



Santé  
publique  
France

# Concevez une application au service de la santé publique

Recommandation Alimentaire  
pour Maladies Héritaire du Métabolisme

Source :



<https://world.openfoodfacts.org/>



Idée d'application



Nettoyage des données



Analyse des données



Faisabilité de l'application



Conclusion



## **Idée d'application**



## Nettoyage des données



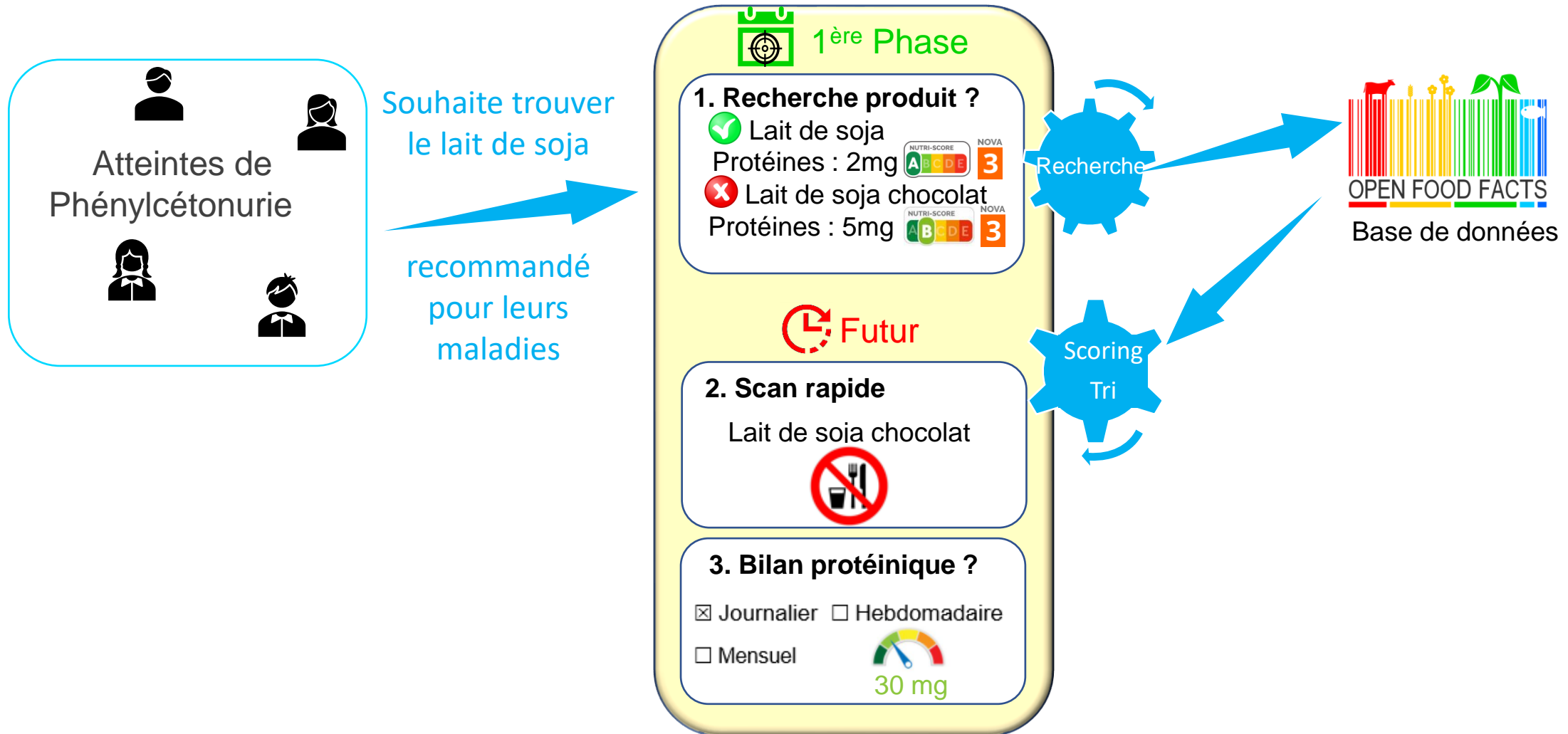
## Analyse des données



## Faisabilité de l'application



## Conclusion



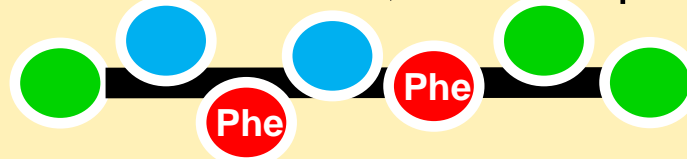
Moteur de **R**ecommandation **A**limentaire pour les personnes atteintes de **P**hénylcétonurie (**M**aladies **M**étaboliques **H**éréditaires)

# Idée d'application – La phénylcétonurie ?

Diagnostic prénatal, test de Guthrie



Alimentation normale, riche en protéines



Molécules de protéine composées d'acides aminés

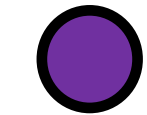


Digestion des aliments



Fragmentation en acides aminés

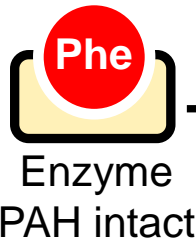
Transformation de la Phe en tyrosine à l'aide de PAH



Tyrosine



Phe



Saine

Personne?

Atteinte PCU

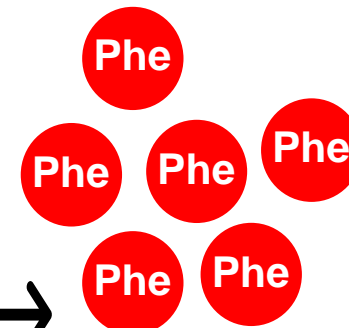


PAH déficiente



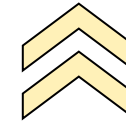
Phe

Phe non transformée en tyrosine, ni assimilée



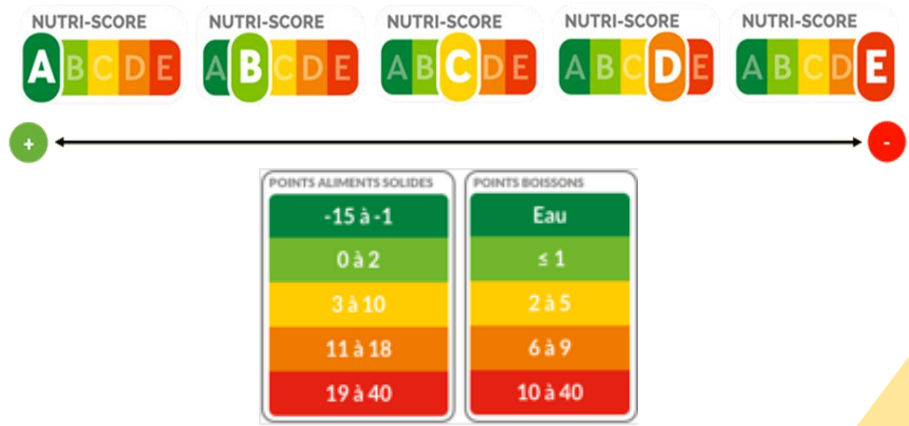
Taux sanguin de Phe élevé : intoxication

Troubles du développement cérébral et dommages physiques et mentaux irréversibles

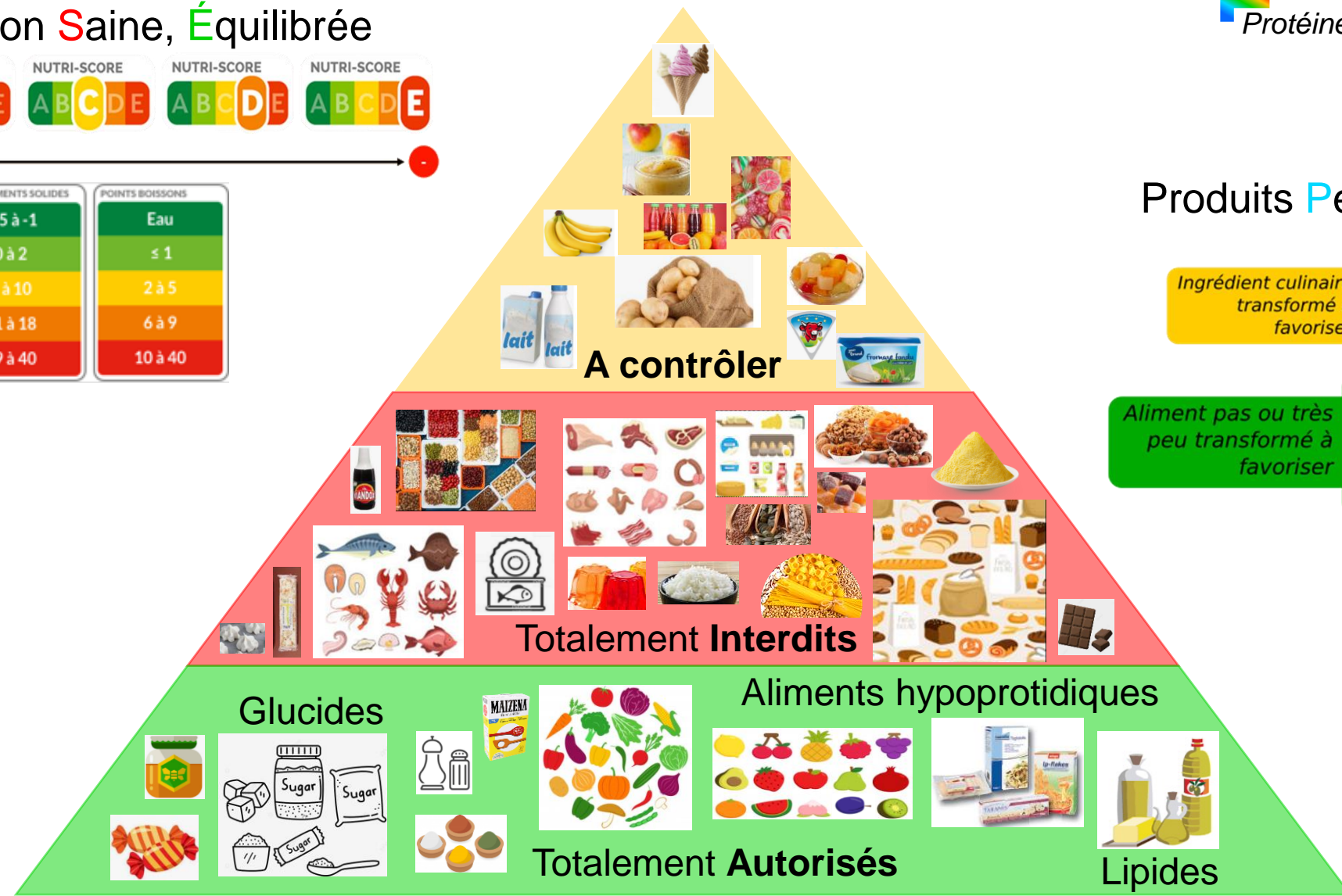


# Idée d'application – Alimentation ?

## Alimentation Saine, Équilibrée



Protéines : 2g/kg/Jour (> 4ans)



## Produits Peu Transformés



## Régime hypoprotidique

- Peu de protéines

## Alimentation saine

- Nutri-score faible
- Peu de sel
- Peu d'additifs
- Pas huile de palme
- Sucre limité
- Peu calorique

## Alimentation équilibrée

- Macro-nutriments
  - Lipides
  - Glucides
- Micro-nutriments
  - Vitamines
  - minéraux

## Peu transformés

- Groupe NOVA

## Produits disponibles

- Nom du produit
- Photo
- Catégorie
- Marque
- Vendus en France



Idée d'application



**Nettoyage des données**



Analyse des données



Faisabilité de l'application

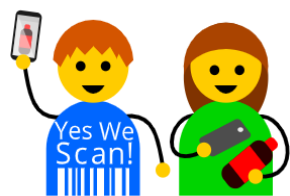


Conclusion

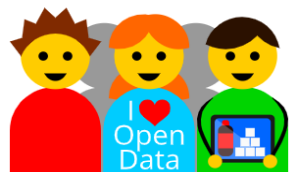




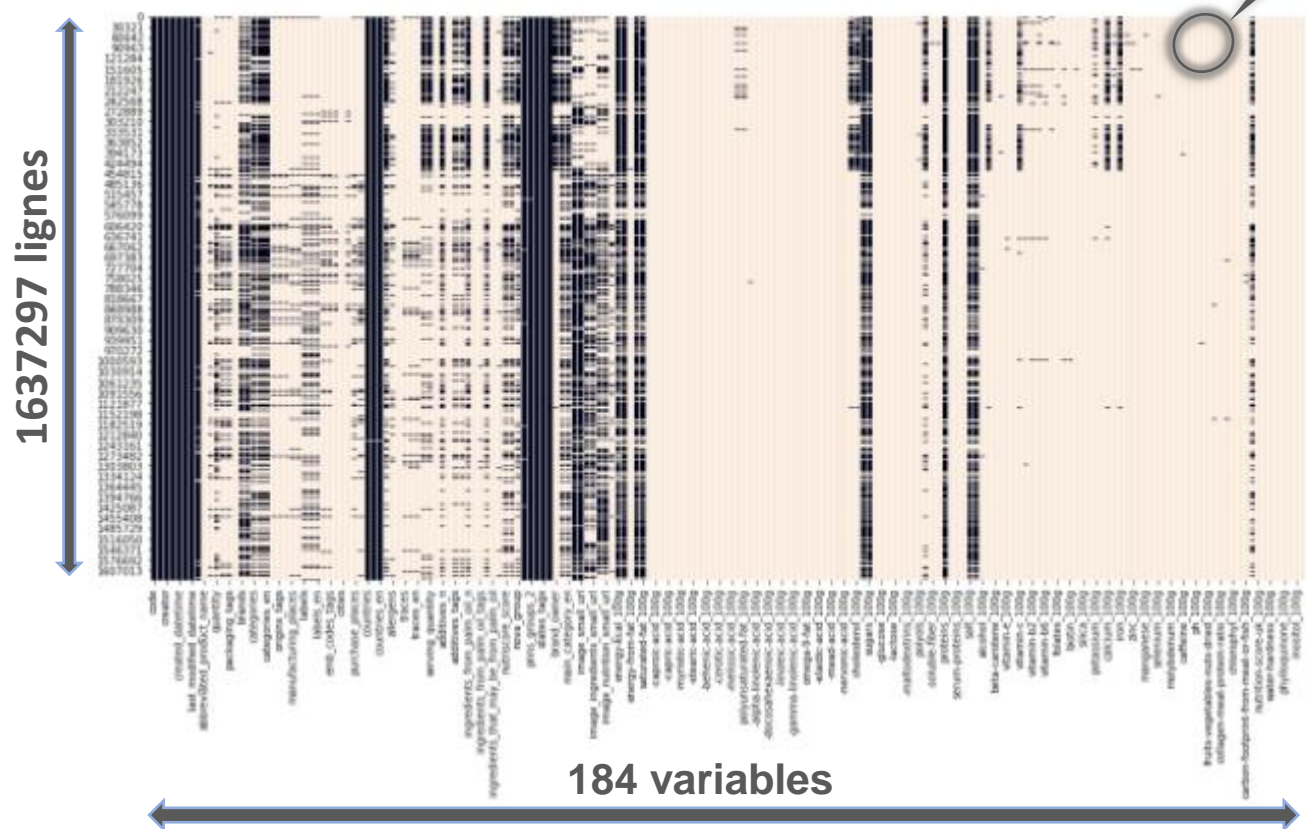
Une base de produits alimentaires



Faite par tout le monde



Pour tout le monde



> 79% de valeurs manquantes  
118 variables > 90% de valeurs manquantes

5 sections

- Informations générales
- tags
- ingrédients
- Données diverses
- Infos nutritionnelles

# Jeu de données - variables utiles

Régime hypoprotidique

proteins\_100g

Alimentation saine

nutriscore\_score  
 nutriscore\_grade  
 sodium\_100g  
 salt\_100g  
 additives\_n  
 additives\_tags  
 additives\_en  
 ingredients\_from\_palm\_oil\_n  
 ingredients\_that\_may\_be\_from\_palm\_oil\_n  
 sugars\_100g  
 energy-kcal\_100g  
 energy\_100g  
 nutrition-score-fr\_100g  
 additives  
 ingredients\_from\_palm\_oil\_tags  
 ingredients\_that\_may\_be\_from\_palm\_oil  
 ingredients\_that\_may\_be\_from\_palm\_oil\_tags  
 energy-kj\_100g  
 energy-from-fat\_100g  
 nutrition-score-uk\_100g  
 allergens  
 ingredients\_from\_palm\_oil  
 allergens\_en

Alimentation équilibrée

fat\_100g  
 saturated-fat\_100g  
 carbohydrates\_100g  
 fiber\_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  
 -nervonic-acid\_100g  
 -sucrose\_100g  
 zinc\_100g  
 copper\_100g  
 manganese\_100g  
 -lactose\_100g  
 -maltose\_100g  
 -maltodextrins\_100g  
 starch\_100g  
 polyols\_100g  
 -soluble-fiber\_100g  
 -insoluble-fiber\_100g  
 casein\_100g  
 serum-proteins\_100g  
 nucleotides\_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  
 magnesium\_100g  
 fluoride\_100g  
 selenium\_100g  
 chromium\_100g  
 molybdenum\_100g  
 iodine\_100g  
 caffeine\_100g  
 taurine\_100g  
 ph\_100g  
 fruits-vegetables-nuts\_100g  
 fruits-vegetables-nuts-dried\_100g  
 fruits-vegetables-nuts-estimate\_100g  
 collagen-meat-protein-ratio\_100g  
 cocoa\_100g  
 chlorophyl\_100g  
 glycemic-index\_100g  
 water-hardness\_100g  
 choline\_100g  
 phylloquinone\_100g  
 beta-glucan\_100g  
 inositol\_100g  
 -glucose\_100g  
 -fructose\_100g  
 iron\_100g  
 magnesium\_100g  
 cholesterol\_100g  
 trans-fat\_100g  
 omega-9-fat\_100g  
 -oleic-acid\_100g  
 -elaidic-acid\_100g  
 -gondoic-acid\_100g  
 -mead-acid\_100g  
 -erucic-acid\_100g  
 carnitine\_100g

Peu transformés

nova\_group

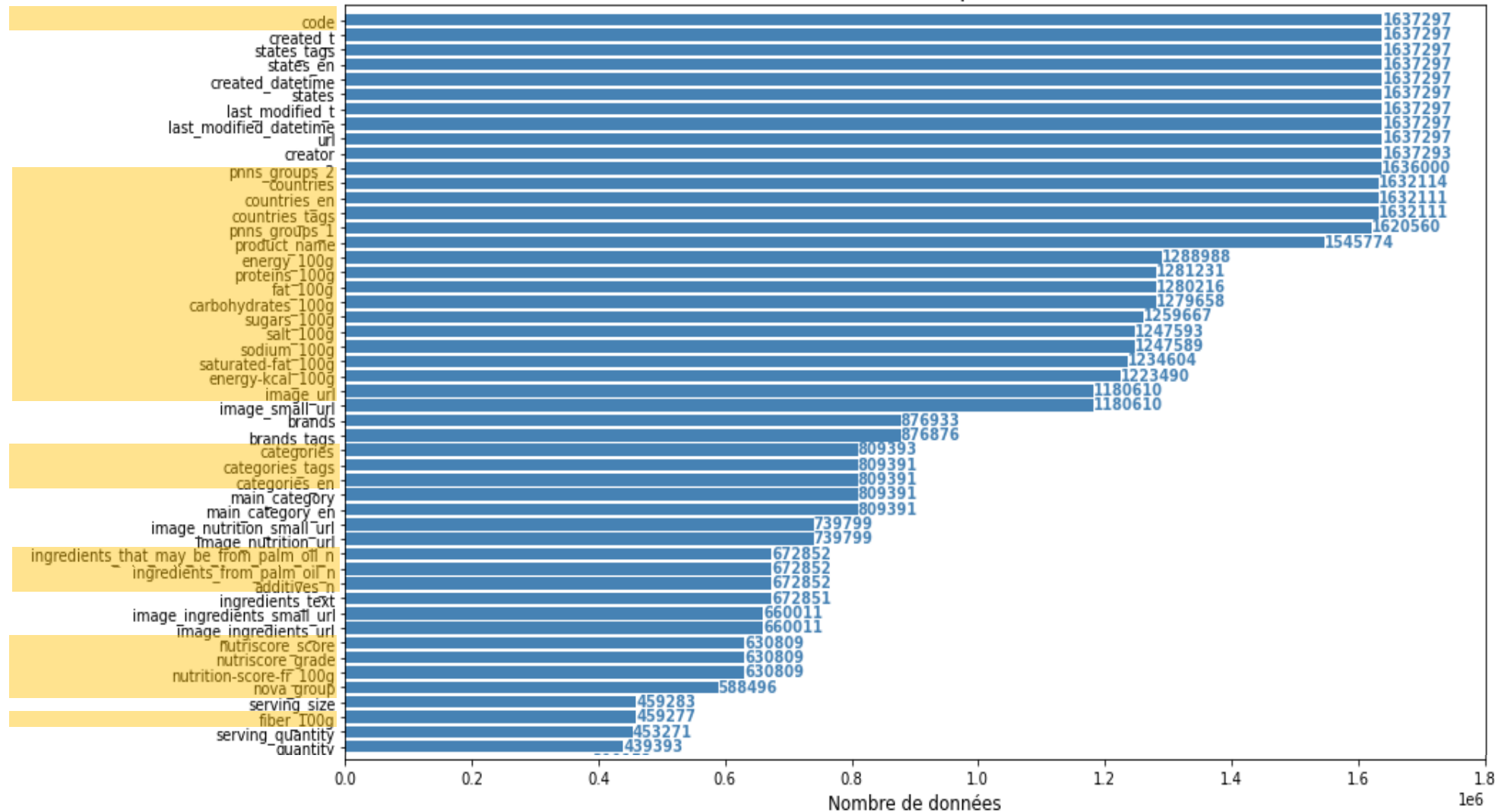
Produits disponibles

code  
 product\_name  
 brands  
 brands\_tags  
 categories  
 categories\_tags  
 categories\_en  
 countries  
 countries\_tags  
 countries\_en  
 pnns\_groups\_1  
 pnns\_groups\_2  
 image\_url  
 generic\_name

Légende :  
 En\_marron : >80% NaN  
 En\_rouge : 100% NaN

# Jeu de données - variables utiles

Nombre de données par variables

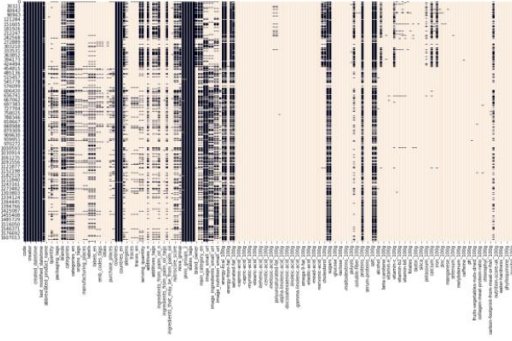


- Régime hypoprotidique ✓
- Alimentation saine ✓
- Alimentation équilibrée ✗
- Peu transformés ✓
- Produits disponibles ✓

# Nettoyage – réduction du jeu de données

De...

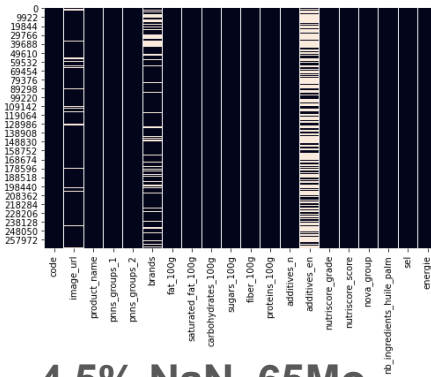
1637297 lignes, 184 variables



79,5% NaN, 3,5Go

... À

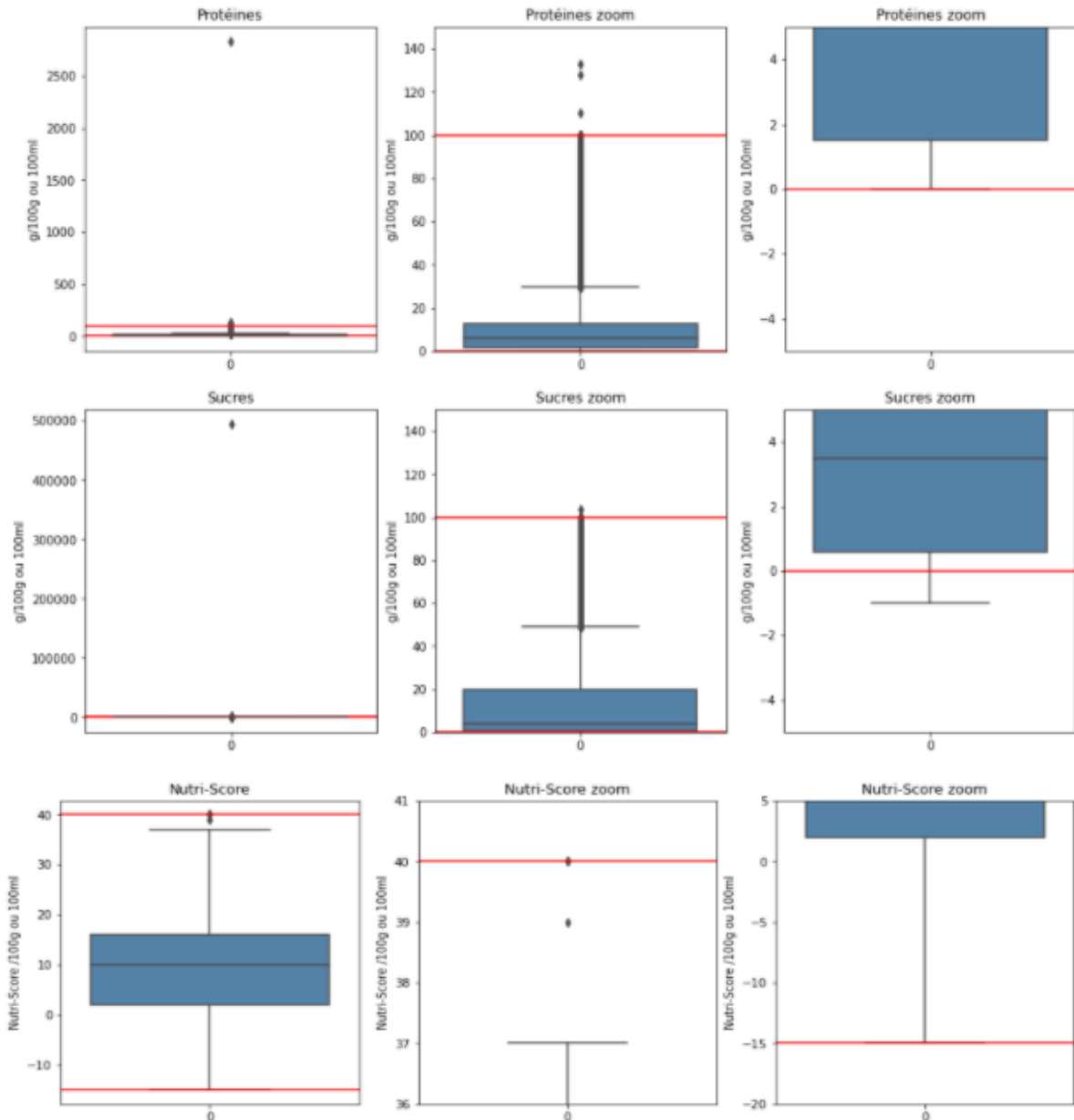
267881 lignes, 20 variables



4,5% NaN, 65Mo

	Nb lignes	Nb var.
Chargement des données	1637297	184
Suppression des 9 variables vides	1637297	175
Remplacer '-' par '_' dans le nom des variables Suppression des 2 doublons sur le code produit	1637295	175
Sélection des variables pertinentes pour l'application	1637295	124
Suppression des variables avec %NaN > 80%	1637295	24
Suppression des lignes sans nom produit	1545773	24
Conservation des produits vendus en France Traduction des groupes et sous-groupes de produit	748948	24
Gestion des valeurs aberrantes 🔍	748703	23
Feature engineering huile de palme	748703	22
Feature engineering sel	748703	21
Feature engineering énergie	748703	20
Gestion des valeurs manquantes 🔍	267881	20

# Nettoyage – Valeurs aberrantes



## Outliers pour 100g ou ml de produit

énergie nutritionnelle > 900 kcal ou > 3766 kJ

nutriments \_100g > 100g (ou 100ml)

nutriments \_100g < 0g (ou 0ml)

le score du nutri-score au dessus de 40

## Macro-nutriment et sous-groupe de nutriment

si la masse totale des glucides inférieure à la masse de sucre

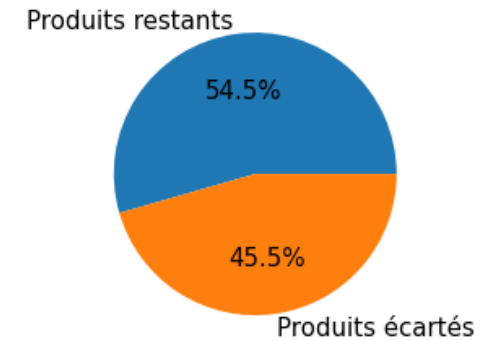
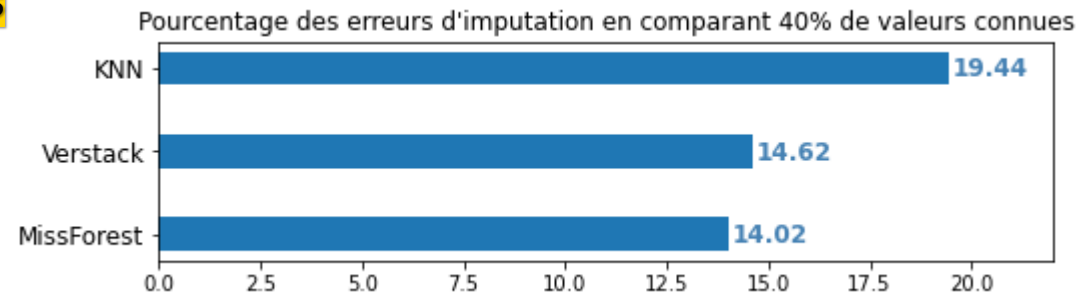
si la masse total des lipides est inférieure à la masse des acides gras saturés.

si la masse de sodium est supérieure à la masse de sel (en g).

# Nettoyage – Valeurs manquantes

1 **Valeurs obligatoires** suppression des produits sans les valeurs obligatoires (application de santé)

2 **Imputations**



Nombres de valeurs manquantes % de valeurs manquantes

additives_en	176027	65.71000
fiber_100g	168919	63.06000
nova_group	119060	44.45000
additives_n	106190	39.64000
nb_ingredients_huile_palm	106190	39.64000
brands	51438	19.20000
image_url	14431	5.39000
saturated_fat_100g	144	0.05000

Imputation par 0

Imputation par MissForest

Pas d'imputation



Idée d'application



Nettoyage des données



**Analyse des données**



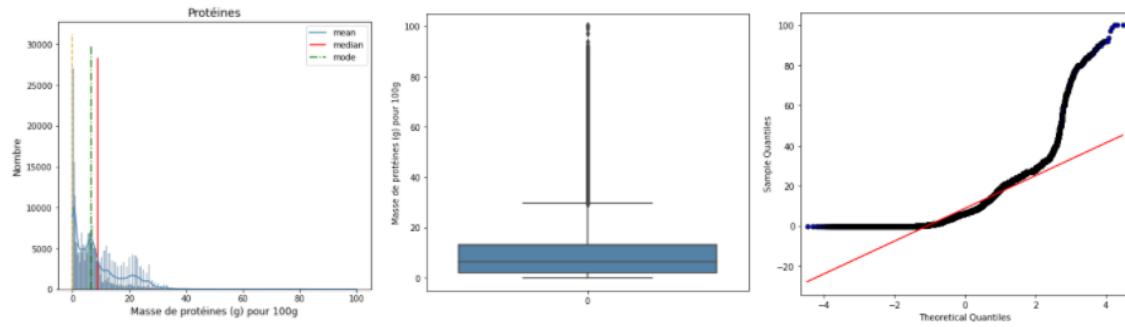
Faisabilité de l'application



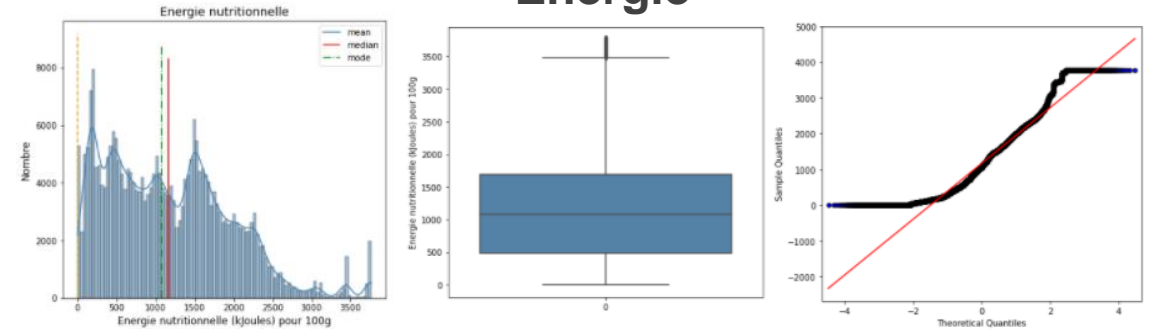
Conclusion



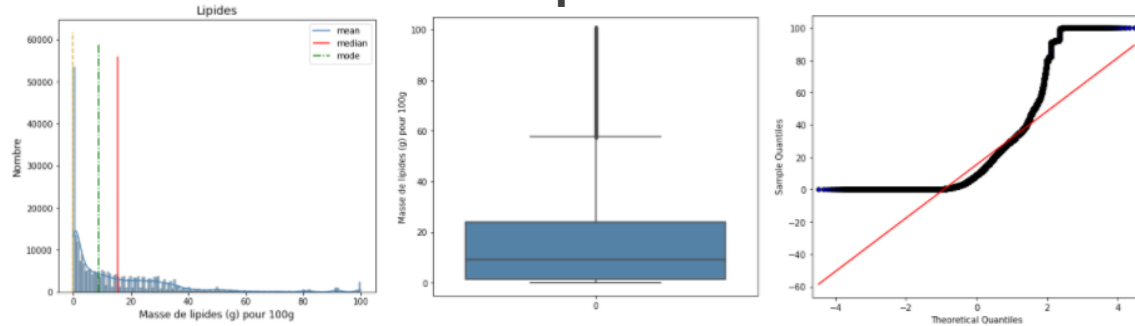
## Protéines



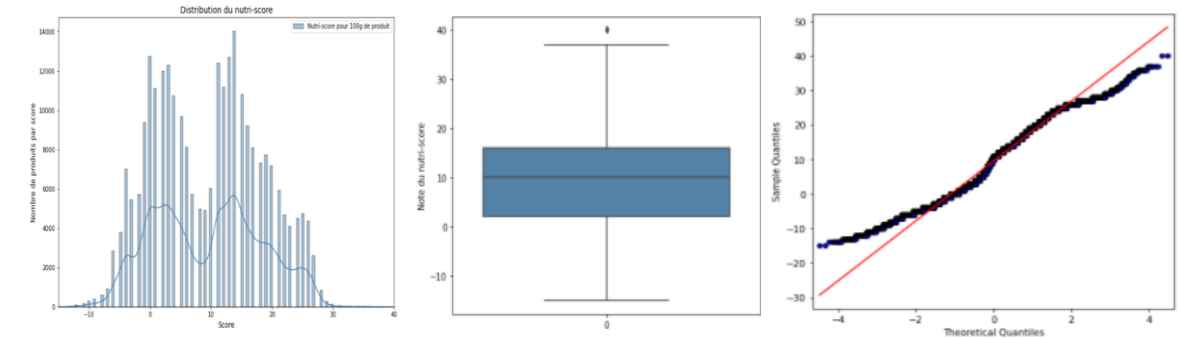
## Energie



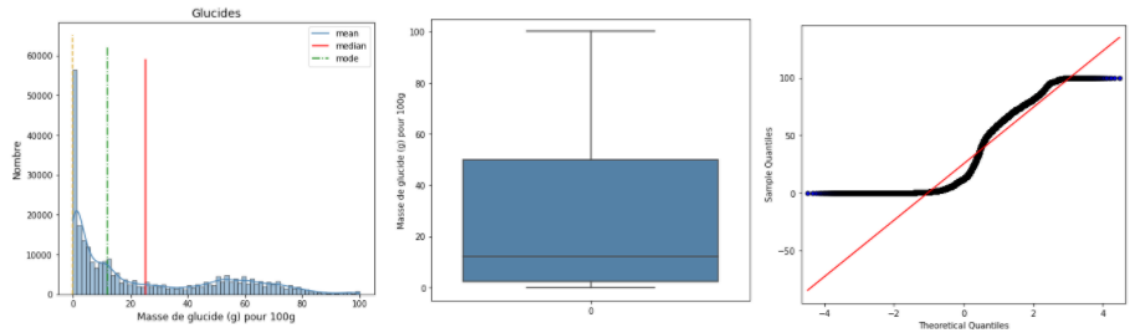
## Lipides



## Nutri-score



## Glucides



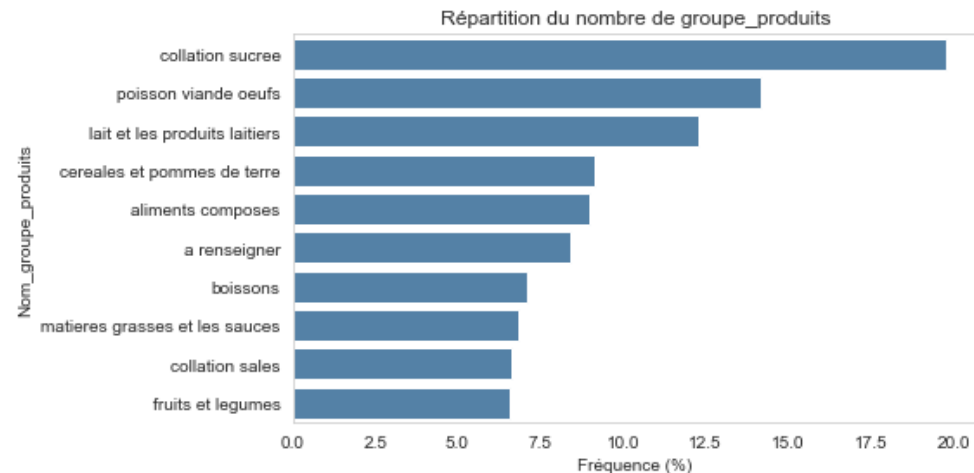
Desc	fat_100g	saturated_fat_100g	carbohydrates_100g	sugars_100g	fiber_100g	proteins_100g	sel	energie	nutriscore_score
mean	15.42544	5.89300	25.20521	12.67856	1.05942	8.81753	1.02137	1168.11084	9.50895
median	9.00000	2.30000	12.00000	3.30000	0.00000	6.40000	0.60000	1079.00000	10.00000
var	352.80405	73.00647	711.97937	347.02606	9.13648	79.57175	5.69719	635797.68750	76.91982
std	18.78308	8.54438	26.68294	18.62864	3.02266	8.92030	2.38688	797.36920	8.77039
skew	2.06164	3.22234	0.81542	1.86204	8.10557	1.92781	17.69135	0.65176	0.10425
kurtosis	5.30530	19.15100	-0.69722	3.04304	130.59885	7.77652	519.88153	0.11710	-0.93296
mode	0 0.0	0 0.0	0 0.0	0 0.0	0 0.0	0 0.0	0 0.0	0 0.0	0 14.0
Min	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	-15.00000
Max	100.00000	100.00000	100.00000	100.00000	100.00000	100.00000	100.00000	3766.00000	40.00000



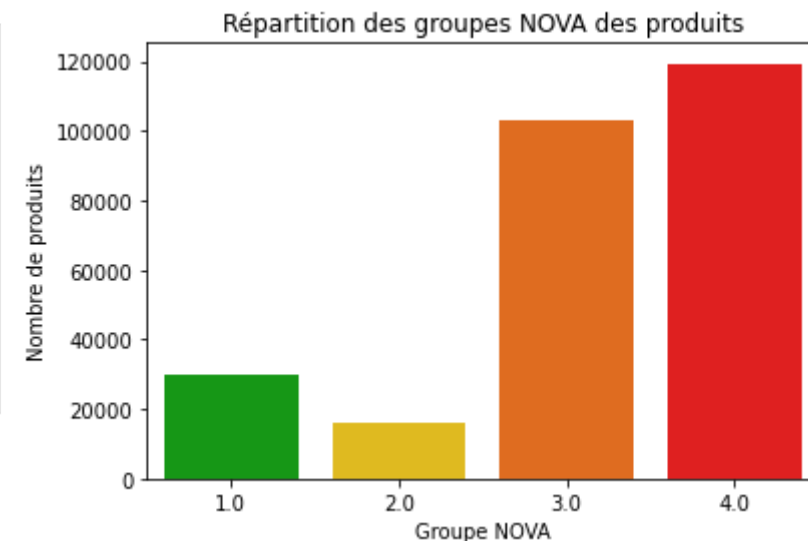
# Analyse univariée – les produits

## Groupe de produits

matieres grasses et les sauces  
poisson viande oeufs  
lait et les produits laitiers  
a renseigner fruits et legumes  
cereales et pommes de terre  
boissons  
aliments composes  
collation sucee

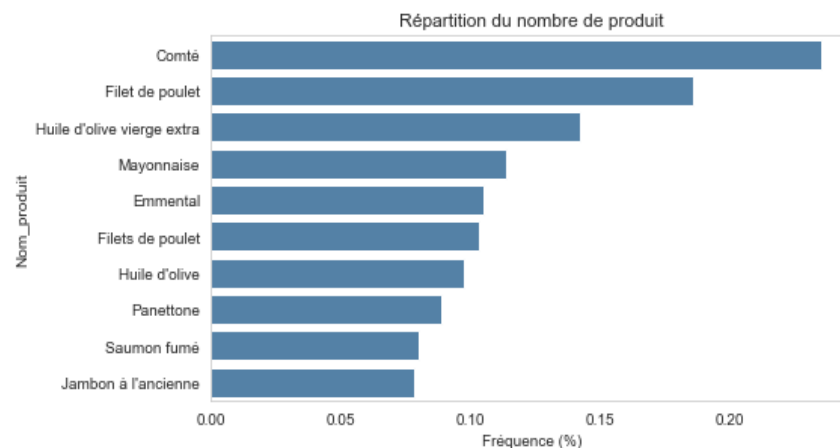


## Groupe Nova

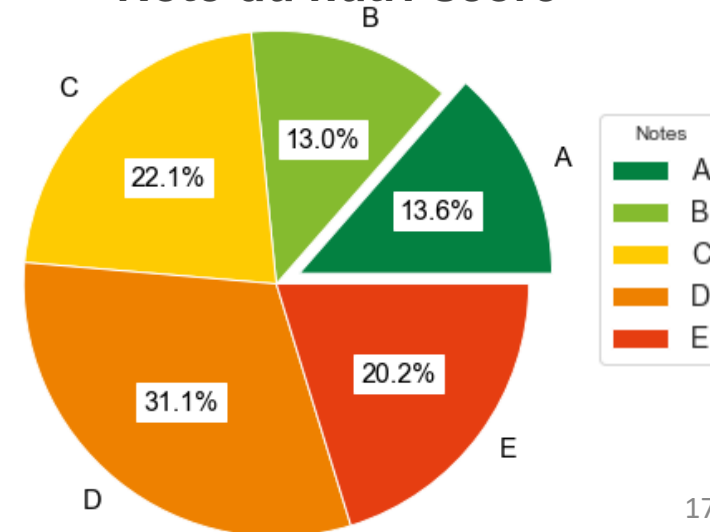


## Produits

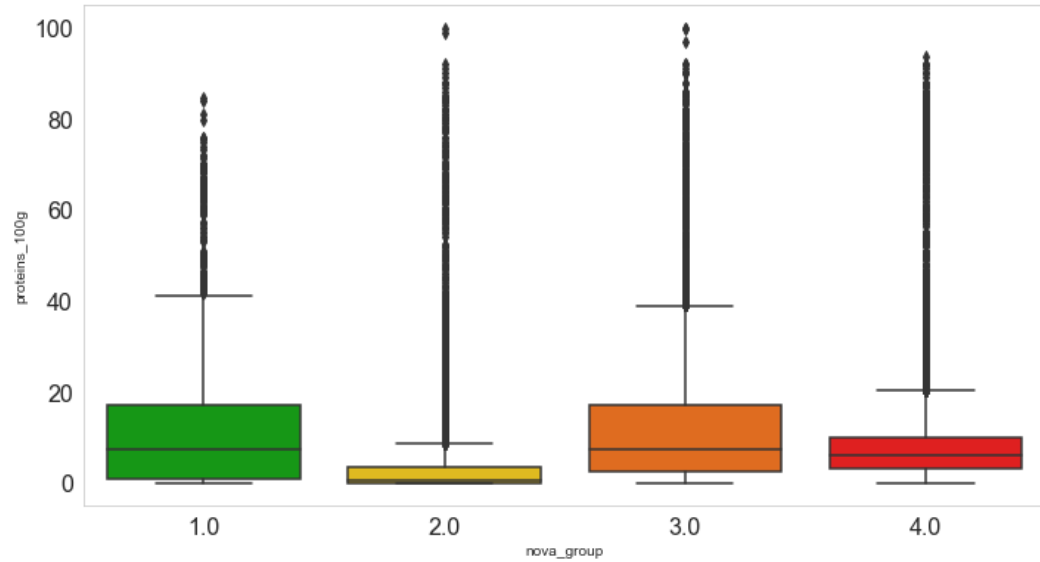
Jus de pomme Mozzarella Panettone Jus d'orange  
Filet de poulet  
Filets de poulet Huile d'olive  
Huile d'olive vierge extra  
Saumon fumé  
Comté  
Jambon à l'ancienne



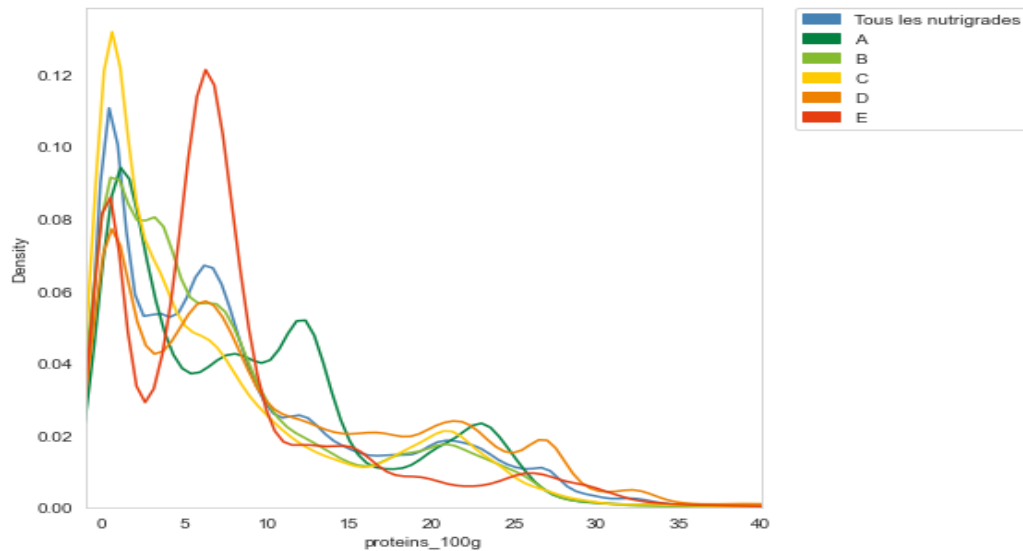
## Note du nutri-score



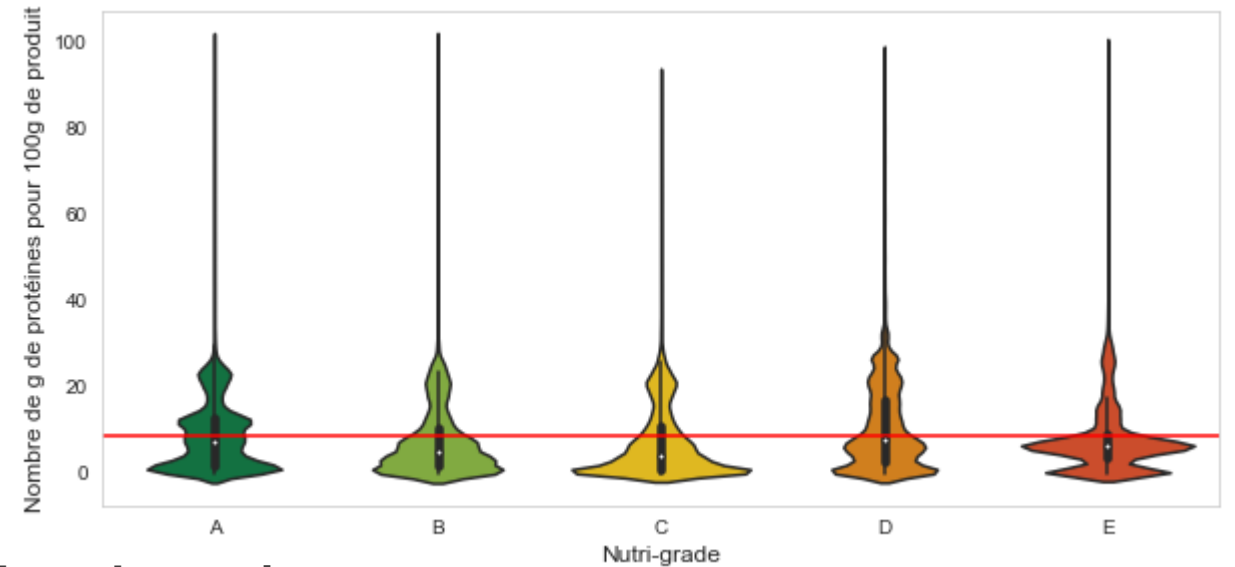
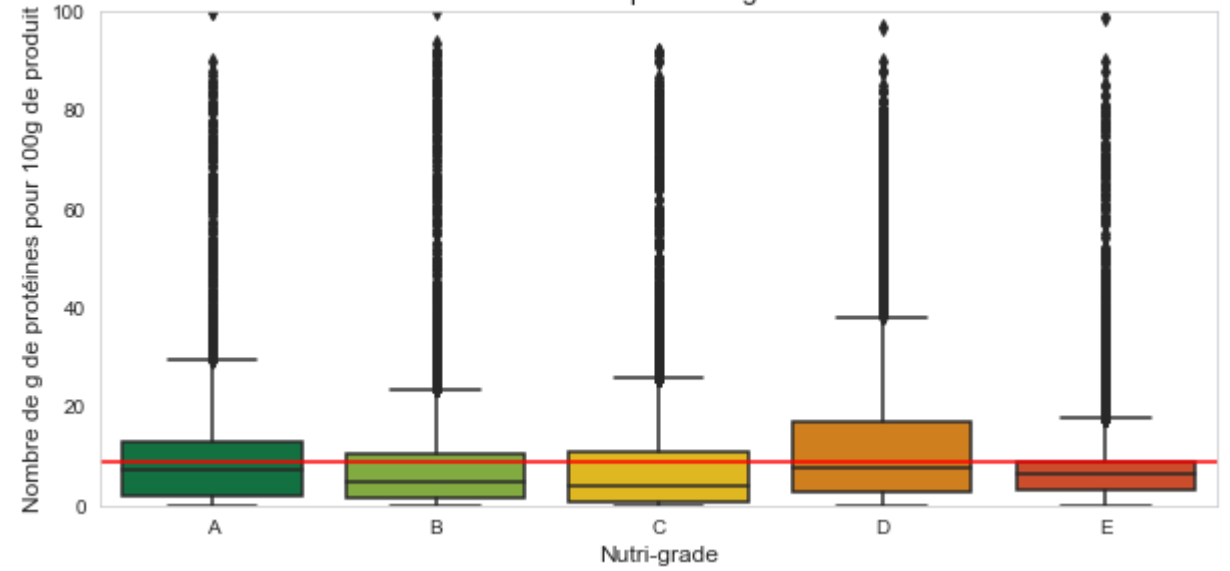
## Groupe Nova



## Distribution des protéines

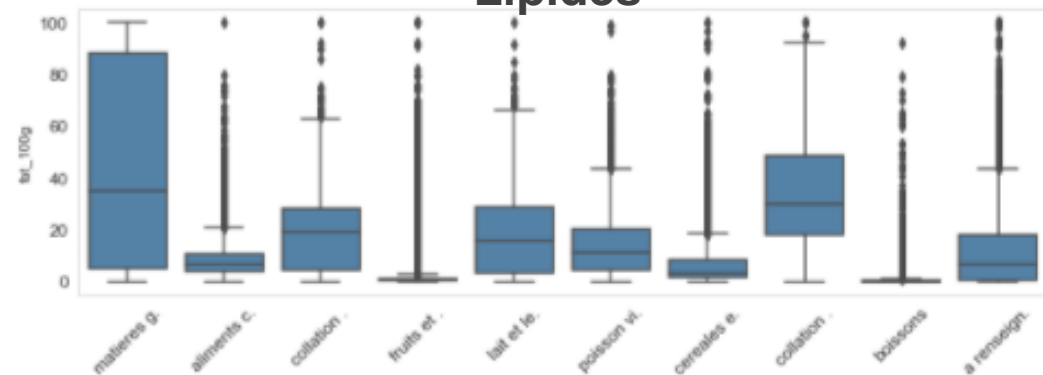


## Protéines par nutri-grade

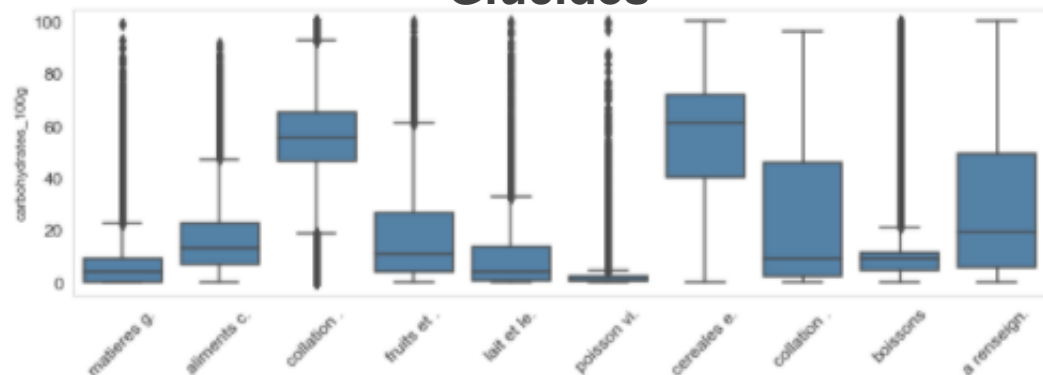


## Note du nutri-score

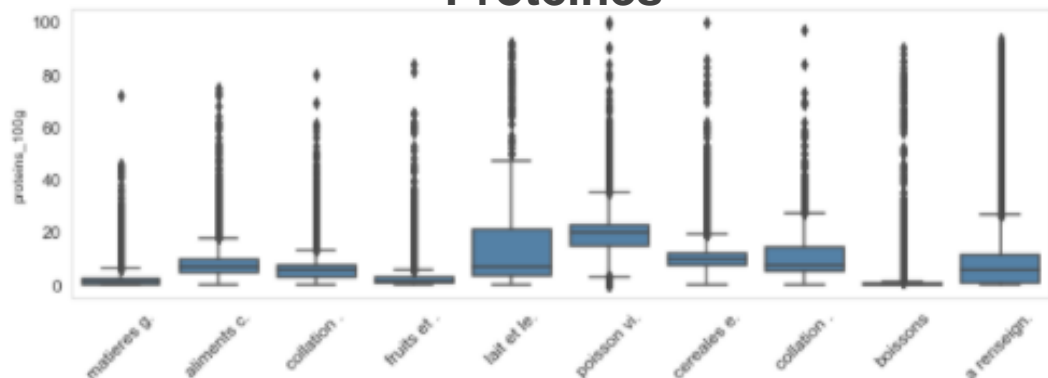
Lipides



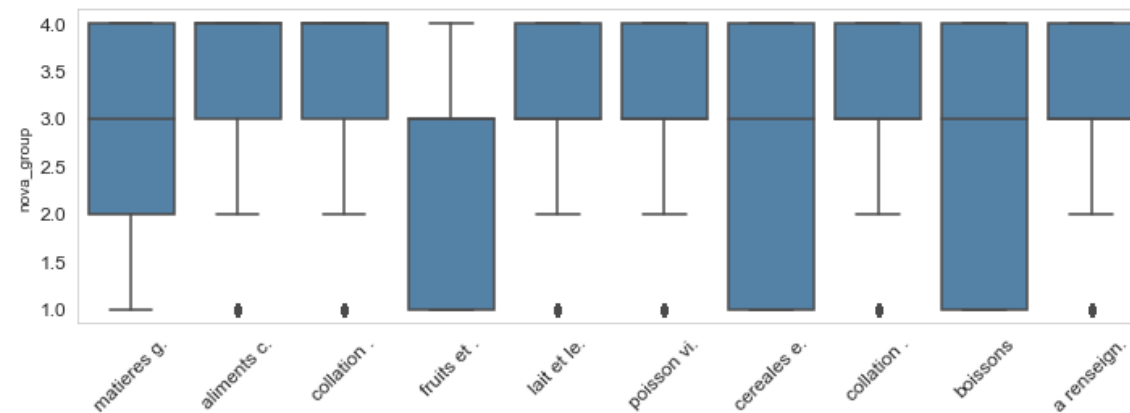
Glucides



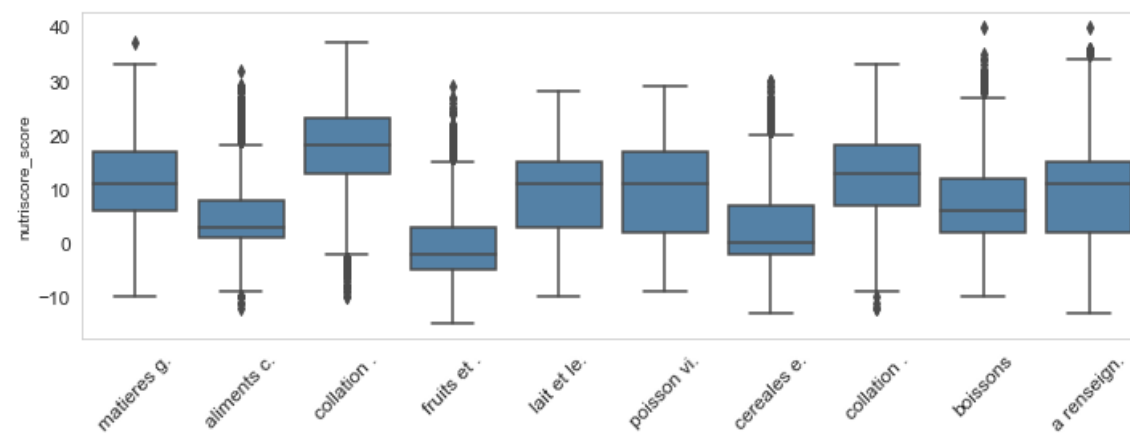
Protéines



Groupe Nova

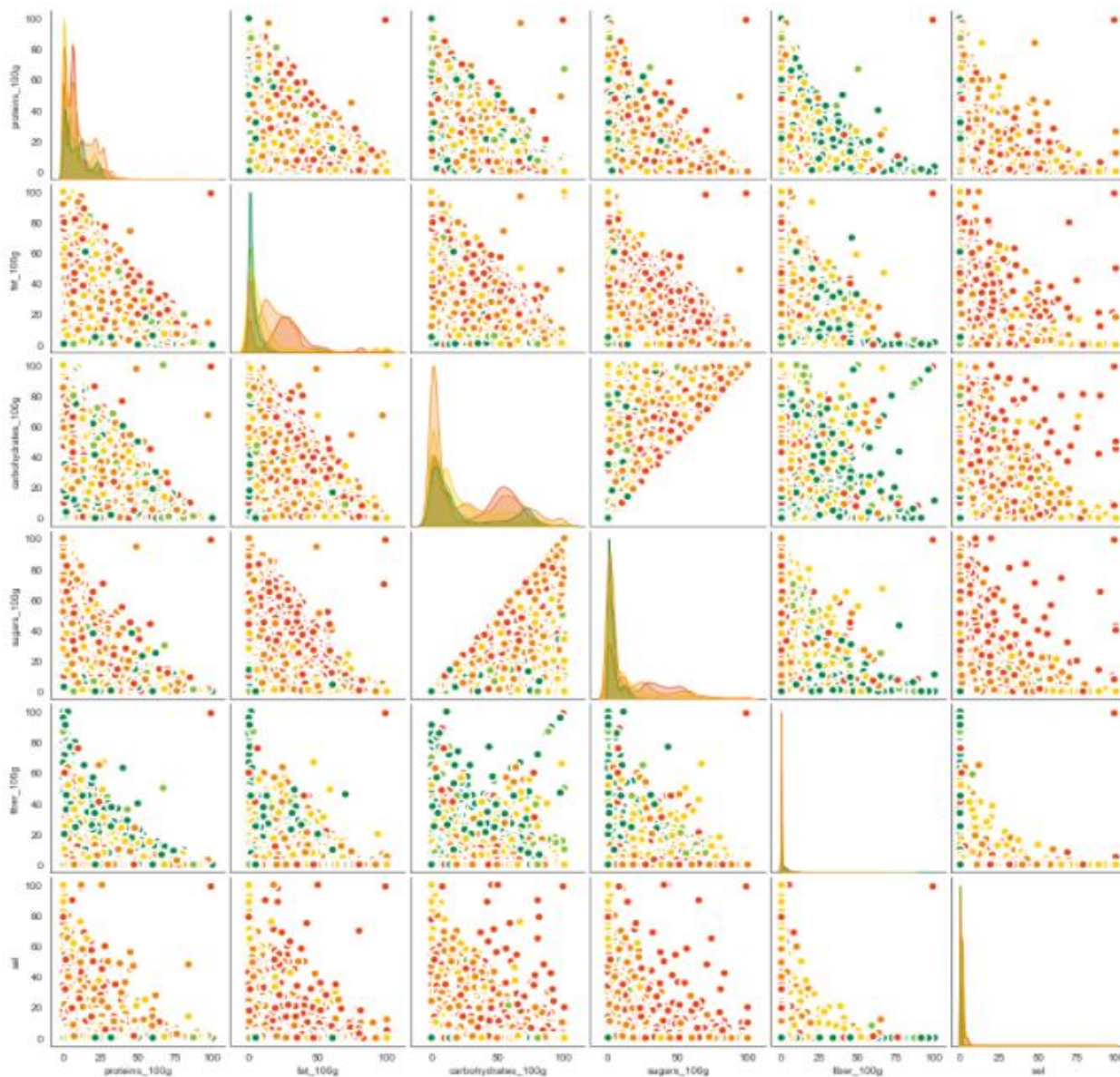


Nutri-score

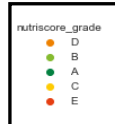
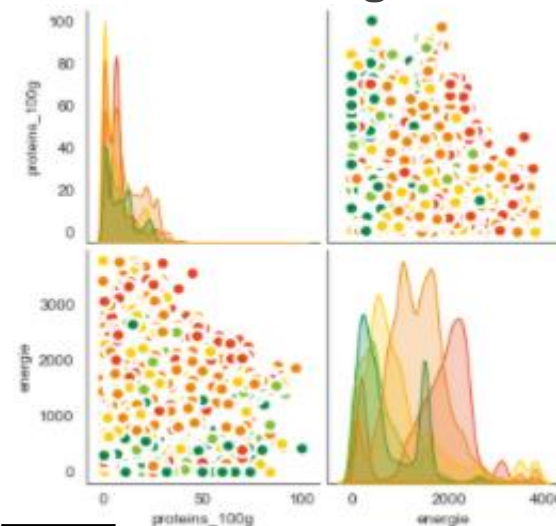


# 3 Analyse multivariée

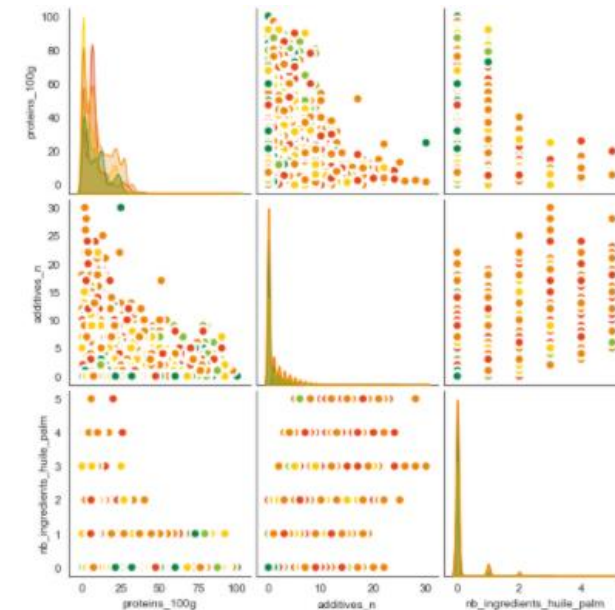
## Nutriment



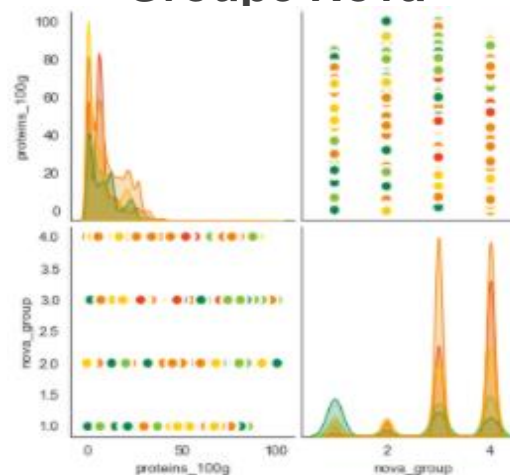
## Énergie



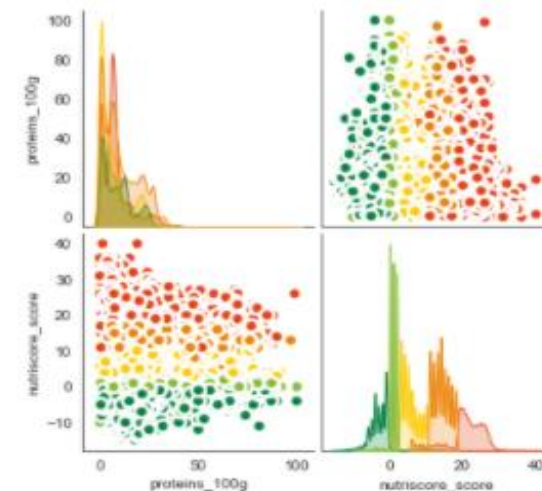
## Additifs



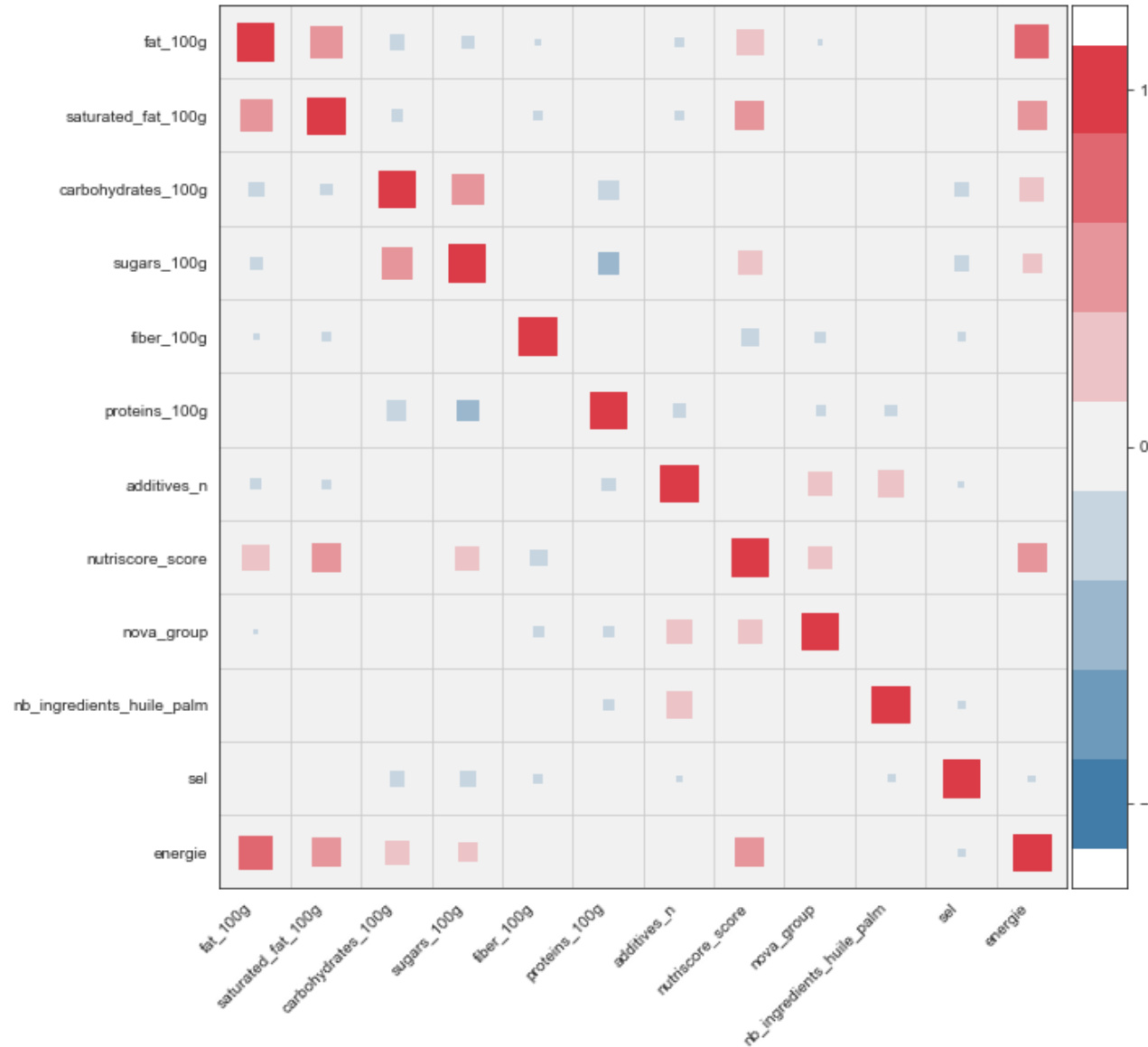
## Groupe Nova



## Nutri-score







Coef. corrélation

	proteins_100g
fat_100g	0.109987
saturated_fat_100g	0.144397
carbohydrates_100g	-0.238518
sugars_100g	-0.308174
fiber_100g	0.061575
proteins_100g	1.000000
additives_n	-0.100708
nutriscore_score	0.065837
nova_group	-0.048588
nb_ingredients_huile_palm	-0.071881
sel	0.177789
energie	0.148193



Idée d'application



Nettoyage des données



Analyse des données



**Faisabilité de l'application**



Conclusion

# 4 Faisabilité de l'application

Le nettoyage et l'analyse exploratoire ont permis de vérifier que les variables sont disponibles en quantité avec qualité

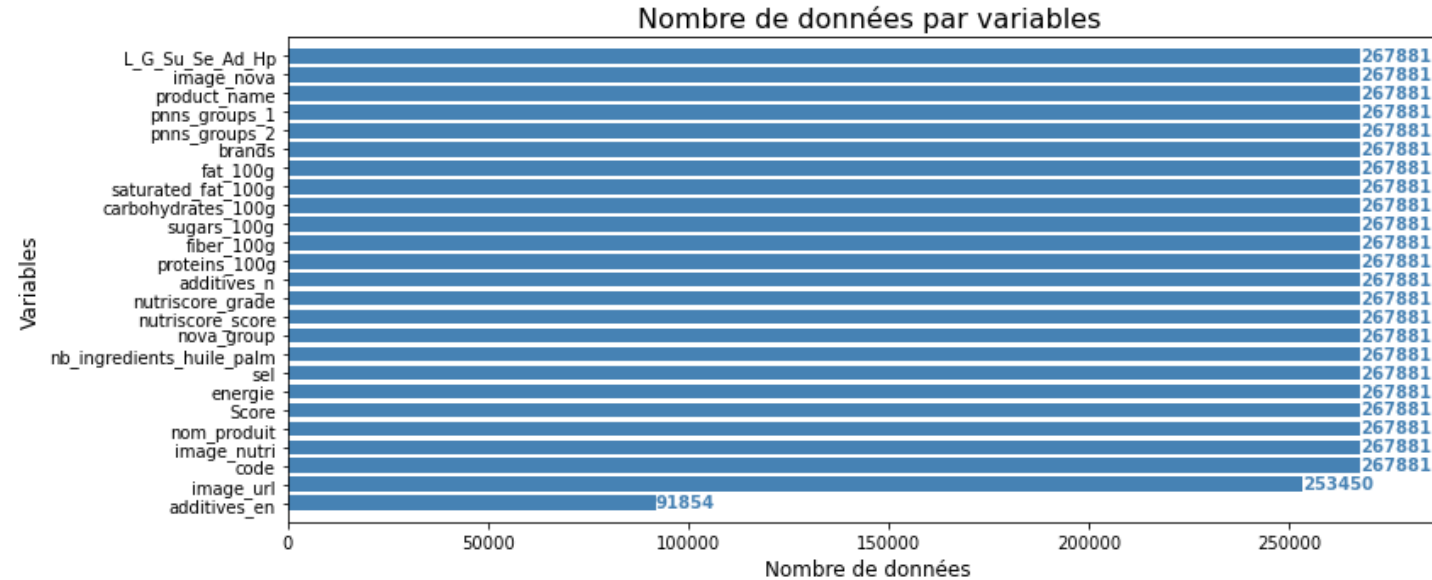
Régime hypoprotidique  
proteins\_100g

Alimentation saine  
nutriscore\_score  
nutriscore\_grade  
sel  
sugars\_100g  
additives\_n  
ingredients\_from\_palm\_oil\_n  
energie

Alimentation équilibrée  
fat\_100g  
carbohydrates\_100g  
fiber\_100g

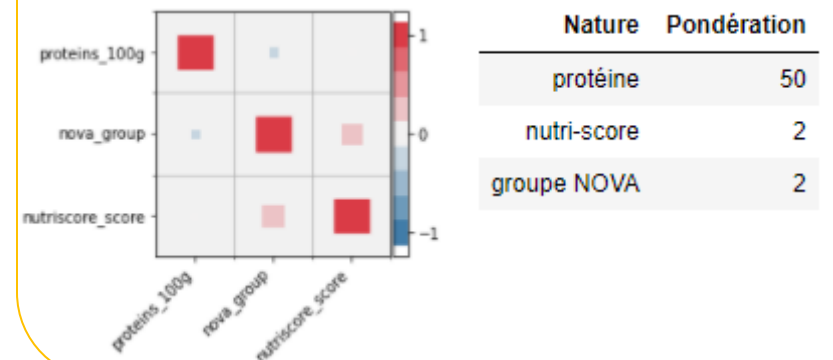
Peu transformés  
nova\_group

Produits disponibles  
product\_name  
brands  
image\_url



Moteur de recommandation réalisable

Variables de scoring



1

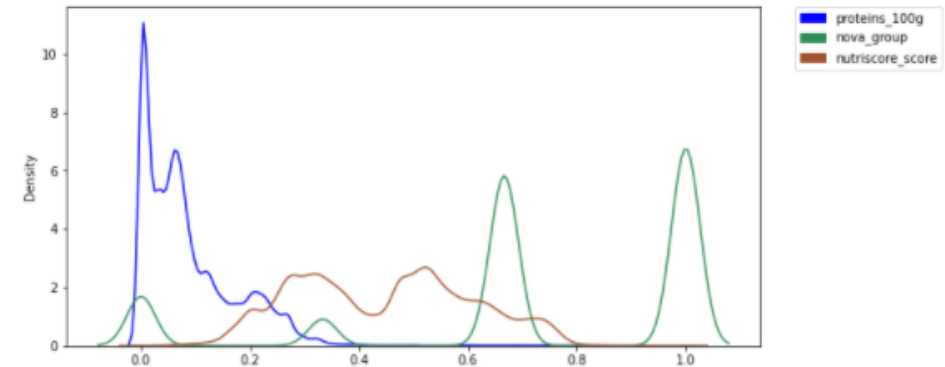
Sélection des variables de scoring :  
protéines, nutri-score et groupe nova

2

Mise à l'échelle :  
MinMaxScaler()

3

Scoring avec pondération :  
Protéines : 50 – Nutri-score et groupe Nova : 2





# 4 Moteur de recommandation

## 1 Pré-traitement des noms de produit:

Pré-processing : clean() de texthero

Suppression stopwords français : NLTK

## 2 Vectorisation :

TfidfVectorizer de scikit-learn

## 3 DataFrame de comparaison :

Features engineering (renommer variables, ajouter images, ajouter variables L\_G\_Su\_Se\_Ad\_Hp)

## 4 Moteur de recommandation :

Vérification de la saisie, suppression des ponctuations, suppression des stopwords, similarité des cosinus scikit-learn, tri des produits similaires et récupération des n produits trier par score descendant

```
# test du pré-moteur
produit = input(str('Quel est le produit recherché ? : '))
rech_produits(produit)
```

Quel est le produit recherché ? : riz

Produits recommandés :

Remarque : L\_G\_Su\_Se\_Ad\_Hp : lipides - Glucides - Sucre - Sel - Additif - huile palme

	Photo	g de protéines/100g prod	Note_nutri_score	Nutri_score	Groupe_NOVA	Marque	L_G_Su_Se_Ad_Hp
--	-------	--------------------------	------------------	-------------	-------------	--------	-----------------

Produit							
---------	--	--	--	--	--	--	--



Riz

2.8



0.0



0.4-27.5-0.2-0.0-0.0-0.0



riz

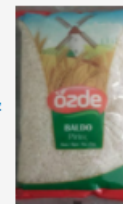
5.0



-2.0



1.0-4.0-2.0-0.0-0.0-0.0



Riz

6.6



0.0



Özde 1.0-78.0-0.4-0.01-0.0-0.0



Riz

7.0



-1.0



Carrefour 0.5-77.0-0.2-0.0-0.0-0.0



3 riz

7.6



-1.0



Lustucru 1.4-78.0-0.4-0.0-0.0-0.0



Idée d'application



Nettoyage des données



Analyse des données



Faisabilité de l'application



**Conclusion**



données nettoyés,  
imputées et analysées

Santé



Basée sur 4  
variables :  
protéines, nutri-  
score, groupe Nova  
et nom du produit

## Limites

- **Pondération** protéines/nutri-score/groupe Nova
- **Moteur de recommandation** : fiabilité de la similarité, tests
- **Fiabilité** de la saisie

## Amélioration

- Produits vendus hors de France (vacances, voyages scolaires)
- Conseils des **médecins/diététiciens**
- Ajout du **nombre de part** de phénylalanine pour chaque aliment
- Magasin où **trouver** le produit **proche de chez soi**
- Consulter les **associations** de personnes atteintes de phénylcétonurie

## Prolongement

- Bilan protéinique journalier/mensuel/ calcul des parts
- Scanner pour les enfants avec logos interdits, autorisés

## Autres MHM

- **Protéines** : aminoacidopathies, troubles du cycle de l'urée, aciduries organiques), tyrosinémie, homocystinurie -> régime hypoprotidique.
- **Sucre** : troubles de la glycolyse et de la néoglycogénèse, glycogénoses) -> régime contrôlé en sucre.
- **Lipides** : déficit de l'oxydation des acides gras -> régime en lipides.
- **Glucides** : déficit en PDH, déficit en Glut1, cytopathies mitochondriales -> régime cétogène (glucides remplacés par lipides).

## Environnement



## Librairies de base



## Outil de bureau



## Visualisation

