

# Etude de faisabilité d'une application basée sur un moteur de recommandation

Tanguy Meyer

# Sommaire

1. Contexte
2. Gestion de projet
3. Analyse du jeu de données
  - 3.1 Etat des lieux
  - 3.2 Classification des données
  - 3.3 Choix des données
4. Nettoyage du jeu de données
  - 4.1 Processus de nettoyage
  - 4.2 Méthodes de nettoyage et de calcul

# 1. Contexte

- Foodflix : application permettant de recommander le meilleur produit à un utilisateur selon un mot clé ou un ensemble de mots clés
- Cahier des charges : Remonter les éléments liés au nutriscore
- Objectifs : Prouver que la donnée récupérée pourra être utilisée pour cette application
- Données : Open Food Facts, une base de données collaborative
- Démarche : Analyse et nettoyage de la donnée



## 2. Gestion de projet : Trello

The screenshot shows a Trello board for the 'open-food-facts' project. The board is organized into five main columns: 'À faire' (To do), 'En cours' (In progress), 'Bloqué' (Blocked), 'Terminé' (Completed), and 'Améliorations futures' (Future improvements). Each column contains a list of tasks or cards, some with progress bars indicating completion status. A sidebar on the left provides additional context with links, images, and documents related to the project.

**Project Resources**

- <https://cdn.discordapp.com/attachments/810826535744831519/822053715845447680/Infographie-calcul-nutriscore-logo-nutritionnel-score-corrigC3A9-scaled.png>
- Kaggle utile pour visualiser les données manquantes
- Plan à suivre
- Analyse ressource
- Exemple présentation
- Calcul energie
- Excel calcul nutriscore

**À faire**

- + Ajouter une carte

**En cours**

- Création de slides pour documenter le raisonnement (15-20 slides, 15 minutes de présentation)
- + Ajouter une autre carte

**Bloqué**

- + Ajouter une carte

**Terminé**

- Création environnement de travail
- Importation des data
- Etat des lieux des data brutes et analyse
- Analyse des notebooks existants sur Kaggle
- Recherche d'informations sur les notions métiers utiles (type nutriscore, yuka, éléments inconnus)
- Mise à jour du Trello
- Définition du besoin client
- Nettoyage des données
- + Ajouter une autre carte

**Améliorations futures**

- Ajouter des colonnes pour afficher les nutriments les + pertinents pour chaque aliment
- + Ajouter une autre carte

## 3. Analyse du jeu de données

### 3.1. Etat des lieux

- Nombre de lignes : 356027
- Nombre de colonnes : 163
- Liste des colonnes :

'code', 'url', 'creator', 'created\_t', 'created\_datetime', 'last\_modified\_t', 'last\_modified\_datetime', 'product\_name', 'generic\_name', 'quantity', 'packaging', 'packaging\_tags', 'brands', 'brands\_tags', 'categories', 'categories\_tags', 'categories\_en', 'origins', 'origins\_tags', 'manufacturing\_places', 'manufacturing\_places\_tags', 'labels', 'labels\_tags', 'labels\_en', 'emb\_codes', 'emb\_codes\_tags', 'first\_packaging\_code\_geo', 'cities', 'cities\_tags', 'purchase\_places', 'stores', 'countries', 'countries\_tags', 'countries\_en', 'ingredients\_text', 'allergens', 'allergens\_en', 'traces', 'traces\_tags', 'traces\_en', 'serving\_size', 'no\_nutriments', 'additives\_n', 'additives', 'additives\_tags', 'additives\_en', 'ingredients\_from\_palm\_oil\_n', 'ingredients\_from\_palm\_oil', 'ingredients\_from\_palm\_oil\_tags', 'ingredients\_that\_may\_be\_from\_palm\_oil\_n', 'ingredients\_that\_may\_be\_from\_palm\_oil', 'ingredients\_that\_may\_be\_from\_palm\_oil\_tags', 'nutrition\_grade\_uk', 'nutrition\_grade\_fr', 'pnns\_groups\_1', 'pnns\_groups\_2', 'states', 'states\_en', 'main\_category', 'main\_category\_en', 'image\_url', 'image\_small\_url', 'energy\_100g', 'energy-from-fat\_100g', 'fat\_100g', 'saturated-fat\_100g', '-butyric-acid\_100g', '-caproic-acid\_100g', '-caprylic-acid\_100g', '-capric-acid\_100g', '-lauric-acid\_100g', '-myristic-acid\_100g', '-palmitic-acid\_100g', '-stearic-acid\_100g', '-arachidic-acid\_100g', '-behenic-acid\_100g', '-lignoceric-acid\_100g', '-cerotic-acid\_100g', '-montanic-acid\_100g', '-melissic-acid\_100g', 'monounsaturated-fat\_100g', 'polyunsaturated-fat\_100g', 'omega-3-fat\_100g', '-alpha-linolenic-acid\_100g', '-eicosapentaenoic-acid\_100g', '-docosahexaenoic-acid\_100g', 'omega-6-fat\_100g', '-linoleic-acid\_100g', '-arachidonic-acid\_100g', '-gamma-linolenic-acid\_100g', '-dihomo-gamma-linolenic-acid\_100g', 'omega-9-fat\_100g', '-oleic-acid\_100g', '-elaidic-acid\_100g', '-gondoic-acid\_100g', '-mead-acid\_100g', '-erucic-acid\_100g', '-nervonic-acid\_100g', 'trans-fat\_100g', 'cholesterol\_100g', 'carbohydrates\_100g', 'sugars\_100g', '-sucrose\_100g', '-glucose\_100g', '-fructose\_100g', '-lactose\_100g', '-maltose\_100g', '-maltodextrins\_100g', 'starch\_100g', 'polyols\_100g', 'fiber\_100g', 'proteins\_100g', 'casein\_100g', 'serum-proteins\_100g', 'nucleotides\_100g', 'salt\_100g', 'sodium\_100g', 'alcohol\_100g', 'vitamin-a\_100g', 'beta-carotene\_100g', 'vitamin-d\_100g', 'vitamin-e\_100g', 'vitamin-k\_100g', 'vitamin-c\_100g', 'vitamin-b1\_100g', 'vitamin-b2\_100g', 'vitamin-pp\_100g', 'vitamin-b6\_100g', 'vitamin-b9\_100g', 'folates\_100g', 'vitamin-b12\_100g', 'biotin\_100g', 'pantothenic-acid\_100g', 'silica\_100g', 'bicarbonate\_100g', 'potassium\_100g', 'chloride\_100g', 'calcium\_100g', 'phosphorus\_100g', 'iron\_100g', 'magnesium\_100g', 'zinc\_100g', 'copper\_100g', 'manganese\_100g', 'fluoride\_100g', 'selenium\_100g', 'chromium\_100g', 'molybdenum\_100g', 'iodine\_100g', 'caffeine\_100g', 'taurine\_100g', 'ph\_100g', 'fruits-vegetables-nuts\_100g', 'fruits-vegetables-nuts-estimate\_100g', 'collagen-meat-protein-ratio\_100g', 'cocoa\_100g', 'chlorophyl\_100g', 'carbon-footprint\_100g', 'nutrition-score-fr\_100g', 'nutrition-score-uk\_100g', 'glycemic-index\_100g', 'water-hardness\_100g'

Nous allons classer les colonnes similaires pour rendre ça plus digeste !

### 3. Analyse du jeu de données

## 3.2 Classification des données

### Les métadonnées

- Donne des informations sur la donnée en elle-même
- Pas utile pour l'application
- Liste des métadonnées :

'code', 'url', 'creator', 'created\_t', 'created\_datetime', 'last\_modified\_t',  
'last\_modified\_datetime'

### 3. Analyse du jeu de données

## 3.2 Classification des données

### Les informations logistiques

- Donne des informations sur l'origine et la dénomination du produit
- Partiellement utile pour l'application afin d'identifier le produit
- Liste des informations logistiques :

'product\_name', 'generic\_name', 'quantity', 'packaging', 'packaging\_tags', 'brands', 'brands\_tags', 'categories', 'categories\_tags', 'categories\_en', 'origins', 'origins\_tags', 'manufacturing\_places', 'manufacturing\_places\_tags', 'labels', 'labels\_tags', 'labels\_en', 'emb\_codes', 'emb\_codes\_tags', 'first\_packaging\_code\_geo', 'cities', 'cities\_tags', 'purchase\_places', 'stores', 'countries', 'countries\_tags', 'countries\_en', 'pnns\_groups\_1', 'pnns\_groups\_2', 'main\_category', 'main\_category\_en'

Les colonnes finissant par \_tags et \_en reprennent les mêmes valeurs que les colonnes du même nom en modifiant la typographie ou la langue.

### 3. Analyse du jeu de données

## 3.2 Classification des données

### Les valeurs nutritionnelles textuelles

- Donne des informations textuelles sur les valeurs nutritionnelles du produit
- Contient des infos importantes pour l'utilisateur de l'application
- Liste des valeurs nutritionnelles textuelles :

'ingredients\_text', 'allergens', 'allergens\_en', 'traces', 'traces\_tags', 'traces\_en', 'serving\_size', 'no\_nutriments', 'additives\_n', 'additives', 'additives\_tags', 'additives\_en', 'ingredients\_from\_palm\_oil\_n', 'ingredients\_from\_palm\_oil', 'ingredients\_from\_palm\_oil\_tags', 'ingredients\_that\_may\_be\_from\_palm\_oil\_n', 'ingredients\_that\_may\_be\_from\_palm\_oil', 'ingredients\_that\_may\_be\_from\_palm\_oil\_tags', 'nutrition\_grade\_uk', 'nutrition\_grade\_fr'

Les colonnes finissant par \_tags et \_en reprennent les mêmes valeurs que les colonnes du même nom en modifiant la typographie ou la langue



### 3. Analyse du jeu de données

## 3.2 Classification des données

### Les valeurs nutritionnelles numériques

- Donne les valeurs nutritionnelles du produit pour 100 grammes
- Contient des infos importantes pour l'utilisateur et le nutriscore
- Liste des valeurs nutritionnelles numériques :

'energy\_100g', 'energy-from-fat\_100g', 'fat\_100g', 'saturated-fat\_100g', 'butyric-acid\_100g', 'caproic-acid\_100g', 'caprylic-acid\_100g', 'capric-acid\_100g', 'lauric-acid\_100g', 'myristic-acid\_100g', 'palmitic-acid\_100g', 'stearic-acid\_100g', 'arachidic-acid\_100g', 'behenic-acid\_100g', 'lignoceric-acid\_100g', 'cerotic-acid\_100g', 'montanic-acid\_100g', 'melissic-acid\_100g', 'monounsaturated-fat\_100g', 'polyunsaturated-fat\_100g', 'omega-3-fat\_100g', 'alpha-linolenic-acid\_100g', 'eicosapentaenoic-acid\_100g', 'docosahexaenoic-acid\_100g', 'omega-6-fat\_100g', 'linoleic-acid\_100g', 'arachidonic-acid\_100g', 'gamma-linolenic-acid\_100g', 'dihomo-gamma-linolenic-acid\_100g', 'omega-9-fat\_100g', 'oleic-acid\_100g', 'elaidic-acid\_100g', 'gondoic-acid\_100g', 'mead-acid\_100g', 'erucic-acid\_100g', 'nervonic-acid\_100g', 'trans-fat\_100g', 'cholesterol\_100g', 'carbohydrates\_100g', 'sugars\_100g', 'sucrose\_100g', 'glucose\_100g', 'fructose\_100g', 'lactose\_100g', 'maltose\_100g', 'maltodextrins\_100g', 'starch\_100g', 'polyols\_100g', 'fiber\_100g', 'proteins\_100g', 'casein\_100g', 'serum-proteins\_100g', 'nucleotides\_100g', 'salt\_100g', 'sodium\_100g', 'alcohol\_100g', 'vitamin-a\_100g', 'beta-carotene\_100g', 'vitamin-d\_100g', 'vitamin-e\_100g', 'vitamin-k\_100g', 'vitamin-c\_100g', 'vitamin-b1\_100g', 'vitamin-b2\_100g', 'vitamin-pp\_100g', 'vitamin-b6\_100g', 'vitamin-b9\_100g', 'folates\_100g', 'vitamin-b12\_100g', 'biotin\_100g', 'pantothenic-acid\_100g', 'silica\_100g', 'bicarbonate\_100g', 'potassium\_100g', 'chloride\_100g', 'calcium\_100g', 'phosphorus\_100g', 'iron\_100g', 'magnesium\_100g', 'zinc\_100g', 'copper\_100g', 'manganese\_100g', 'fluoride\_100g', 'selenium\_100g', 'chromium\_100g', 'molybdenum\_100g', 'iodine\_100g', 'caffeine\_100g', 'taurine\_100g', 'ph\_100g', 'fruits-vegetables-nuts\_100g', 'fruits-vegetables-nuts-estimate\_100g', 'collagen-meat-protein-ratio\_100g', 'cocoa\_100g', 'chlorophyll\_100g', 'carbon-footprint\_100g', 'nutrition-score-fr\_100g', 'nutrition-score-uk\_100g', 'glycemic-index\_100g', 'water-hardness\_100g'

### 3. Analyse du jeu de données

## 3.2 Classification des données

#### Les informations additionnelles

- Donne le statut de la ligne en question et un lien renvoyant vers une image du produit
- Le lien vers l'image peut être utile pour l'application
- Liste des informations additionnelles :

`'states', 'states_tags', 'states_en', 'image_url', 'image_small_url'`

### 3. Analyse du jeu de données

## 3.3 Choix des données

- Les données choisies pour être utilisées dans le nettoyage doivent :
  - Etre pertinentes pour l'utilisateur et l'application :

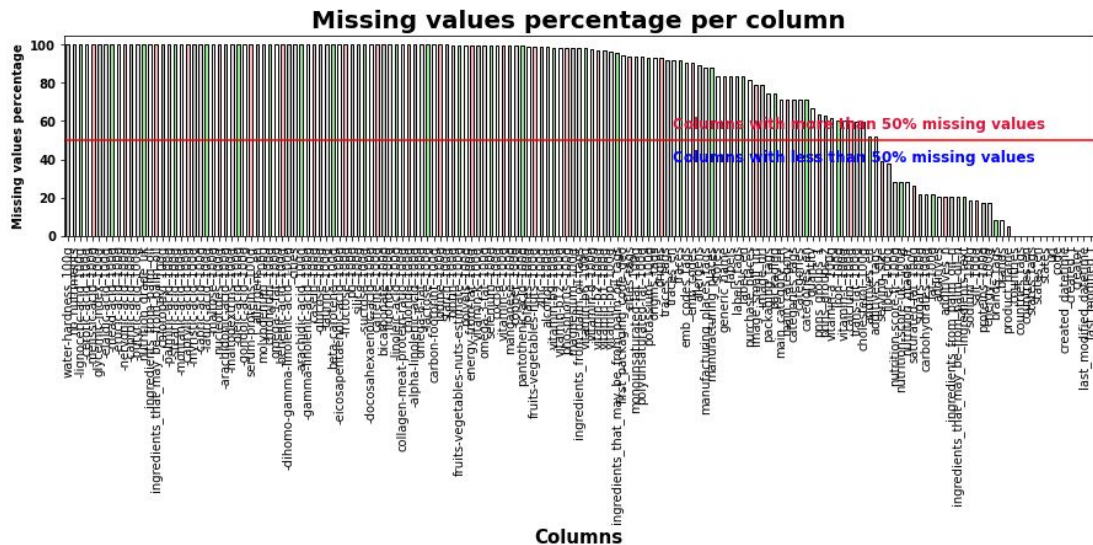
Valeurs nutritionnelles les plus répandues, allergènes, présence d'huile de palme, nutriscore

- Permettre de calculer le nutriscore
- Permettre de compléter des valeurs manquantes (ces colonnes pourront être supprimées à la fin du nettoyage)

# 3. Analyse du jeu de données

## 3.3 Choix des données

Valeurs manquantes :



Beaucoup de valeurs pas renseignées

100% de valeurs manquantes pour :

'ingredients\_that\_may\_be\_from\_palm\_oil', 'ingredients\_from\_palm\_oil', 'nutrition\_grade\_uk', '-butyric-acid\_100g', '-caproic-acid\_100g', '-nervonic-acid\_100g', '-erucic-acid\_100g', '-mead-acid\_100g', '-elaidic-acid\_100g', 'glycemic-index\_100g', '-melissic-acid\_100g', '-cerotic-acid\_100g', '-lignoceric-acid\_100g', 'no\_nutriments', 'water-hardness\_100g'

### 3. Analyse du jeu de données

## 3.3 Choix des données

#	Code	Type	Description	Raison
1	<b>product_name</b>	Object	Nom du produit	Pour l'application
2	<b>generic_name</b>	Object	Nom alternatif	Pour le nettoyage
3	<b>brands</b>	Object	Marque du produit	Pour l'application
4	<b>countries_en</b>	Object	Pays de distribution	Pour l'application et le nettoyage
5	<b>pnns_groups_1</b>	Object	Catégorie du produit	Pour l'application et le nettoyage
6	<b>ingredients_text</b>	Object	Liste des ingrédients	Pour l'application
7	<b>allergens</b>	Object	Liste des allergènes	Pour l'application
8	<b>traces</b>	Object	Liste des traces d'allergènes	Pour l'application
9	<b>additives_en</b>	Object	Liste des additifs	Pour l'application

### 3. Analyse du jeu de données

## 3.3 Choix des données

#	Code	Type	Description	Raison
10	<b>ingredients_from_palm_oil_n</b>	Object	Présence d'huile de palme	Pour l'application
11	<b>nutrition_grade_fr</b>	Object	Nutrigrade	Pour le nettoyage
12	<b>fat_100g</b>	Float	Graisses pour 100g	Pour l'application
13	<b>carbohydrates_100g</b>	Float	Glucides pour 100g	Pour l'application
14	<b>energy_100g</b>	Float	Energie pour 100g	Pour le nutriscore
15	<b>saturated-fat_100g</b>	Float	Graisses saturées pour 100g	Pour le nutriscore
16	<b>sugars_100g</b>	Float	Sucre pour 100g	Pour le nutriscore
17	<b>fiber_100g</b>	Float	Fibres pour 100g	Pour le nutriscore
18	<b>salt_100g</b>	Float	Sel pour 100g	Pour le nutriscore

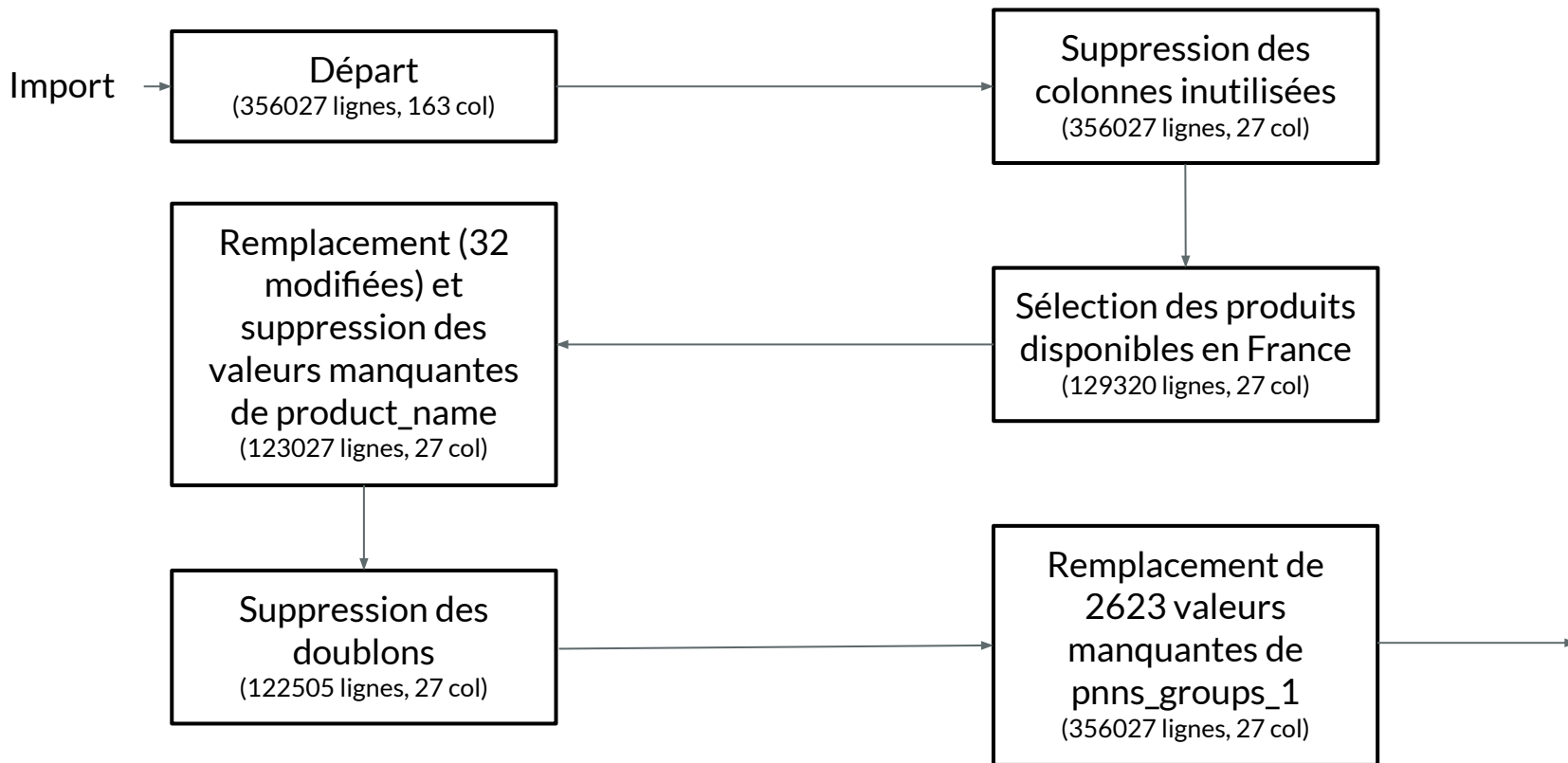
### 3. Analyse du jeu de données

## 3.3 Choix des données

#	Code	Type	Description	Raison
19	proteins_100g	Float	Protéines pour 100g	Pour le nutriscore
20	nutrition-score-fr_100g	Float	Valeur du nutriscore	Pour l'application
21	fruits-vegetables-nuts_100g	Float	Taux de fruits, légumes et noix	Pour le nutriscore
22	vitamin-d_100g	Float	Vitamine D pour 100g	Pour l'application
23	vitamin-c_100g	Float	Vitamine C pour 100g	Pour l'application
24	calcium_100g	Float	Calcium pour 100g	Pour l'application
25	iron_100g	Float	Fer pour 100g	Pour l'application
26	image_url	Object	Lien de l'image du produit	Pour l'application
27	pnns_groups_2	Object	Catégorie du produit alternative	Pour le nettoyage

## 4. Nettoyage du jeu de données

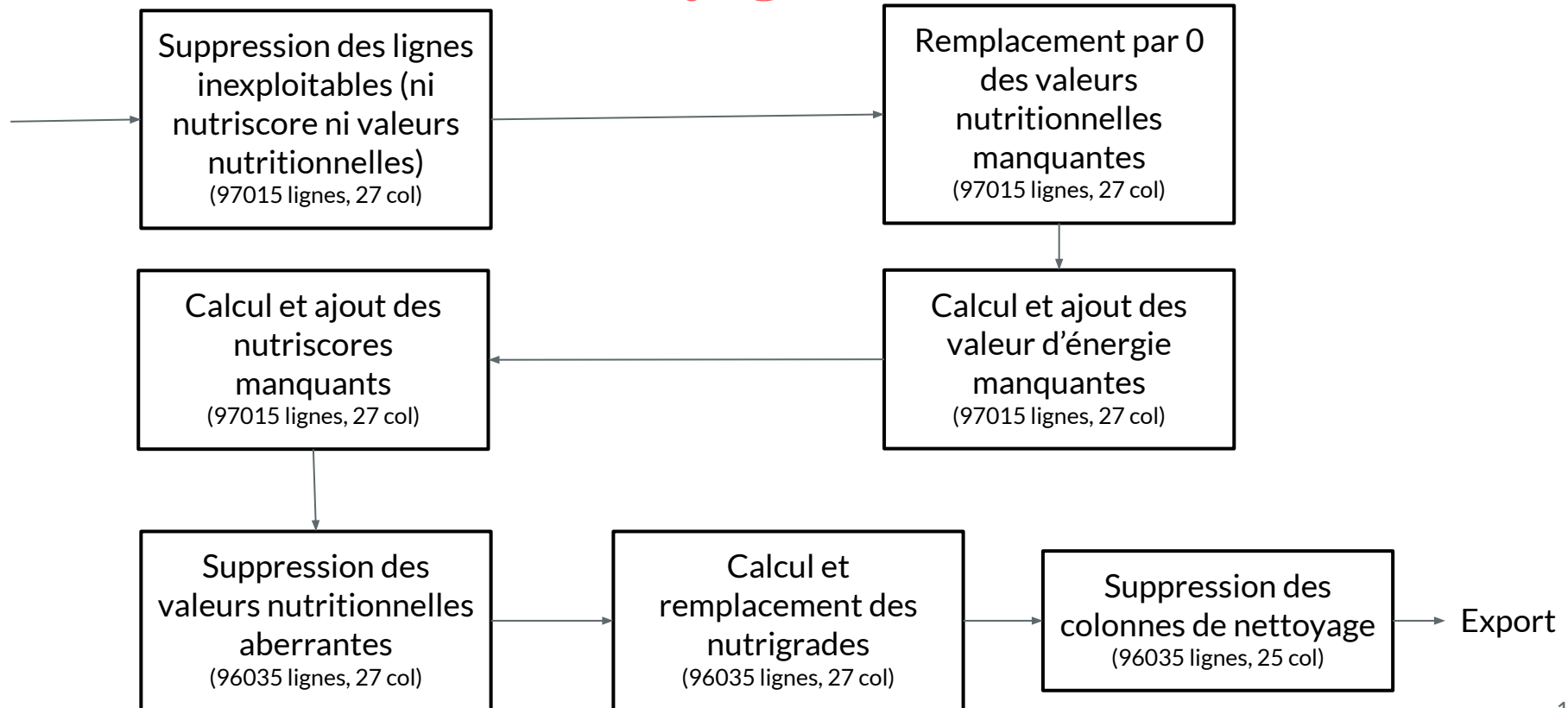
### 4.1 Processus de nettoyage





## 4. Nettoyage du jeu de données

### 4.1 Processus de nettoyage

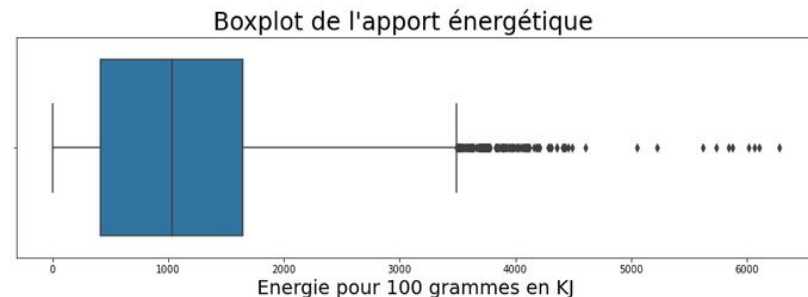


## 4. Nettoyage du jeu de données

### 4.2 Méthodes de nettoyage et de calcul

#### - Energie :

```
# E KJ = (37 x lipides) + (17 x protéines) + (17 x glucides) + (8 x fibres)
# On calcule l'énergie suivant la formule ci-dessus quand la valeur est manquante
df["energy_100g"] = df.apply(
    lambda row: 37*row['fat_100g'] + 17*row['proteins_100g'] + 17*row['carbohydrates_100g'] + 8*row['fiber_100g']
    if np.isnan(row["energy_100g"])
    else row["energy_100g"],
    axis=1)
```



#### - Nutrigrade :

```
# Liste des valeurs de nutrigrade
grad=['a', 'b', 'c', 'd', 'e']
# Valeurs limites de nutriscore pour les solides
scoS=[-1, 2, 10, 19]
# Valeurs limites de nutriscore pour les liquides
scoL=[1, 5, 9]
```

```
def nutrigrade(catg, nutriscore):
    """Calcule le nutrigrade à partir du nutriscore"""
    # Les critères ne sont pas les mêmes pour les liquides et les solides
    if "beverage" in str(catg).lower():
        # L'eau a toujours un nutrigrade de a
        if "eau" in str(catg).lower():
            return grad[0]
        return grad[calcul_points(nutriscore, scoL) + 1]
    return grad[calcul_points(nutriscore, scoS)]
```

```
# On calcule le nutrigrade pour chaque valeur manquante
df["nutrition_grade_fr"] = df.apply(
    lambda row: nutrigrade(row["pnns_groups_1"], row["nutrition-score-fr_100g"]),
    axis=1)
```

## 4. Nettoyage du jeu de données

# 4.2 Méthodes de nettoyage et de calcul

### - Nutriscore :

```
: # Mise en place des listes contenant les limites pour chaque valeur nutritionnelle
# Limites spécifiques aux solides
ptsNrjS=[335, 670, 1005, 1340, 1675, 2010, 2345, 2680, 3015, 3350]
ptsSucS=[4.5, 9, 13.5, 18, 22.5, 27, 31, 36, 40, 45]
ptsLegS=[40,60,80,80,80]

# Limites partagées entre solides et liquides
ptsSod=[0.090, 0.180, 0.270, 0.360, 0.450, 0.540, 0.630, 0.720, 0.810, 0.900]
ptsGra=[1, 2.1, 3.2, 4.3, 5.4, 6.5, 7.6, 8.7, 9.8, 10.9]
ptsFib=[0.9, 1.9, 2.8, 3.7, 4.7]
ptsPro=[1.6, 3.2, 4.8, 6.4, 8.0]

# Limites spécifiques aux liquides
ptsNrjL=[0, 30, 60, 90, 120, 150, 180, 210, 240, 270]
ptsSucL=[0, 1.5, 3, 4.5, 6, 7.5, 9, 10.5, 12, 13.5]
ptsLegL=[40, 40, 60, 80, 80, 80, 80, 80, 80, 80]

def calcul_points(val, pts):
    """Donne le nombre de points de la valeur val par rapport à son placement relatif à la liste pts"""
    points=0
    # On parcourt la liste des limites jusqu'à trouver une limite majorant notre valeur nutritionnelle
    for points in range(len(pts)):
        if val <= pts[points]:
            return points
    return points+1

def nutriscore(catg, nrj, suc, gra, sel, leg, fib, pro):
    """Calcule le nutriscore à partir des valeurs nutritionnelles"""
    # Compteur de points positifs
    p=0
    # Compteur de points négatifs
    n=0
    p = calcul_points(fib, ptsFib)
    n = n + calcul_points(gra, ptsGra) + calcul_points(sel/2.54, ptsSod)
    # Cas où le produit est un liquide
    if "beverage" in str(catg).lower():
        p = p + calcul_points(leg, ptsLegL)
        n = n + calcul_points(nrj, ptsNrjL) + calcul_points(suc, ptsSucL)
    # Cas où le produit est un solide
    else:
        p = p + calcul_points(leg, ptsLegS)
        n = n + calcul_points(nrj, ptsNrjS) + calcul_points(suc, ptsSucS)
    # Cas où les points des protéines ne sont pas pris en compte
    if n >= 11 and leg < 100 and "fromage" not in str(catg).lower():
        return(n-p)
    # Cas où les points des protéines sont pris en compte
    p = p + calcul_points(pro, ptsPro)
    return (n-p)

: # On calcule les nutriscores qui ne sont pas renseignés
df["nutrition-score-fr_100g"] = df.apply(
    lambda row: nutriscore(row["pnns-groups_1"], row["energy_100g"],
                           row["sugars_100g"], row["saturated-fat_100g"],
                           row["salt_100g"], row["fruits-vegetables-nuts_100g"],
                           row["fiber_100g"], row["proteins_100g"]),
    axis=1)
if np.isnan(row["nutrition-score-fr_100g"]):
    else row["nutrition-score-fr_100g"]
axis=1)
```