

Objectifs :

Traiter les données alimentaires de manière à pouvoir :

- Prédire le nutriscore d'un plat ou d'un aliment en fonction de la proportion de ce qu'ils contiennent.

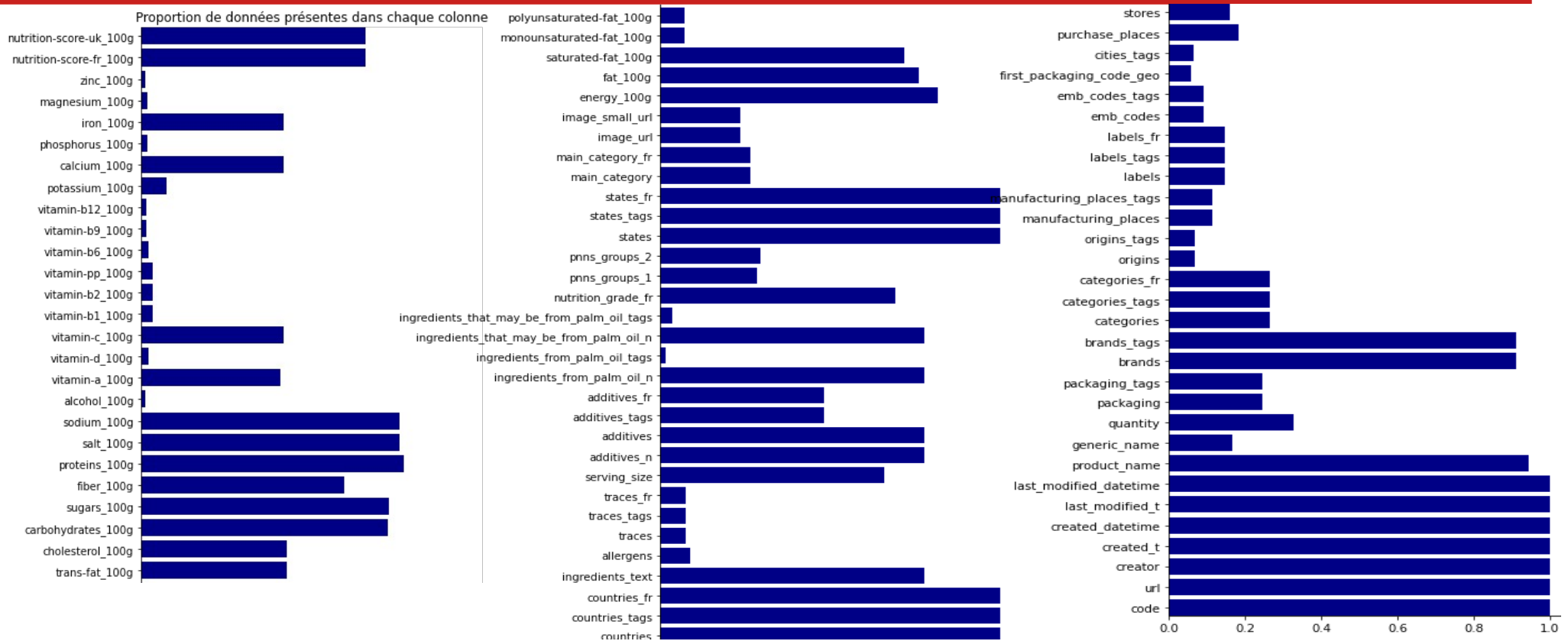
Description des données :

**Base de donnée : - information sur des produits alimentaires
([en.openfoodfacts.org.products](https://en.openfoodfacts.org/products))**

- taille (320 772, 162)

Taille base de donnée après suppression des doublons : - (320 750 x 162)

Visualisation des données



Nettoyage des données

Conservation uniquement des variables dont le remplissage est supérieur à 1 %

Application de règles métiers.

Objectifs : retirer les valeurs aberrantes

$\text{energy_100g} \leq 3700$

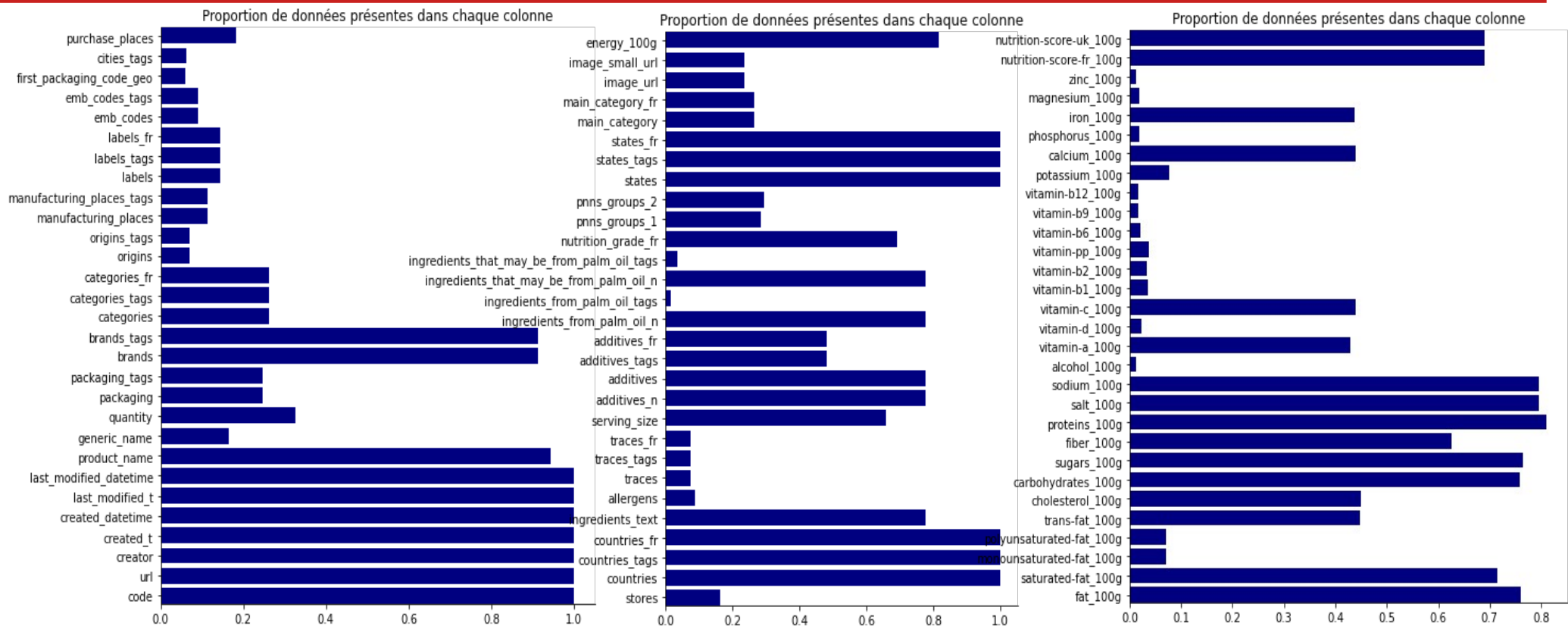
$\text{saturated_fat_100g} \leq \text{fat_100g}$

$\text{sugars_100g} \leq \text{carbohydrates_100g}$

$\text{sodium_100g} \leq \text{salt_100g}$

$(\text{fat_100g} + \text{carbohydrates_100g} + \text{proteins_100g} + \text{salt_100g}) \leq 100$

2 Données après nettoyage



Algorithme du K-NN

Classe que l'algorithme doit déterminer :

nutriscore_grade_fr

Variables sur lesquelles il peut se baser :

energy_100g, fat_100g, saturated-fat_100g, salt_100g, sodium_100g, proteins_100g, carbohydrates_100g, sugars_100g, iron_100g

Échantillonnage : on prend 15 000 aliments aléatoires

Sur ces 15 000 aliments, 80 % servent à entraîner l'algorithme, 20 % à tester la précision de l'algorithme

pourcentage d'erreur en fonction du nombre de plus proche voisin retenu pour l'algorithme avec la méthode suivante : uniform

