

# Projet 3 : Application en lien avec l'alimentation pour Santé publique France

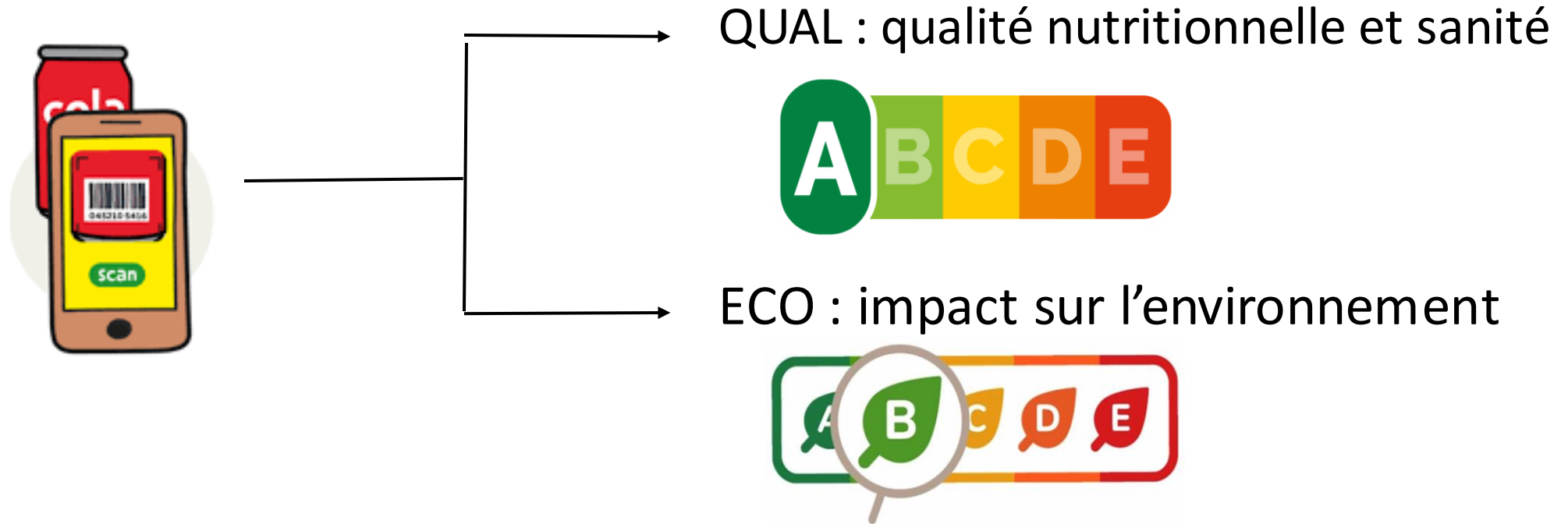
Fanjamalala Rajaonalison  
14/06/2021

# Plan

1. Présentation de l'application
2. Opérations de nettoyage
3. Description et Analyse des variables
4. Analyse multivariée
5. Pertinence et Faisabilité de l'application

## QUALECO

Application, qui après avoir identifié le produit, retourne deux scores



## DEUX SCORES

### QUAL

- Fournit une note allant de A à E



- Calculé à partir de :
  - nutriscore
  - nombre d'additifs
  - nombre d'ingrédients issus de l'huile de palme

### ECO

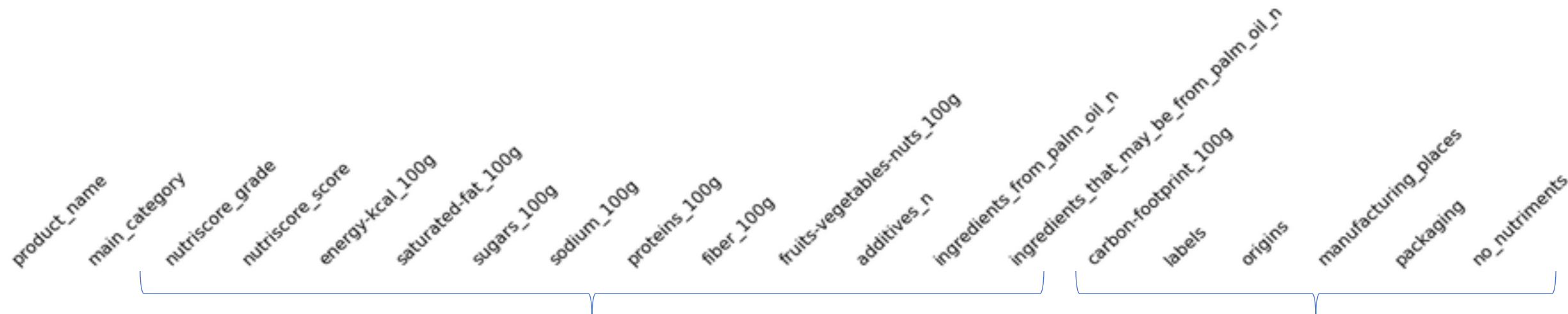
- Fournit une note allant de A à E



- Calculé à partir de :
  - présence de label bio
  - ingrédients locaux ou non
  - produits transformés/fabriqués localement ou non
  - empreinte carbone
  - type d'emballage

## 1. Sélection de colonnes

186 colonnes  $\longrightarrow$  20 colonnes



### Produit

2 colonnes :  
2 var. quali

### QUAL

12 colonnes :  
- 1 var. quali  
- 11 var. quanti

### ECO

6 colonnes :  
- 5 var. quali  
- 1 var. quanti

1. Sélection de colonnes

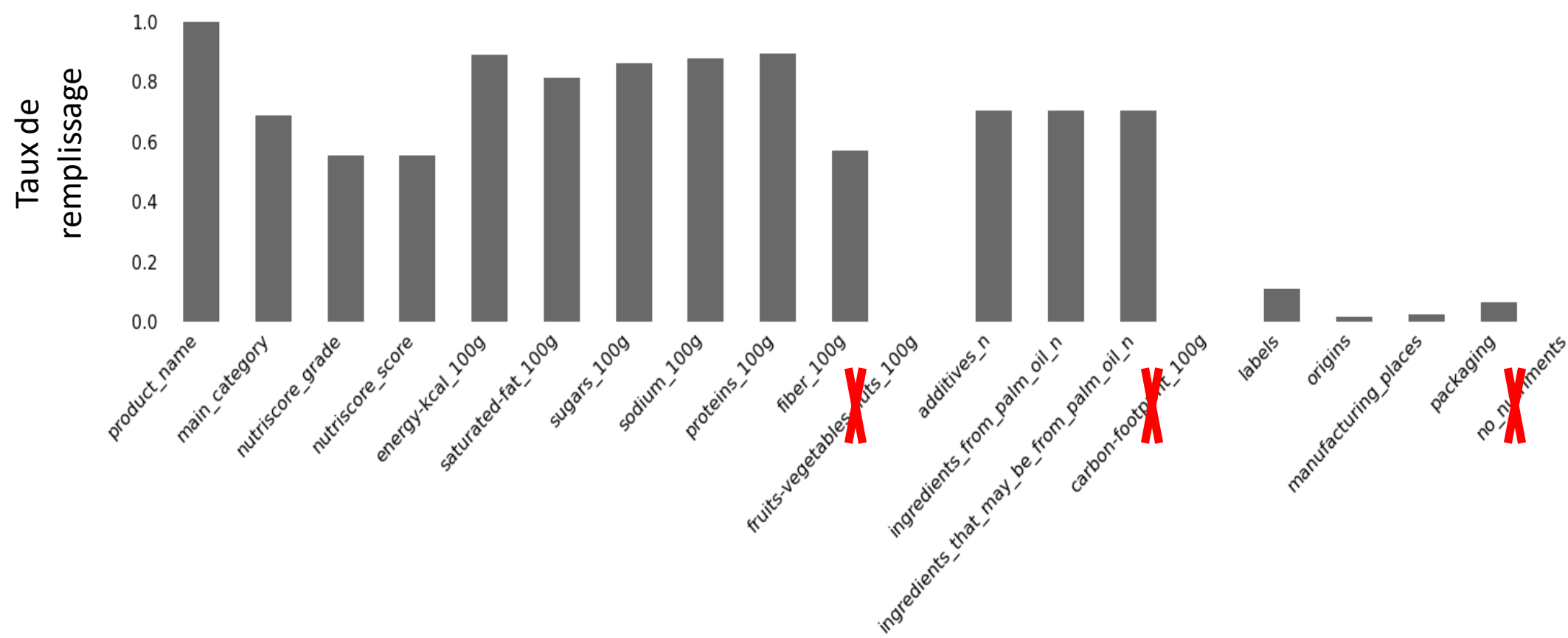
## 2. Suppression des données non pertinentes

- Lignes : Données dupliquées & Lignes vides

1. Sélection de colonnes

## 2. Suppression des données non pertinentes

- Lignes : Données dupliquées & Lignes vides
- Colonnes : Données avec des taux de remplissage < 5%



1. Sélection de colonnes
2. Suppression des données non pertinentes

## 3. Traitement des outliers

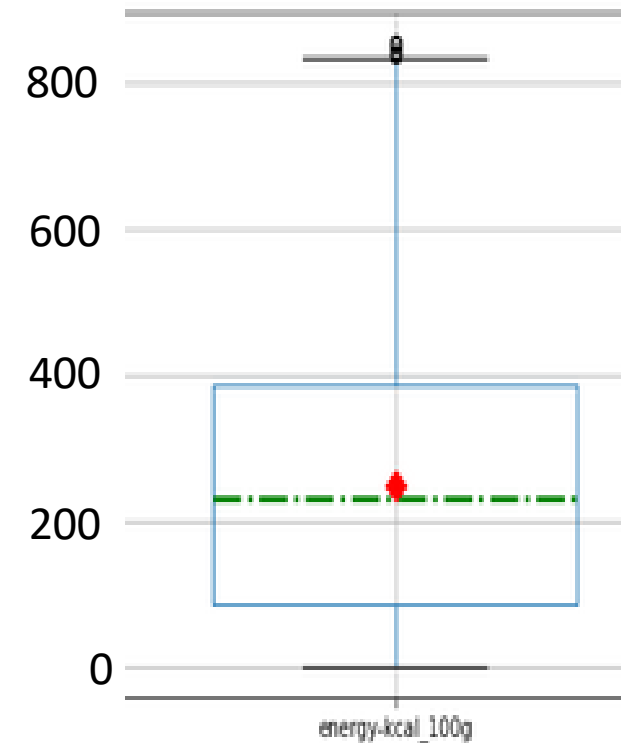
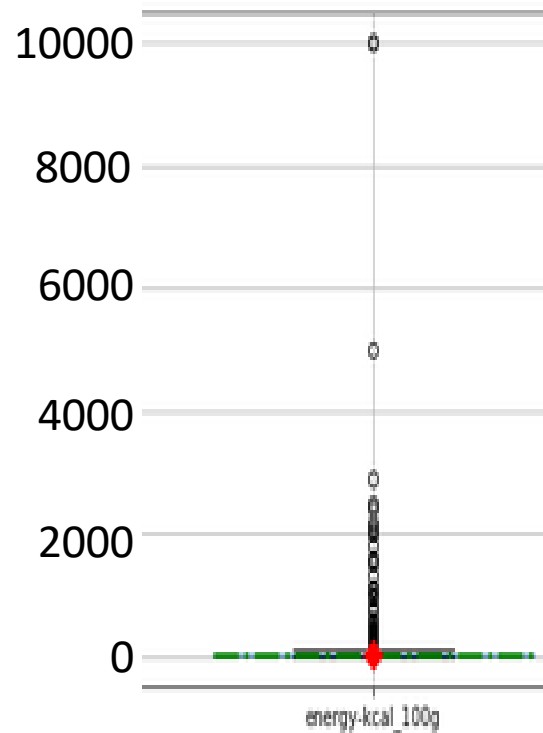
- Identification des outliers : valeurs distantes des extrémités des moustaches



1. Sélection de colonnes
2. Suppression des données non pertinentes

## 3. Traitement des outliers

- Identification des outliers : valeurs distantes des extrémités des moustaches
- Remplacement des outliers par NaN



1. Sélection de colonnes
2. Suppression des données non pertinentes
3. Traitement des outliers

## 4. Traitement des valeurs manquantes

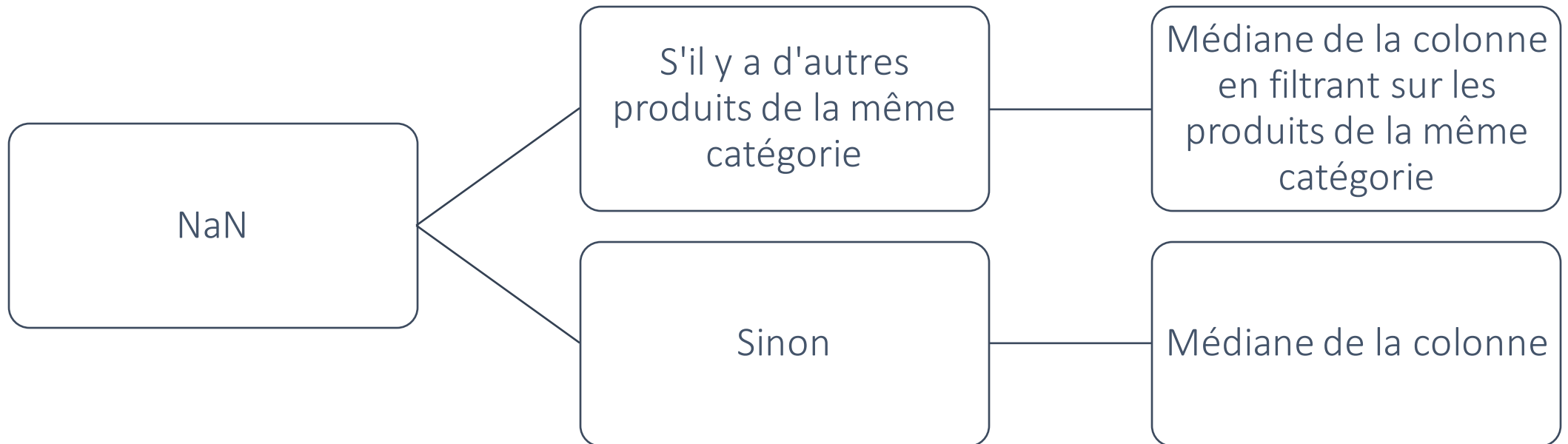
- Pour les variables qualitatives



1. Sélection de colonnes
2. Suppression des données non pertinentes
3. Traitement des outliers

## 4. Traitement des valeurs manquantes

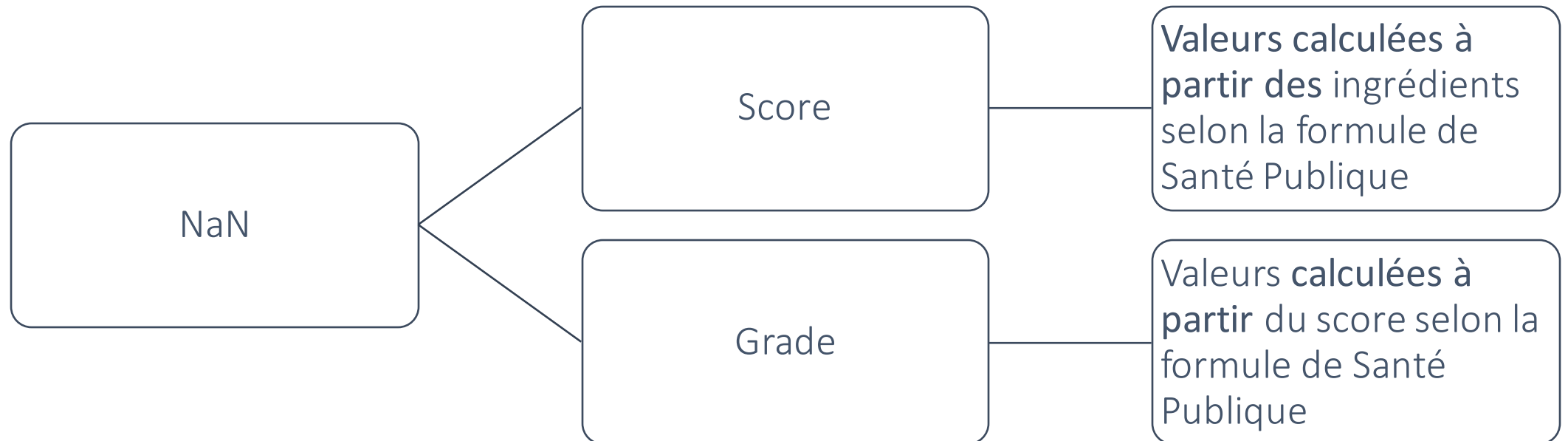
- Pour les variables qualitatives
- Pour les variables quantitatives



1. Sélection de colonnes
2. Suppression des données non pertinentes
3. Traitement des outliers

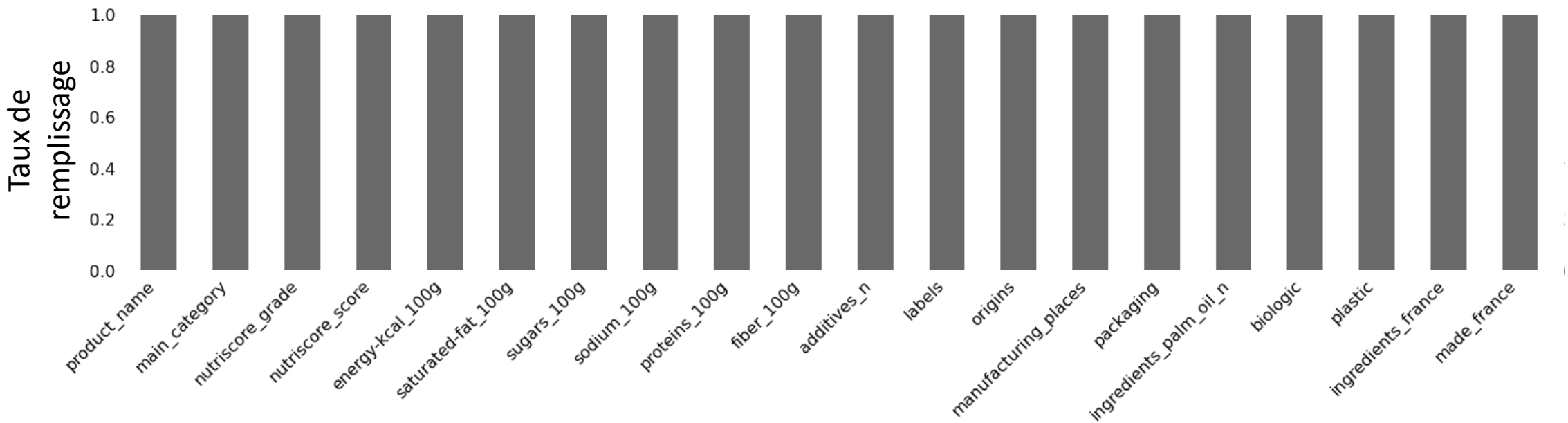
## 4. Traitement des valeurs manquantes

- Pour les variables qualitatives
- Pour les variables quantitatives
- Pour les variables nutriscores



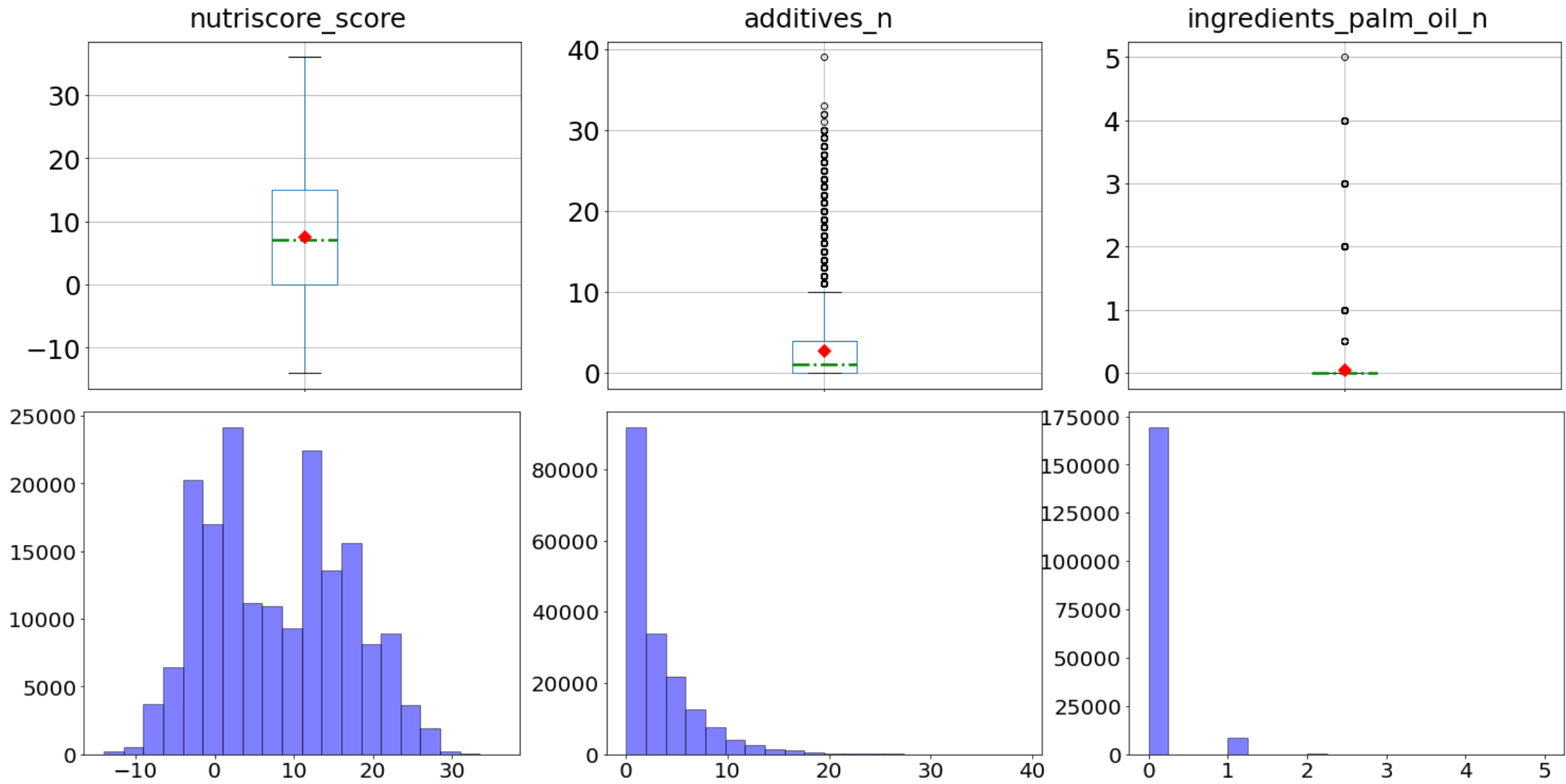
- 1. Sélection de colonnes
- 2. Suppression des données non pertinentes
- 3. Traitement des outliers
- 4. Traitement des valeurs manquantes

## 5. Vérification des données



→ Le jeu de donnée est prêt pour les étapes d'analyses.

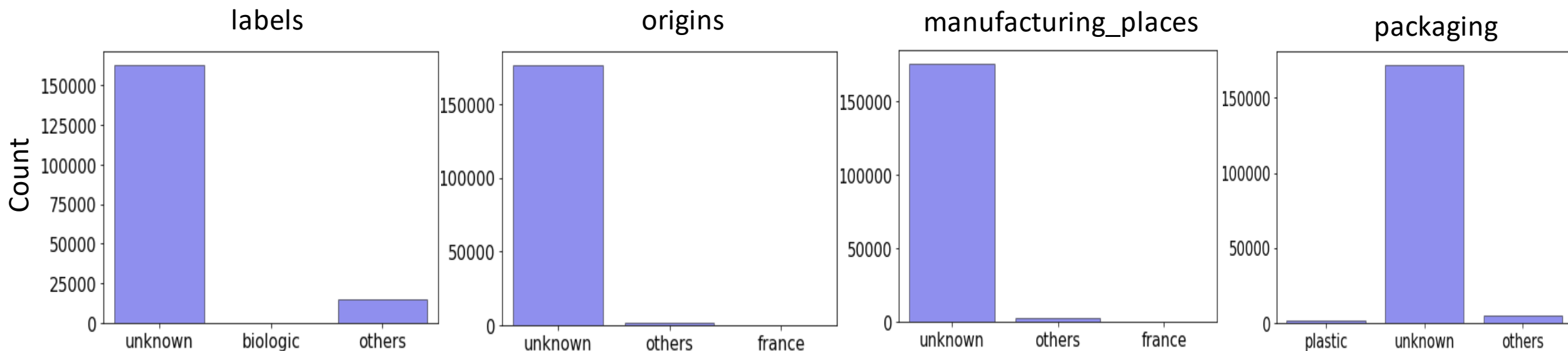
# 1. Variables quantitatives



→ La distribution des variables ne suivent pas la loi normale

1. Variables quantitatives

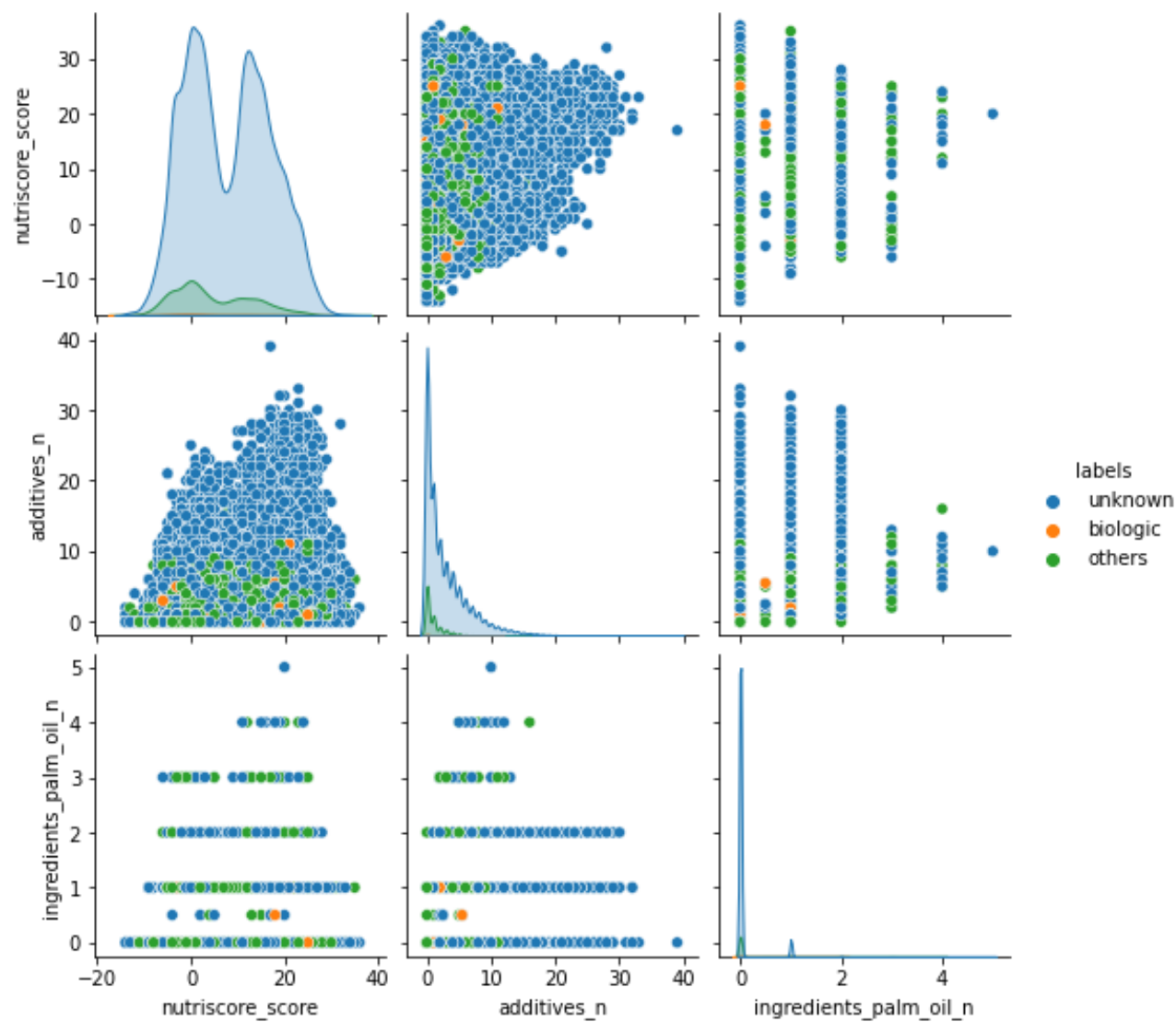
## 2. Variables qualitatives



→ Il y a beaucoup de données "inconnus" pour les 4 variables, largement plus que les données renseignées

→ Moins bonne qualité des données : réelle manque d'informations

# 1. Relation entre variables



→ les produits ayant des labels inconnus ont un plus fort nutriscore, un plus grand nombre d'additifs et d'ingrédients provenant de l'huile de palme

→ on ne note pas de relation flagrante entre les variables nutriscore et le nombre d'additifs

→ la variable N ingrédients d'huile de palme ne montre aucune relation avec les autres variables



1. Relation entre variables

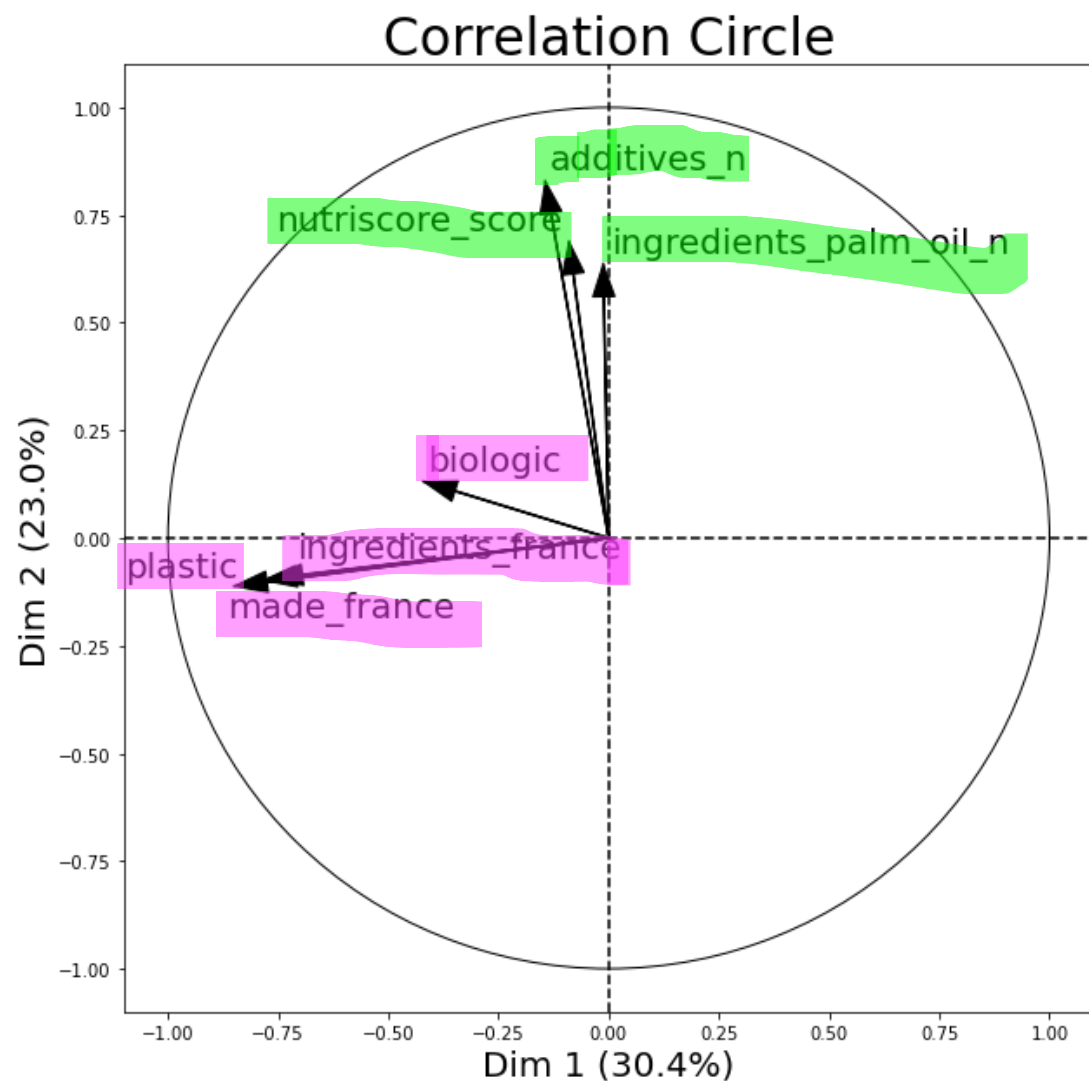
2. Corrélation

2. Correlation

| QUAL |                        |                  |             | ECO                    |          |         |                    |             |
|------|------------------------|------------------|-------------|------------------------|----------|---------|--------------------|-------------|
|      |                        | nutriscore_score | additives_n | ingredients_palm_oil_n | biologic | plastic | ingredients_france | made_france |
| QUAL | nutriscore_score       | 1.00             | 0.37        | 0.12                   | 0.10     | 0.03    | 0.03               | 0.02        |
|      | additives_n            | 0.37             | 1.00        | 0.31                   | 0.14     | 0.05    | 0.04               | 0.04        |
|      | ingredients_palm_oil_n | 0.12             | 0.31        | 1.00                   | 0.02     | -0.03   | -0.01              | -0.02       |
| ECO  | biologic               | 0.10             | 0.14        | 0.02                   | 1.00     | 0.26    | 0.15               | 0.21        |
|      | plastic                | 0.03             | 0.05        | -0.03                  | 0.26     | 1.00    | 0.46               | 0.58        |
|      | ingredients_france     | 0.03             | 0.04        | -0.01                  | 0.15     | 0.46    | 1.00               | 0.60        |
|      | made_france            | 0.02             | 0.04        | -0.02                  | 0.21     | 0.58    | 0.60               | 1.00        |

- on note une relation faible pour les variables définissant QUAL
- les relations sont plus significatives pour les variables définissant ECO
- il semble ne pas avoir de relation entre les variables définissant QUAL et celles définissant ECO

# 1. Existence de 2 groupes de variables



→ on note 2 groupes de variables corrélées entre elles

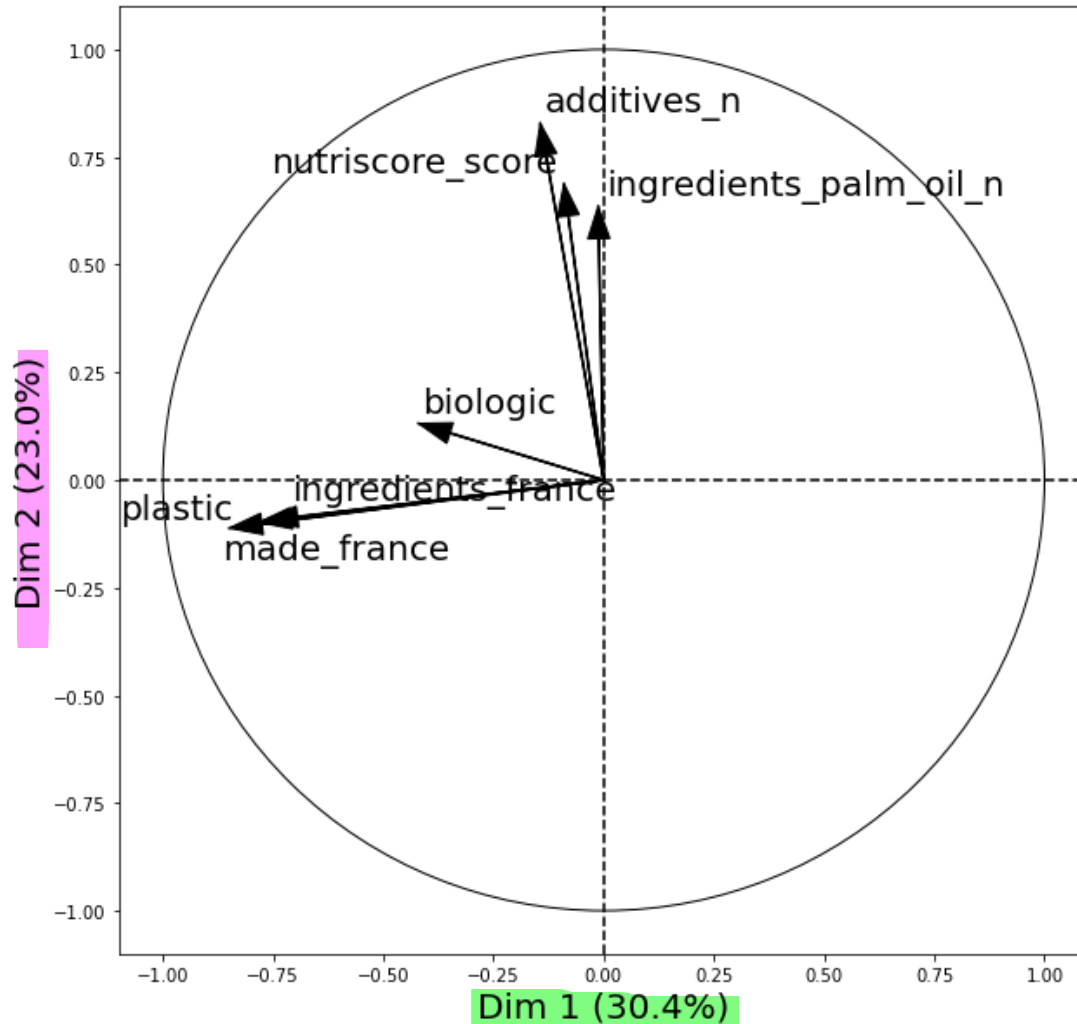
→ possibilité de résumer ces 2 groupes de variables par 2 variables synthétiques qui pourront être les 2 scores de l'application : **QUAL** et **ECO**

(1<sup>ère</sup> possibilité de définir de QUAL et ECO)

1. Existence de 2 groupes de variables

## 2. Pertinence de 2 axes principales

### Correlation Circle



→ les 2 groupes de variables identifiés précédemment sont respectivement corrélés aux 2 axes principaux

→ ces 2 axes expliquent plus de la moitié (53,4%) de la variance des données

→ ces 2 axes pourraient être les 2scores de l'application : **QUAL** et **ECO**

(2<sup>ème</sup> possibilité de définir de QUAL et ECO)

1. Existence de 2 groupes de variables
2. Pertinence de 2 axes principales

### 3. Contribution des variables à ces scores

|                               | ECO       | QUAL      |
|-------------------------------|-----------|-----------|
| <b>nutriscore_score</b>       | -0.107672 | 0.536689  |
| <b>additives_n</b>            | -0.130236 | 0.639554  |
| <b>ingredients_palm_oil_n</b> | -0.026127 | 0.487845  |
| <b>biologic</b>               | -0.294329 | 0.144872  |
| <b>plastic</b>                | -0.529990 | -0.103531 |
| <b>ingredients_france</b>     | -0.528233 | -0.123589 |
| <b>made_france</b>            | -0.569398 | -0.134021 |

→ les variables qui contribuent le plus aux scores ECO et QUAL sont ceux attendus

→ ces 2 scores ECO et QUAL peuvent être définis par la combinaison linéaire des variables y participant

$$\begin{aligned}\text{QUAL} = & 0.54 * \text{nutriscore} \\ & + 0.64 * \text{N additifs} \\ & + 0.49 * \text{N ingrédients huile de palme}\end{aligned}$$

$$\begin{aligned}\text{ECO} = & -0.29 * \text{bio} \\ & -0.53 * \text{plastic} \\ & -0.53 * \text{ingrédients de France} \\ & -0.57 * \text{produit en France}\end{aligned}$$

# Synthèse & Conclusion



Conception de l'application QUALECO permettant de calculer 2 scores QUAL et ECO

# Synthèse & Conclusion

Conception de l'application QUALECO permettant de calculer 2 scores QUAL et ECO



Faisabilité à partir des données de Santé Publique?

# Synthèse & Conclusion

Conception de l'application QUALECO permettant de calculer 2 scores QUAL et ECO

Faisabilité à partir des données de Santé Publique?



- Variables pertinentes et liées en 2 groupes
- Définition de 2 axes possible
- Contribution des variables aux 2 axes

# Synthèse & Conclusion

Conception de l'application QUALECO permettant de calculer 2 scores QUAL et ECO

Faisabilité à partir des données de Santé Publique?

- Variables pertinentes et liées en 2 groupes
- Définition de 2 axes
- Contribution des variables aux 2 axes



**Calcul des 2 scores de l'application QUALECO faisable  
...mais peut encore être amélioré**





# MERCI

---

Questions



Incompréhensions

