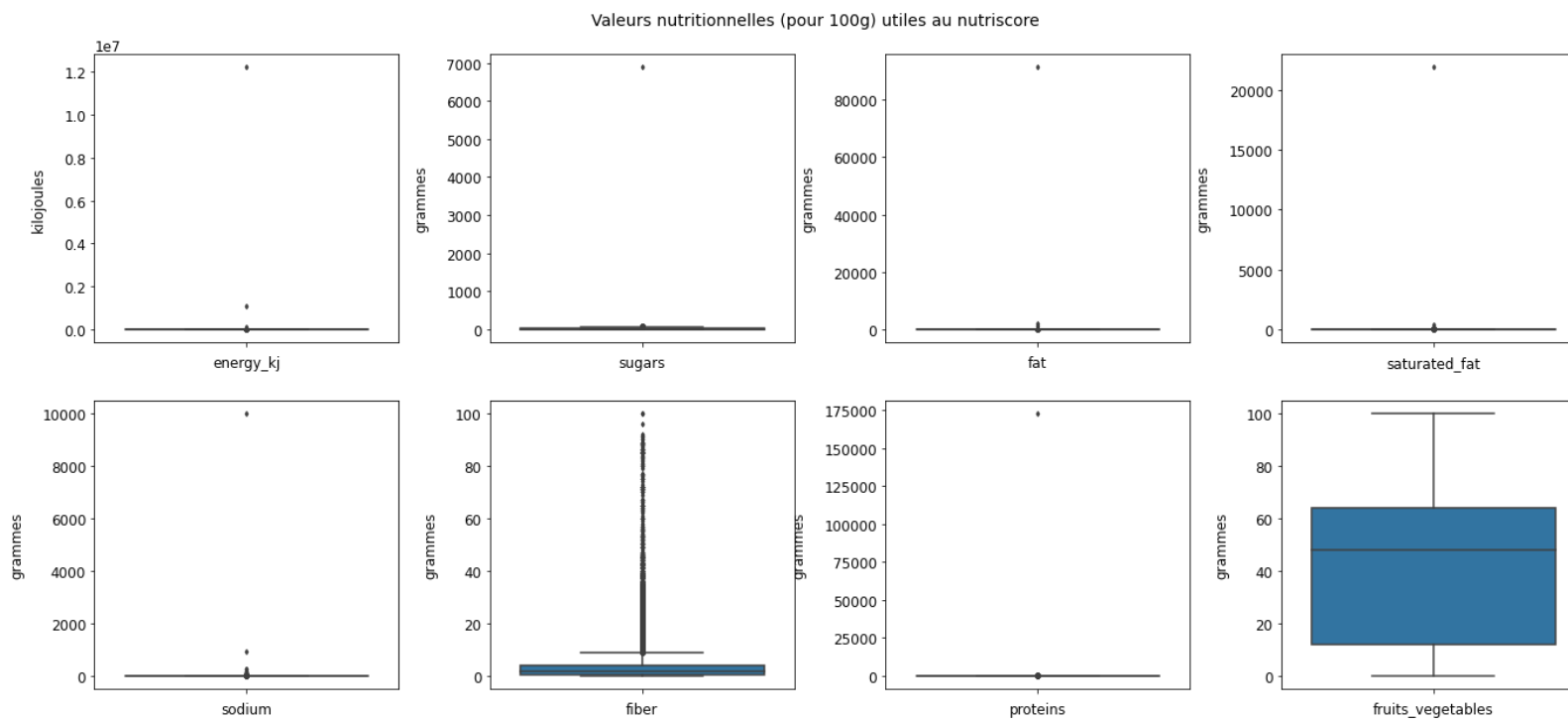
# Analyse univariée

# Analyse univariée

Détection derniers outliers sur 8 features



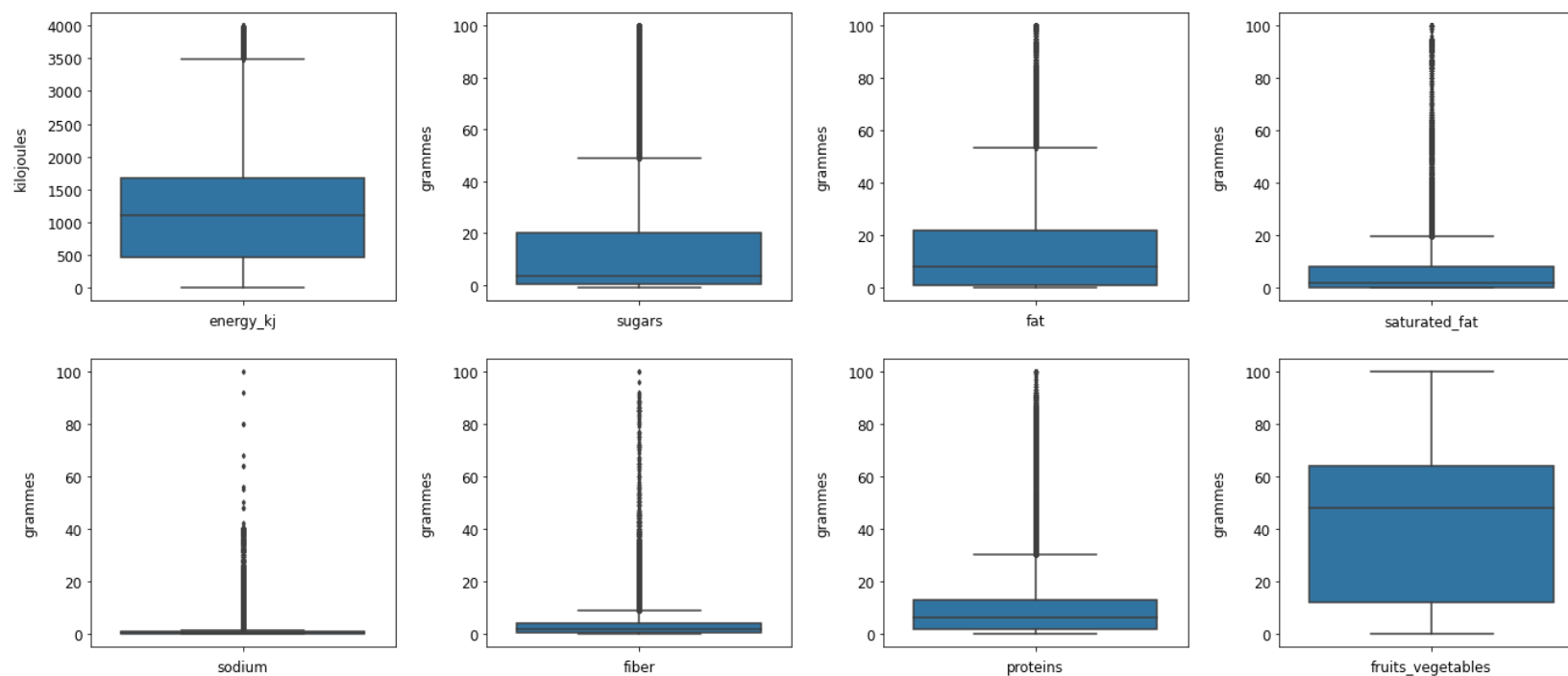
Elimination de 1 800 produits



# Analyse univariée

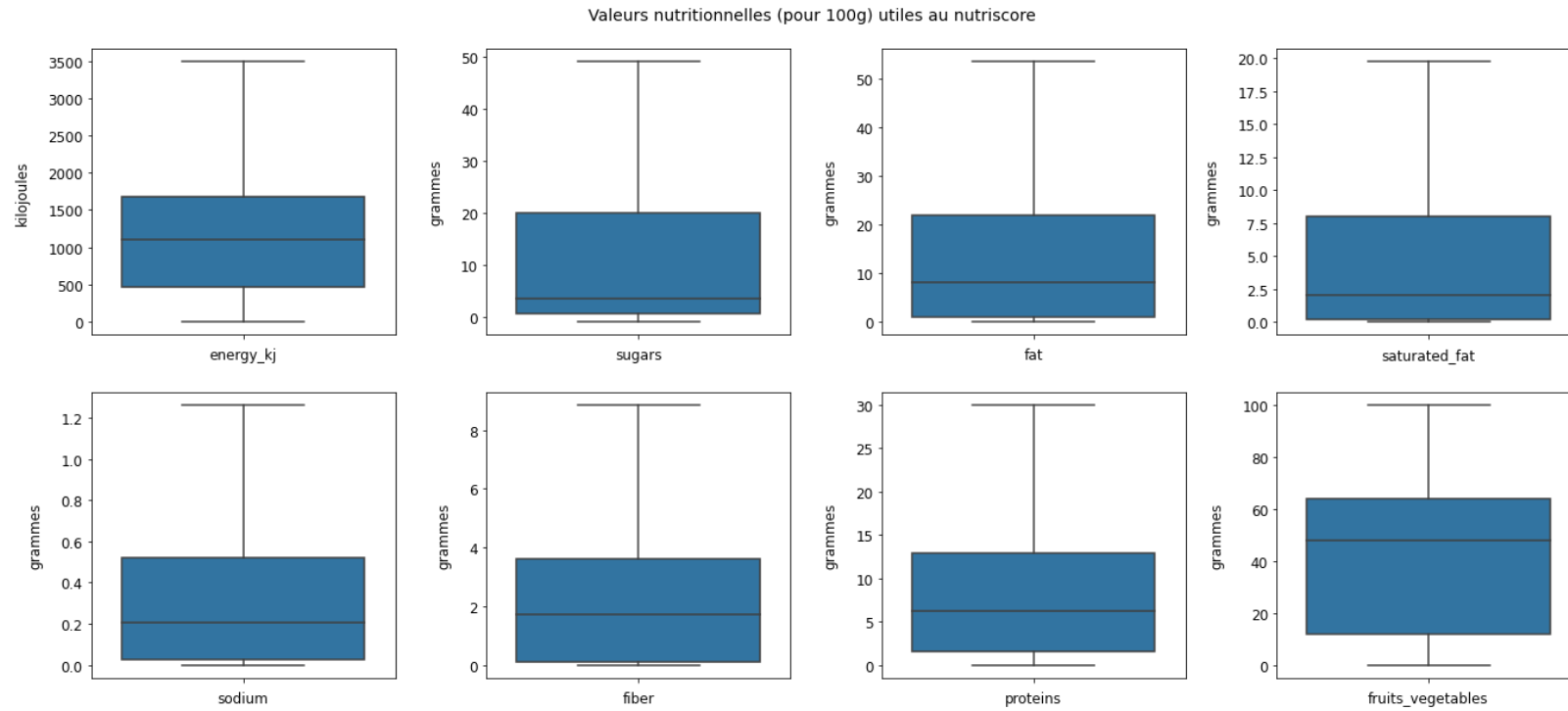
## Valeurs nutritionnelles

Valeurs nutritionnelles (pour 100g) utiles au nutriscore



# Analyse univariée

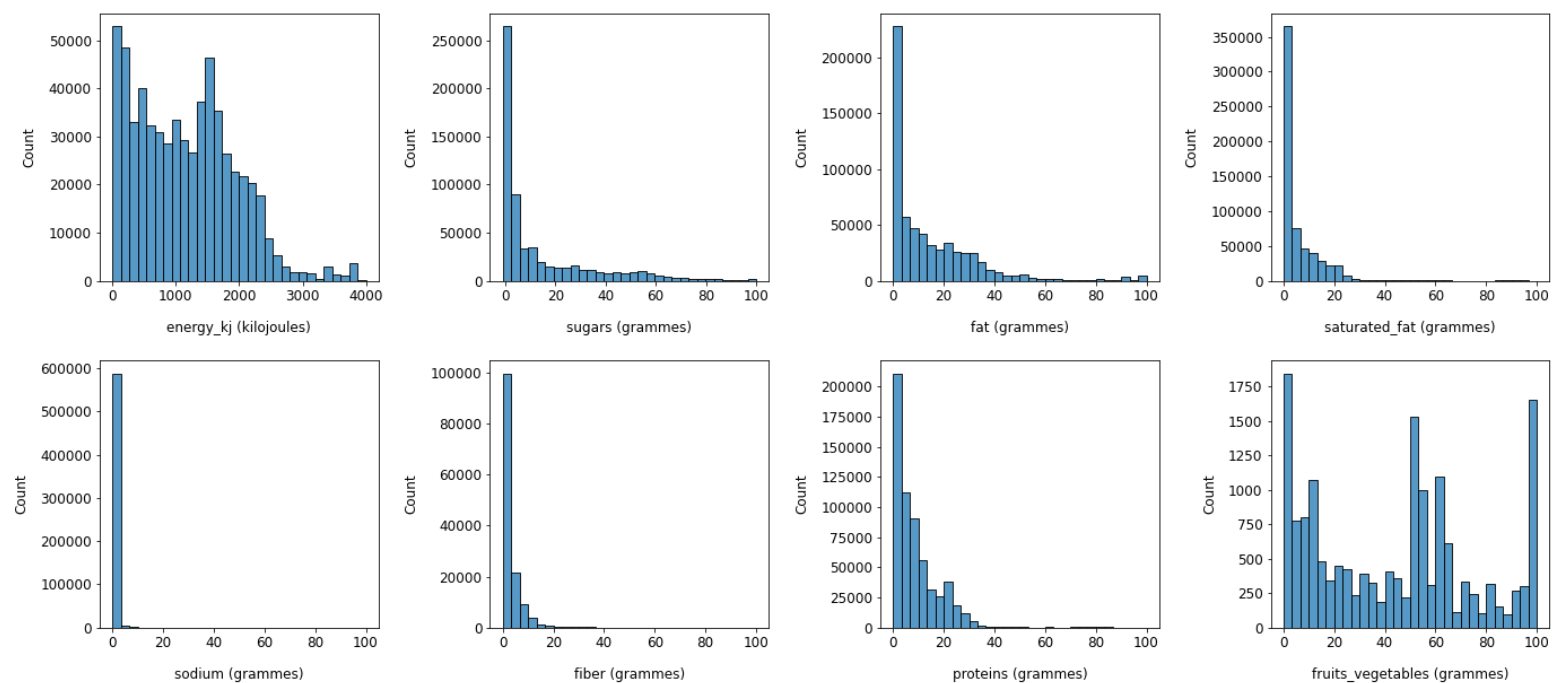
## Valeurs nutritionnelles



# Analyse univariée

## Valeurs nutritionnelles

Distribution des valeurs nutritionnelles (pour 100g) utiles au nutriscore

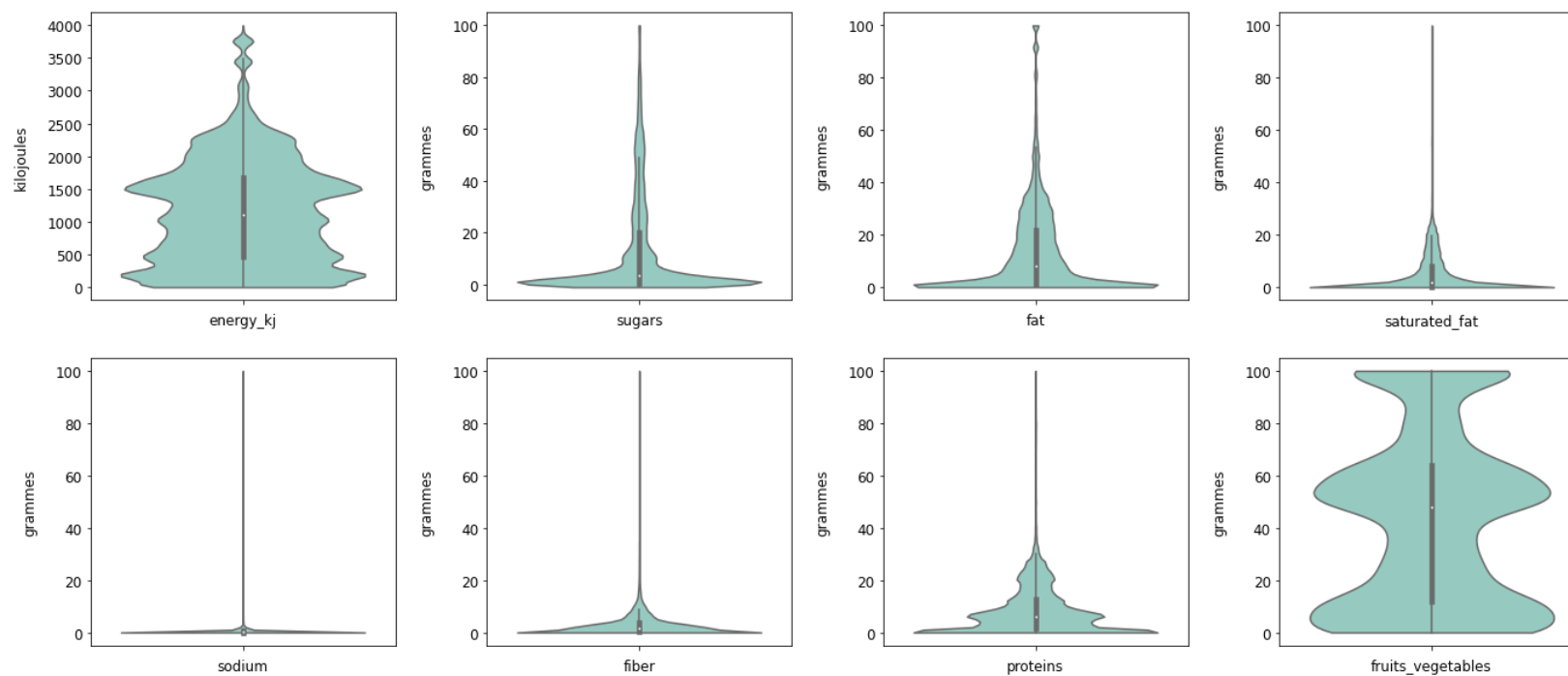


3 nodes > possibilité de discretiser

# Analyse univariée

## Valeurs nutritionnelles

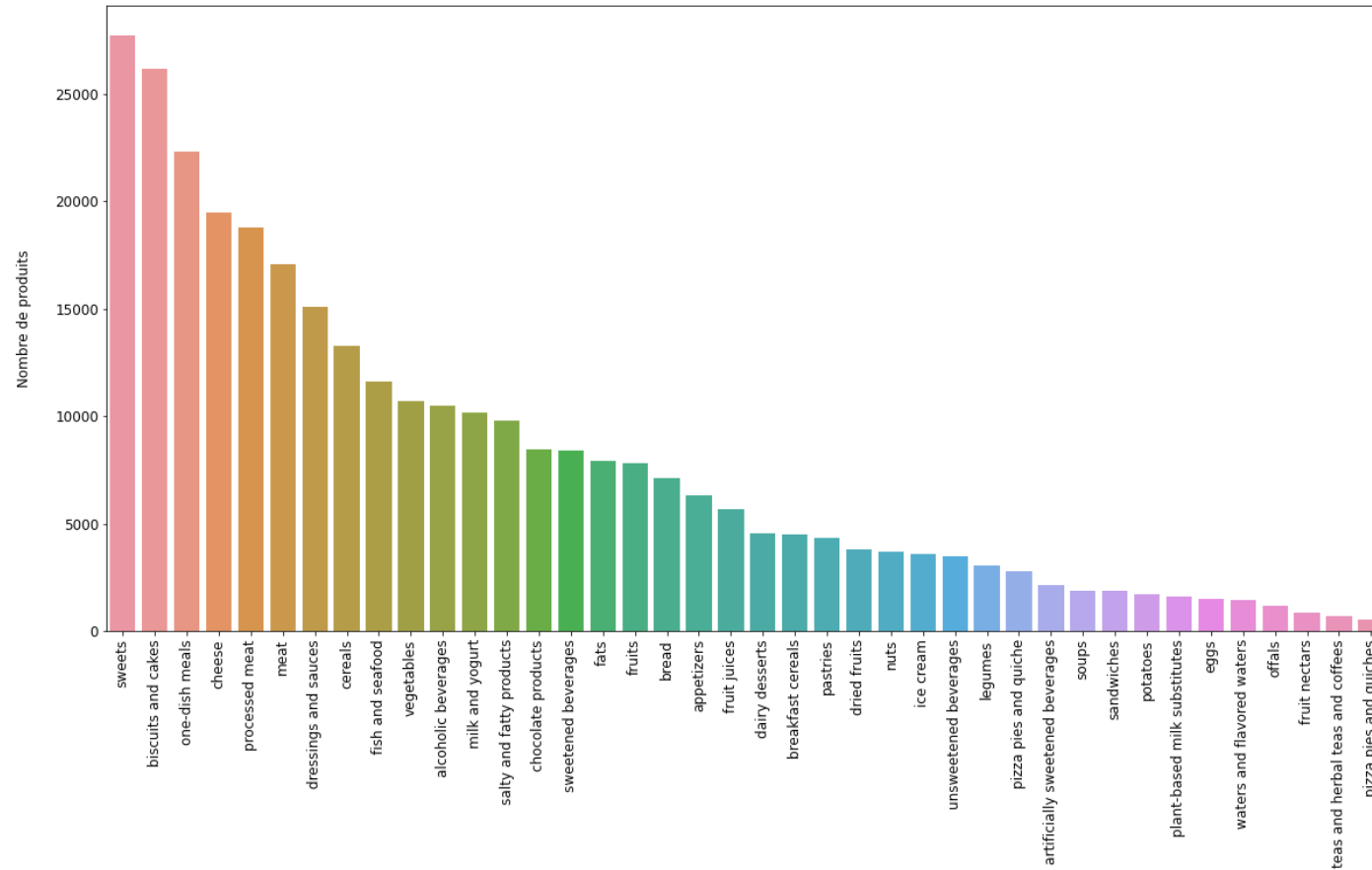
Valeurs nutritionnelles (pour 100g) utiles au nutriscore



# Analyse univariée

## Catégories de produits

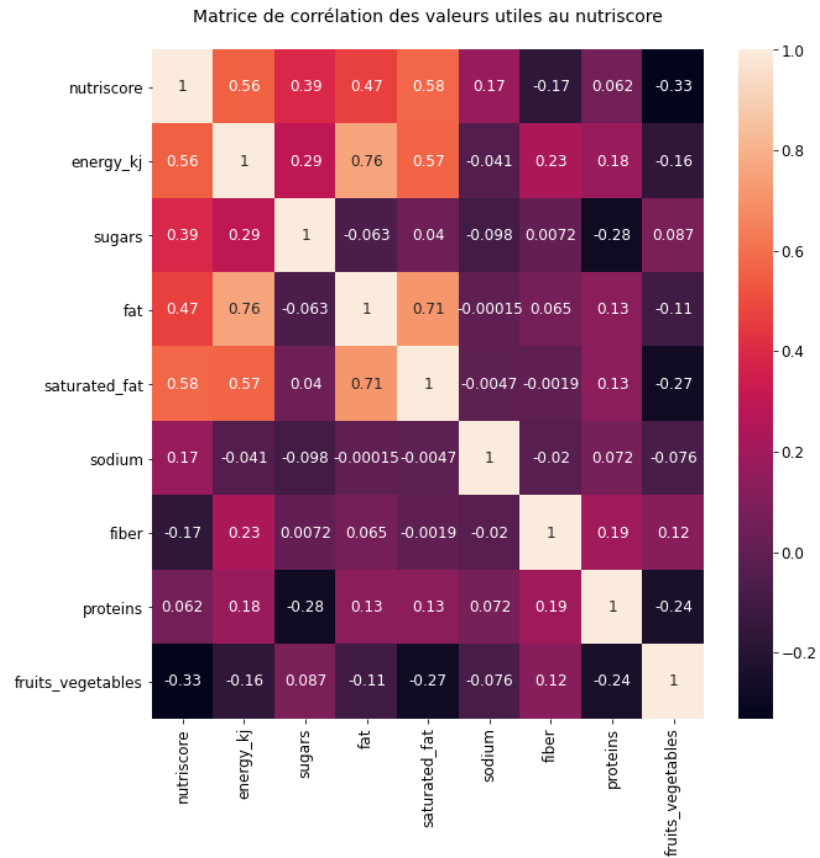
nombre de produits par categorie pnns





# Analyse multivariée

# Analyse multivariée



## Corrélation

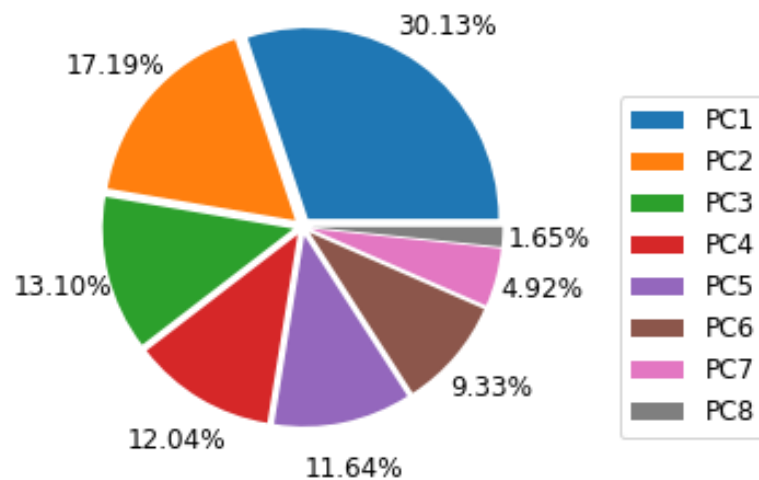
Limitation aux produits "données complètes" ~ 600 000

Remplissage par 0

Positive/Négative  
Mauvais/bon pour nutriscore

# Analyse multivariée

Ratios d'explication de la variance selon les composantes principales



## Reduction dimensionnelle

8 features d'entrée

Réduction à 6 pc pour 93%  
d'explication de variation

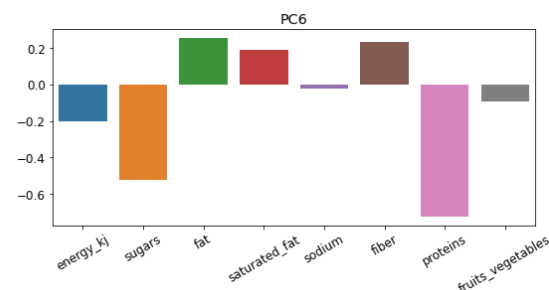
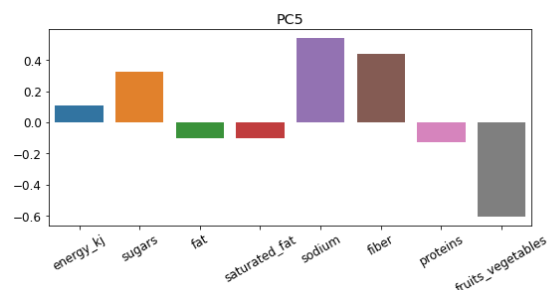
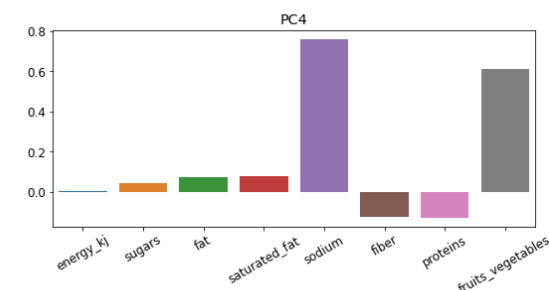
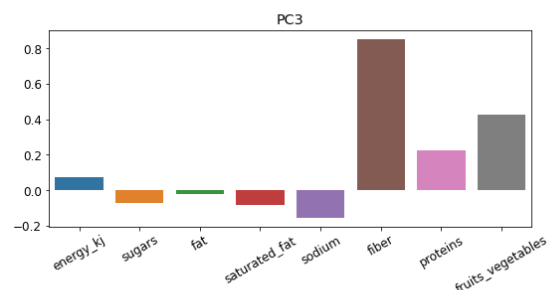
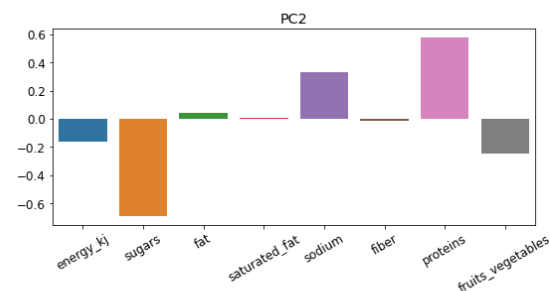
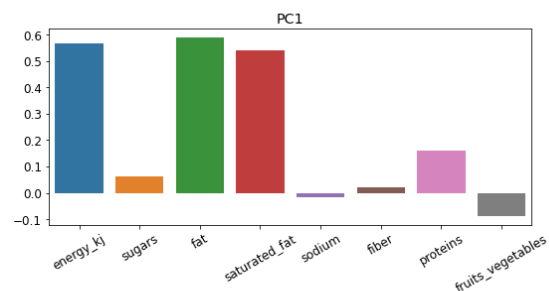
ou

Réduction à 5 pc pour 84%  
d'explication de variation

# Analyse multivariée

## Composantes nouvelles dimensions

Features des composantes principales (PC)



Pc 1

Energy / Fat : corrélation 0.76

Pc 2

Sucre / Prot : opposés

Pc 3

Fibre / Prot / Fruits : amélioration nutriscore

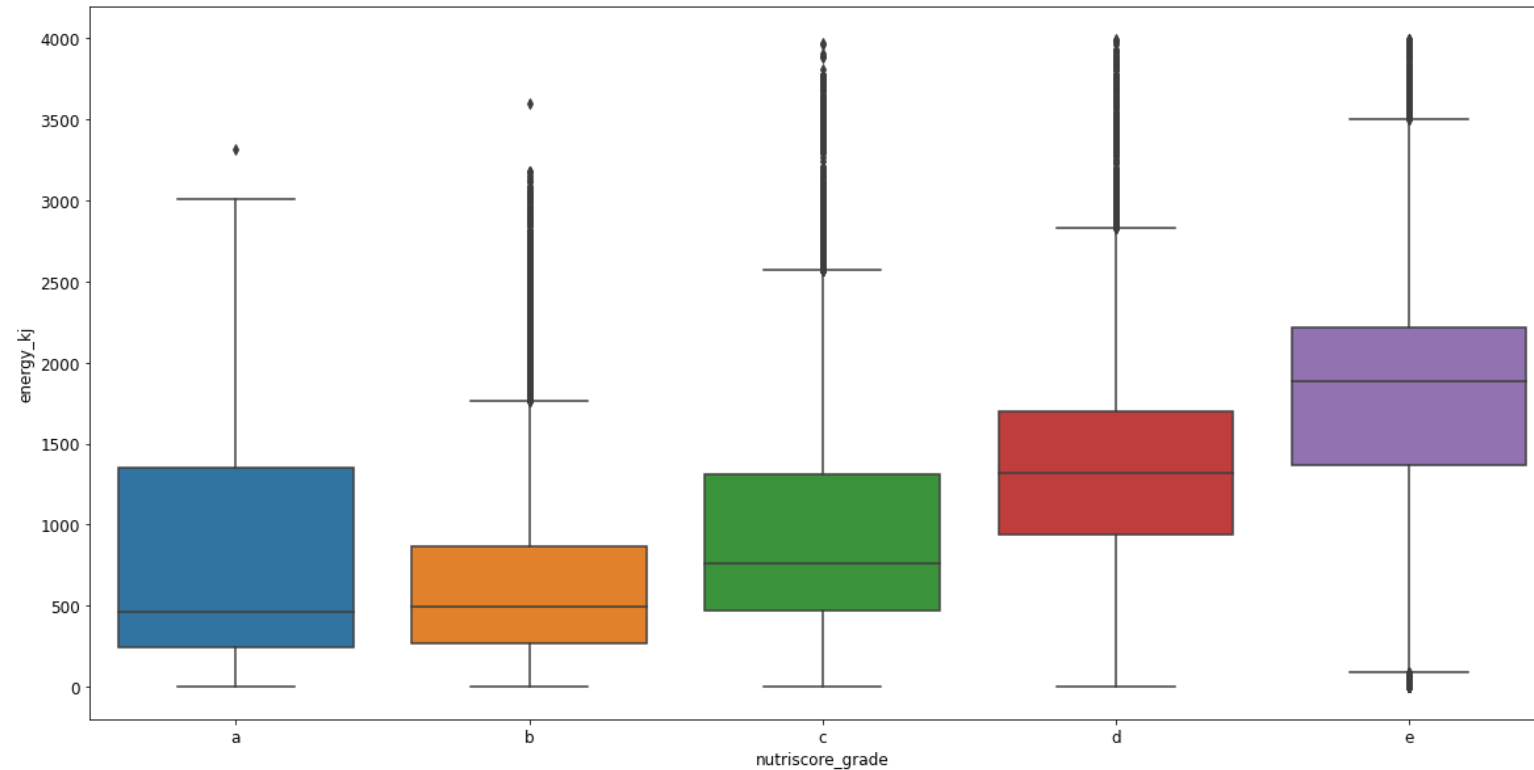
# Analyse multivariée

ANOVA : Comparaison des groupes



Limite aux produits qui ont un nutriscore : 275 000

Repartition des valeurs de energy\_kj selon les groupes formés par nutriscore\_grade



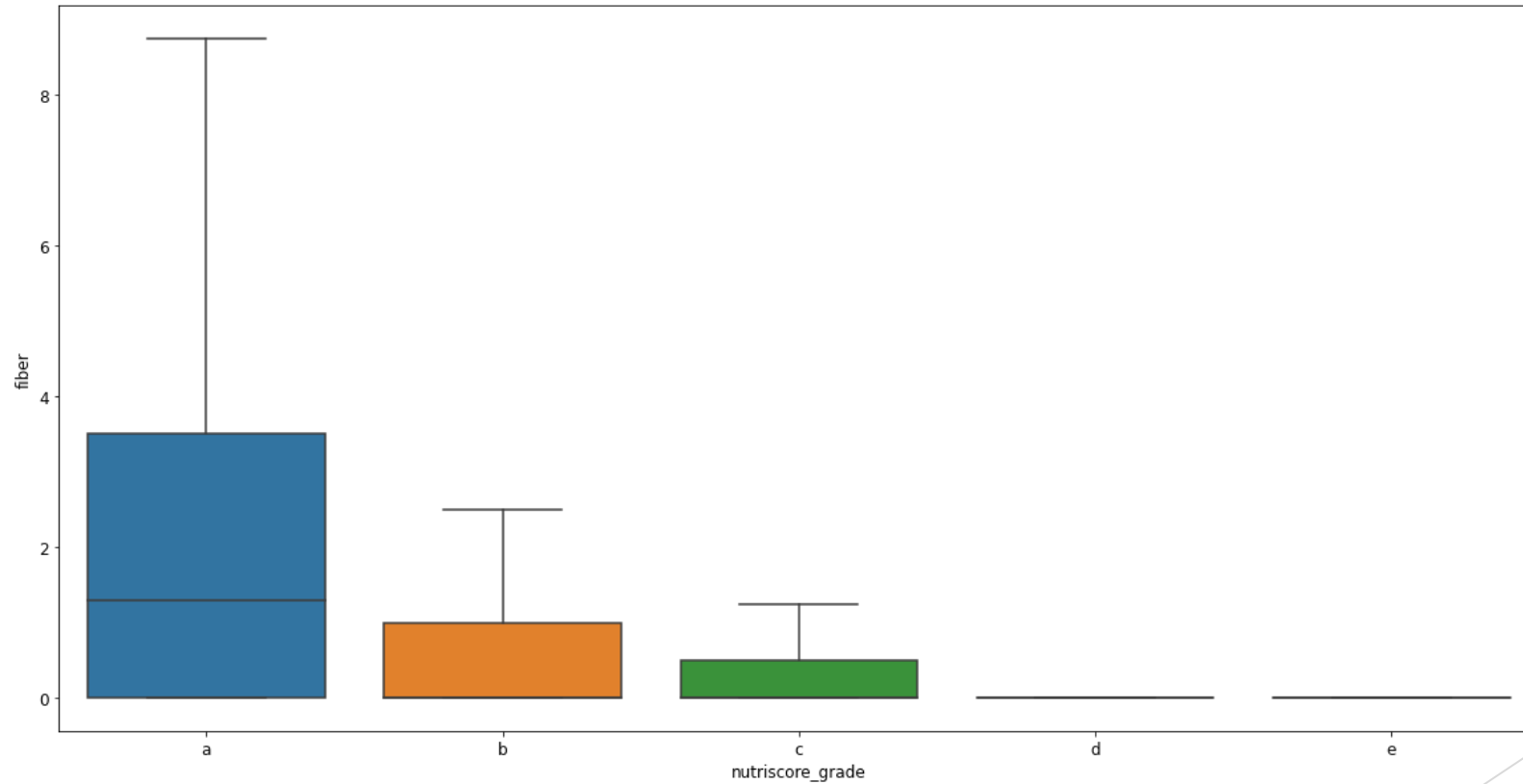
# Analyse multivariée

ANOVA : Comparaison des groupes



Limite aux produits qui ont un nutriscore : 275 000

Repartition des valeurs de fiber selon les groupes formés par nutriscore\_grade



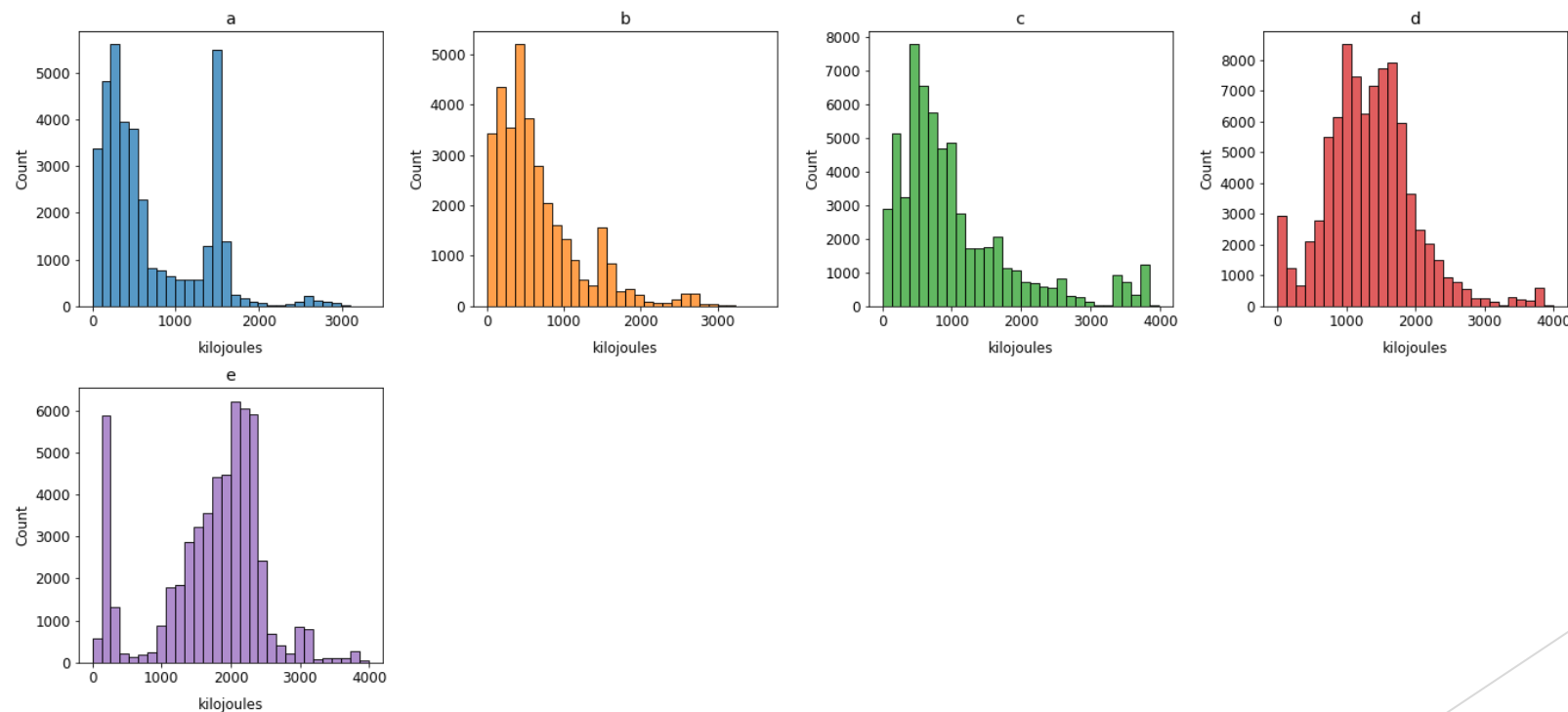
# Analyse multivariée

ANOVA : Comparaison des groupes



Limite aux produits qui ont un nutriscore : 275 000

distribution de energy\_kj en fonction des groupes formés par nutriscore\_grade



# Analyse multivariée

ANOVA : Comparaison des groupes



Limite aux produits qui ont un nutriscore : 275 000

Variable	Resultats du test
Energy_kj	F = 18 900 / Pvalue < 0.05
Sugars	F = 8 900 / Pvalue < 0.05
Fat	F = 14 300 / Pvalue < 0.05
Saturated_fat	F = 26 300 / Pvalue < 0.05
Sodium	F = 2 160 / Pvalue < 0.05
Fibers	F = 4 360 / Pvalue < 0.05
Proteins	F = 1 720 / Pvalue < 0.05
Fruits	F = 650 / Pvalue < 0.05

Groupes par nutriscore



# Analyse multivariée

ANOVA : Comparaison des groupes



Limite aux produits qui ont un pnns-group : 270 000

Variable	Resultats du test
Energy_kj	F = 12 500 / Pvalue < 0.05
Sugars	F = 13 400 / Pvalue < 0.05
Fat	F = 12 900 / Pvalue < 0.05
Saturated_fat	F = 6 700 / Pvalue < 0.05
Sodium	F = 1 400 / Pvalue < 0.05
Fibers	F = 1 200 / Pvalue < 0.05
Proteins	F = 12 400 / Pvalue < 0.05
Fruits	F = 550 / Pvalue < 0.05

Groupes par pnns-groups

The background features abstract, overlapping green geometric shapes, primarily triangles and polygons, in various shades of green, creating a modern and dynamic visual effect.

Pertinence/Faisabilité

# Pertinence / Faisabilité

## ANOVA :

Des groupes plutôt bien séparés, on peut donc estimer le groupe en fonction des variables retenues

## Modélisation KNN :

Algorithme de classification utilisé

Prédiction correcte du nutriscore à ~ 73%

Prédiction correcte de catégorie à ~ 69%

## Qualité du dataset:

Nombre de valeurs et de produits « erronées »

Beaucoup d'informations manquantes