# Health Application Project: Alimentation Médicale Préventive

Adonija ZIO

October 2022

# Executive summary

**Data mining**

- Data selection
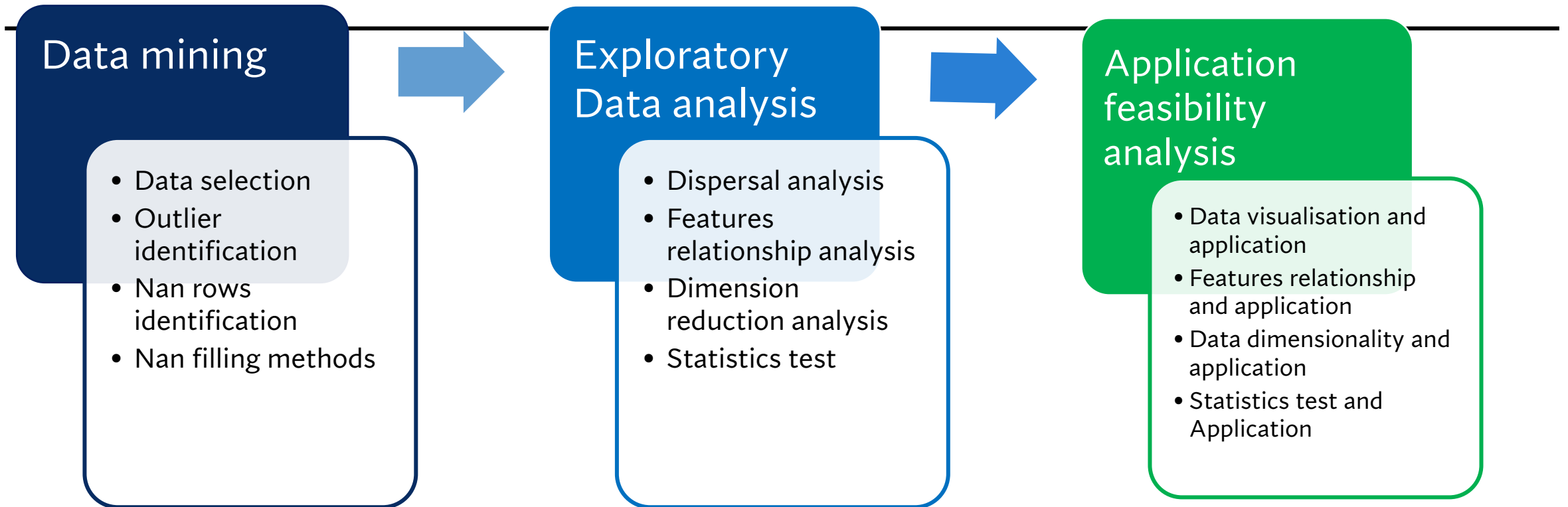- Outlier identification
- Nan rows identification
- Nan filling methods

**Exploratory Data analysis**

- Dispersal analysis
- Features relationship analysis
- Dimension reduction analysis
- Statistics test

**Application feasibility analysis**

- Data visualisation and application
- Features relationship and application
- Data dimensionality and application
- Statistics test and Application

# Introduction

- Steady increase in diabetes cases

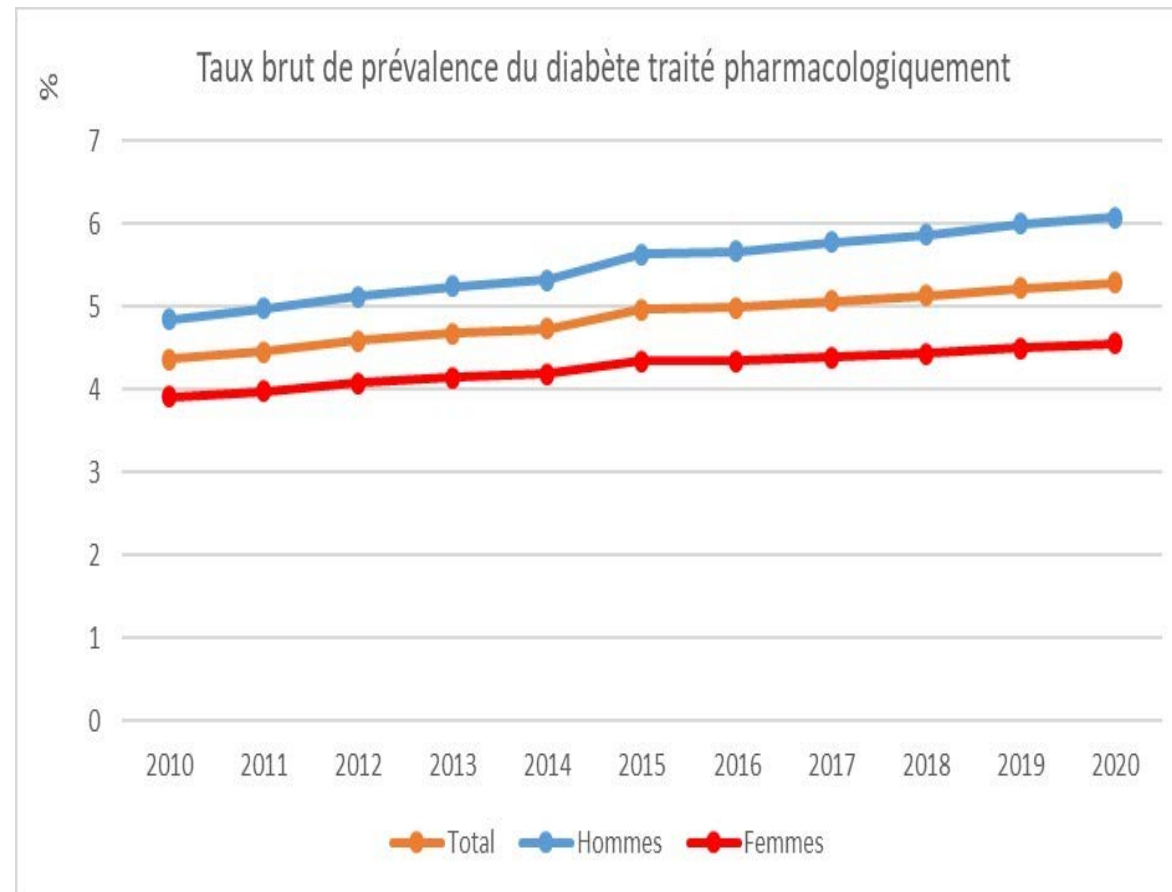- **463 million adults** (20-79 years) with diabetes worldwide in 2019

- This is expected to rise to **700 million by 2045**

- Over **4 million people** identified in France in 2019

- **3,5 millions** under treatment in France in 2020

- 10% des diabétiques sont de type 1

- These treatments are very high cost for the French health system (**19 billion** per year according to CNAM)

Taux brut de prévalence du diabète traité pharmacologiquement

**Source:** Santé Publique France

# Introduction

High need for prevention

The National Health Strategy 2018-2022 and the National Priority Prevention Plan set the framework for diabetes prevention policy
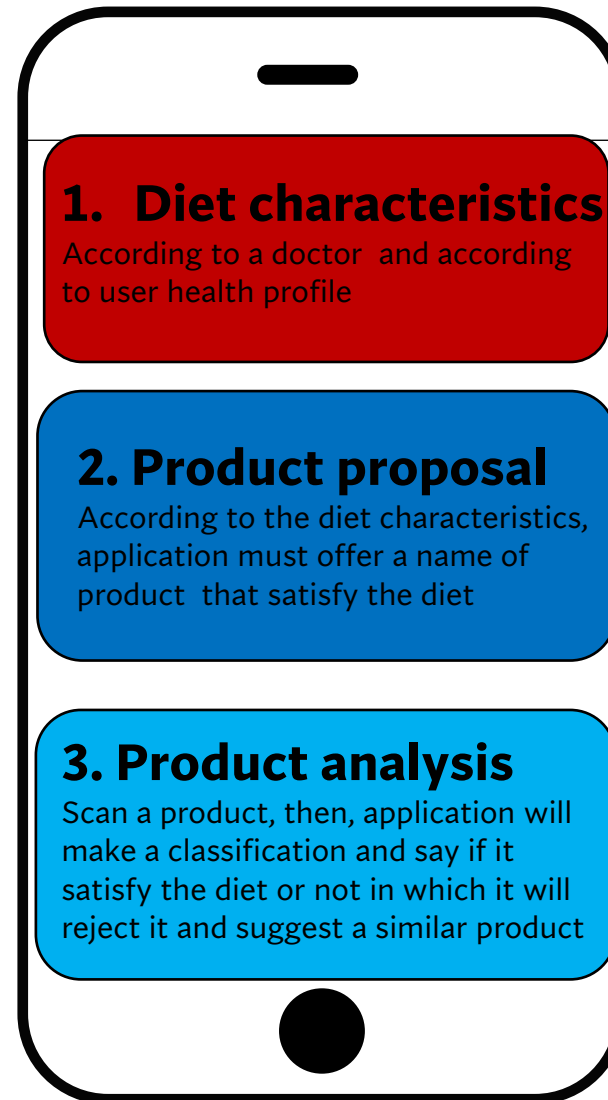
Diet is recognised as an essential preventive element in the most common type of diabetes (Type 2) in the world and in France (90%)
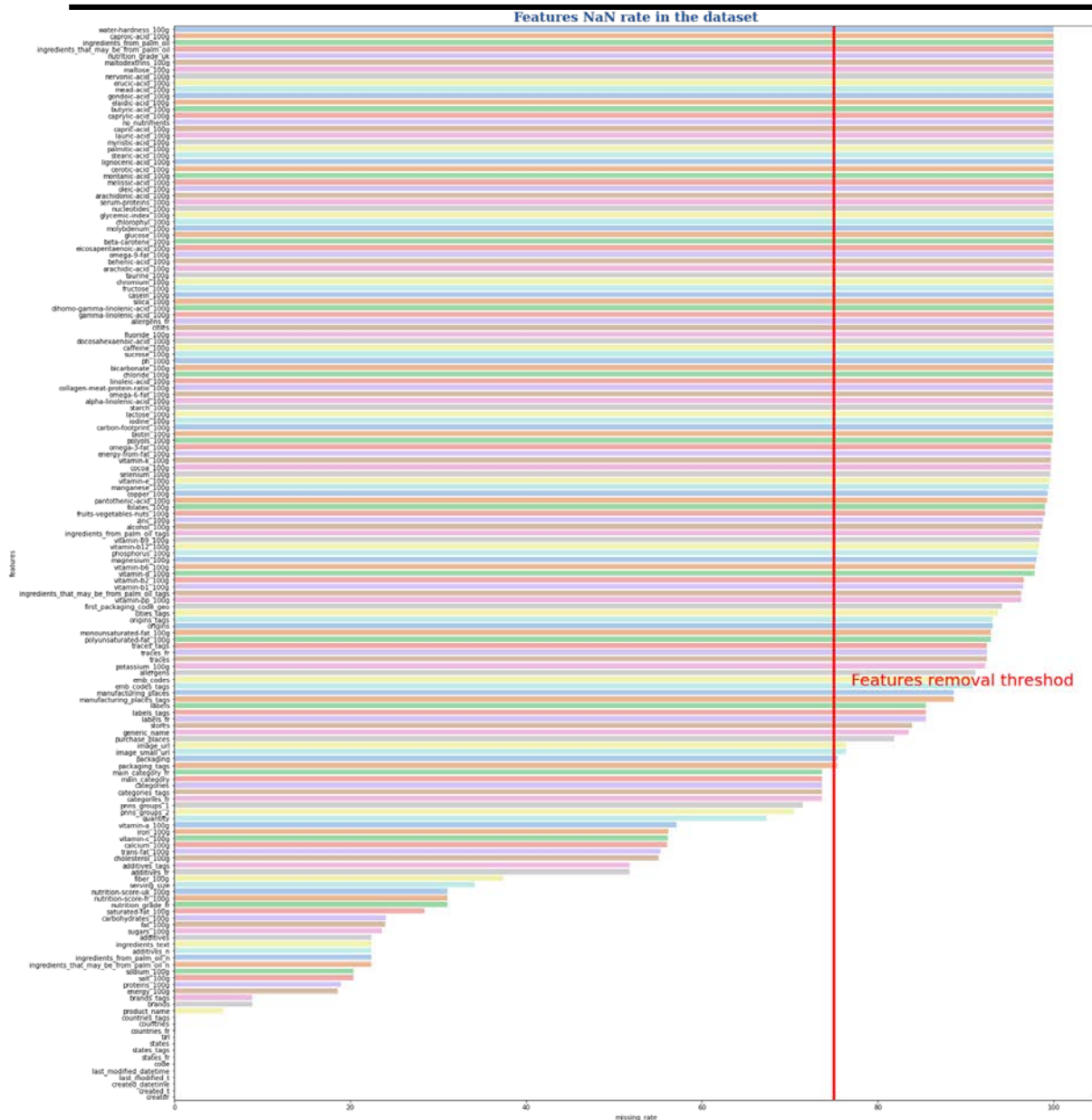
How to support the population in implementing a healthy diet to prevent type 2 diabetes?

# Application Idea

**1. Diet characteristics**
According to a doctor and according to user health profile

**2. Product proposal**
According to the diet characteristics, application must offer a name of product that satisfy the diet

**3. Product analysis**
Scan a product, then, application will make a classification and say if it satisfy the diet or not in which it will reject it and suggest a similar product

# Data Mining

# Technical selection

Our dataset contains 162 columns and 320772 rows

106 numerics columns

50 feature out of 162 have at least 25% data available

There's 30403 unusable rows

# Value treatment

- There are 44459 rows which contain absurd values

Traitement des valeurs aberrantes

- Il existe 31565 lignes dupliquées

Traitement des valeurs dupliquées

- There are 89591 rows which contain Interquartile Outliers

Traitement des outliers

| |
|---|
| Inter Quartile Range for energy_100g is 1284.0 |
| Inter Quartile Range for proteins_100g is 9.41 |
| Inter Quartile Range for salt_100g is 1.40602 |
| Inter Quartile Range for sodium_100g is 0.480811023622047 |
| Inter Quartile Range for additives_n is 3.0 |
| Inter Quartile Range for sugars_100g is 27.41 |
| Inter Quartile Range for fat_100g is 22.02 |
| Inter Quartile Range for carbohydrates_100g is 50.33 |
| Inter Quartile Range for saturated-fat_100g is 7.875 |
| Inter Quartile Range for nutrition-score-fr_100g is 13.0 |
| Inter Quartile Range for nutrition-score-uk_100g is 14.0 |
| Inter Quartile Range for fiber_100g is 3.3 |
| Inter Quartile Range for cholesterol_100g is 0.025 |
| Inter Quartile Range for trans-fat_100g is 0.0 |
| Inter Quartile Range for calcium_100g is 0.089 |
| Inter Quartile Range for vitamin-c_100g is 0.005 |
| Inter Quartile Range for iron_100g is 0.00203 |
| Inter Quartile Range for vitamin-a_100g is 8.34e-05 |

# Nan Filling

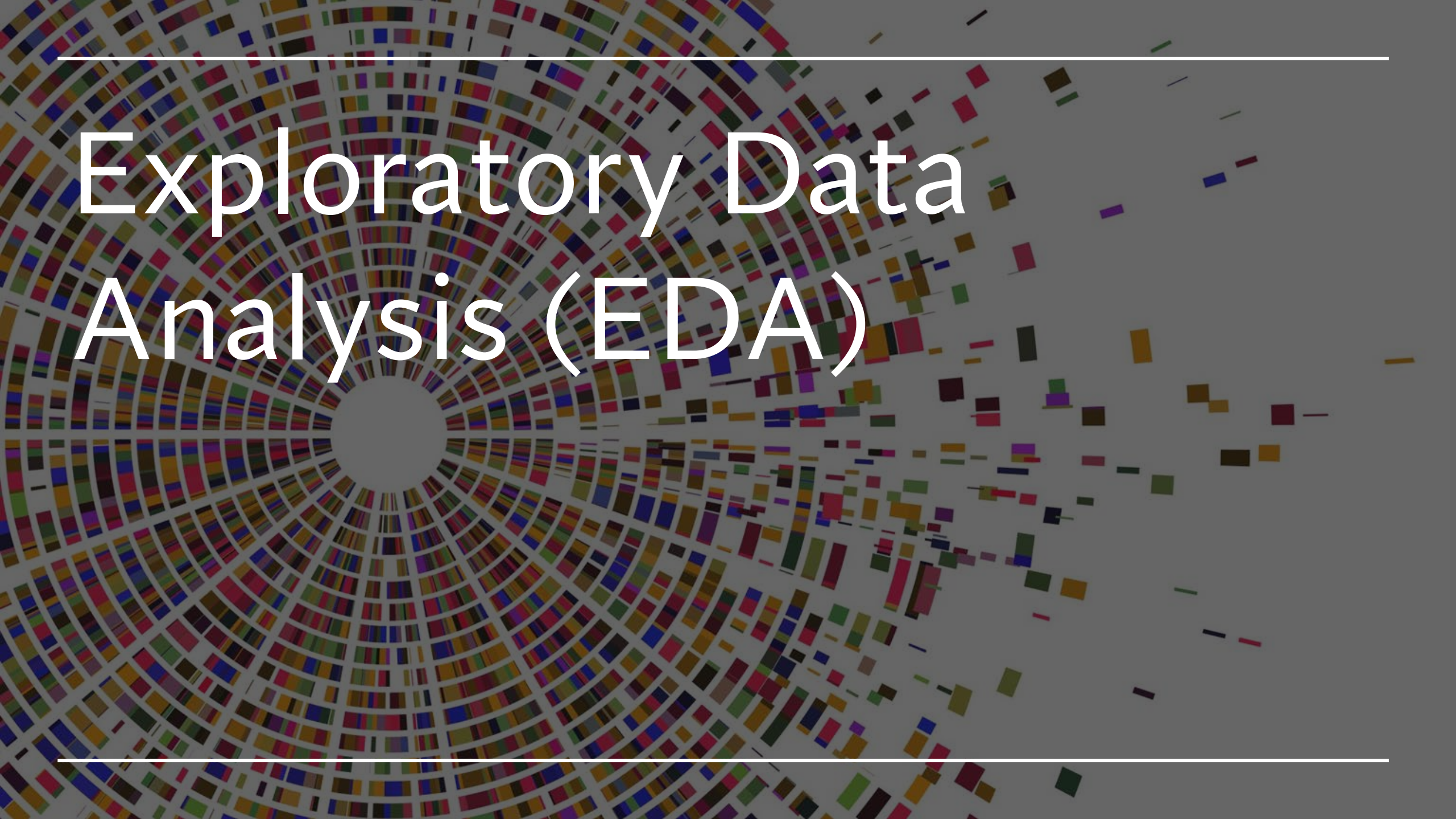**01** Filling discrete variables with 0

**02** Filling in the secondary variables with the median

**03** KNN imputation with 5 neighbourhoods for the key variable

**04** Iterative imputation with sklearn for all the others numeric variables.
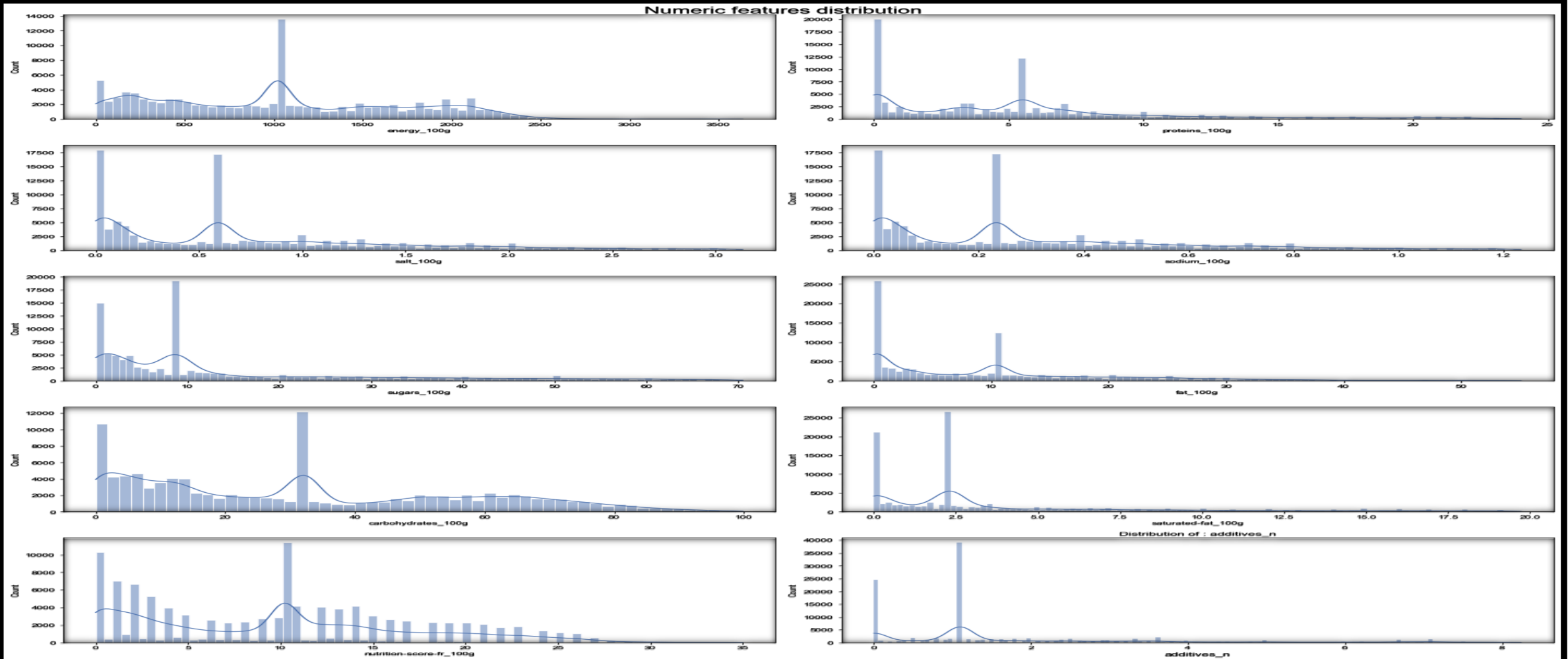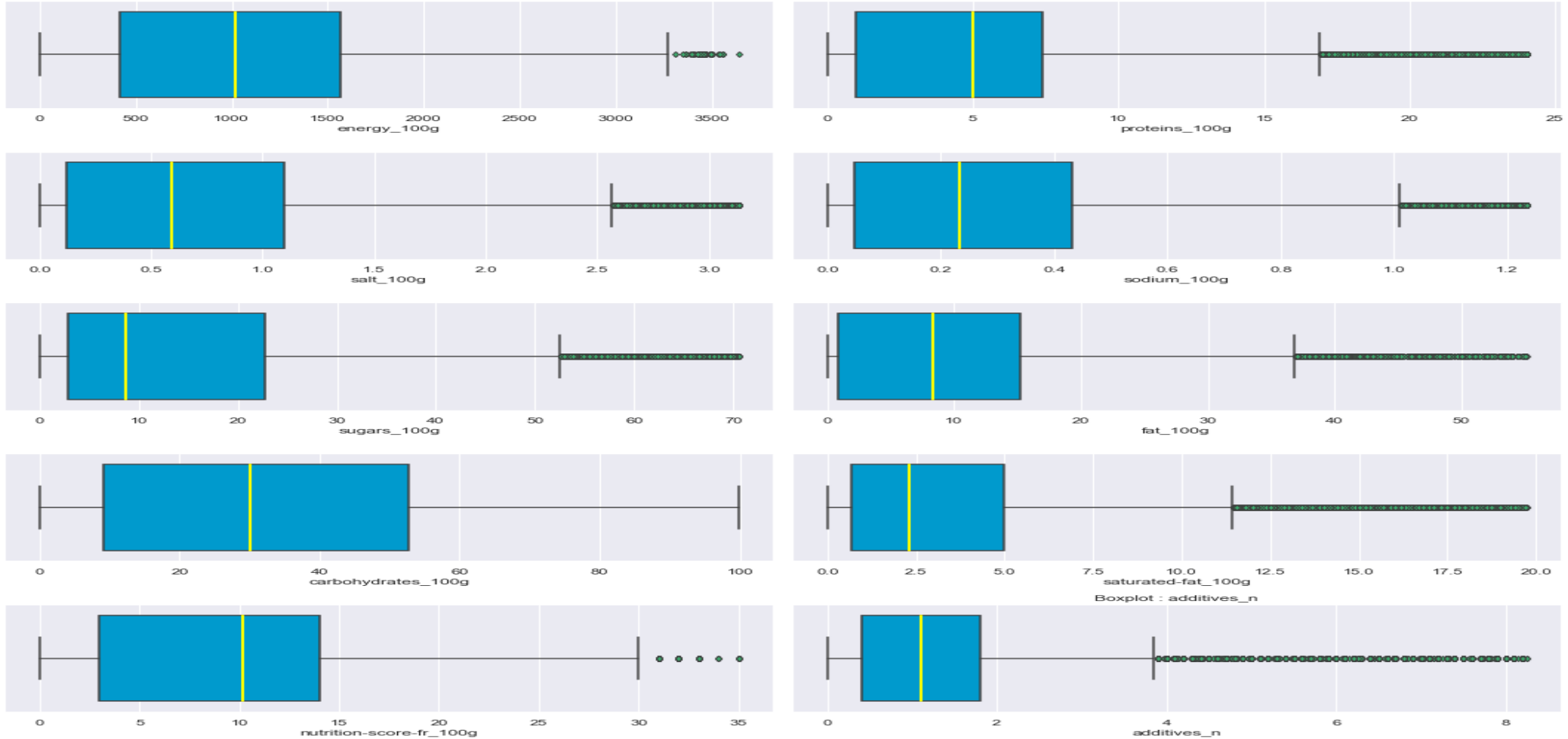
# Exploratory Data Analysis (EDA)

Univariate analysis
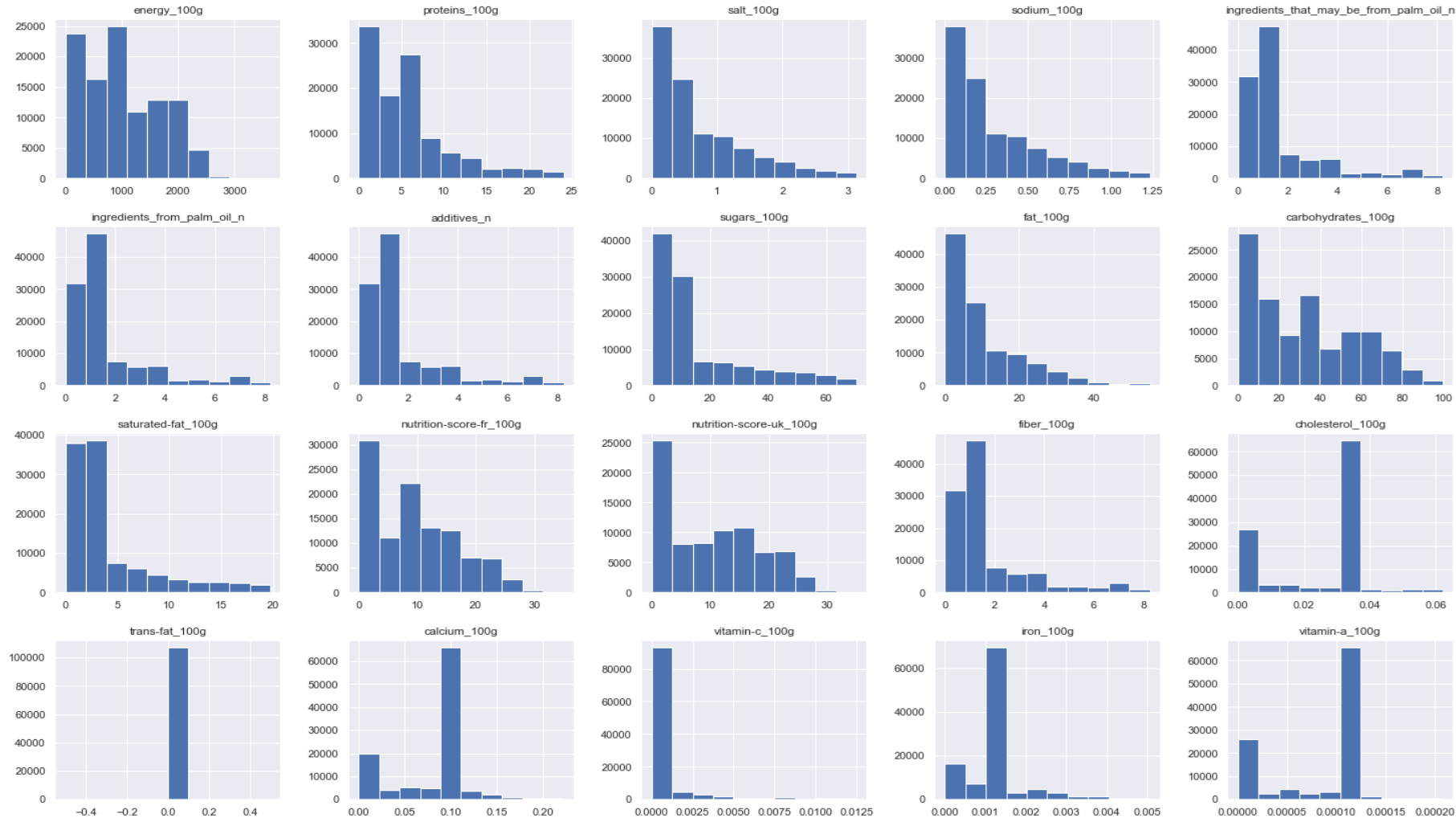
# Displot

# Boxplot



Numerical features Boxplot

# Histogram

# Categorical features visualisation
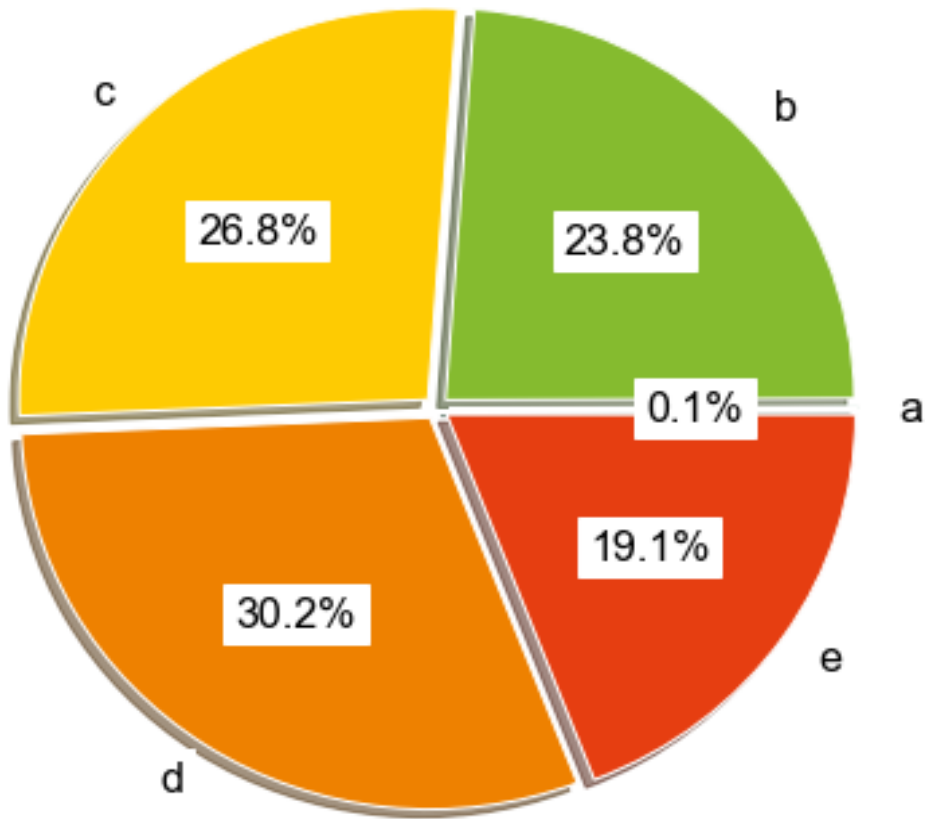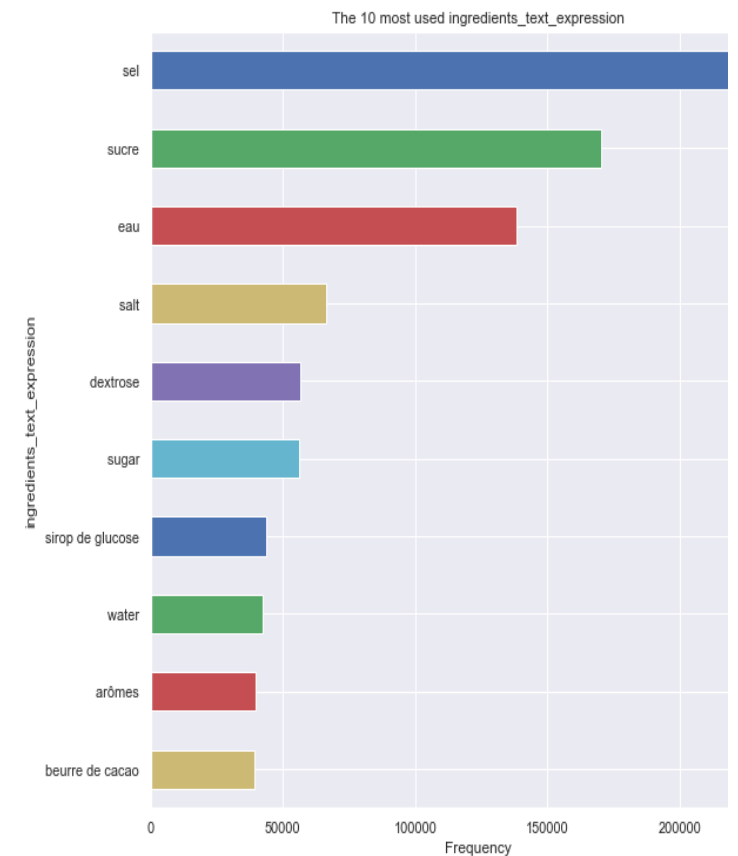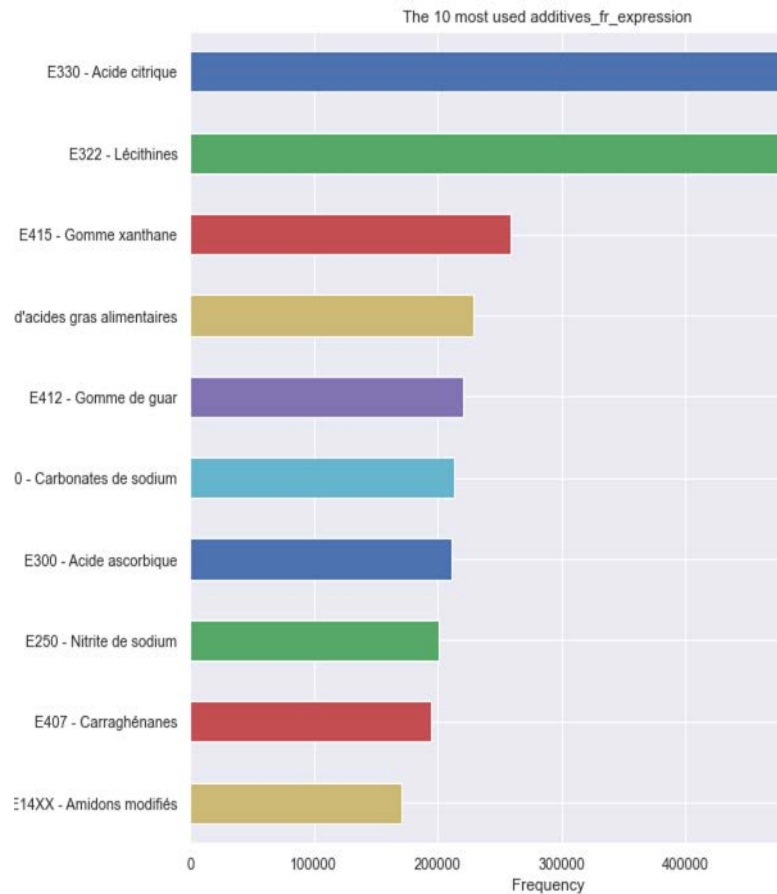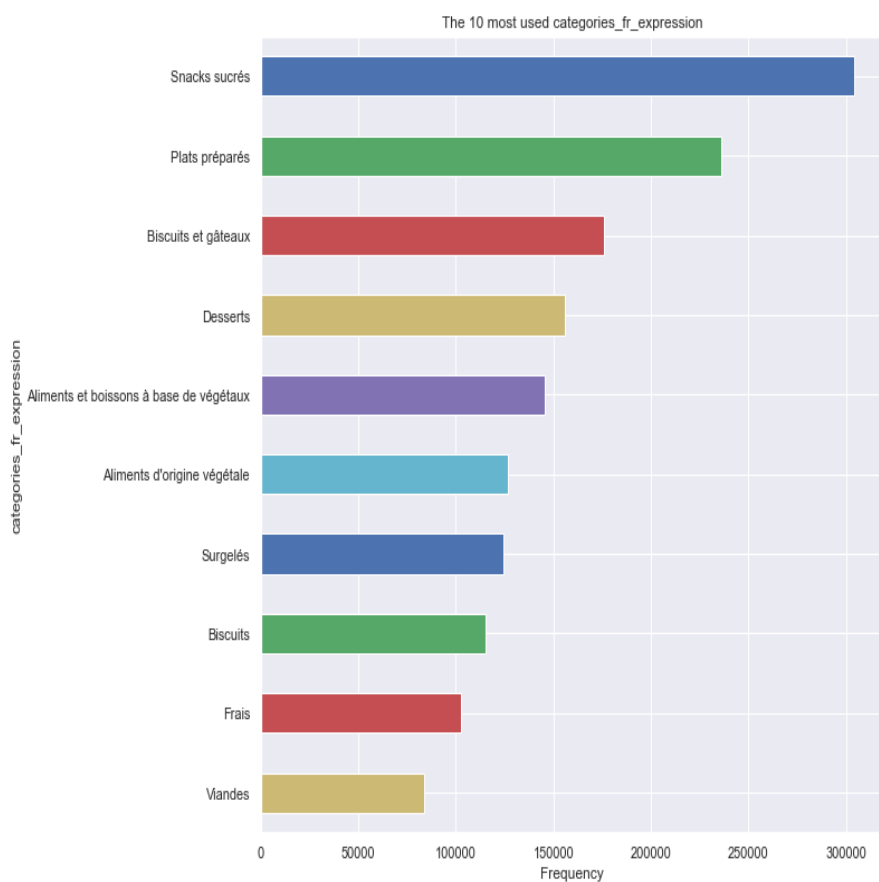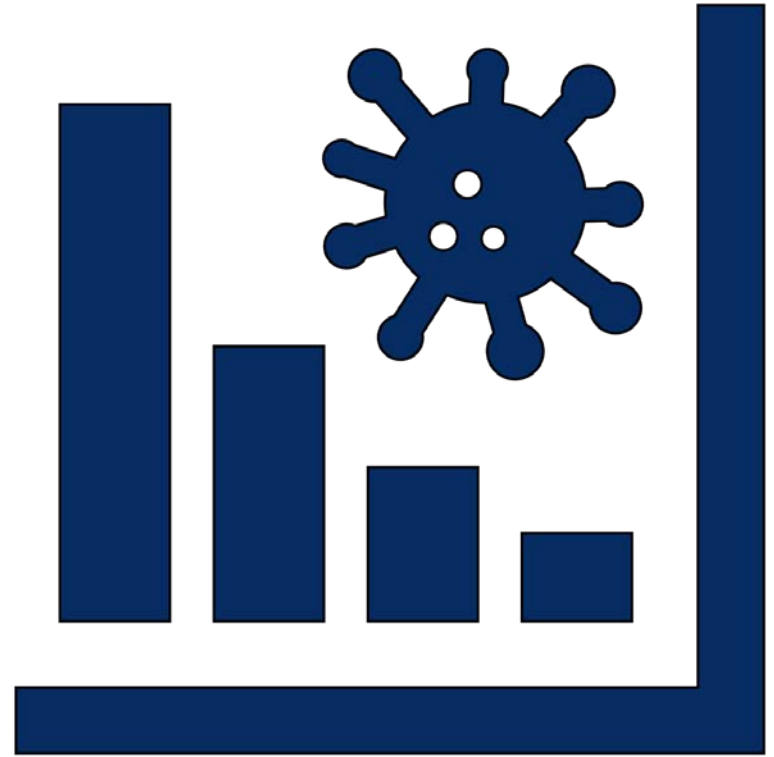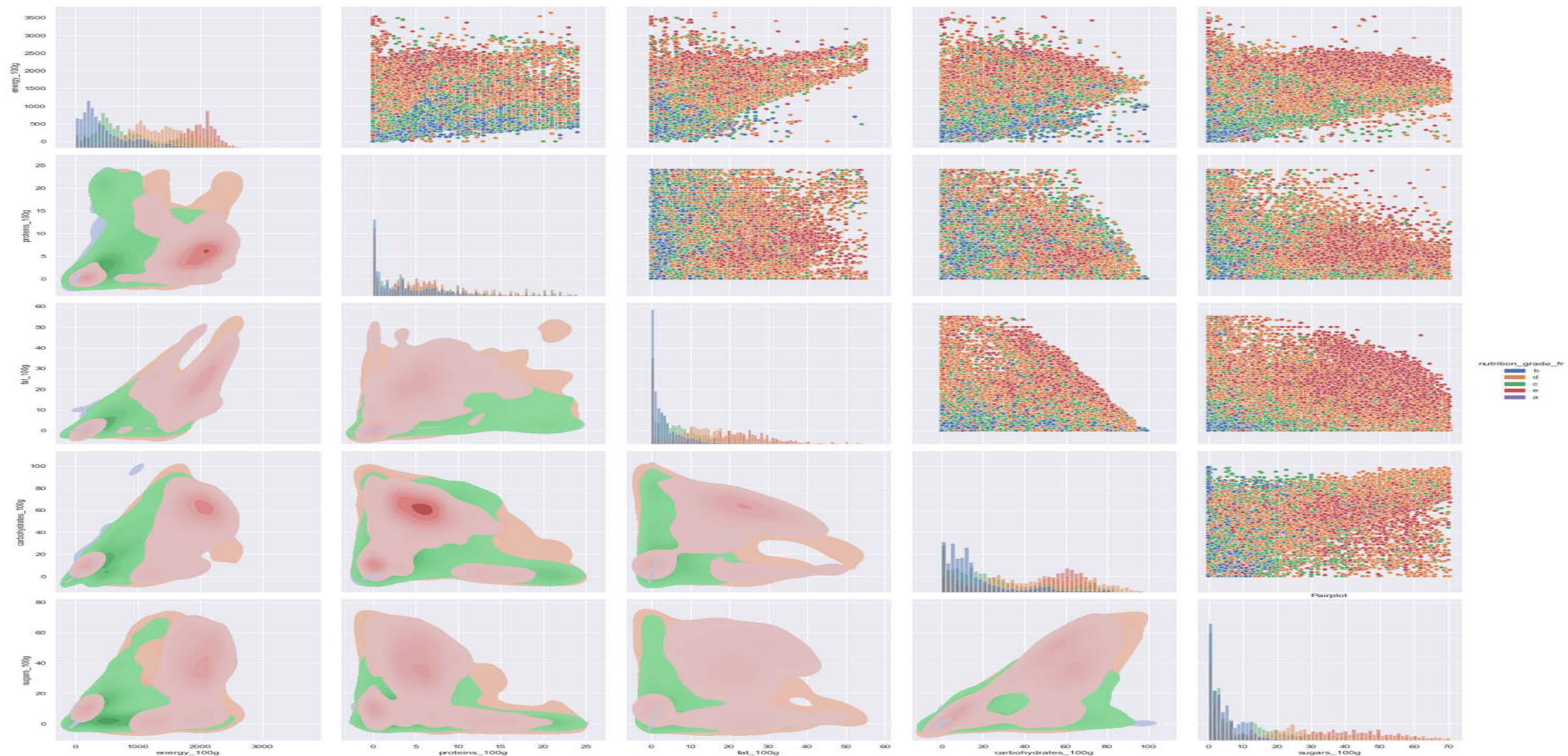


Distribution of nutrition grade
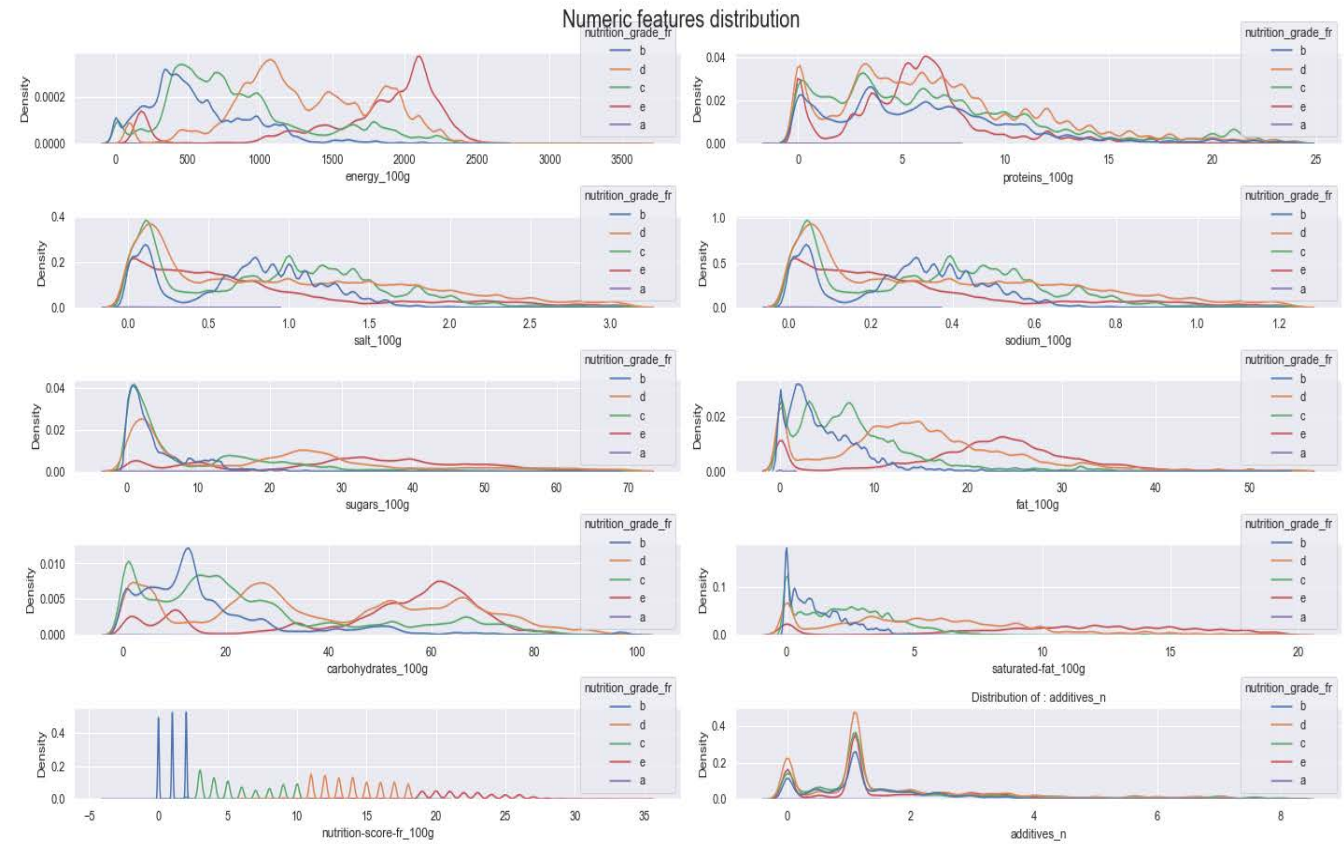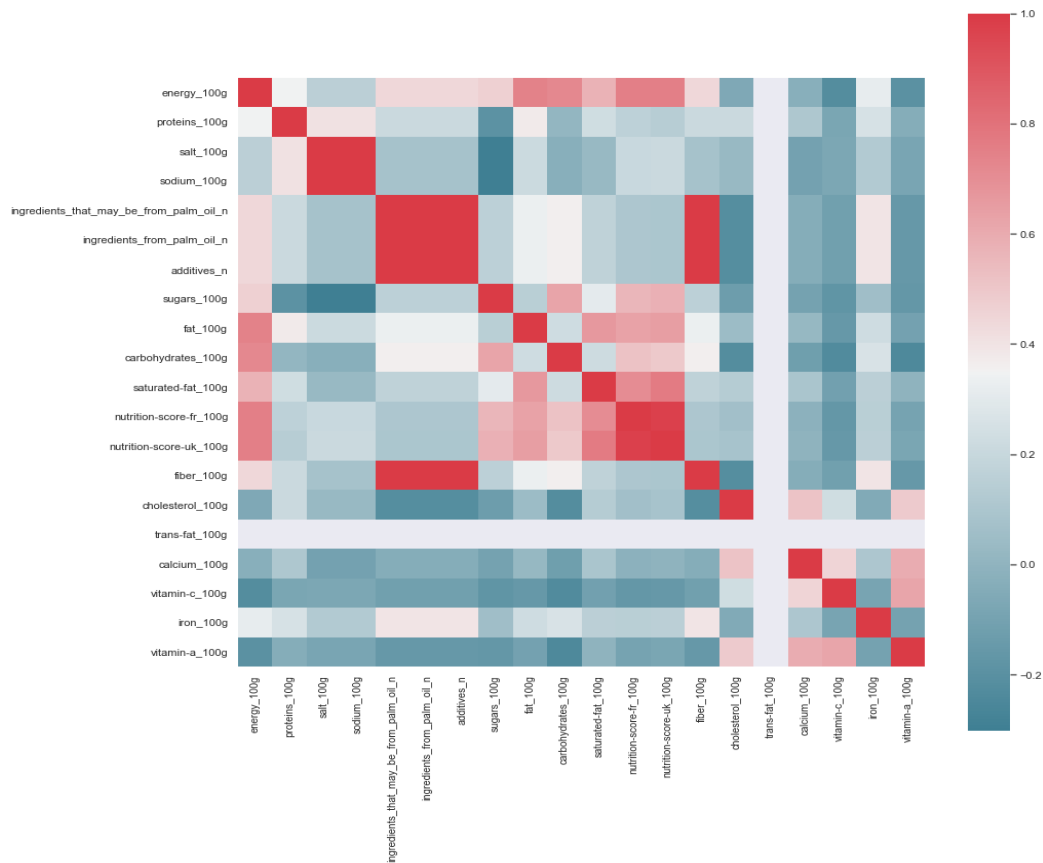
# Categorical features visualisation
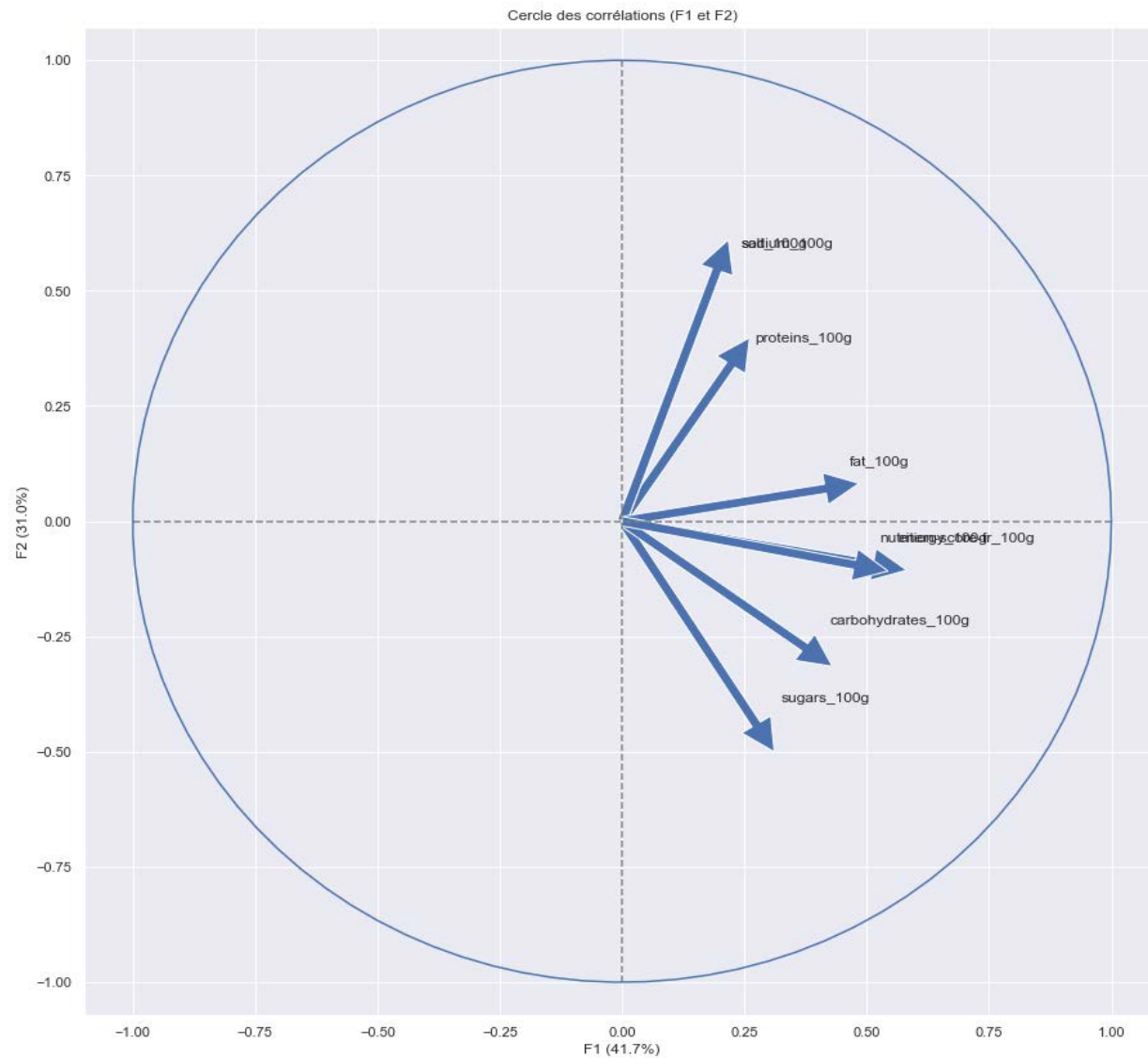
# Bivariate analysis

# Features relationship
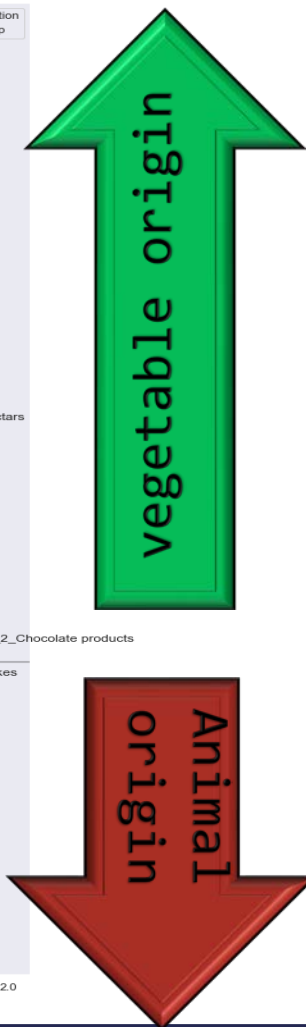
# Correlation analysis

# Multivariate analysis

# PCA



Cercle des corrélations (F1 et F2)

Row principal coordinates

# MCA

Healthy

Unhealthy

vegetable origin

Animal origin

Row and column principal coordinates

- nutrition
- group

group_2_cereals-and-potatoes

group_2_fruits
group_2_vegetables-fresh vegetables

nutrition_grade_fr_a

group_2_Soups
group_21_Vegetables
group_2_Dried fruits
group_2_Fruits

group_2_Non-sugared beverages
group_1_Beverages
group_2_Fruit juices
group_2_Artificially sweetened beverages

group_2_Fruit nectars
group_2_Sweetened beverages

group_2_One-dish meals
group_1_Composite foods

nutrition_grade_fr_b
group_2_Pizza pies and quiche

group_2_Potatoes-Legumes
group_1_Fresh
group_1_Cereals and potatoes

nutrition_grade_fr_e

group_2_Chocolate products

group_2_Breakfast cereals
group_2_Milk and yogurt
group_2_Eggs
nutrition_grade_fr_c

group_2_unknown

group_2_Sweets
group_2_Biscuits and cakes
group_2_pastries

group_1_Milk and dairy products
group_2_Dairy desserts
group_2_Fish and seafood
group_2_Meat
group_2_Eggs
nutrition_grade_fr_d
group_2_Cheese
Processed meat

group_2_Dressings and sauces
group_2_Fats

group21_Appetizers
group_2_Salty snacks
group_2_Salty and fatty products

Component 1 (4.09% inertia)

Component 0 (4.69% inertia)

# Statistic tests

| | Chi-square | Degree of freedom | Critical value | p_value |
|---|---|---|---|---|
| **pnns_groups_1** | 11171.42 | 44. | 60.48 | 0.0 |
| **pnns_groups_2** | 18517.63 | 156 | 186.15 | 0.0 |
| countries | 5030.93 | 2180 | 2289.74 | $5.92.10^{-226}$ |
| categories_fr | 77966.73 | 34084 | 34514.59 | 0.0 |
| additives_fr | 104693.14 | 63472 | 64059.18 | 0.0 |

| | df | sum_sq | mean_sq | F | PR(>F) |
|---|---|---|---|---|---|
| **pnns_groups_1** | 11.0 | 5.48e+05 | 49809.45 | 11308.08 | 0.0 |
| **nutrition_grade_fr** | 4.0 | 1.06e+06 | 263771.33 | 59883.17 | 0.0 |
| **pnns_groups_1:nutrition_grade_fr** | 44.0 | 2.07e+04 | 469.62 | 106.62 | 0.0 |
| **Residual** | 29297.0 | 1.29e+05 | 4.40 | NaN | NaN |

- Our features not satisfy normality test

# **Application feasibility analysis**

# Thank You for your attention!