

# **Projet 3 : Préparez des données pour un organisme de santé publique**

**11/06/2024**

**Soukaina GUAOUA ELJADDI**

**Parcours Data Scientist  
OpenClassrooms**

# Plan:

- ❑ Problématique
- ❑ Nettoyage et filtration des features et produits
- ❑ Identification et traitement des valeurs aberrantes
- ❑ Identification et traitement des valeurs manquantes
- ❑ Analyses uni-variée et bi-variée
- ❑ Analyse multi-variée
- ❑ Conclusion

# Problématique

**Contexte** : Projet d'amélioration de la base de données **Open Food Facts** de l'agence **santé publique france**.

**Objectif** : Création d'un système de **suggestion** ou d'**auto-complétion** pour aider les usagers à remplir plus efficacement la base de données.

**Mission** : **nettoyage** et **exploration** des données, afin de déterminer la faisabilité de l'idée d'application de Santé publique France.

# Étape 1 : Nettoyez et filtrez des features et produits

# Étape 1 : Nettoyez et filtrez des features et produits

## Analyse exploratoire du jeu de données:

[fr.openfoodfacts.org.products.csv](https://fr.openfoodfacts.org/products.csv) : 320772 observations, 162 variables

Fiche produit	Tags	Ingrédients et additifs	Informations nutritionnels
<ul style="list-style-type: none"><li>- code</li><li>- url</li><li>- creator</li><li>- created_t</li><li>- product_name</li></ul>	<ul style="list-style-type: none"><li>- packaging_tags</li><li>- brand_tags</li><li>- categories_tags</li><li>- origin_tags</li></ul>	<ul style="list-style-type: none"><li>- ingredients_text</li><li>- allergens</li><li>- additives</li></ul>	<ul style="list-style-type: none"><li>- fat_100g</li><li>- sugars_100g</li><li>- fiber_100g</li><li>- salt_100g</li></ul>

# Étape 1 : Nettoyez et filtrez des features et produits

## Identification de la cible et des features pertinentes:

Données originales : 320772 lignes, 162 colonnes

Cible = 'nutrition\_grade\_fr'

Features pertinentes : Features avec plus de 50% de valeurs présentes ⇒ 43 colonnes

Suppression des produits en double ⇒ 221210 lignes

# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

Suppression des colonnes redondantes: 'additives\_n', 'additives', 'additives\_tags'

'additives\_fr' plus informative  $\Rightarrow$  221210 lignes  $\times$  40 colonnes

### Exemple :

product_name	Peanut Butter Power Chews
additives_n	3.0
additives	[ peanut-butter -> en:peanut-butter ] [ but...
additives_tags	en:e170,en:e322,en:e410
additives_fr	E170 - Carbonate de calcium,E322 - Lécithines,...

# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

Suppression des colonnes Temps sous format `_t` (UNIX timestamp format): '`created_t`' et '`last_modified_t`'

`_datetime` (format ISO 8601) : '`created_datetime`',  
'`last_modified_datetime`' ⇒ 221210 lignes × 38 colonnes

### Exemple:

product_name	Peanut Butter Power Chews
<code>created_t</code>	1489138486
<code>created_datetime</code>	2017-03-10T09:34:46Z
<code>last_modified_t</code>	1489138486
<code>last_modified_datetime</code>	2017-03-10T09:34:46Z



# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

Suppression des colonnes redondantes \_tags : 'brands\_tags', 'countries\_tags', 'states\_tags'

'brands', 'countries', 'states' ⇒ 221210 lignes × 35 colonnes

### Exemple:

product_name	Peanut Butter Power Chews
brands_tags	sunridge
brands	Sunridge
countries_tags	en:united-states
countries	US
states_tags	en:to-be-completed,en:nutrition-facts-complete...
states	en:to-be-completed, en:nutrition-facts-complet...

# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

Suppression de la colonne redondante : 'nutrition-score-uk\_100g'

'nutrition-score-fr\_100g'  $\Rightarrow$  221210 lignes  $\times$  34 colonnes

### Exemple :

product_name	Peanut Butter Power Chews
nutrition-score-uk_100g	9.0
nutrition-score-fr_100g	9.0

# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

Suppression des colonnes redondantes : 'countries', 'states'

'countries\_fr' et 'states\_fr'  $\Rightarrow$  221210 lignes  $\times$  32 colonnes

### Exemple :

product_name	Peanut Butter Power Chews
countries	US
countries_fr	États-Unis
states	en:to-be-completed, en:nutrition-facts-complet...
states_fr	A compléter, Informations nutritionnelles compl...

# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

Suppression des colonnes inutiles pour l'analyse : 'code', 'url', 'creator', 'created\_datetime', 'last\_modified\_datetime', 'brands', 'origins', 'countries\_fr', 'ingredients\_text', 'serving\_size', 'additives\_fr', 'states\_fr'

⇒ 221210 lignes × 21 colonnes

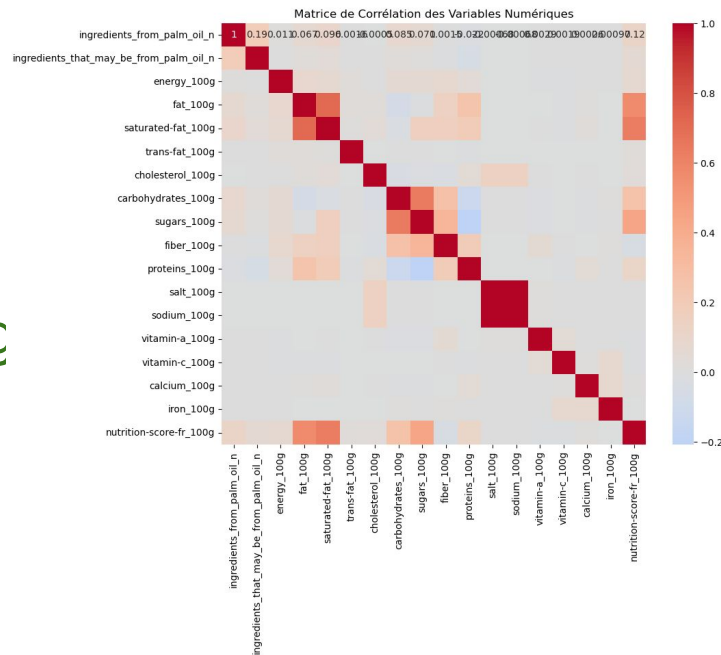
# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

Suppression d'une des 2 variables corrélées

- 'sodium\_100g' == 'salt\_100g'
- 'sugars\_100g' == 'carbohydrates\_100g'
- 'fat\_100g' == 'saturated-fat\_100g'

⇒ 221210 rows × 18 columns



# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

Suppression des 2 colonnes mal renseignés :

'ingredients\_from\_palm\_oil\_n' : 0 et NaN

et 'ingredients\_that\_may\_be\_from\_palm\_oil\_n'

⇒ 221210 lignes × 16 colonnes

# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

Suppression de la colonne : 'nutrition-score-fr\_100g' correspondante à 'nutrition\_grade\_fr'

⇒ 221210 lignes × 15 colonnes

<b>nutrition_grade</b>	<b>nutrition-score</b>
<b>A</b>	<b>-15 et -2</b>
<b>B</b>	<b>-1 à +3</b>
<b>C</b>	<b>+4 à +11</b>
<b>D</b>	<b>+12 à +16</b>
<b>E</b>	<b>+17 à +40</b>

Source :Wikipedia

# Étape 1 : Nettoyez et filtrez des features et produits

## Filtrage et nettoyage des données :

- Suppression des lignes avec '**variables \_100g**' > **100**, à l'exception du variable energy\_100g (en kj )
- Remplacement des valeurs négatives par 0  $\Rightarrow$  **221133** lignes  $\times$  **15** colonnes

### 12 variables numériques

energy\_100g, saturated-fat\_100g,  
trans-fat\_100g, cholesterol\_100g,  
carbohydrates\_100g, fiber\_100g,  
proteins\_100g, salt\_100g,  
vitamin-a\_100g, vitamin-c\_100g,  
calcium\_100g, iron\_100g.

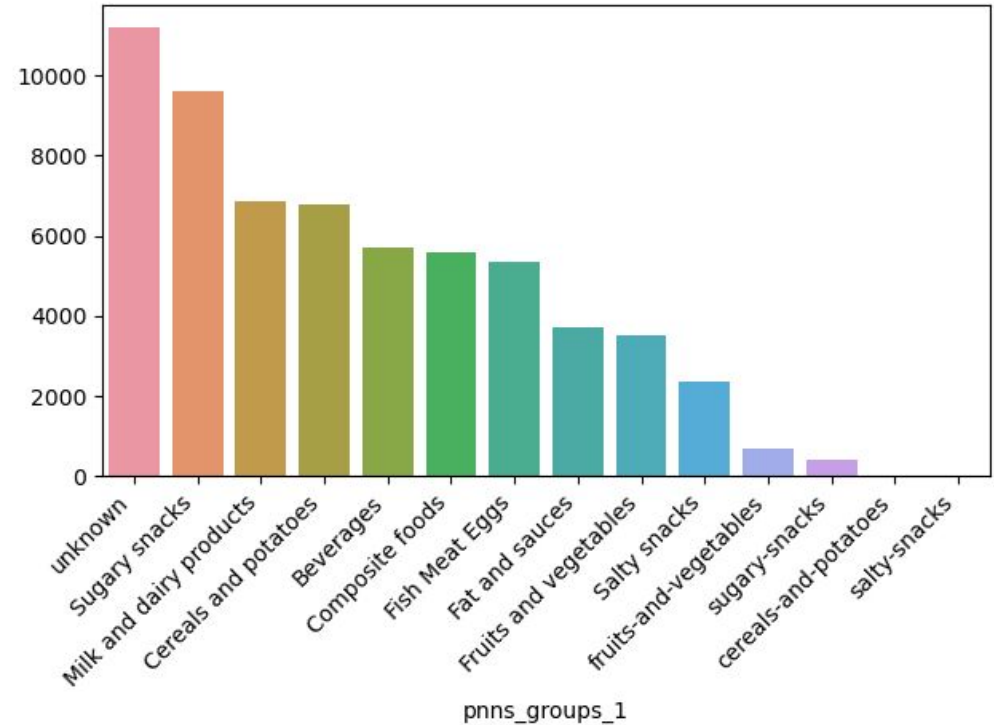
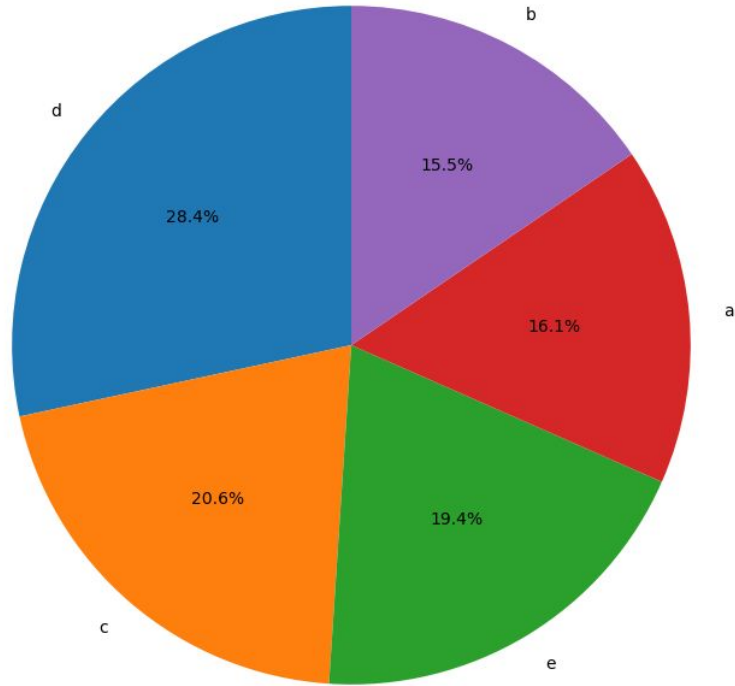
### 3 catégorielles

product\_name,  
pnns\_groups\_1,  
**nutrition\_grade\_fr** (cible)



# Étape 1 : Nettoyez et filtrez des features et produits

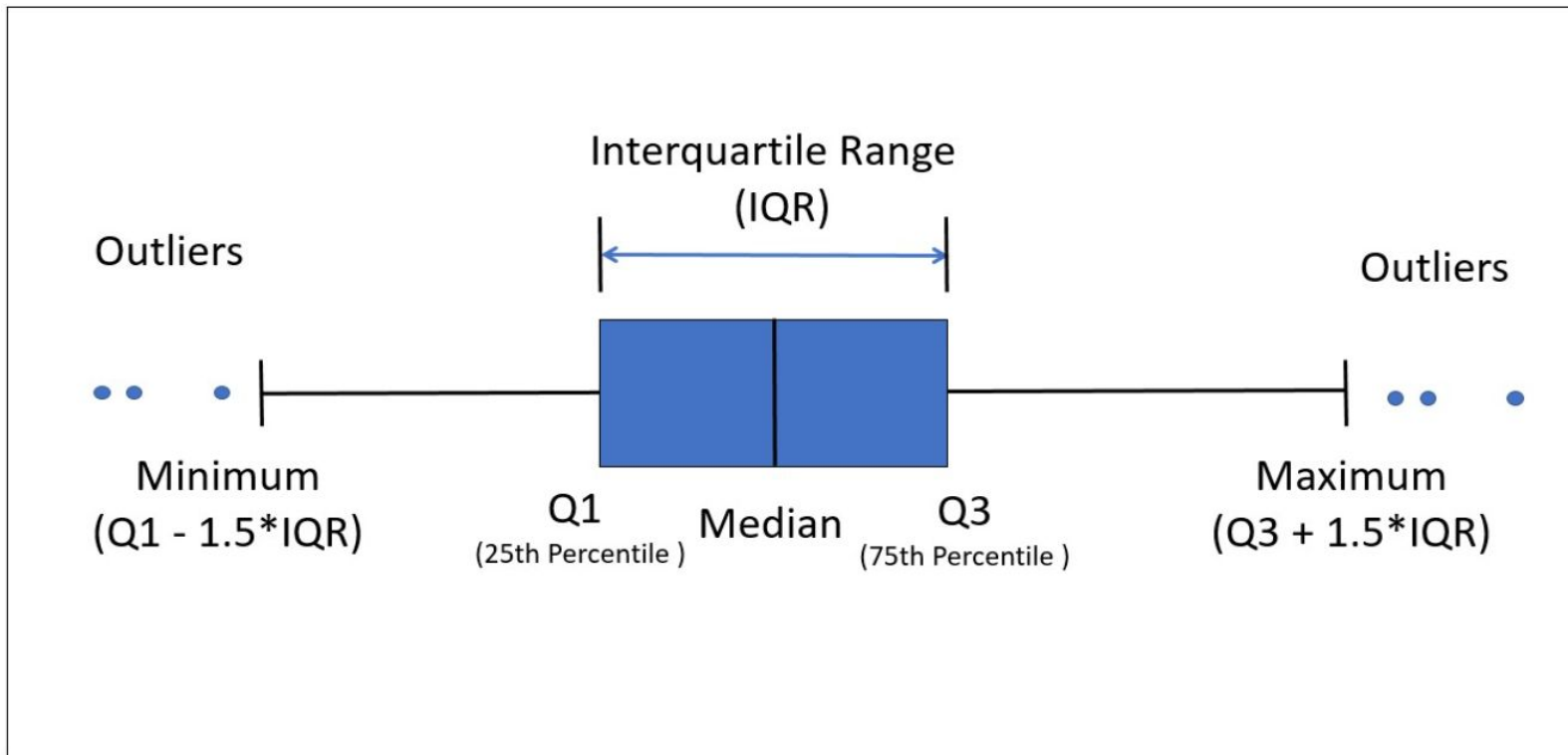
Répartition des catégories de 'nutrition\_grade\_fr'



# **Étape 2 : Identifiez et traitez les valeurs aberrantes**

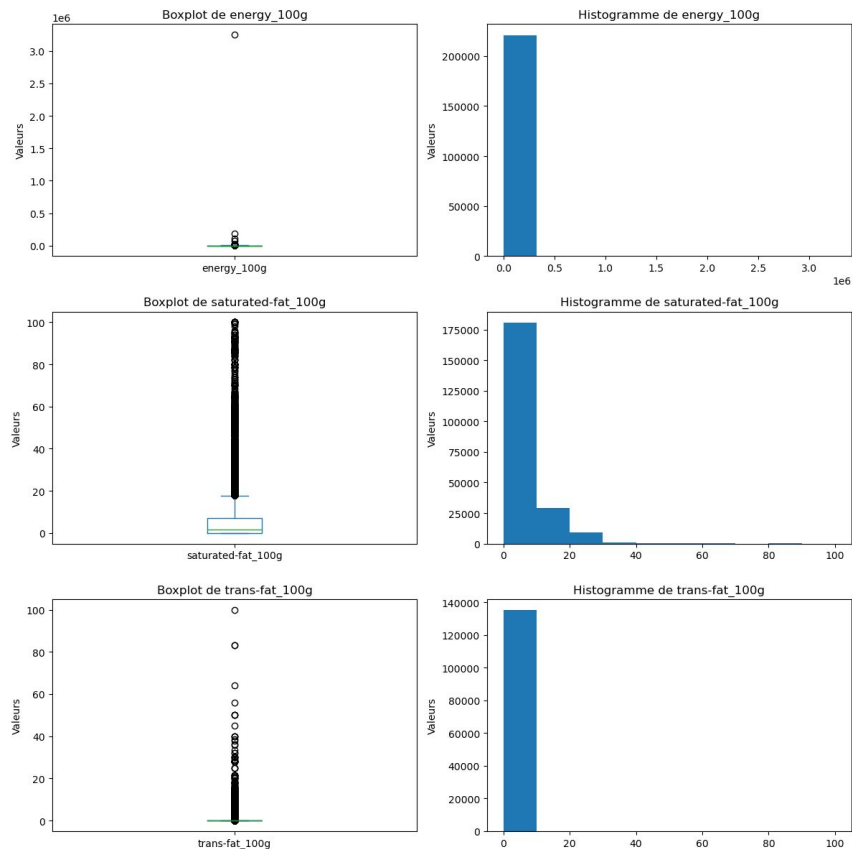
# Étape 2 : Identifiez et traitez les valeurs aberrantes

## Méthode des plages interquartiles

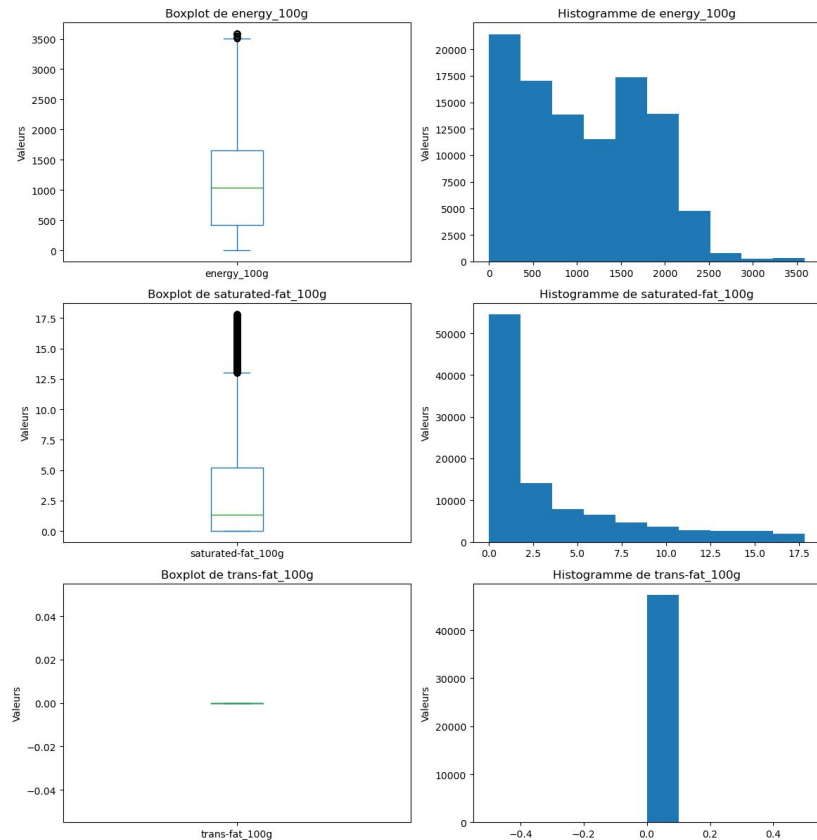


# Étape 2 : Identifiez et traitez les valeurs aberrantes

Avant suppression des valeurs aberrantes

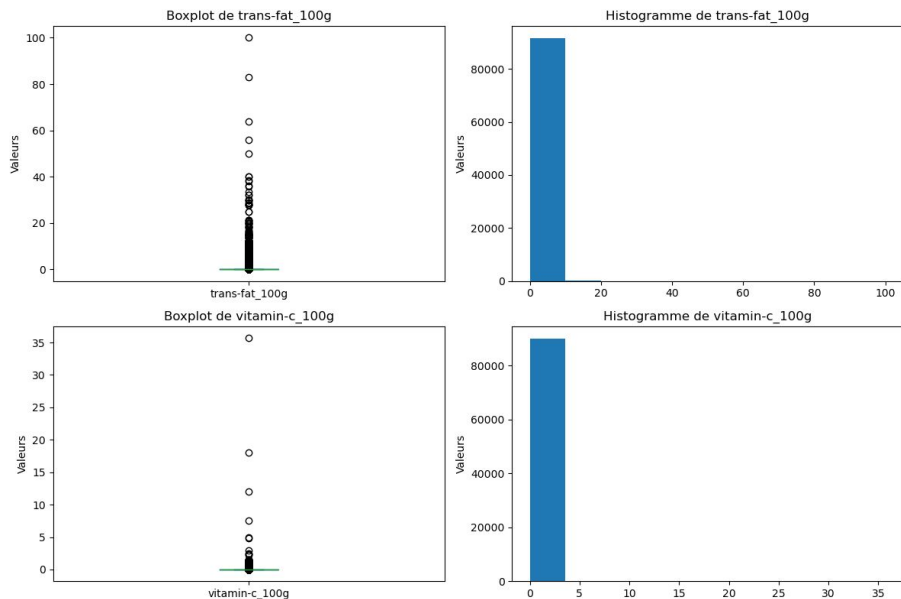


Après suppression des valeurs aberrantes

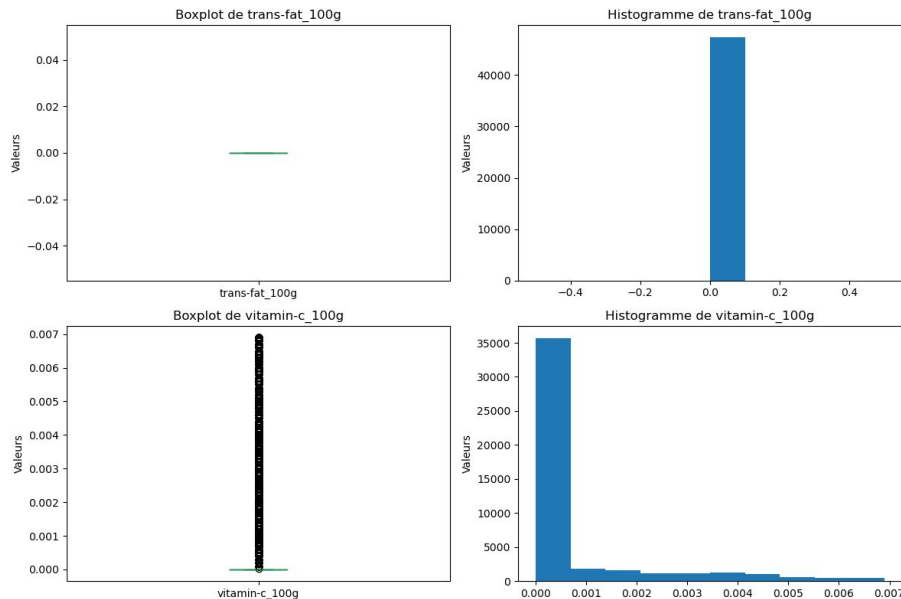


# Étape 2 : Identifiez et traitez les valeurs aberrantes

## Avant suppression des valeurs aberrantes



## Après suppression des valeurs aberrantes

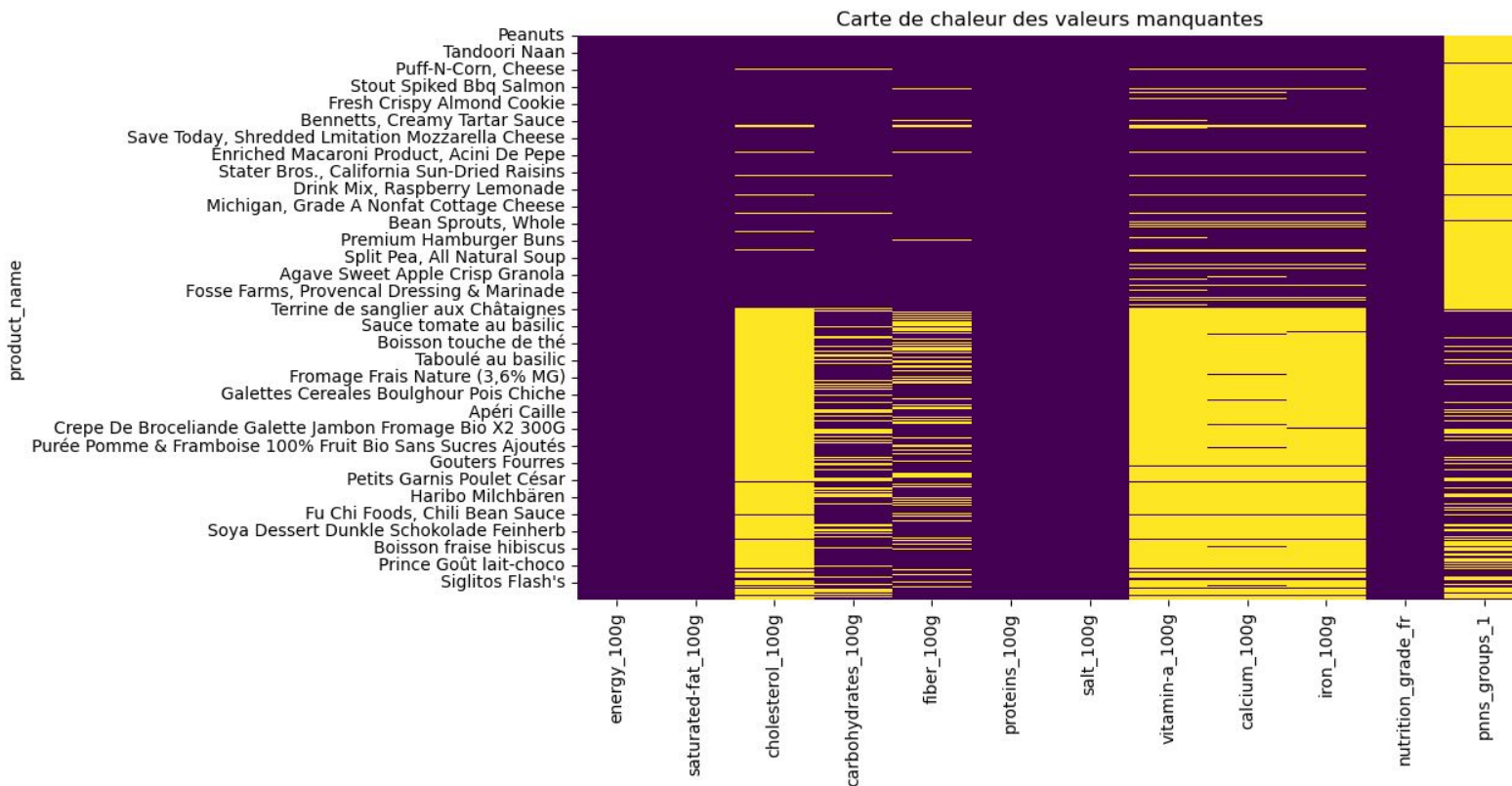


⇒ Suppression de la variable nulle 'trans-fat\_100g', et la variable 'vitamin-c\_100g' presque nulle avec que des valeurs aberrantes ⇒ 10 variables restantes

# **Étape 3 : Identifiez et traitez les valeurs manquantes**

# Étape 3 : Identifiez et traitez les valeurs manquantes

## Identification des valeurs manquantes



# Étape 3 : Identifiez et traitez les valeurs manquantes

## Traitement des valeurs manquantes

### Méthode 1

Après suppression des valeurs aberrantes : 101444 lignes, 12 colonnes

Suppression des lignes avec que des NaN  $\Rightarrow$  101281 lignes  $\times$  12 colonnes

### Méthode 2

Suppression des lignes avec que des 0  $\Rightarrow$  101080 lignes  $\times$  12 colonnes

### Méthode 3

Imputation des valeurs manquantes qui ont une catégorie pnns\_groups\_1 avec la méthode KNNI: 101080 lignes  $\times$  12 colonnes

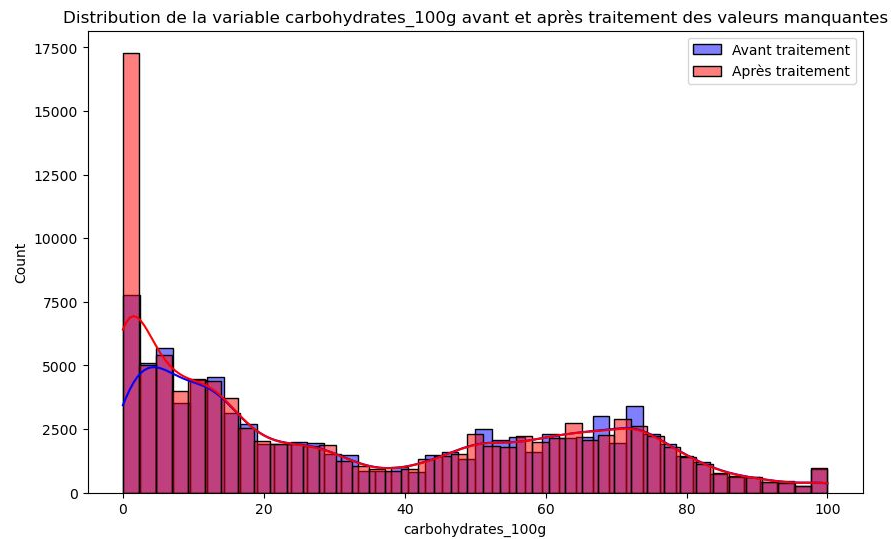
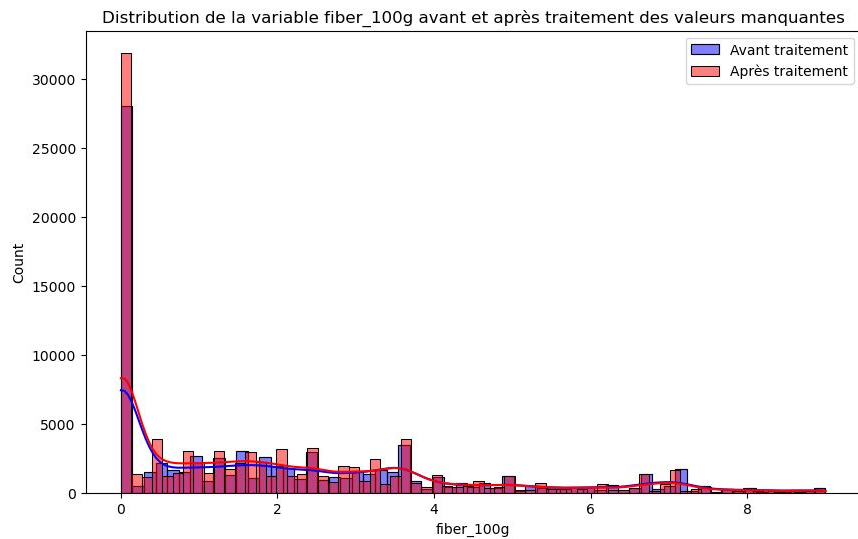
### Méthode 4

Remplissage des valeurs manquantes restantes par 0  $\Rightarrow$  101080 lignes  $\times$  12 colonnes



# Étape 3 : Identifiez et traitez les valeurs manquantes

## Visualisation des variables avant et après traitement des valeurs manquantes



# Étape 4 : Effectuez les analyses uni-variée et bi-variée

# Étape 4 : Effectuez les analyses uni-variée et bi-variée

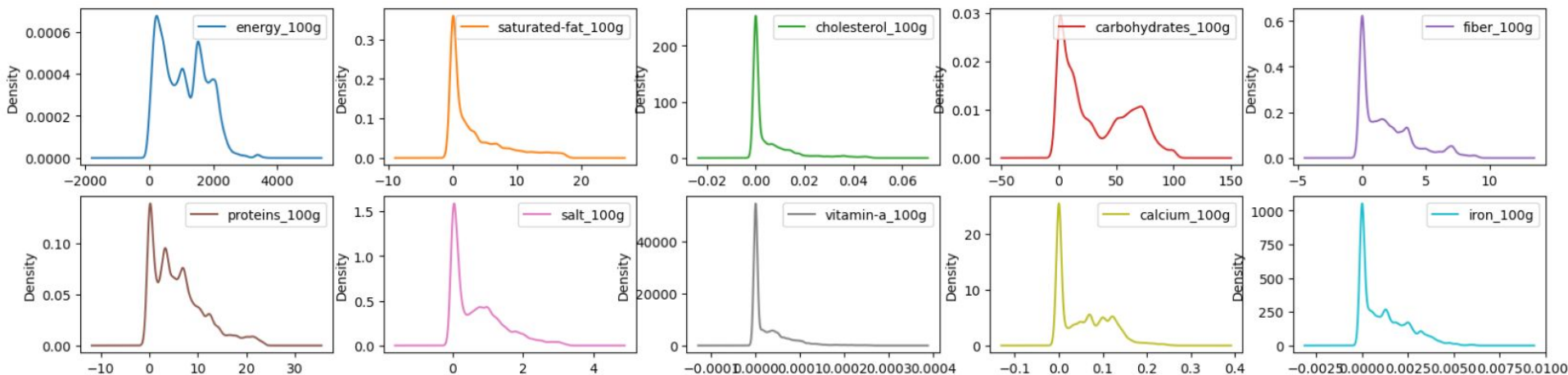
## Analyse uni-variée : statistiques descriptives

	energy_100g	saturated-fat_100g	cholesterol_100g	carbohydrates_100g	fiber_100g	proteins_100g	salt_100g	vitamin-a_100g	calcium_100g	iron_100g
count	101080.000000	101080.000000	101080.000000	101080.000000	101080.000000	101080.000000	101080.000000	101080.000000	101080.000000	101080.000000
mean	1088.438248	3.385542	0.005141	32.909335	1.931544	6.028504	0.707040	0.000025	0.057136	0.001196
std	710.960931	4.455492	0.009197	29.384576	2.131785	5.327604	0.725666	0.000041	0.057313	0.001279
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	423.000000	0.000000	0.000000	5.830000	0.000000	1.670000	0.091440	0.000000	0.000000	0.000000
50%	1046.000000	1.400000	0.000000	23.000000	1.340000	5.000000	0.508000	0.000000	0.048000	0.000840
75%	1653.000000	5.200000	0.007000	60.470000	3.100000	8.850000	1.115060	0.000040	0.102445	0.002030
max	3590.000000	17.840000	0.047000	100.000000	9.000000	23.640000	3.251200	0.000259	0.261000	0.006260

# Étape 4 : Effectuez les analyses uni-variée et bi-variée

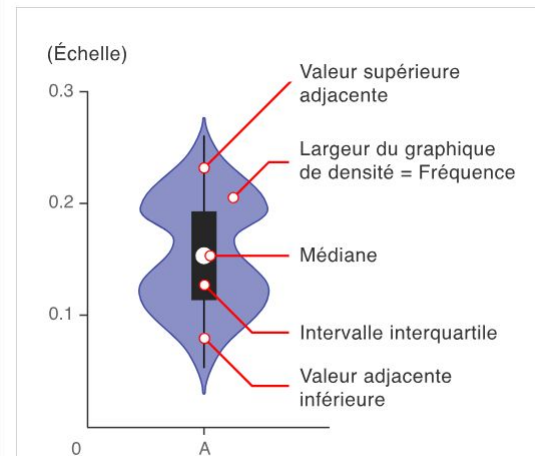
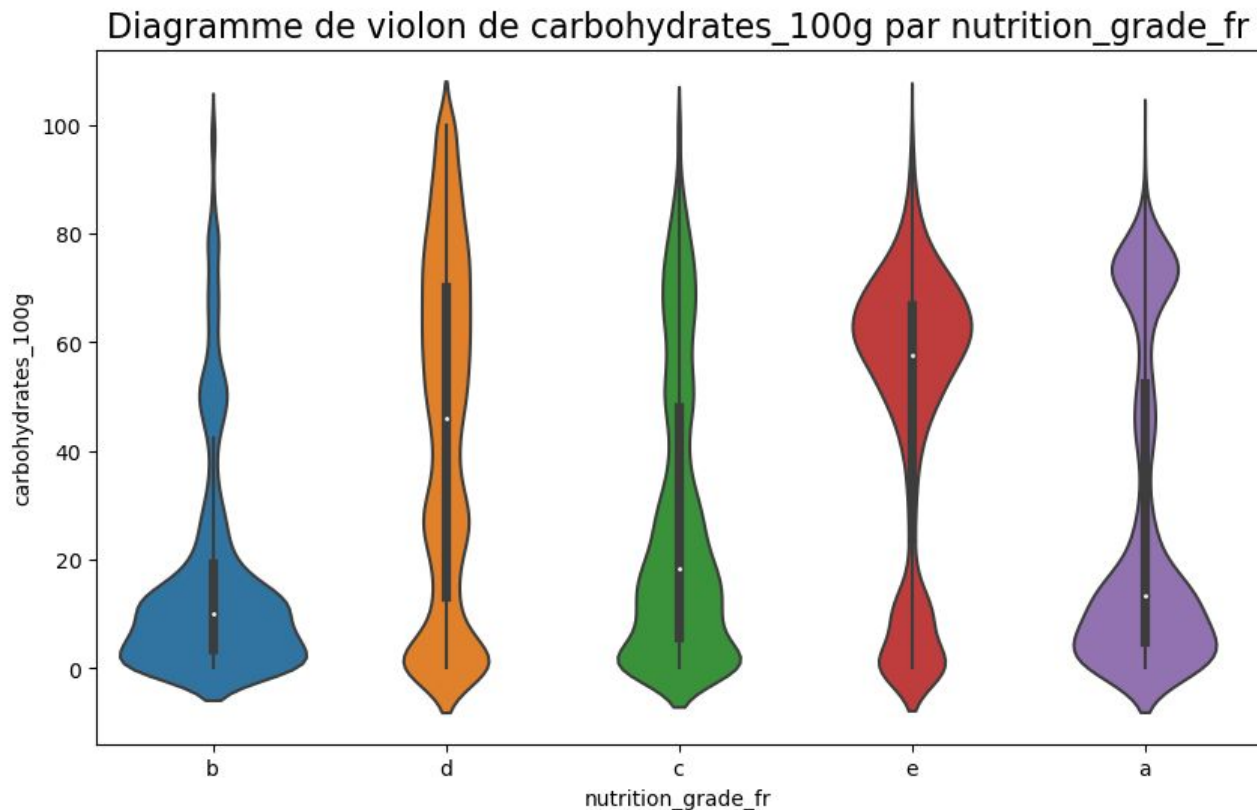
## Visualisation uni-variée : graphiques de densité

Graphiques de densité des variables



# Étape 4 : Effectuez les analyses uni-variée et bi-variée

## Visualisation bi-variée : diagrammes de violon

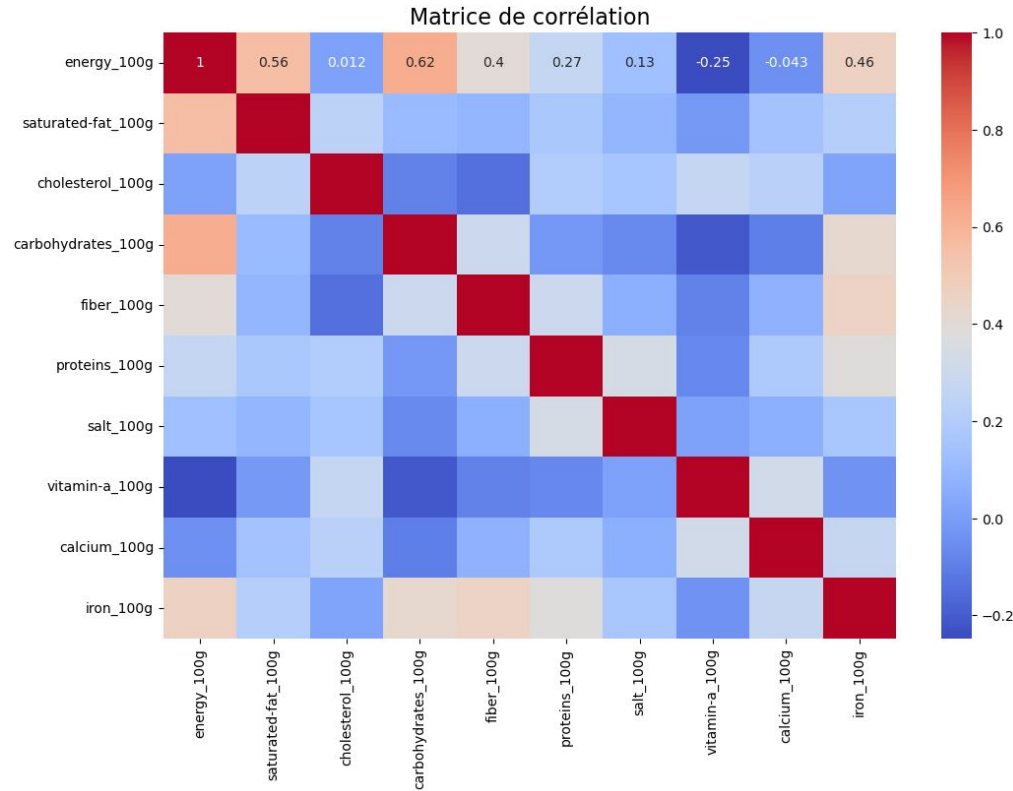


<https://datavizcatalogue.com/>

# Étape 5 : Réalisez une analyse multi-variée

# Étape 5 : Réalisez une analyse multi-variée

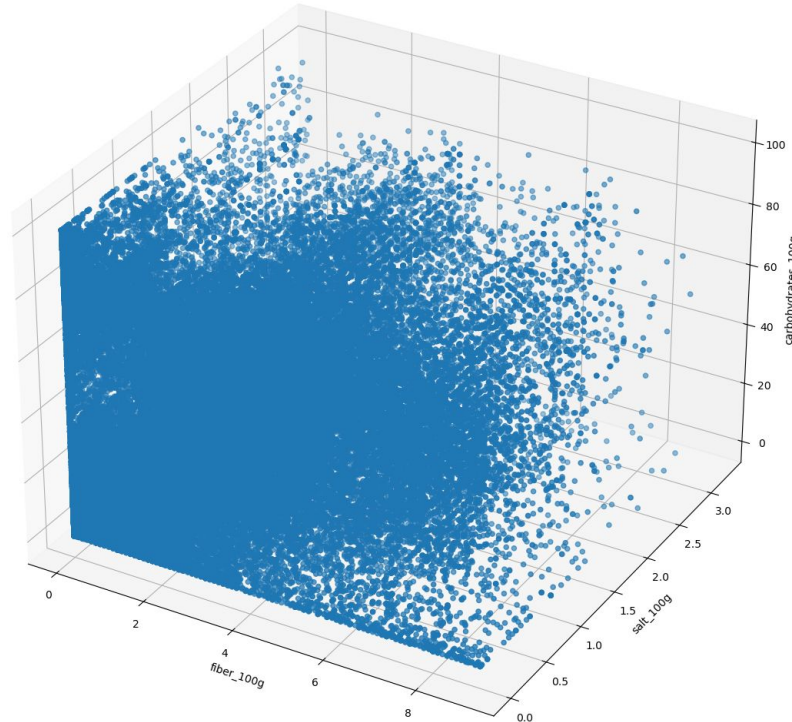
## Visualisation avancée: carte de chaleur



# Étape 5 : Réalisez une analyse multi-variée

## Visualisation avancée: graphique en 3D

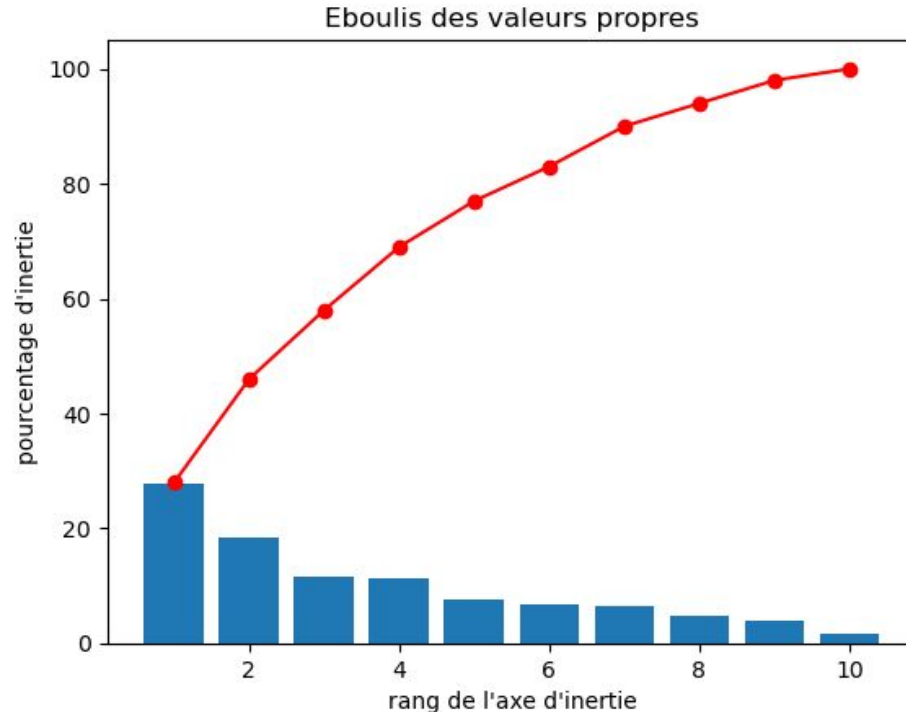
Nuage de points en 3D





# Étape 5 : Réalisez une analyse multi-variée

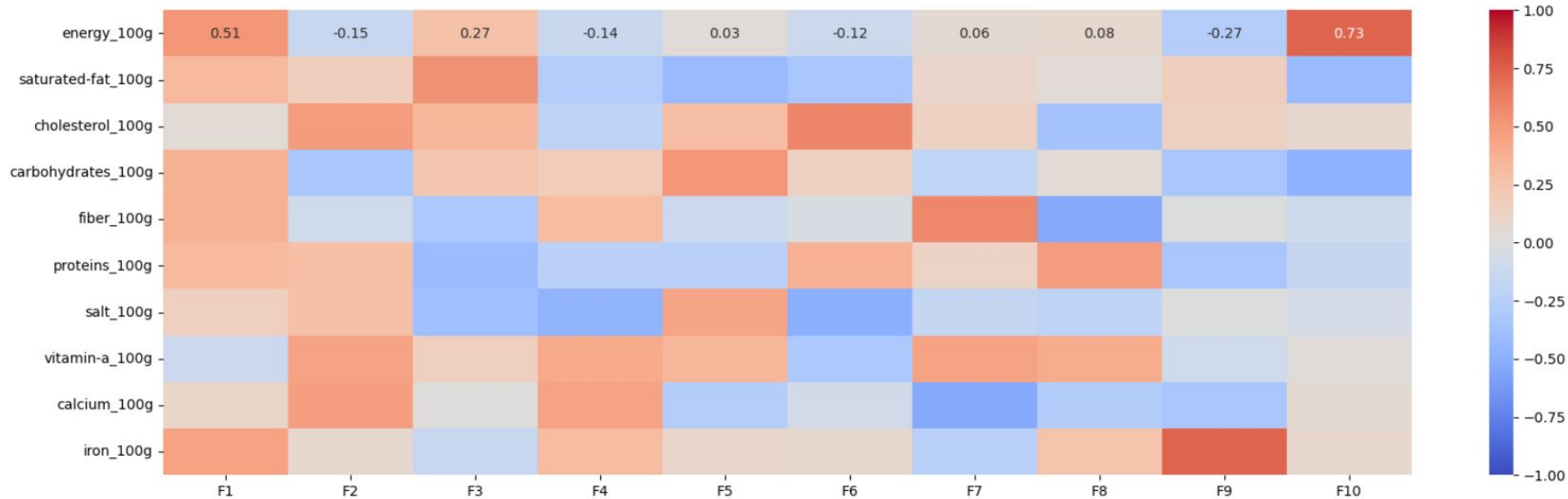
## Analyse statistique : analyse en composante principale (ACP)



# Étape 5 : Réalisez une analyse multi-variée

## Analyse statistique : analyse en composante principale (ACP)

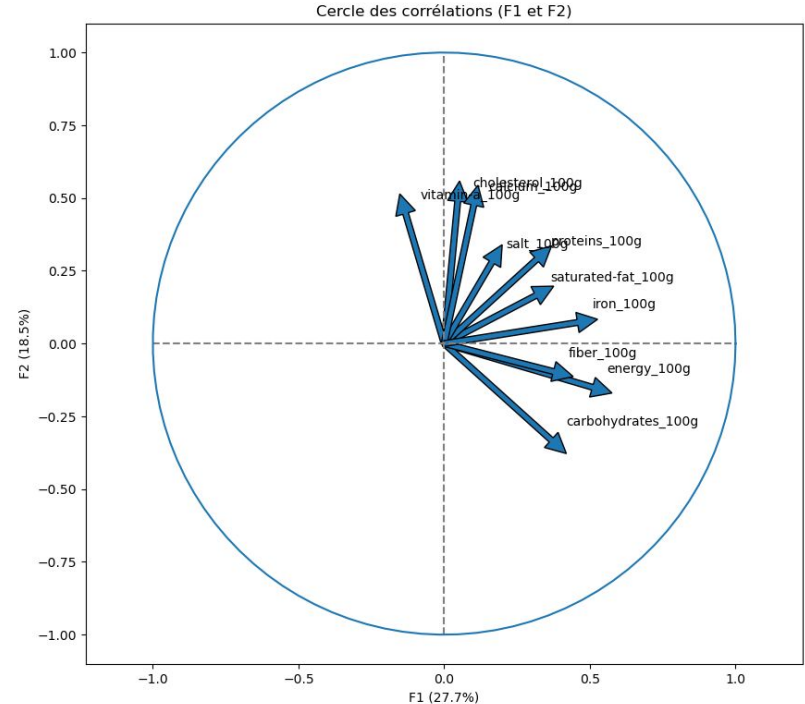
Carte de chaleur des variables en fonction des 10 composantes



# Étape 5 : Réalisez une analyse multi-variée

## Analyse statistique : analyse en composante principale (ACP)

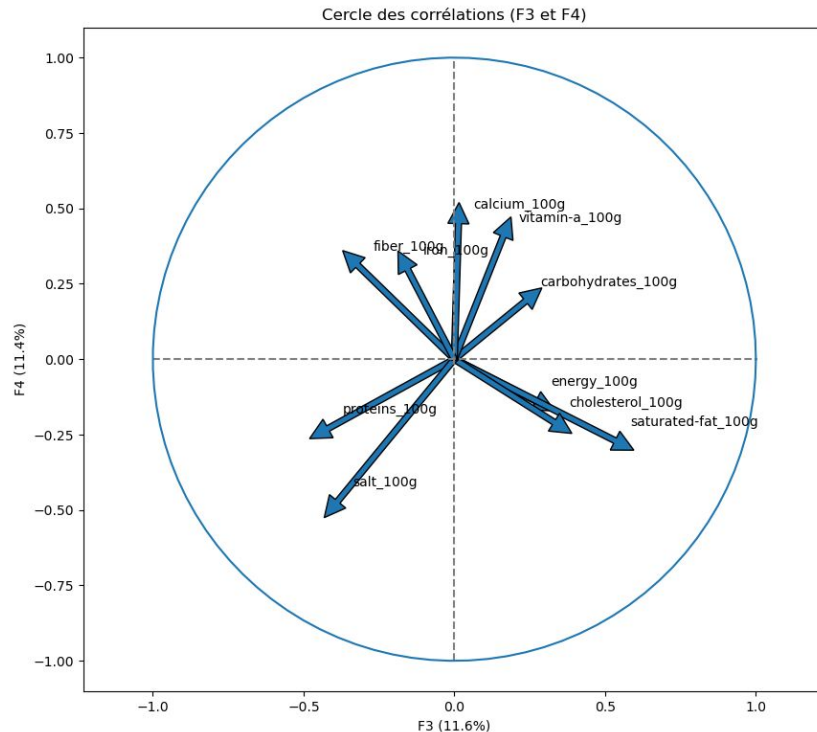
- **energy\_100g** et **iron\_100g** : les plus corrélées positivement à **F1**
- **vitamin-a\_100g**, **cholesterol\_100g** et **calcium\_100g** : les plus corrélées positivement à **F2**
- **energy\_100g** et **fiber\_100g** : très corrélées entre eux
- **energy\_100g** et **cholesterol\_100g** : indépendantes



# Étape 5 : Réalisez une analyse multi-variée

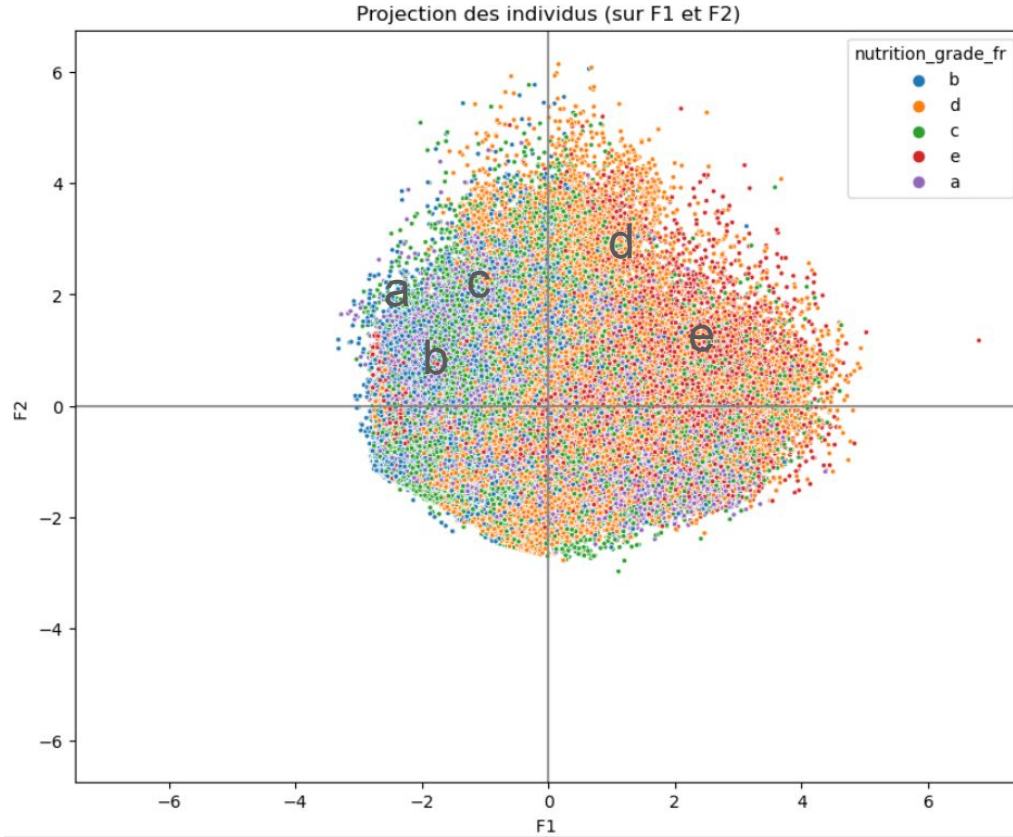
## Analyse statistique : analyse en composante principale (ACP)

- **saturated-fat\_100g** : la plus corrélée positivement à **F3**
- **salt\_100g** : la plus corrélée négativement à **F4**
- **calcium\_100g** est à 0.5 sur **F4** et à 0 sur **F3** : bien corrélée avec **F4** et elle n'a pas d'impact sur **F3**
- **vitamin-a\_100g**, **energy\_100g** : indépendantes



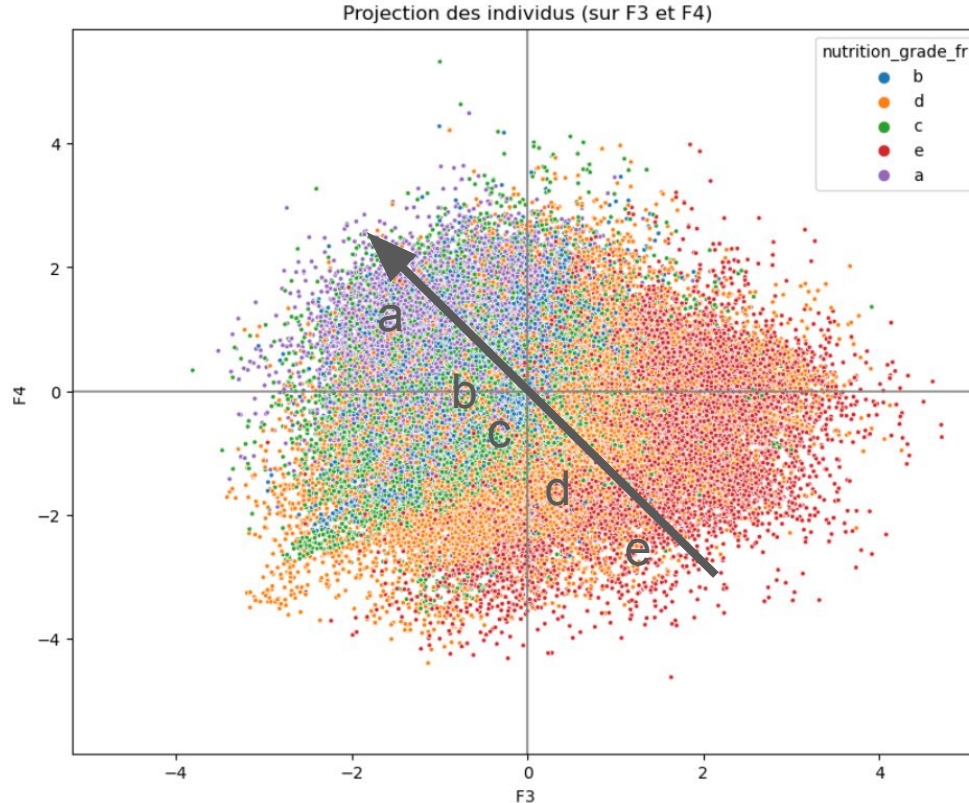
# Étape 5 : Réalisez une analyse multi-variée

## Analyse statistique : analyse en composante principale (ACP)



## Étape 5 : Réalisez une analyse multi-variée

### Analyse statistique : analyse en composante principale (ACP)



# Conclusion

# Conclusion :

Le nettoyage et l'exploration des données montre que ces 10 variables peuvent contribuer pour prédire un Nutri-Score automatique:

- saturated-fat\_100g



- energy\_100g



- cholesterol\_100g



- proteins\_100g



- vitamin-a\_100g

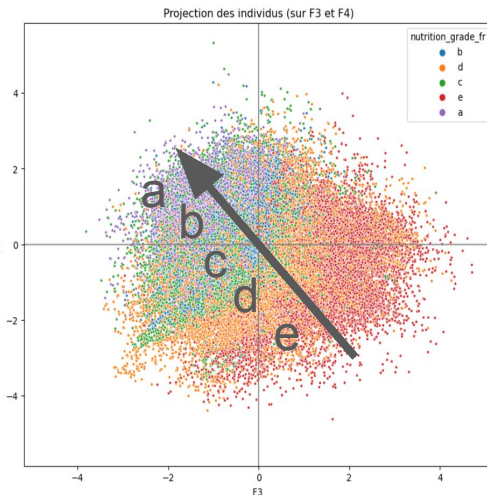
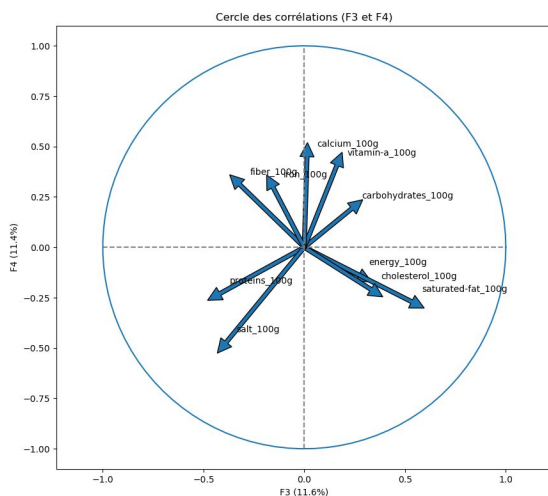


- carbohydrates\_100g

- calcium\_100g

- iron\_100g

- fiber\_100g





# Conclusion :

Le Règlement Général sur la Protection des Données (**RGPD**) est basé sur **cinq grands principes** qui encadrent la protection des données personnelles au sein de l'Union Européenne. Voici ces principes :

1. Licéité, loyauté et transparence
2. Limitation des finalités
3. Minimisation des données
4. Exactitude
5. Limitation de la conservation