

October 30,  
2018

# Analysez des données nutritionnelles

Maurice clère

# Objectifs

1. Problématique & OpenFoodFacts
2. Présentation du nettoyage effectué
3. Analyse exploratoire
4. Pistes de modélisation
5. Q&A

# Problématique

- Utilisé la base de données OpenFoodFact comme socle informatif pour générer des recettes saines.



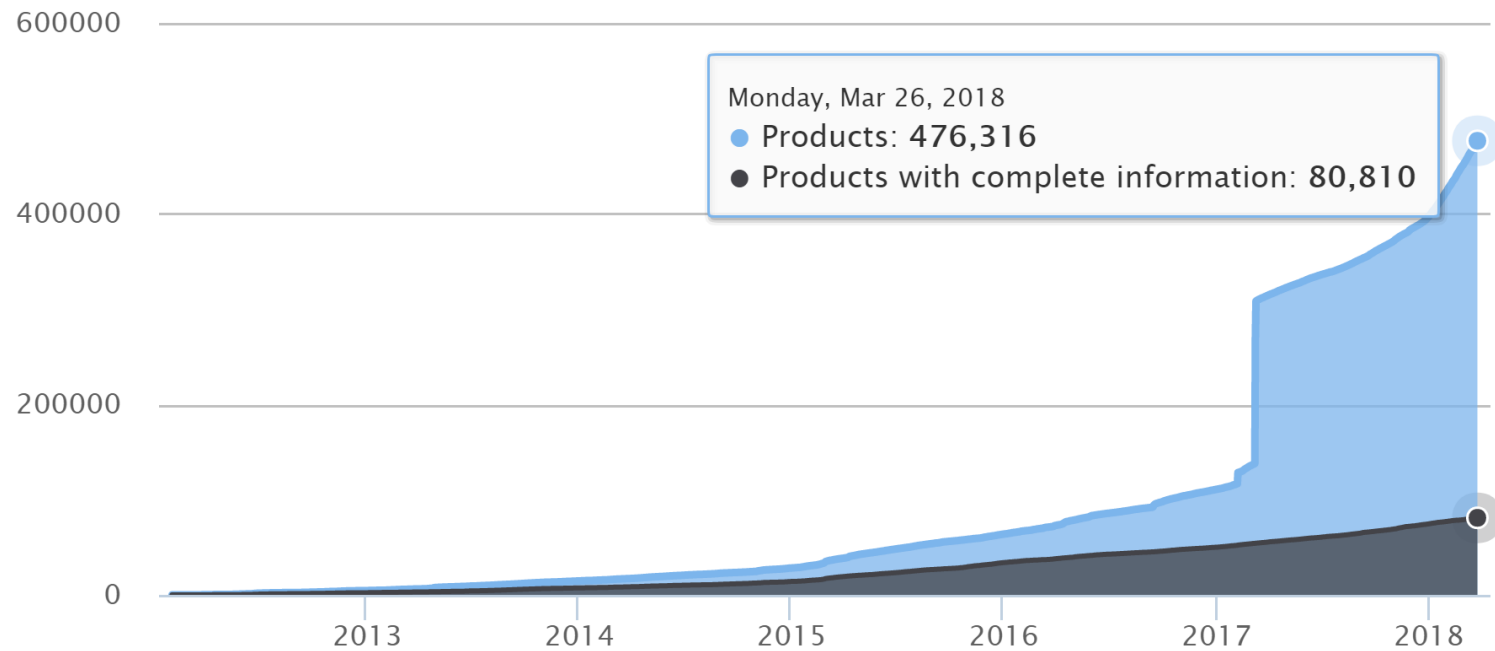
# Open Food Facts

- le projet a débuté en mai 2012 sous l'impulsion de Stéphane Gigandet
- Une base de produits alimentaires
- Faite par tout le monde
- Pour tout le monde
- **Utilisation de licences libres**



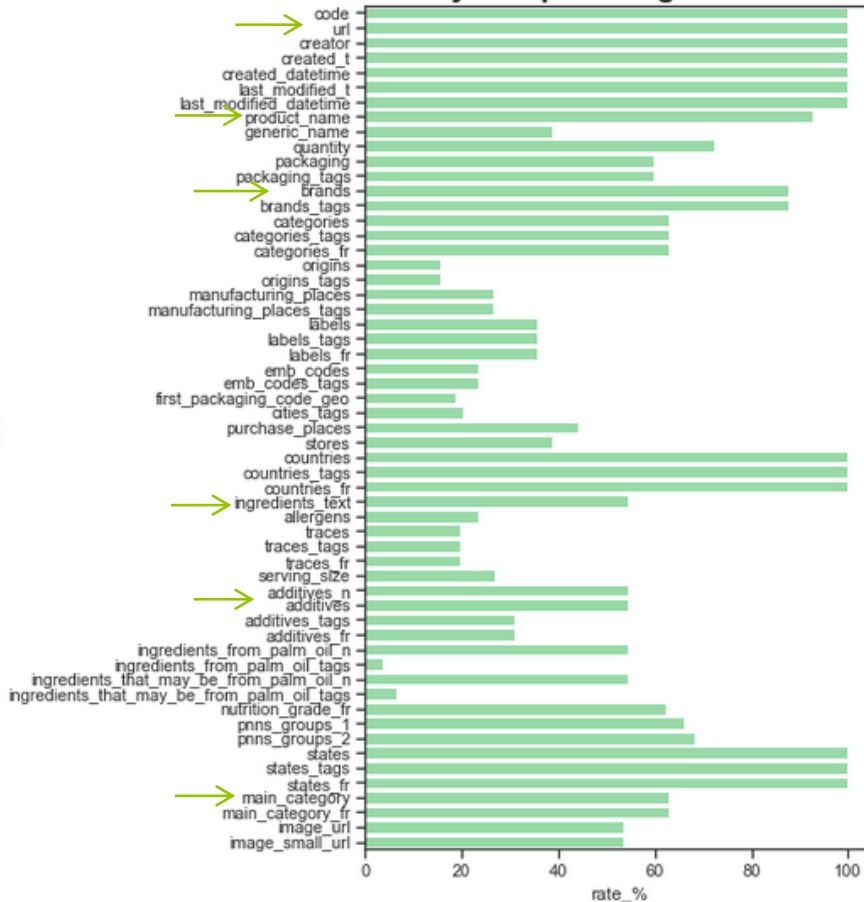
# La base de donnée

- Fichier: [fr.openfoodfacts.org.products.csv](http://fr.openfoodfacts.org/products.csv)
- 162 variables x 98440 produits (en France)

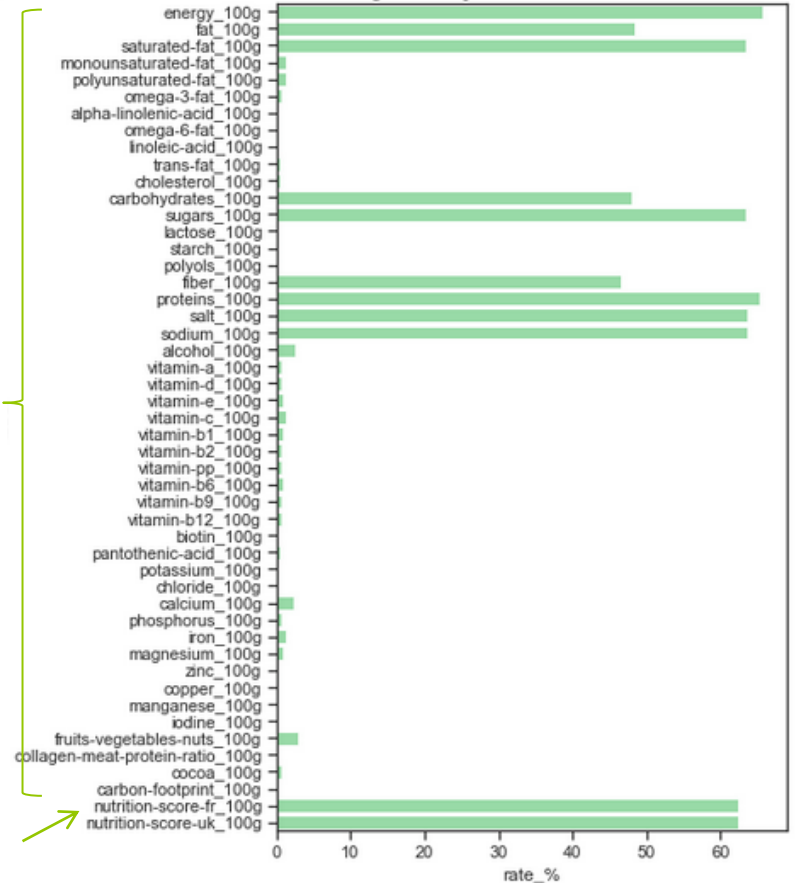


# Les Variables

recovery rate per categorical variables



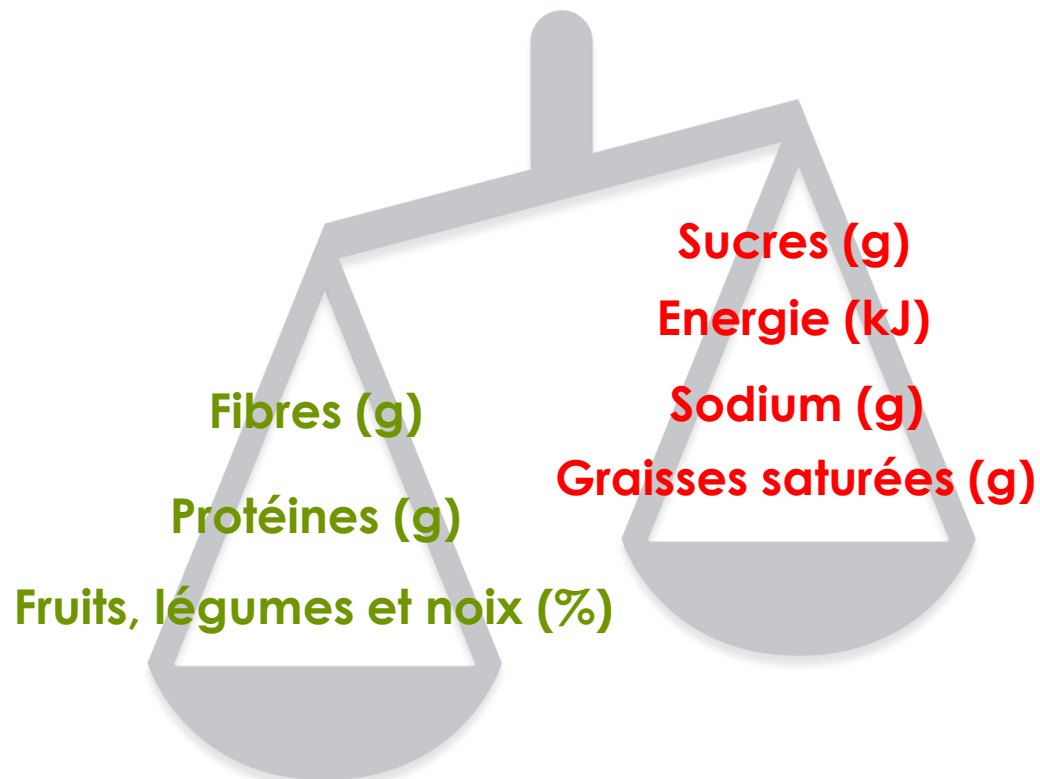
recovery rate per numerical variables



# Sélection des variables

- Numériques:
  - toutes les variables avec un taux de recouvrement supérieur à 0.1%.
- Catégorical:
  1. **code**: indentification du produit (unique)
  2. **url**: lien sur la fiche produit du site officiel d'openFoodFacts
  3. **product\_name**: nom du produit
  4. **Brands**: marque
  5. **categories\_fr**: liste de mots clé (soupes, biscuit, etc...)
  6. **countries\_fr**: pays dans lequel le produit est répertorié
  7. **ingredients\_text**: liste d'ingrédients
  8. **ingredients\_from\_palm\_oil\_n**: huile de palme
  9. **additives\_fr**: liste des additifs
  10. **main\_category\_fr**: catégorie qui représente le mieux le produit
  11. **nutrition\_grade\_fr**: notation (A, B, C, D, E)

# Nutri-score





# Nettoyage

- Outliers
- Valeurs incohérentes
- Valeurs manquantes

	min	max
energy_100g	0.0	3251373.0
fat_100g	0.0	380.0
saturated-fat_100g	0.0	210.0
carbohydrates_100g	0.0	190.0
sugars_100g	-0.1	105.0
fiber_100g	0.0	178.0
salt_100g	0.0	211.0
vitamin-b1_100g	0.0	161.0
carbon-footprint_100g	0.0	2520.0
nutrition-score-fr_100g	-15.0	40.0
nutrition-score-uk_100g	-15.0	36.0

# Nettoyage

type	variables	Min	max	action	nombre d'occurrences
Outliers	toutes les valeurs pour 100g	0	110*	value set to None	10
Outliers	Energy_100g	0	3700**	value set to None	246
Incoherence	somme des macronutriments sur 100g	-	110*	produit retiré	69
Incoherence	somme des Graisses / 100g	-	110*total Graisse	value set to None	89
Incoherence	somme des Glucides / 100g	-	110*total Glucide	valeur set to None	3758
Valeur absentes	nutri-score	-	-	produit retiré	37020
Valeur absentes	product_name	-	-	produit retiré	343
Doublon	Code	-	-	un doublon retiré	0

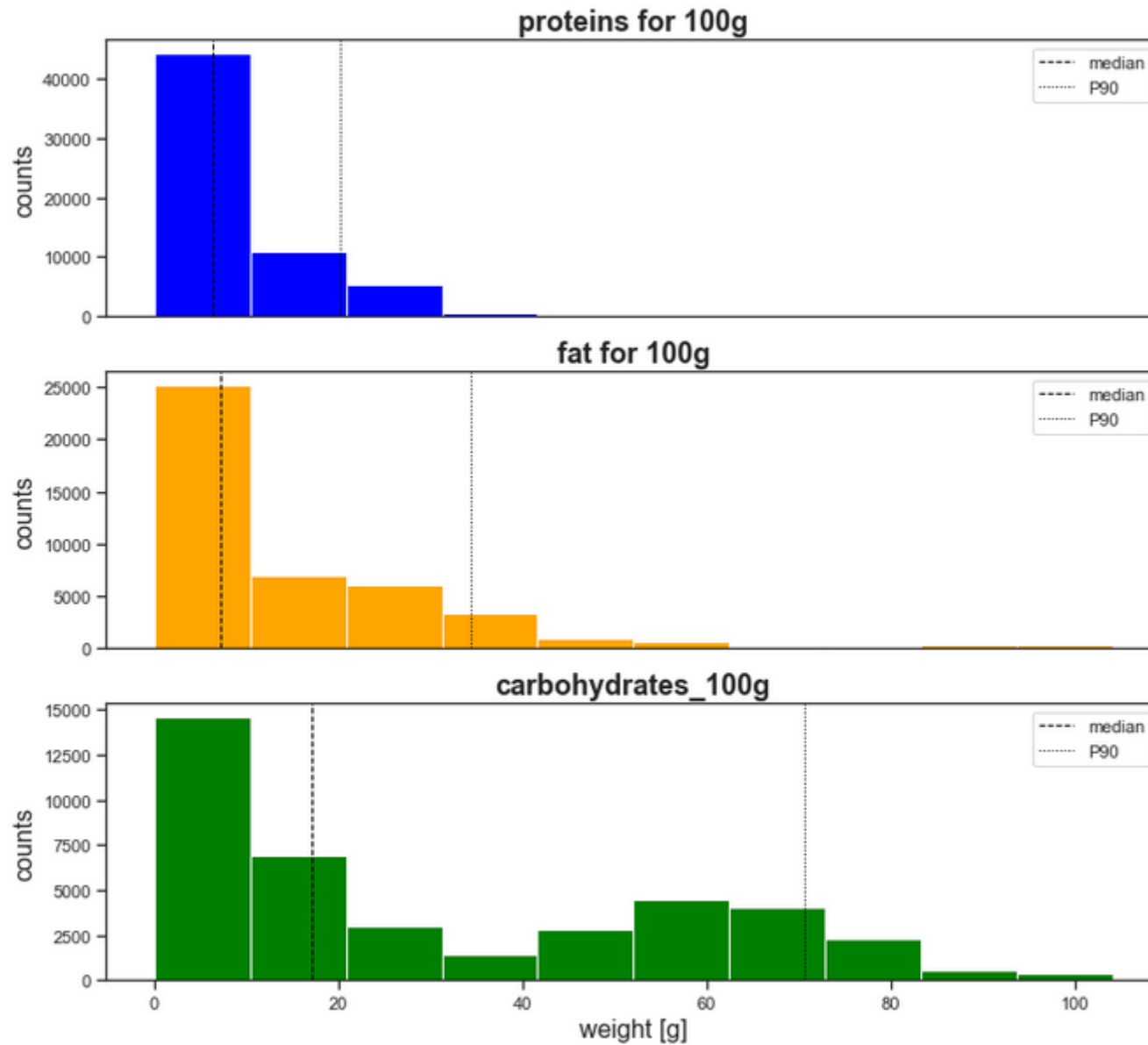
\* 110 = limite haute + 10% car certains produits sont en ml

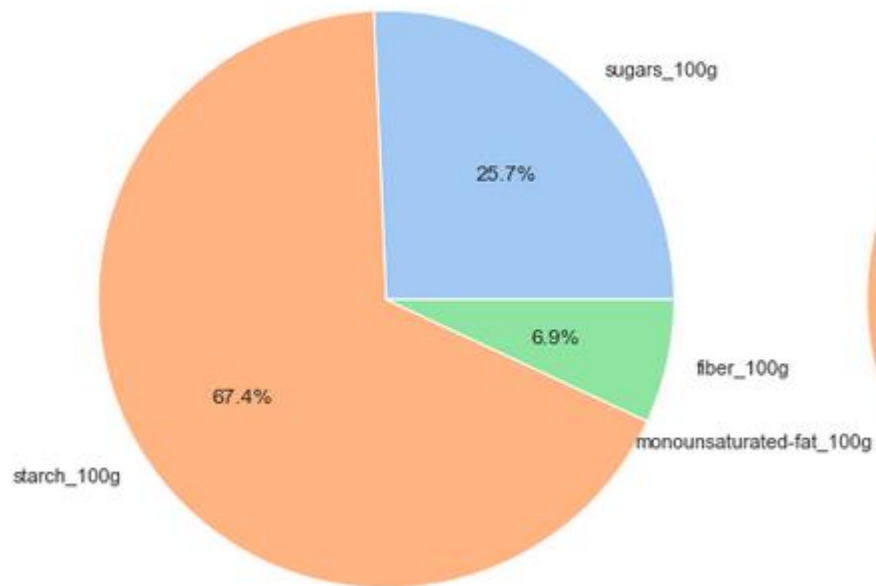
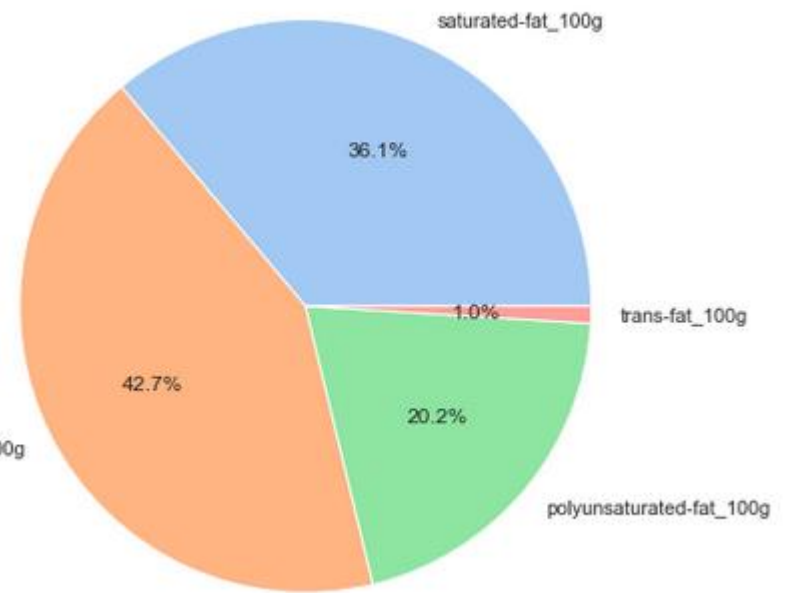
\*\* 3700 = 37kj/g\*100g (teneur maximum sur un produit composé à 100% de graisse)

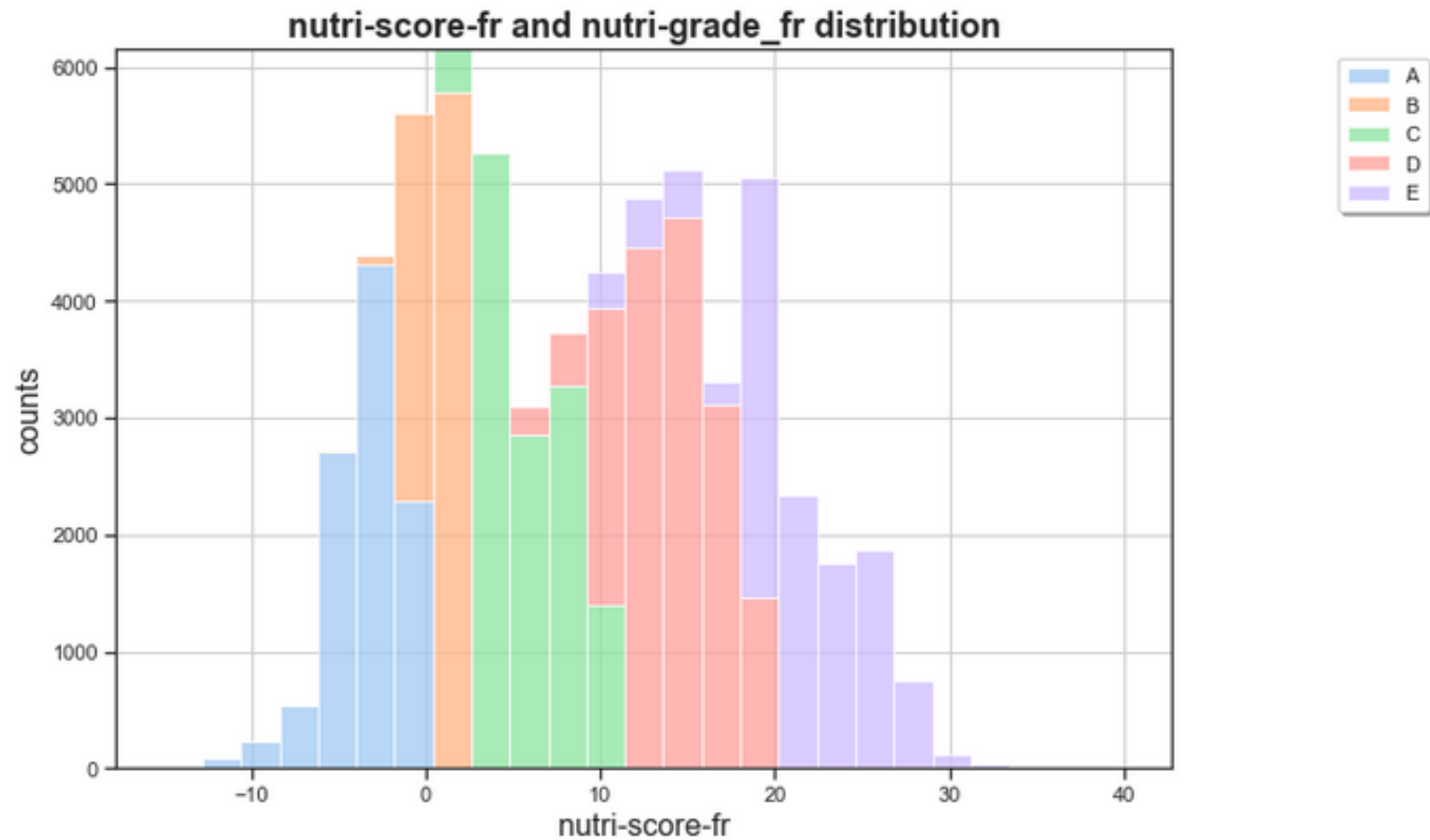
# Autres

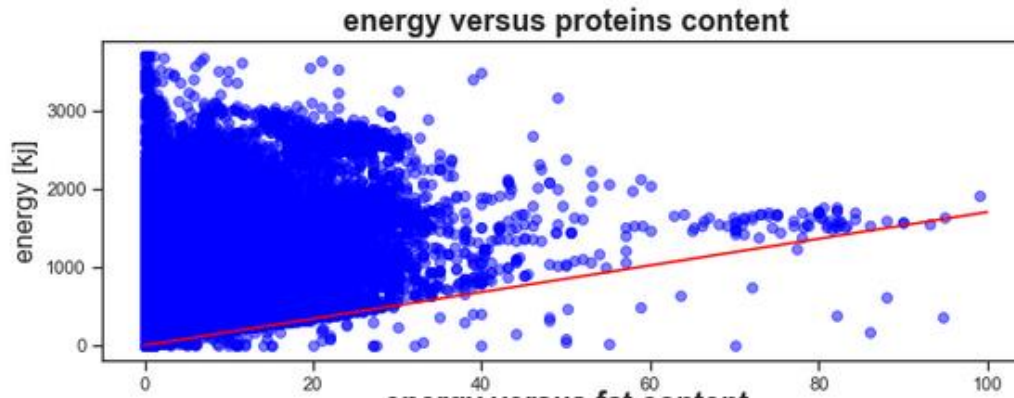
type	variables	action
formatage	Toutes les valeurs des variables catégorical	Changer en minuscule
formatage	Nom de la variable 'energy_100g'	Renomme 'energy_kj_100g'
feature engineering	'energy_Kcal_100g'	Energy convertie en Kcal
feature engineering	'check_energy_kj_100g'	Energy recalculer a partir des macronutriments

# Analyse exploratoire



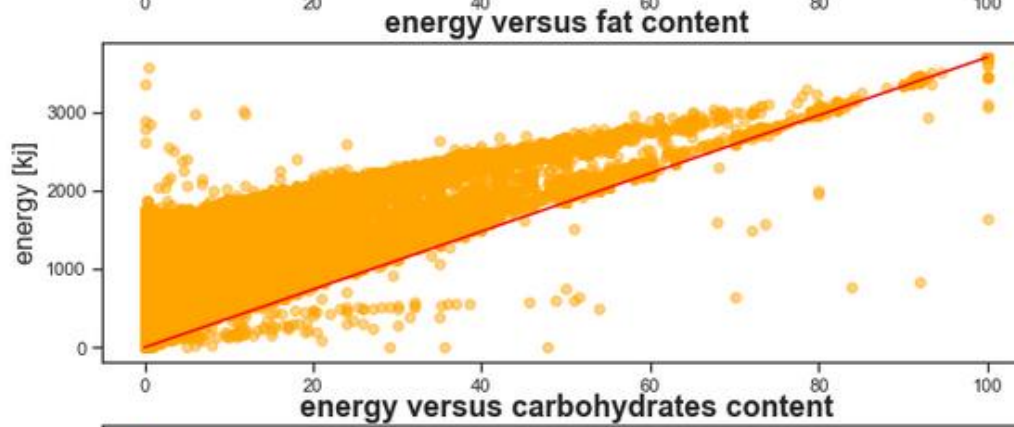
**carbohydrates repartition****fats repartition**





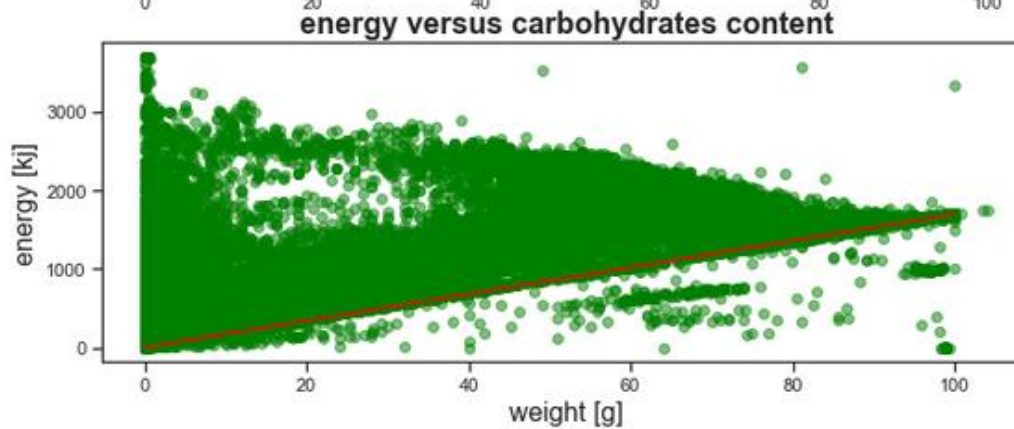
$$y = 17x$$

(1g protéines = 17kj)



$$y = 37x$$

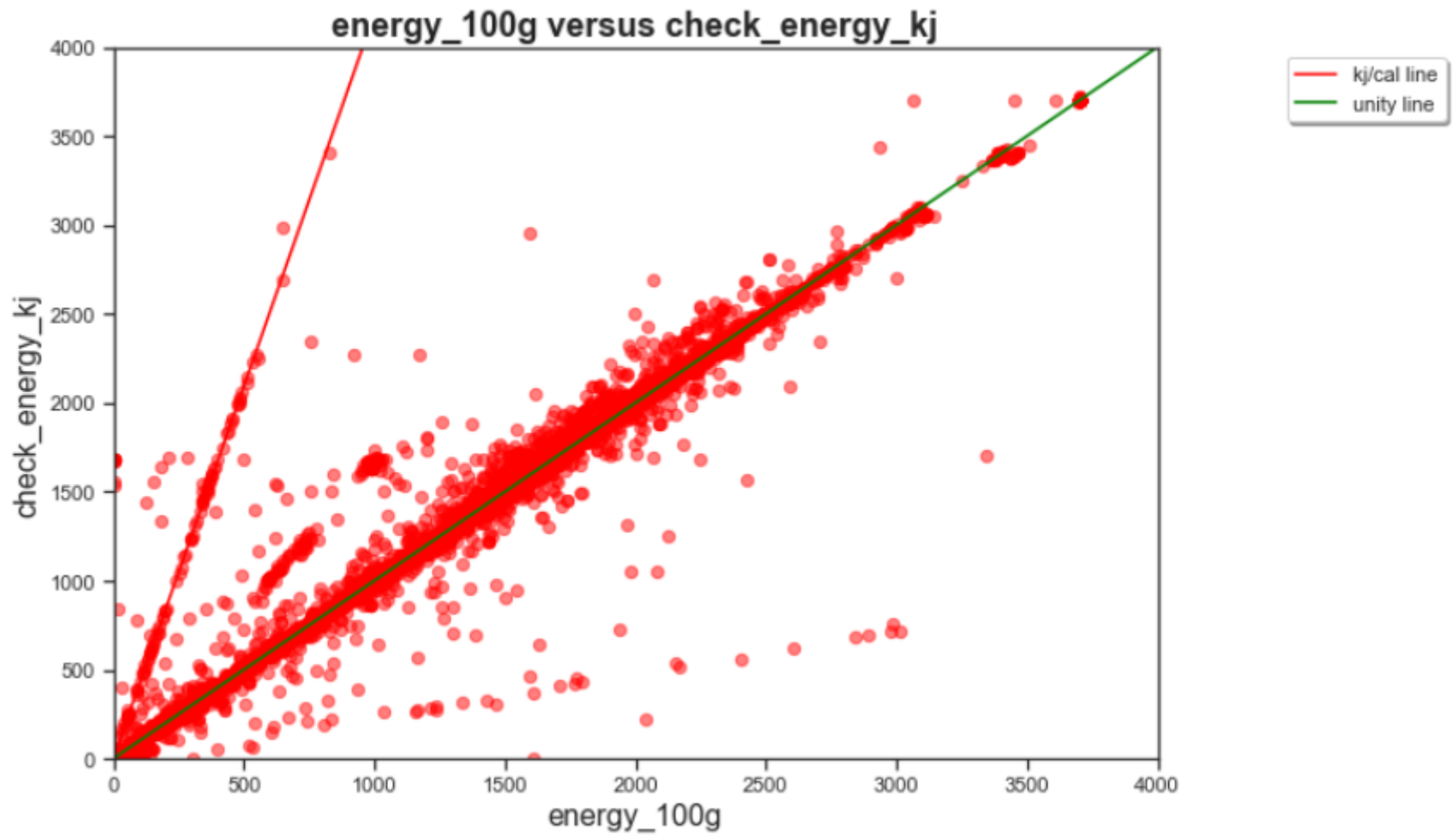
(1g graisse = 37kj)

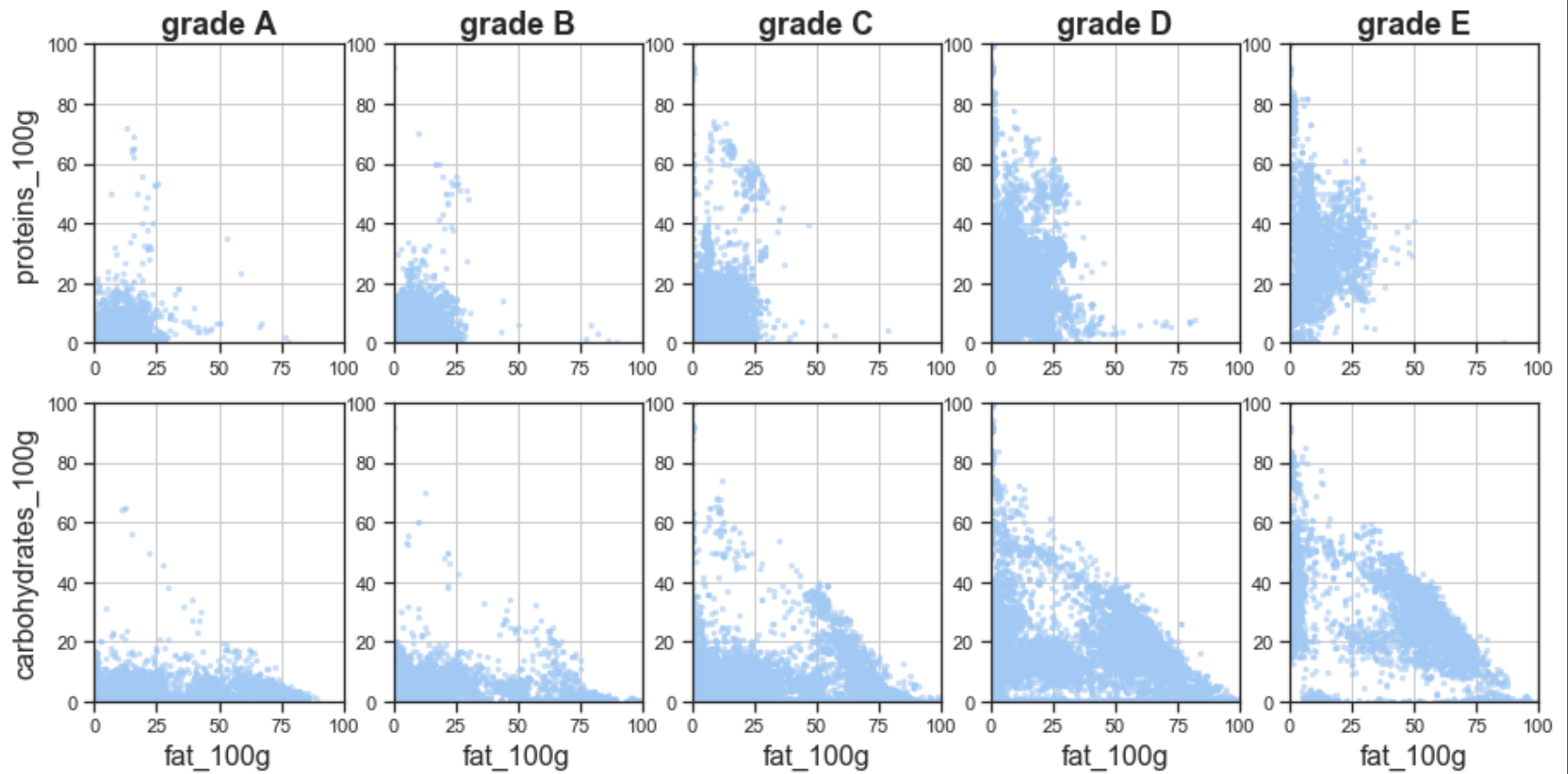


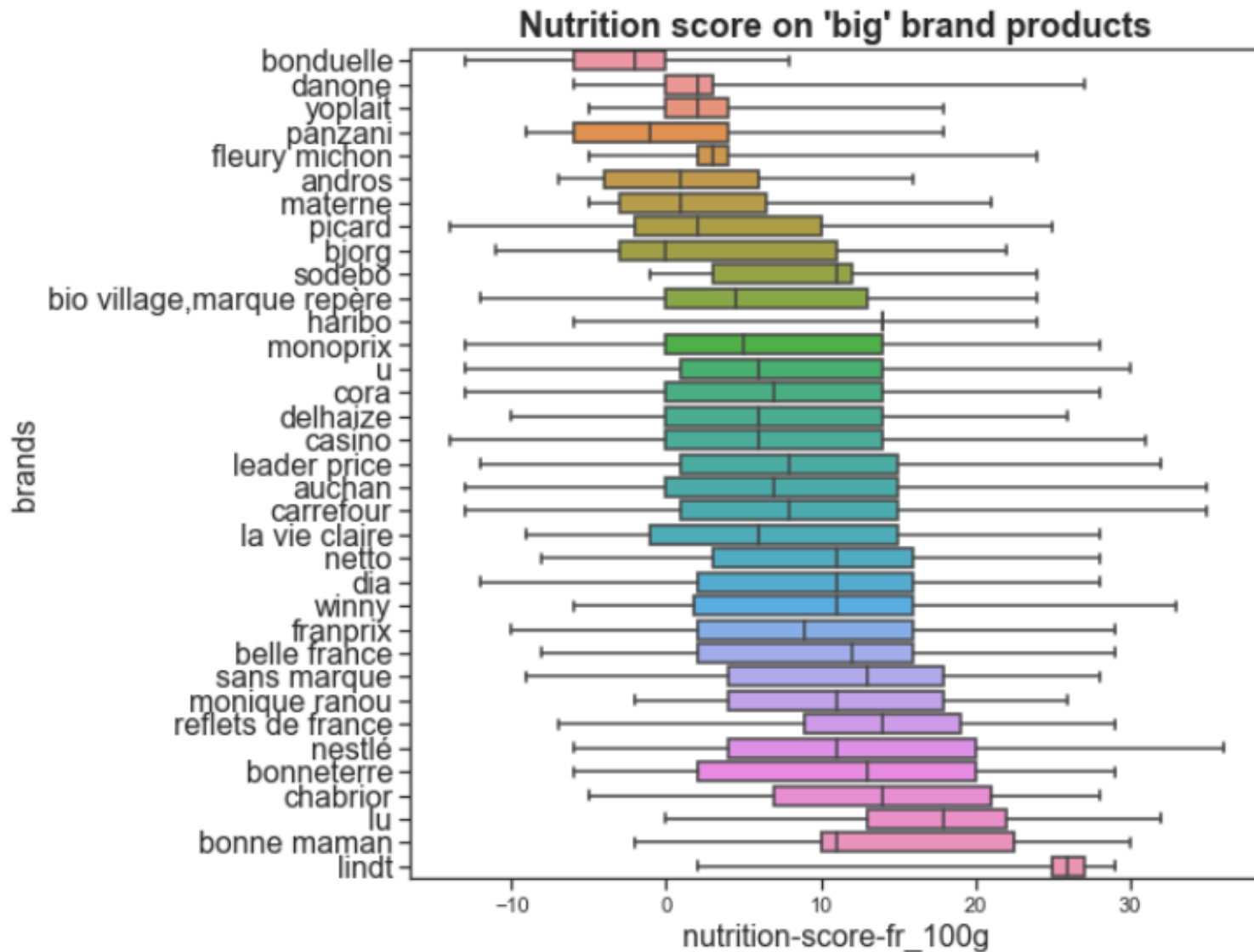
$$y = 17x$$

(1g glucide = 17kj)

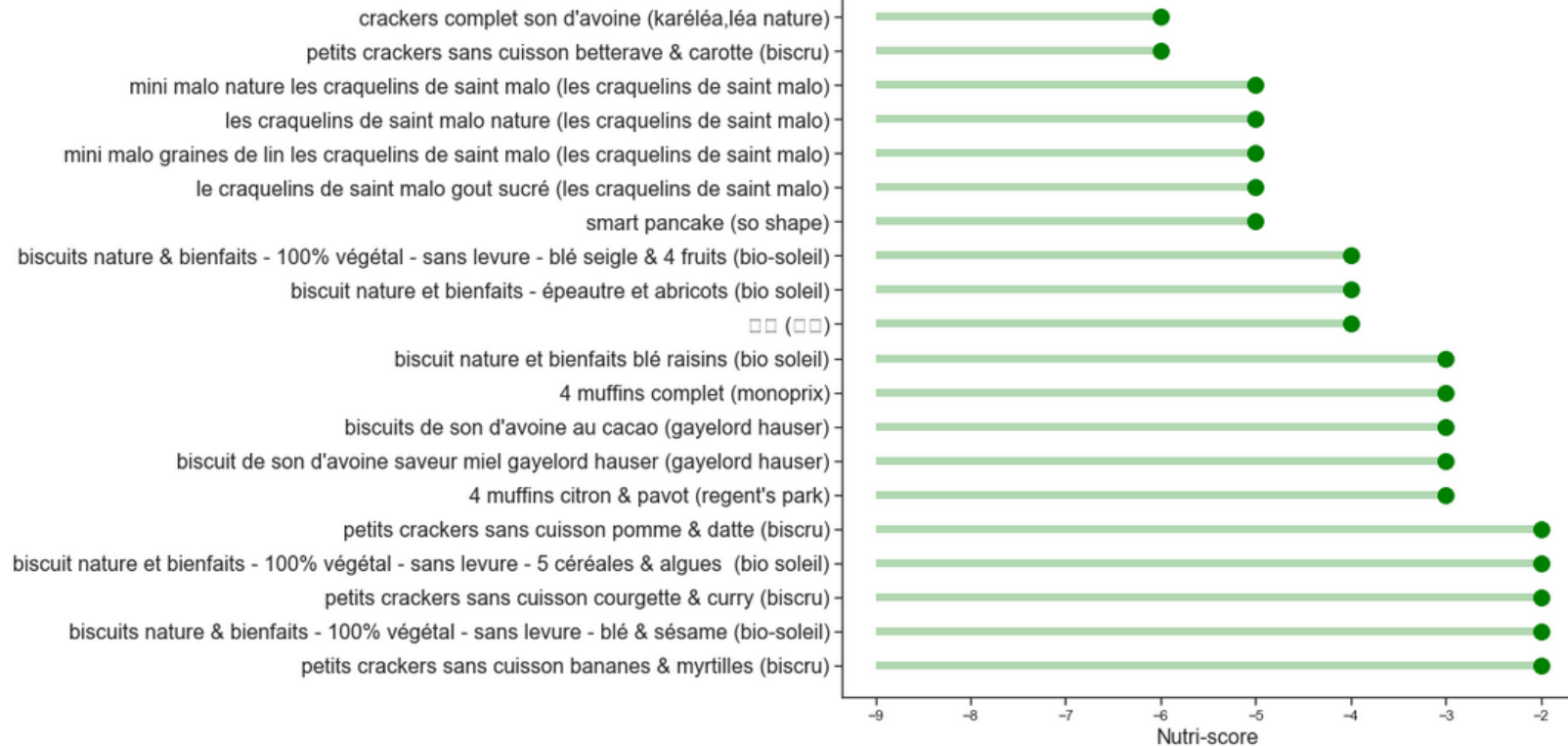








### Top healthy products for biscuit



# Pistes de modélisation

- Recherche automatique des meilleurs produits par mots clés
- Sauvegarde des produits (ajoutés à la recette)
- Bilan nutritionnel par recette
- Ajout d'une colonne favori (suivant les disponibilités des produits)
- Update automatique de la base de donnée
- Redéfinition du nutri-score / nutri-grade
  - Intégration des additifs et de l'huile de palme

# Q&A