

# Conception d'une application au service de la santé publique



## Contexte :

Santé publique France lance un appel à projets pour trouver des idées innovantes d'applications en lien avec l'alimentation.



## Données :

Jeu de données sur les produits alimentaires du monde entier



## Mission :

- Proposer une idée d'application pour répondre à l'appel à projet
- Effectuer une analyse exploratoire des données Open Food Facts afin déterminer la pertinence et la faisabilité de l'application.

## I) Présentation de l'idée d'application

## II) Présentation et préparation des données

- A. Présentation des données
- B. Nettoyage des données
- C. Analyse et prédiction des valeurs manquantes

## III) Analyse exploratoire des données

- A. Analyse multivariée des variables nutritionnelles, du Nutri-score et des groupes de produits
- B. Répartition des Grades Nutri-score dans les différentes catégories de produits
- C. Analyse des additifs
- D. Analyses des produits d'origine France
- E. Analyses des produits Bio

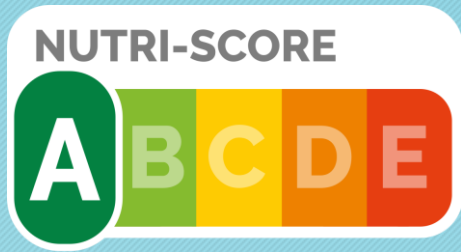


# Partie I : Présentation de l'application

3

# Idée d'application : 1) manger plus sain

## 1) MANGER MOINS GRAS, MOINS SUCRÉS ET MOINS SALÉS :



Valeurs nutritives pour 100g :

+ Fibres, protéines, fruits légumes, légumineuses, fruits à coques, huile de colza, de noix et d'olive.

- Energie, acides gras saturés, sucres, sel.

## 2) LIMITER LES ADDITIFS ALIMENTAIRES

Ingrédients d'origine naturel ou de synthèse

- Qualité sanitaire (conservateurs, antioxydants)
- Aspect et Goût (colorants, édulcorants, exhausteurs de goût)
- Texture (épaississants, gélifiants)
- Stabilité (émulsifiants, antiagglomérants, stabilisants).

Qui peuvent être nocifs pour la santé :

- Nitrites et nitrates (E249, E250, E251 et E252) <sup>1</sup>
- Dioxyde de titane (E171) <sup>2</sup>
- Carboxyméthylcellulose (E466) et Polysorbate-80 (E433) (Emulsifiants) <sup>3</sup>

Sources :

1) Évaluation des risques liés à la consommation de nitrates et nitrites. Anses, Rapport d'expertise collective. Juillet 2022

2) Re-evaluation of titanium dioxide (E 171) as a food additive. EFSA Journal, 2016

3) Dietary emulsifiers impact the mouse gut microbiota promoting colitis and metabolic syndrome. Chassaing et al. Nature, 2015

# Idée d'application : 2) manger plus responsable



## 1) PRODUITS D'ORIGINE FRANCE <sup>1</sup> :

- L'alimentation en France = 140 Mt CO<sub>2</sub> eq (2,1 t par personne), soit 22% de l'empreinte carbone totale.
- 46% de ces émissions sont associées aux importations.



## 2) PRODUITS BIO <sup>2</sup> :

- Préservation de la biodiversité et fertilité des sols (- pesticides de synthèse, utilisation d'engrais vert riche en matière organique)
- Préservation de la qualité de l'eau (- de rejets azoté)
- Limitation des émissions de GEF (- de rejets, + de séquestration par fixation du carbone dans le sol)

### Sources :

1) Ministère de la transition Ecologique et de la Cohésion des Territoires, DOCUMENT DE TRAVAIL N° 59, Juillet 2022.

2) Agence française pour le développement et la promotion de l'agriculture biologique, basée sur le rapport : Sautereau N., Benoit M., 2016. Quantification et chiffrage des externalités de l'agriculture biologique, Rapport d'étude ITAB, 136 p



# Concept de l'application

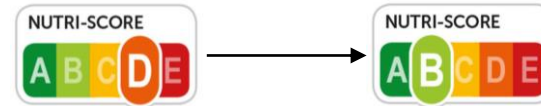
6



Proposition d'une liste de  
produits de la même famille

## Manger plus sain

- Nutri-score



- - d'additifs



## Manger plus responsable

- Origine France
- Label Biologique



# Partie II : Présentation et préparation des données

7

- A) Présentation des données
- B) Nettoyage des données
- C) Analyse et prédiction des valeurs manquantes



- Export du 26/08/2022
- 2,5 millions de produits
- 186 variables décrivant les produits (infos générales, tags, ingrédients, valeurs nutritives)



INFOS GÉNÉRALES	ADDITIFS	TAGS	ORIGINE DU PRODUIT	GROUPE DE PRODUITS	ELEMENTS CONSTITUANT LE NUTRIScore	NUTRI-SCORE
<ul style="list-style-type: none"><li>• Code</li><li>• Product_name</li><li>• Brands</li></ul>	<ul style="list-style-type: none"><li>• additives_n</li><li>• additives_fr</li></ul>	<ul style="list-style-type: none"><li>• Labels_fr</li></ul> <div>↓</div> <div>Bio</div>	<ul style="list-style-type: none"><li>• Origins_fr</li><li>• Manufacturing_places</li></ul> <div>↓</div> <div>Origine France</div>	<ul style="list-style-type: none"><li>• Pnns_groups_1</li><li>• Pnns_groups_2</li></ul>	<ul style="list-style-type: none"><li>• fibers_100g</li><li>• Saturated-fat_100g</li><li>• Fat_100g</li><li>• Proteins_100g</li><li>• Sugars_100g</li><li>• Sodium_100g</li><li>• Fruits-vegetables-nuts-estimate-from-ingredients_100g</li></ul>	<ul style="list-style-type: none"><li>• Nutriscore_score</li><li>• Nutriscore_grade</li></ul>

# B) Nettoyage du jeu de données

- Etape 1 : Nettoyage des données numériques

- ELEMENTS CONSTITUANT LE NUTRISCORE [0;100] (Sinon : NaN)
- ENERGIE\_100g : [0;3700] car les lipides (fat) ont une valeur de 37(kJ/g) (sinon : NaN)

- Etape 2 : Nettoyage des données textuelles

- Création de la colonne « Origine France »



Origine ou manufacturing places

« Fra » = 1

- Création de la colonne « Bio »



Labels

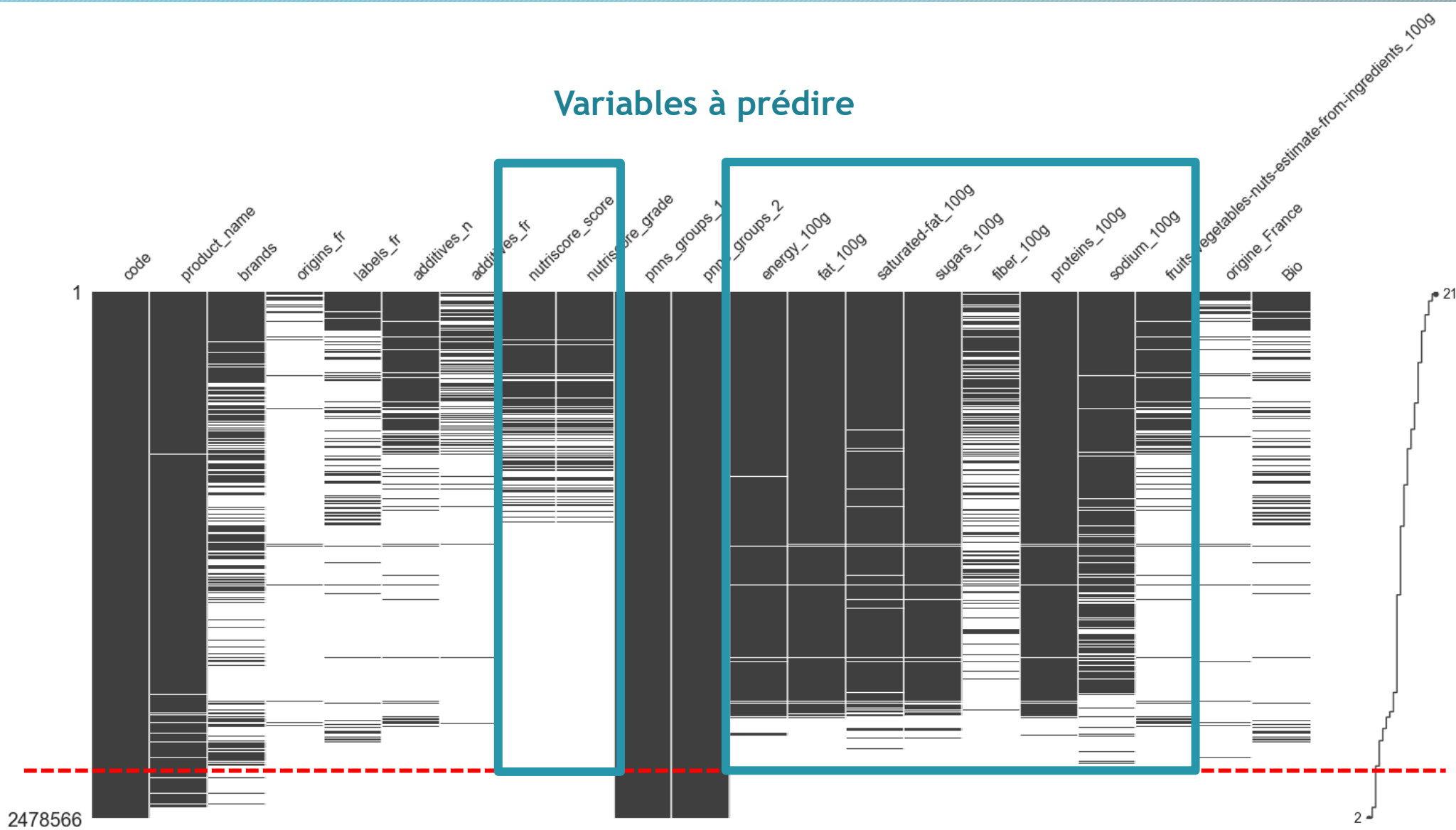
« Bio » ou « Organic » = 1

- Vérification et modification des groupes de produits (exemple : 'sugary-snacks' > 'Sugary snacks')

- Etape 3 : Suppression des doublons

# C) Analyse et prédiction des valeurs manquantes

10



Retrait des produits avec très peu d'informations (moins de 5 variables connus).

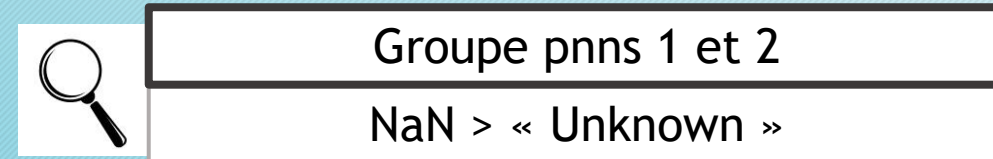
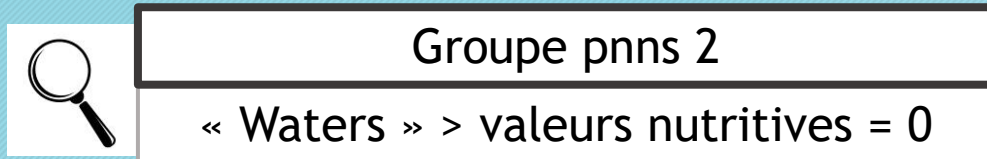
➤ 245 783 produits retirées (~10%)



# C) Analyse et prédiction des valeurs manquantes

## Objectif : Prédire le Grade Nutri-score manquant sur certains produits

- Etape 1 : Remplacement des valeurs manquantes



- Etape 2 : Imputation des valeurs manquantes de Nutriscore-Score, fibers\_100g, Saturated-fat\_100g, Fat\_100g, Proteins\_100g, Sugars\_100g, Sodium\_100g, Fruits-vegetables-nuts-estimate-from-ingredients\_100g

### 3 formats de données en entrée :

- Données Transformées (log naturel +1), Normalisées (MinMax scaler)
- Données Transformées (log naturel +1), Normalisées (MinMax scaler) + pnns\_groups\_2 encodées
- Données Normalisées (MinMax scaler)

5 méthodes : Régression linéaire; KNN; Random Forest; Moyenne; Médiane

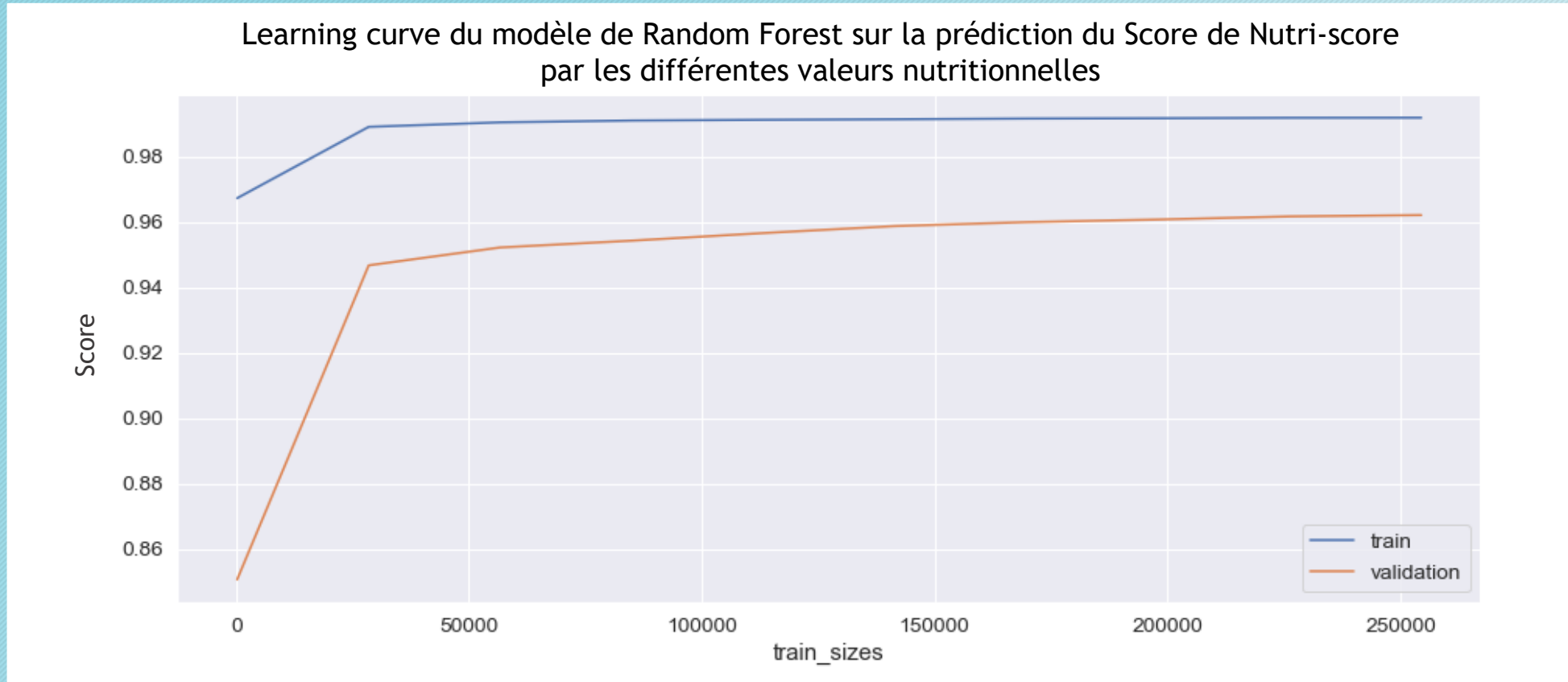
### 3 Critères de décision :

- Erreur moyenne absolue (Moyenne (Valeurs prédites - Valeurs du test set))
- Score du modèle (Coefficient de détermination  $R^2$ )
- Facilité à imputer les données par la suite.

# C) Analyse et prédiction des valeurs manquantes

12

- Optimisation (gridsearchCV et Learning curve) et entraînement du modèle

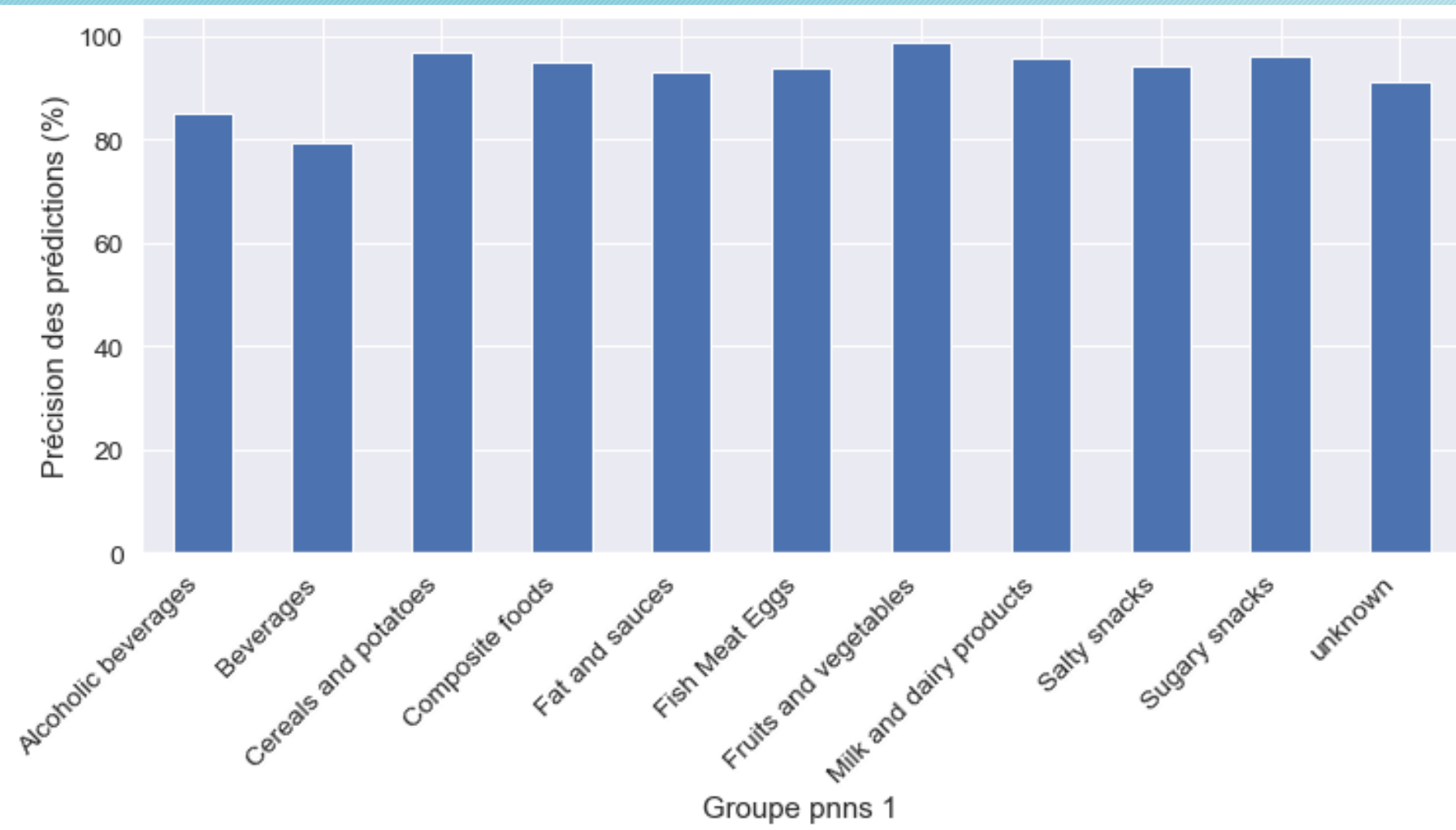


- Normalisation des données de base (minMax Scaler)
- Imputation des valeurs manquantes (Iterative Imputer utilisant le modèle entraîné sur 100 000 produits)
- Restitution des dimensions d'origines (Inverse Scaling)

# C) Analyse et prédiction des valeurs manquantes

13

- Prédiction du Grade Nutri-score en fonction du Score de Nutri-score prédit (et du type d'aliment).



- Taux d'erreur des prédictions : 6,22%
- Les prédictions sont moins précises pour les boissons (alcoolisées ou non)



## Partie III : Analyse Exploratoire des données

14

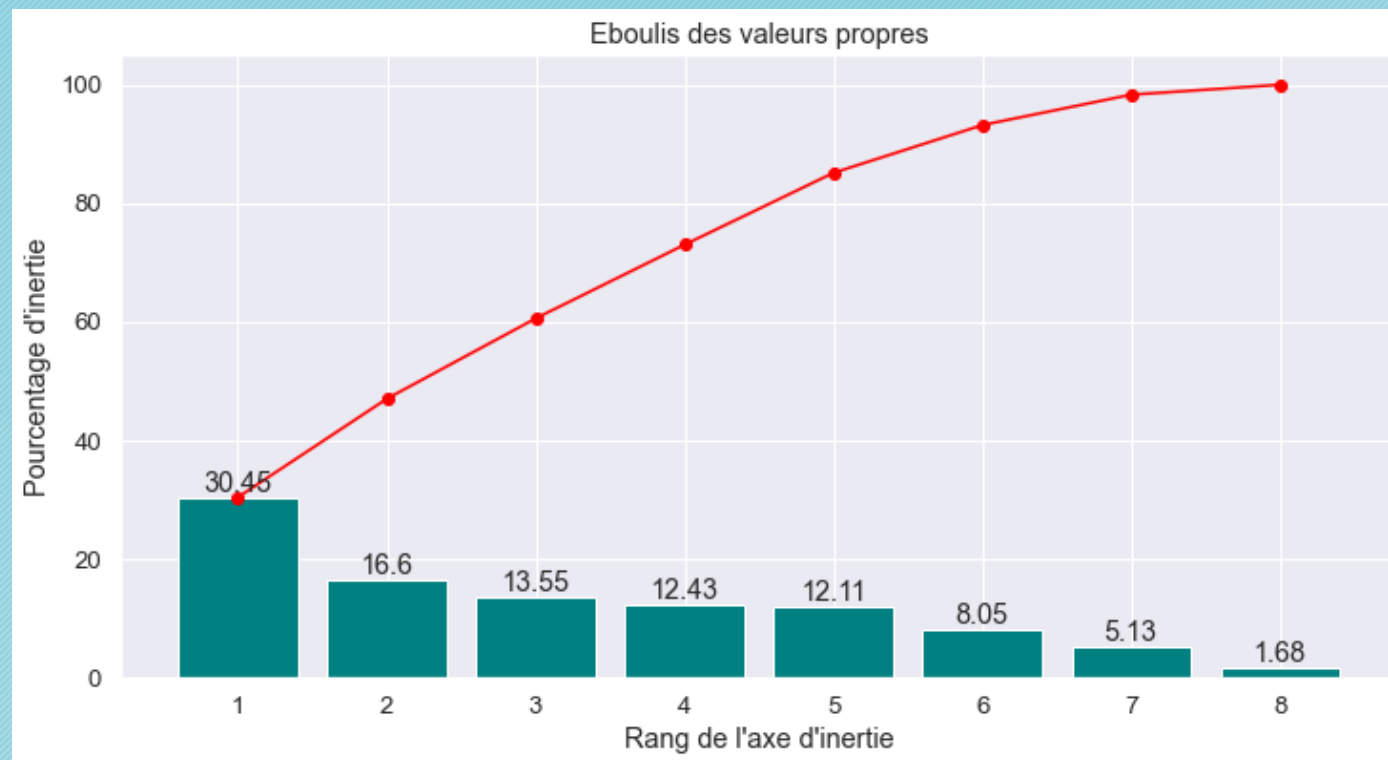
- A) Analyse multivariée des variables nutritionnelles, du Nutri-score et des groupes de produits
- B) Répartition des Grades Nutri-score dans les différentes catégories de produits
- C) Analyse des additifs
- D) Analyses des produits d'origine France
- E) Analyses des produits Bio

# A) Analyse multivariée des variables nutritionnelles, du Nutri-score et des groupes de produits

15

**Principe :** Résumer l'information qui est contenue dans de nombreuses variables en un certain nombre d'axes synthétiques (Composantes principales) en gardant le plus d'information possible.

**Détermination du nombre d'axes d'intérêts :**



**Critère de Kayser :** on ne garde que les composantes  $> (100/p)\%$  où  $p$  est le nombre de variables.

➤  $(100/p)\% = 12,5\% > \text{Axe 1 à 3}$

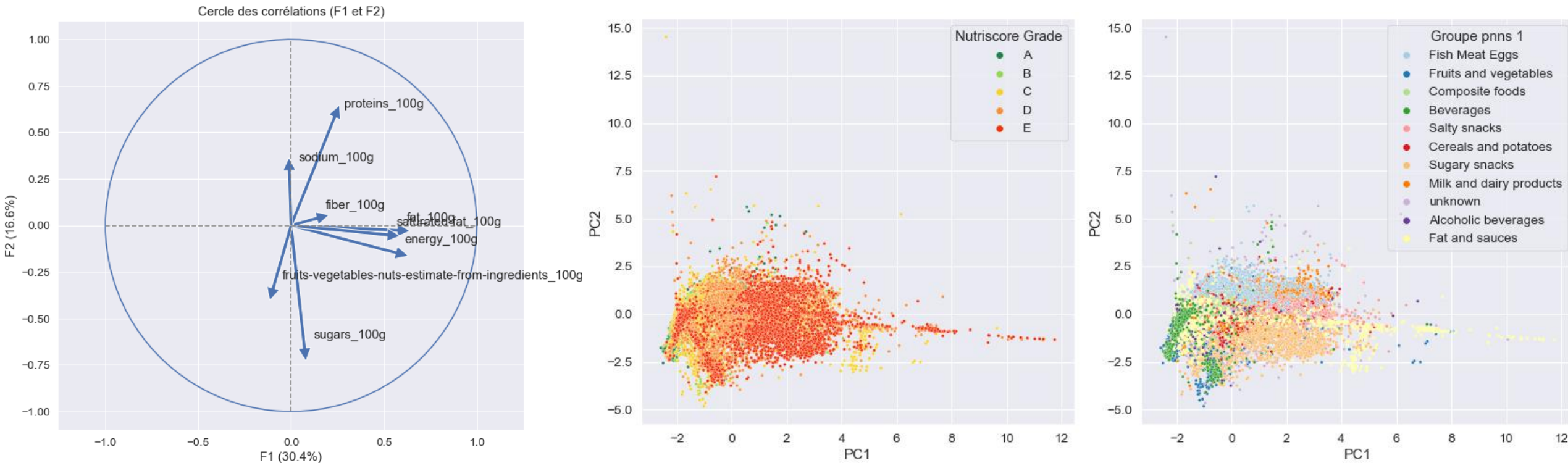
**Méthode du coude :** méthode visuelle

➤ Axe 1 à 5

Nous analyserons les **3 premières dimensions**, expliquant 61% de la variance totale.

# A) Analyse multivariée des variables nutritionnelles, du Nutri-score et des groupes de produits

16

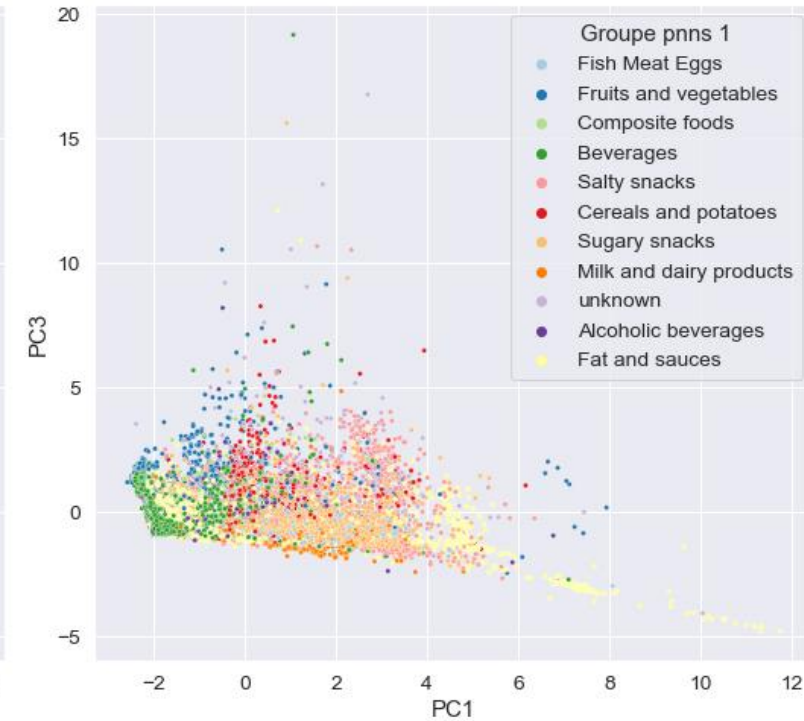
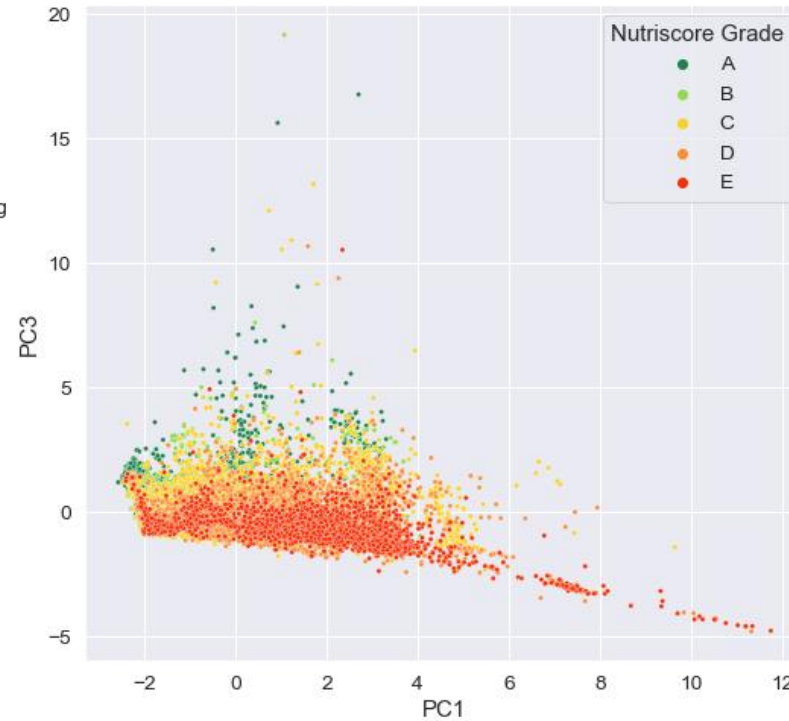
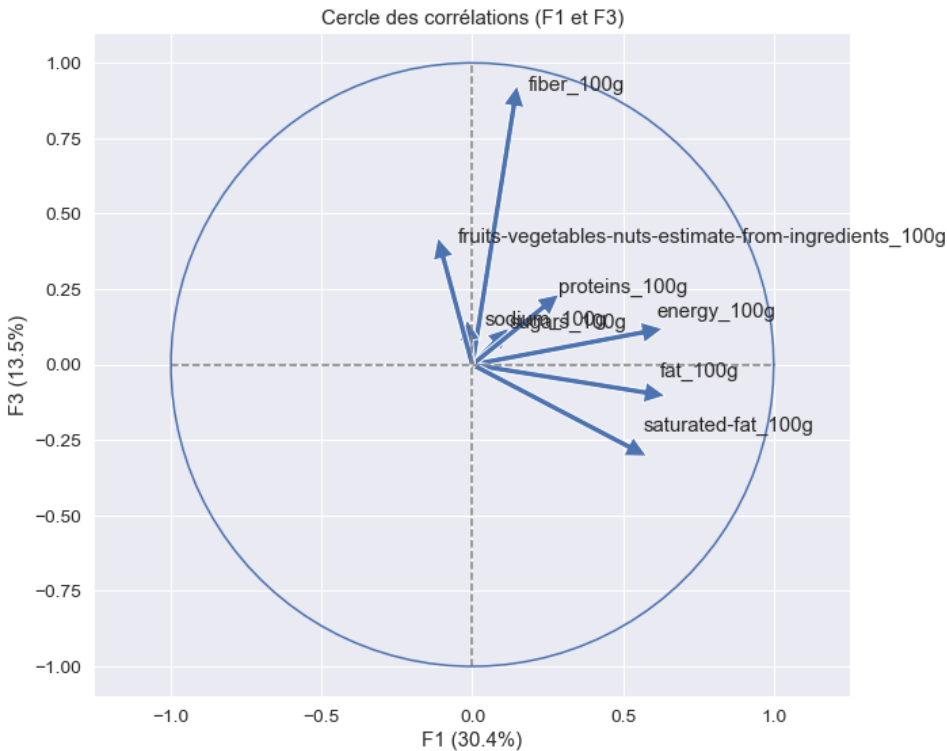


- Les produits de Nutri-score D et E sont plus riches en énergie, lipides et sucres
- Les produits « graisses et sauces » semblent très riches en énergie et lipides
- Les produits « biscuits sucrés » semblent riches en sucre
- Les groupes Poisson-Viande-Œufs et les produits laitiers apparaissent bien comme riches en protéine.



# A) Analyse multivariée des variables nutritionnelles, du Nutri-score et des groupes de produits

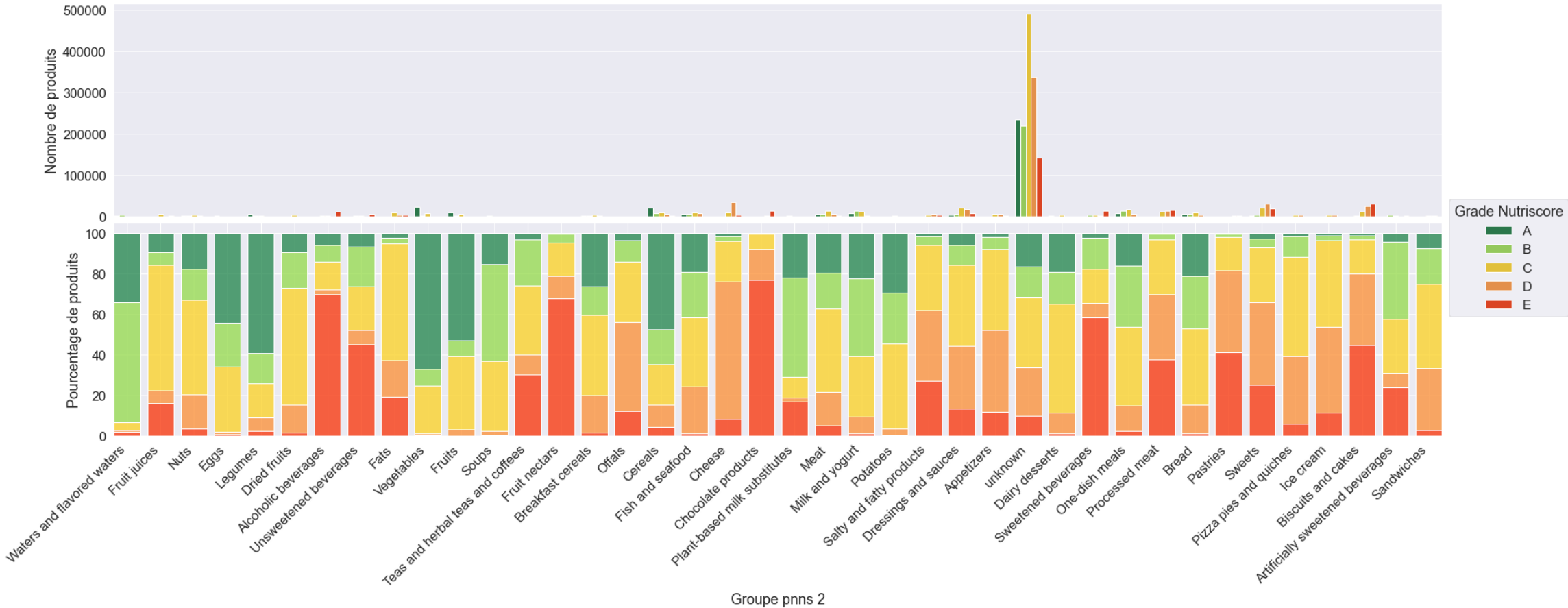
17



- Les produits de Nutri-score A et B sont plus riches en fibre
- Produits de catégorie Fruits et légumes

# B) Répartition des Grades Nutri-score dans les différentes catégories de produits

18



**Le Grade Nutri-score est-il réparti équitablement entre les différentes catégories de produits ?**

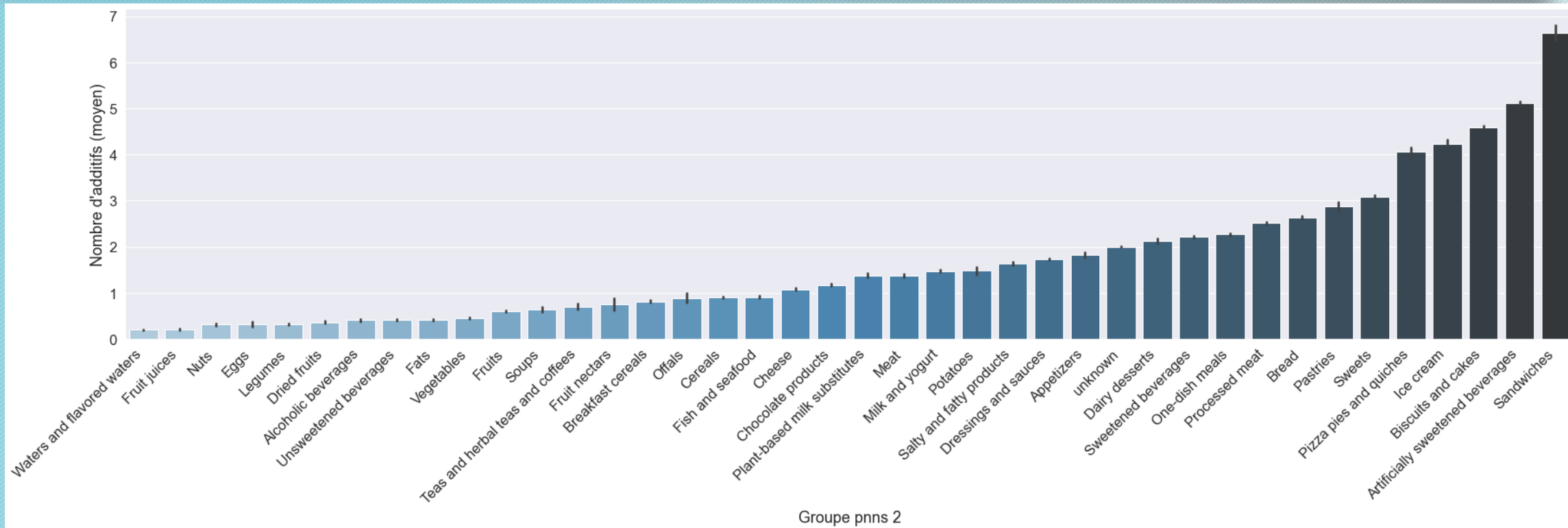
➤ Non, Il y a une différence significative de grade Nutri-score entre les différents groupes ( $\chi^2 = 678746.7$  , p-value = 0)

**Existe-t-il pour chaque produit une alternative de Grade Nutri-score plus faible ?**

➤ Pour chaque catégorie, il existe des alternatives de produits à Nutri-score plus faible

# C) Analyse des additifs : Groupe de produits

19



## ➤ Quelle est l'effet du groupe de produits sur le nombre d'additifs ?

La distribution du nombre d'additifs s'approche d'une Loi de Poisson (données de comptage) : Modèle linéaire généralisé et comparaison d'AIC avec un modèle nul.

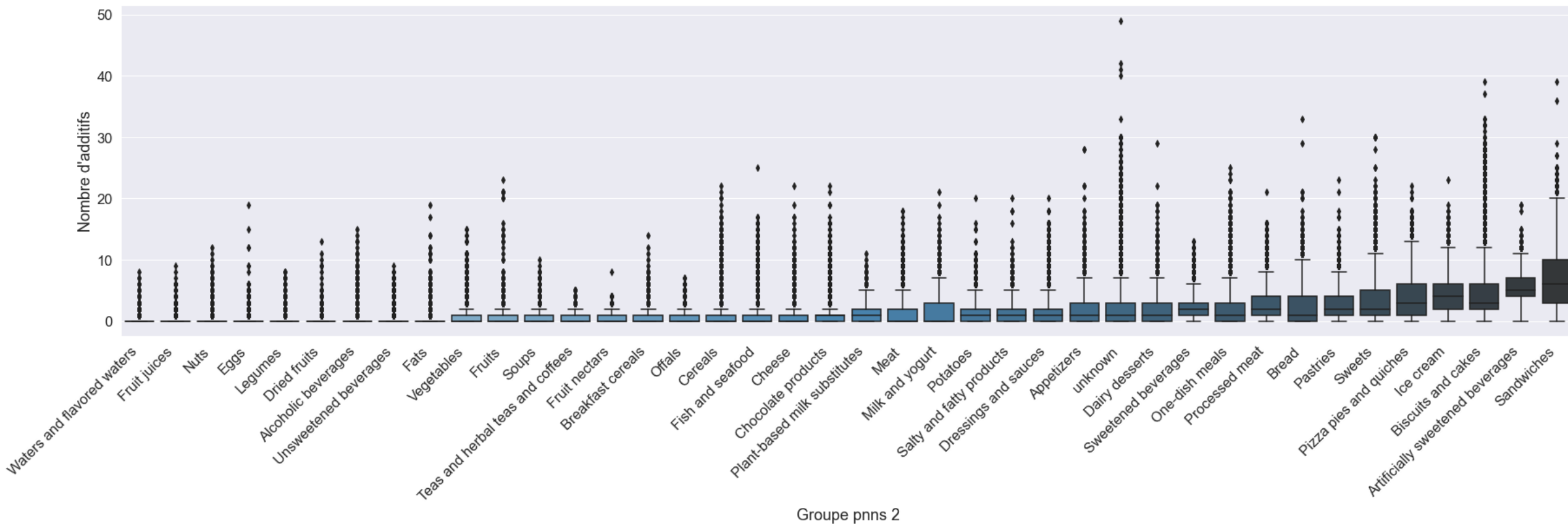
- Modèle nul : Nombre d'additifs  $\sim 1$  > AIC = 3 951 255.46
- Modèle Grade : Nombre d'additifs  $\sim$  Groupe pnns 2 > AIC = 3 460 217.96

➤ Il y a un nombre moyen d'additifs différents entre les différents groupes de produits (validé par Test Post hoc de Dunn)



# C) Analyse des additifs : Groupe de produits

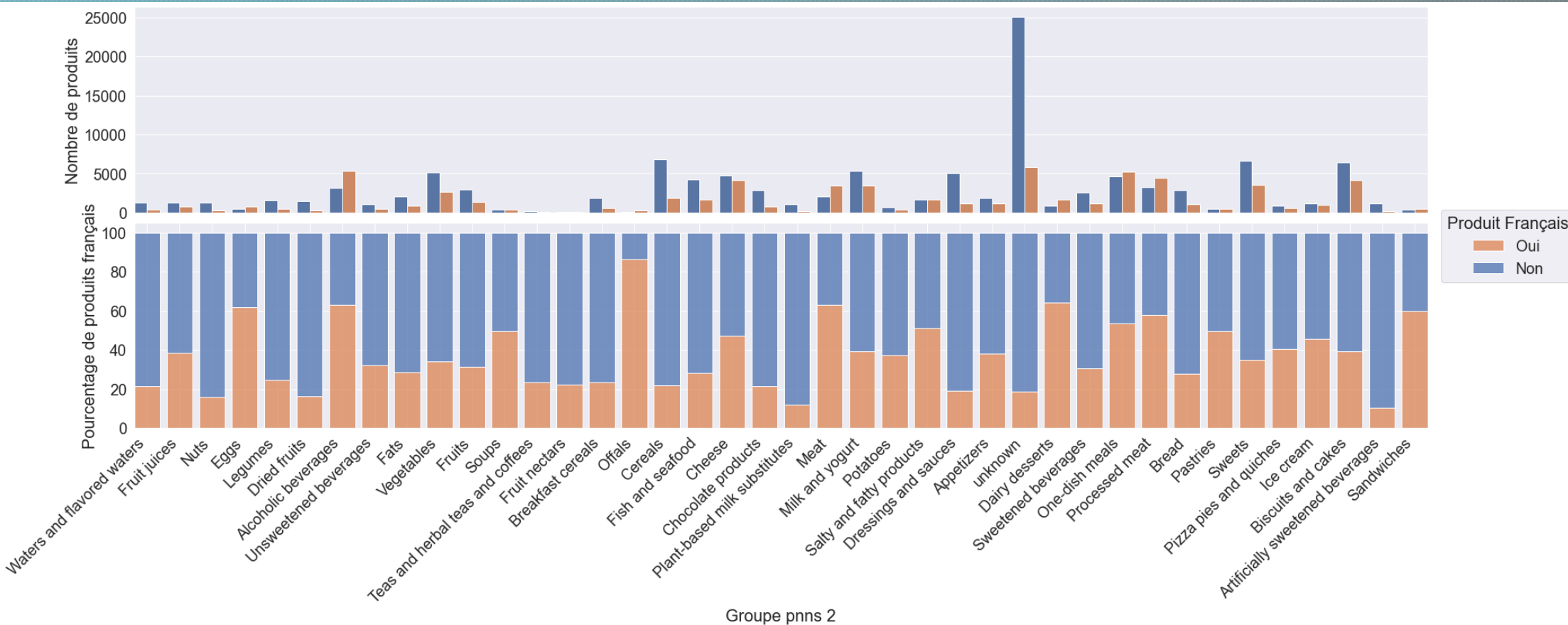
20



**Existe-t-il des alternatives de produits avec moins d'additifs ?**

➤ Oui, pour chaque groupe de produit, il existe des alternatives avec moins (ou sans) additifs.

# D) Analyse des produits d'origine France



Groupe pnns 2

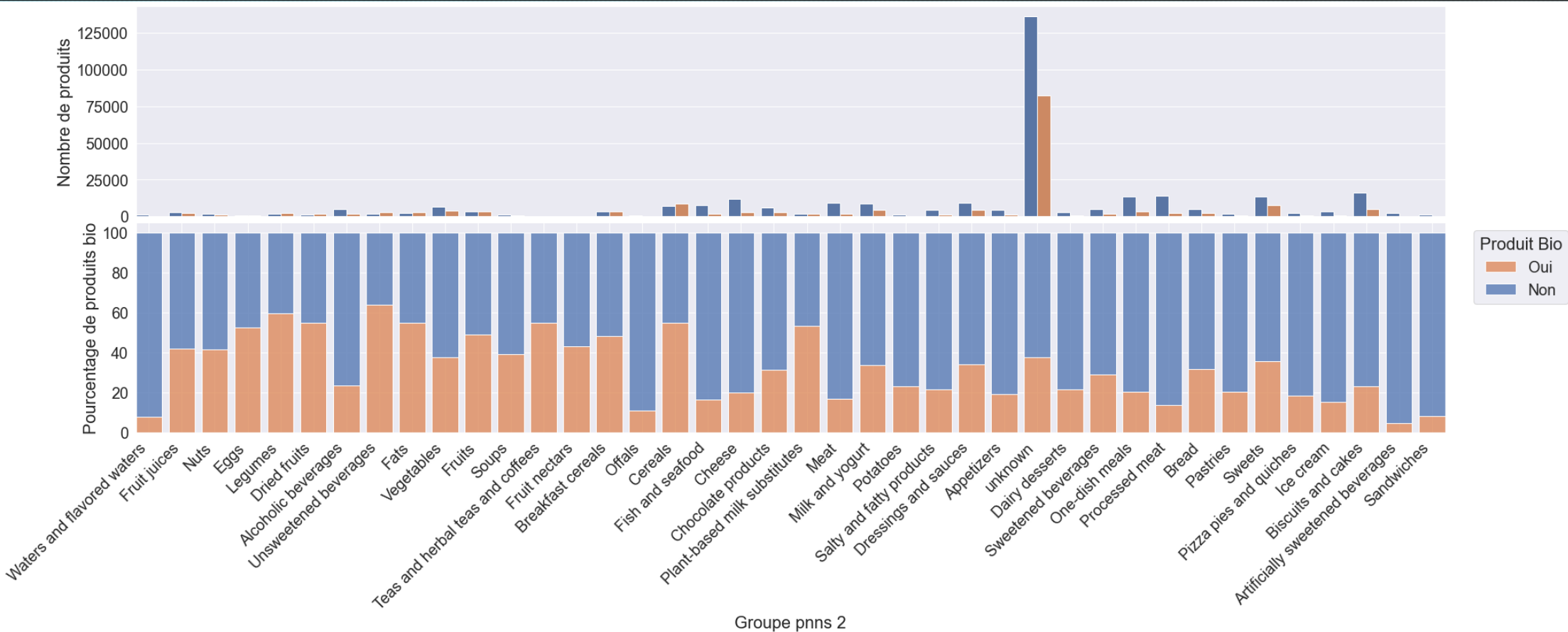
**Les produits Français sont-ils répartis équitablement entre les différentes catégories de produits ?**

➤ Non, Il y a une différence significative d'origine des produits entre les différents groupes ( $\chi^2 = 18221.4$  , p-value = 0)

**Existe-t-il des alternatives d'origine France pour chaque groupe de produits ?**

➤ Oui, pour chaque catégorie, il existe des alternatives françaises

# E) Analyse des produits bio



**Les produits Bio sont-ils répartis équitablement entre les différentes catégories de produits ?**

➤ Non, Il y a une différence significative de répartition des produits bio entre les différents groupes ( $\chi^2 = 25416.9$  ,  $p\text{-value} = 0$ )

**Existe-t-il des alternatives Bio pour chaque groupe de produits ?**

➤ Oui, pour chaque catégorie, il existe des alternatives bio.

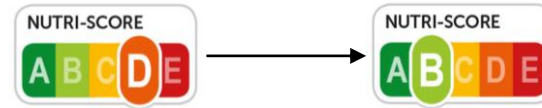




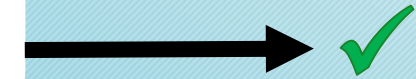
Proposition d'aliments  
de la même famille

## Manger plus sain

- Nutri-score



- - d'additifs



## Manger plus responsable

- Origine France



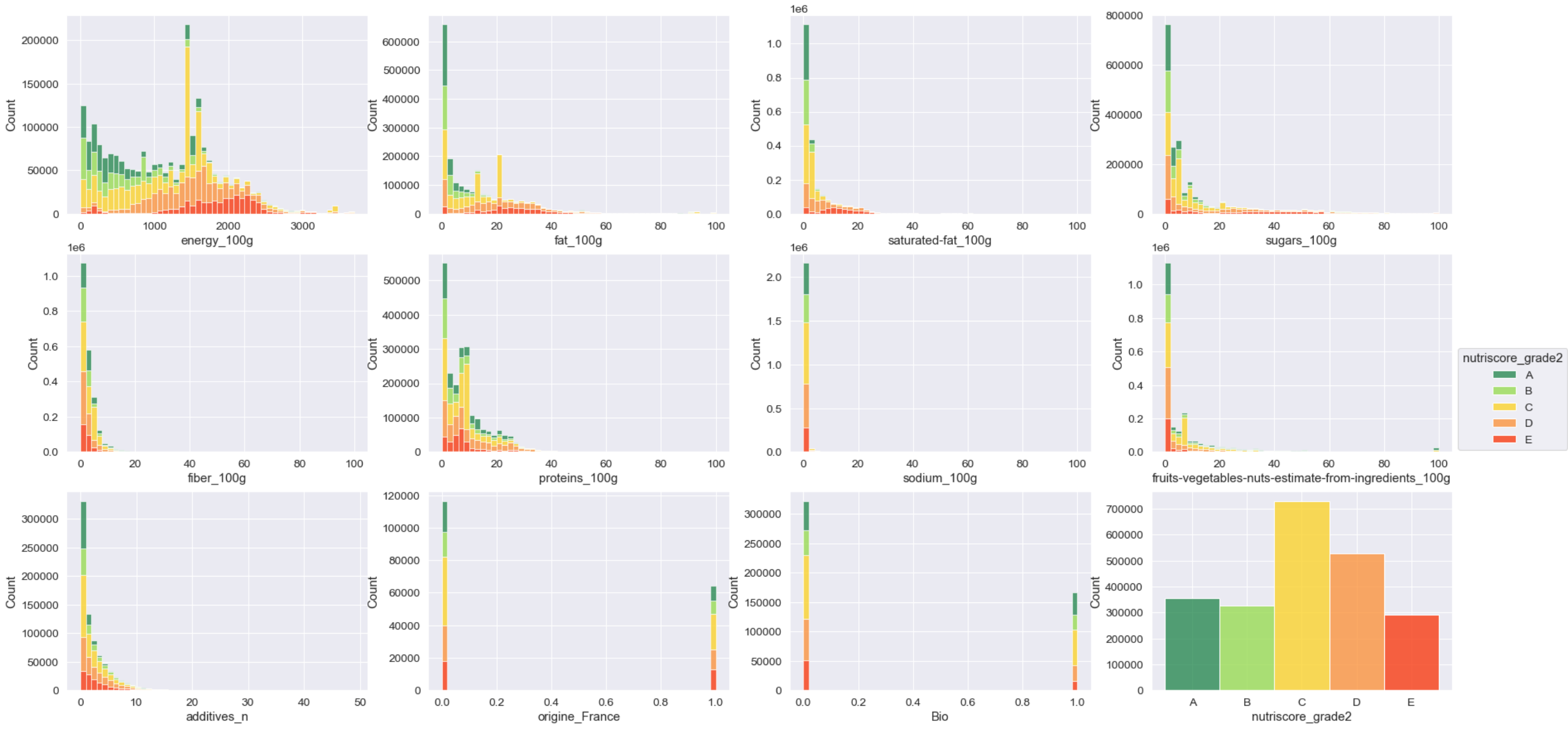
- Label Biologique





# Annexe 1 : Répartition empiriques des données

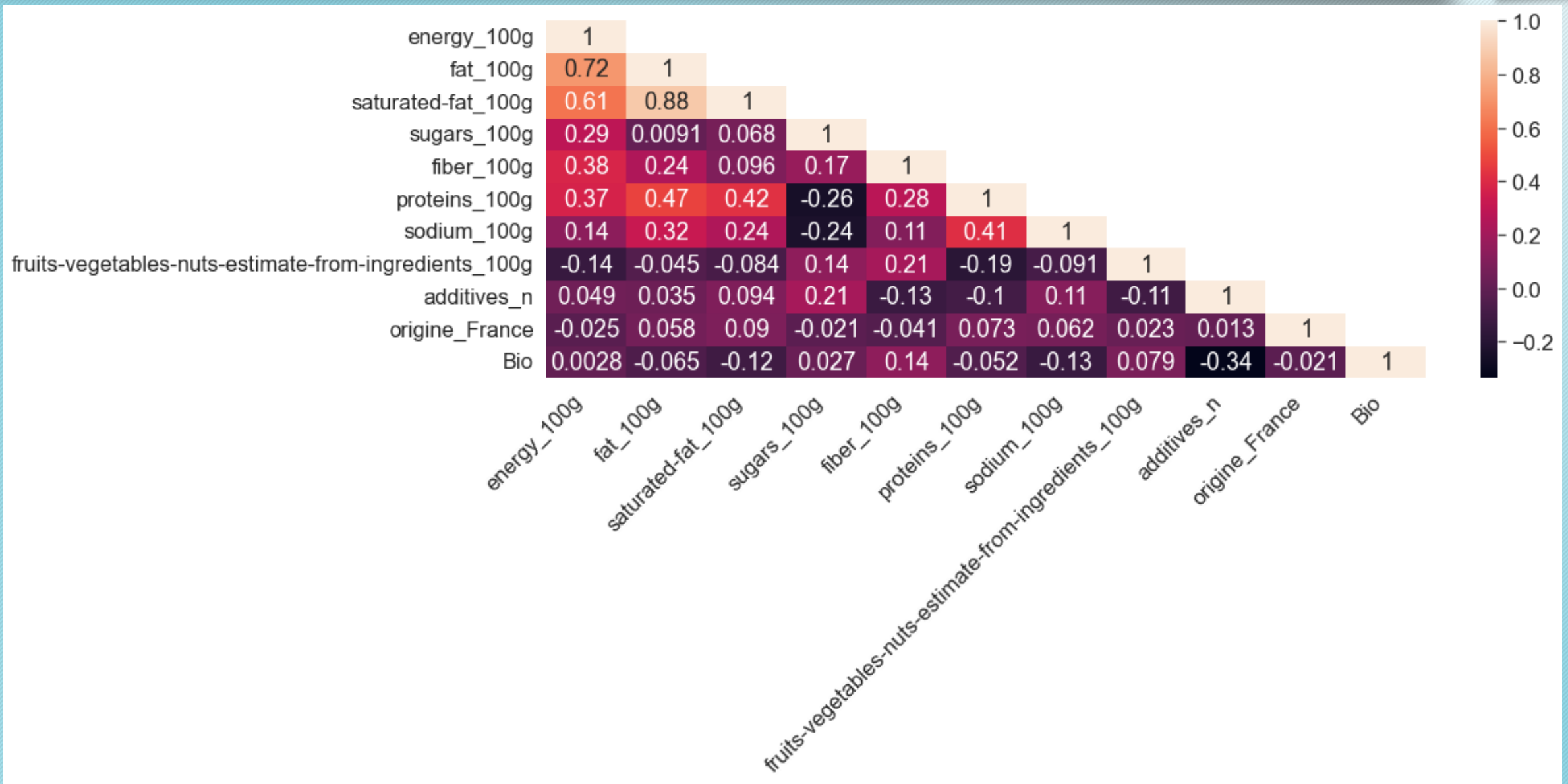
24





# Annexe 2 : Coefficients de corrélation (Spearman) entre les différentes variables

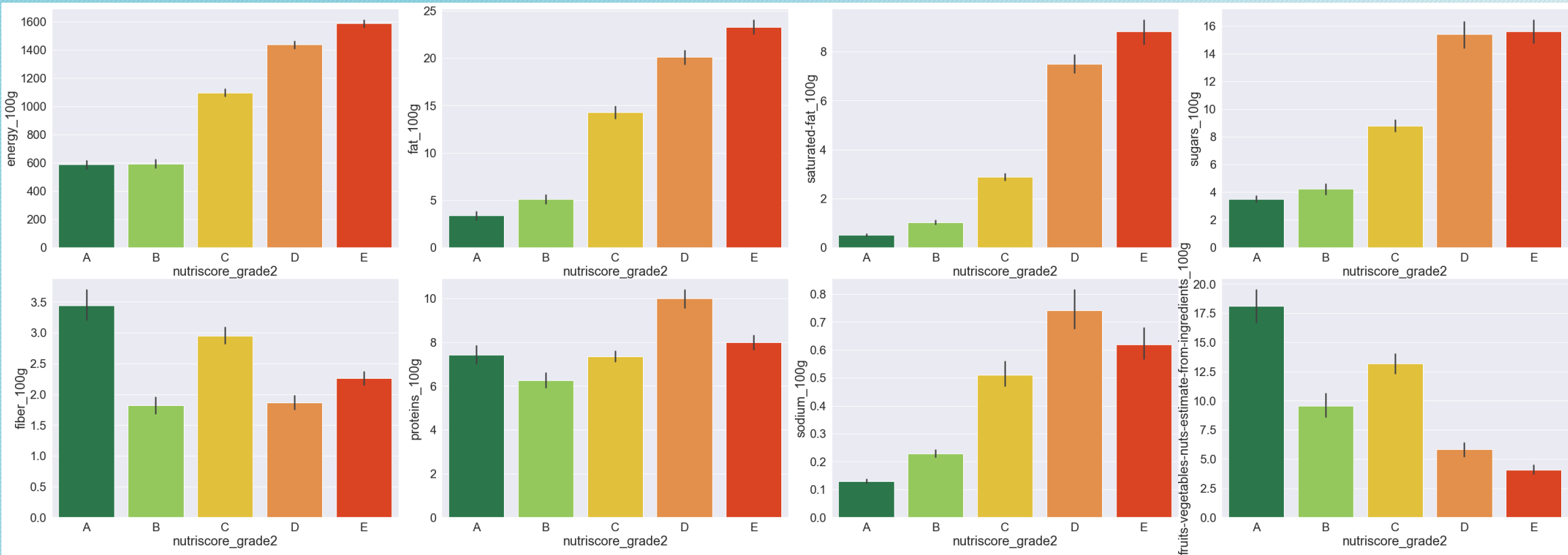
25



# Annexe 3 : Coefficients de corrélation (Spearman) entre les différentes variables

26

➤ Quelle est la relation entre chaque variable nutritionnelle et le grade Nutri-score ?



Pour chaque variable, les conditions d'applications de l'ANOVA ne sont pas respectées

- Normalité des résidus : Test de Kolmogorov Smirnov = **X**
- Homoscédasticité : Test de Bartlett = **X**

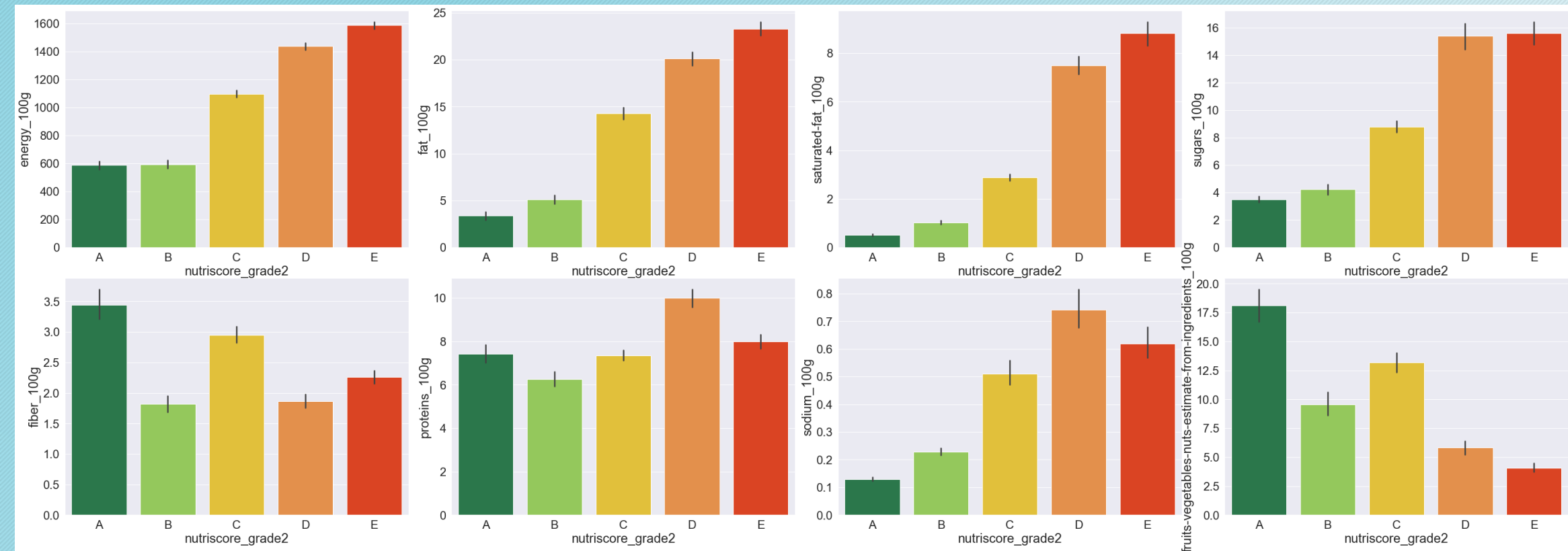
➤ Equivalent de l'ANOVA en non paramétrique : Test de Kruskal Wallis

➤ Comparaisons 2 à 2 : Test post-hoc de Dunn

# Annexe 3 : Coefficients de corrélation (Spearman) entre les différentes variables

26

➤ Quelle est la relation entre chaque variable nutritionnelle et le grade Nutri-score ?

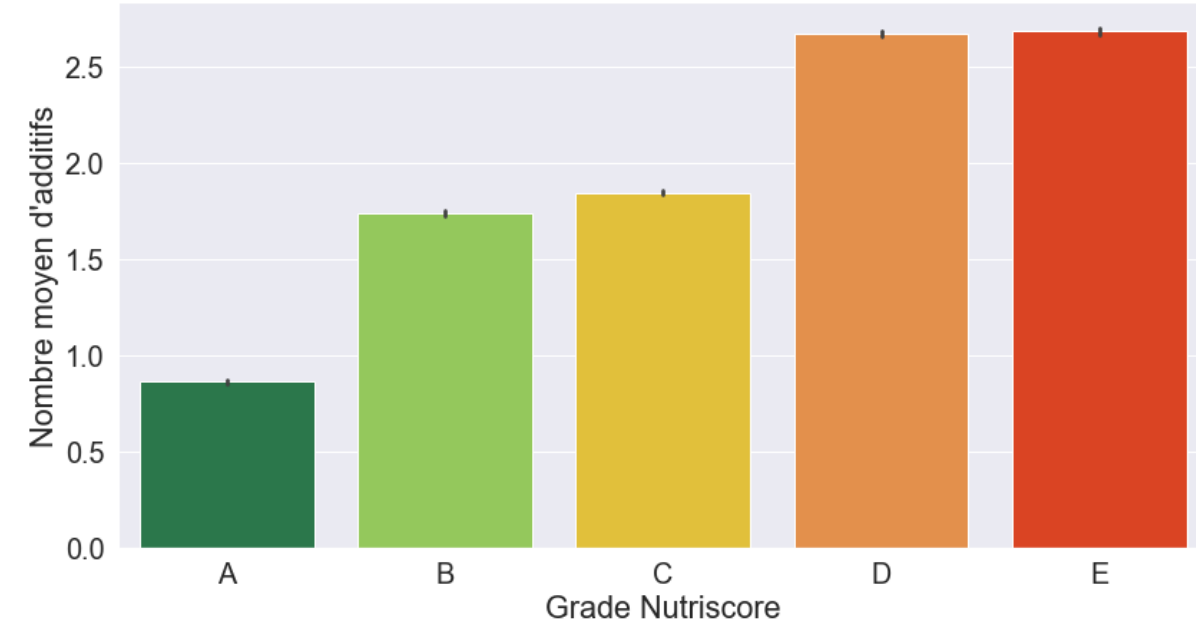
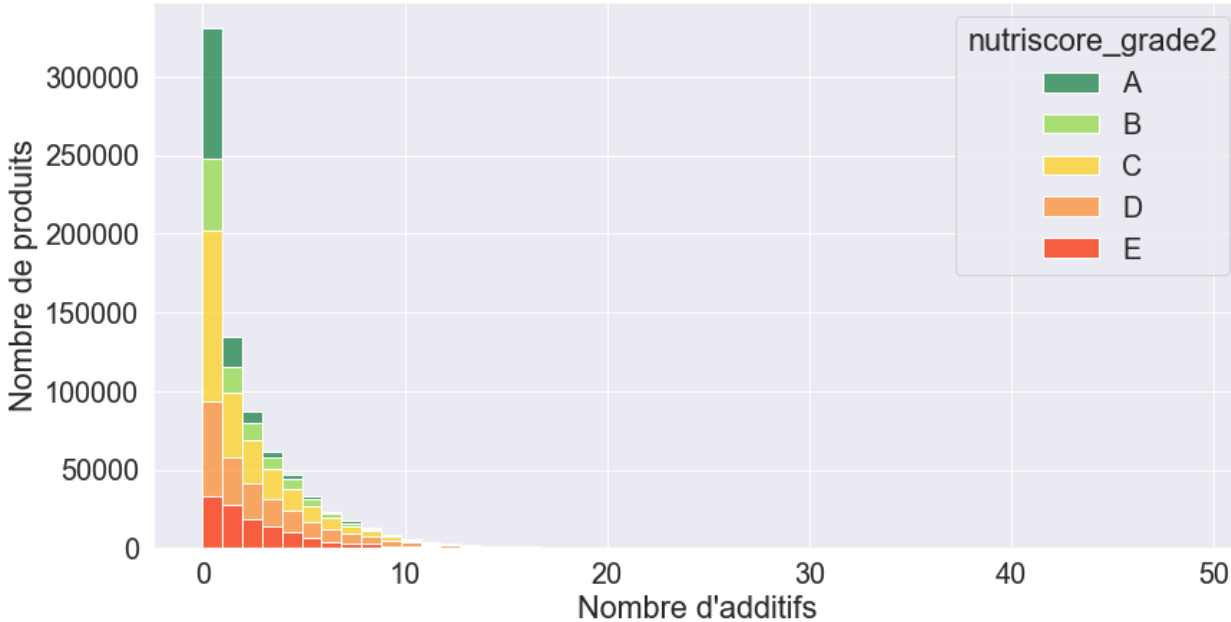


- L'effet de chaque variable de valeur nutritive sur le Grade Nutri-score est statistiquement significatif (p-value = 0)
- Les aliments de grade Nutri-score D et E sont plus riches en énergie, acides gras ( saturés + insaturés ) sucres et sel.
- Les aliments de grade Nutri-score A sont plus riches en fibres, fruits et légumes.



# Annexe 4 : Analyse des additifs et du Grade Nutri-score

27



## Quelle est la relation entre le nombre d'additifs et le Grade Nutri-score ?

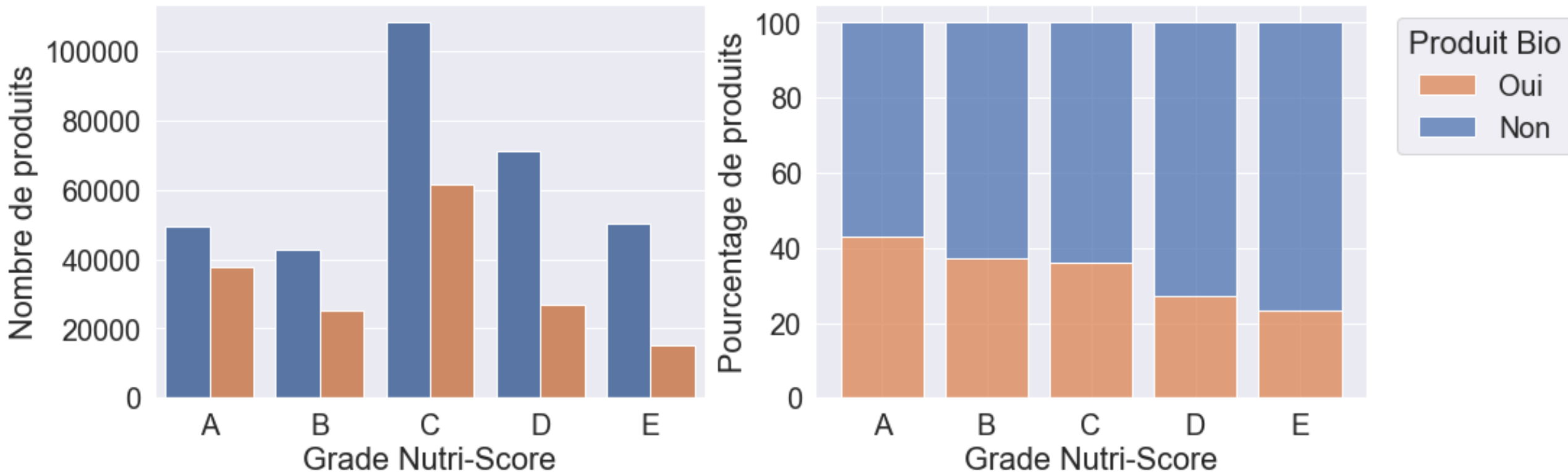
La distribution du nombre d'additifs s'approche d'une Loi de Poisson (données de comptage) : Modèle linéaire généralisé et comparaison d'AIC avec un modèle nul.

- Modèle nul : Nombre d'additifs  $\sim 1$  > AIC = 3 951 255.46
  - Modèle Grade : Nombre d'additifs  $\sim$  Grades Nutri-scores > AIC = 3 779 254.99
- Il y a un nombre d'additifs différents entre les grades Nutri-score (validé par Test Post hoc de Dunn)

# Annexe 5 : Analyse des produits bio et du Nutri-score

28

- Quelle est la relation entre les produits Bio et le Grade Nutri-score ?

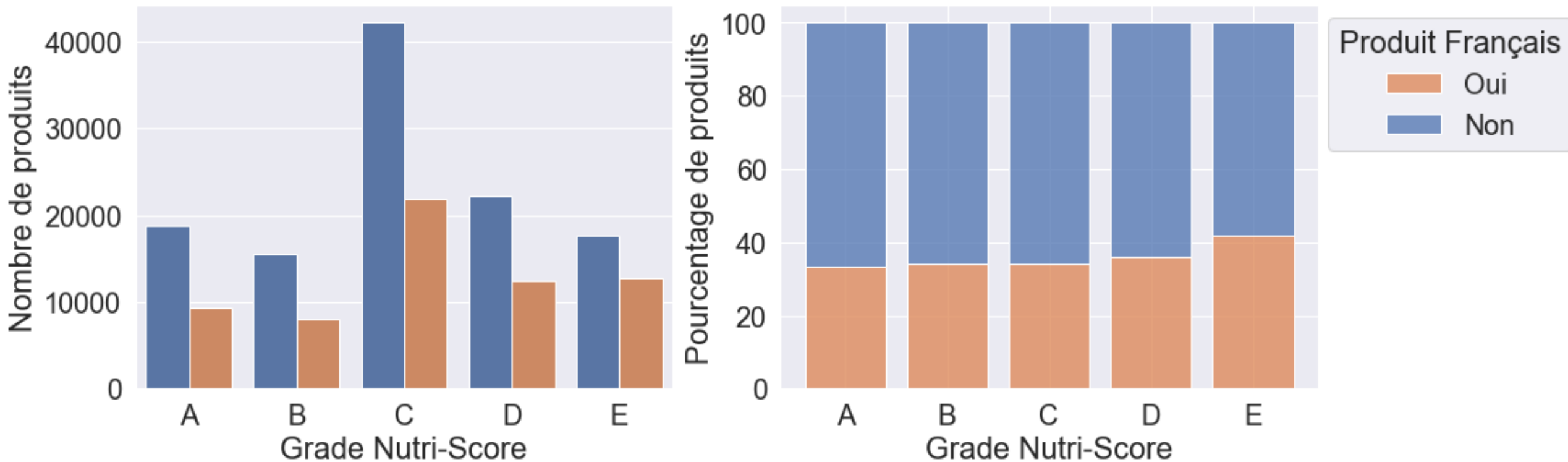


- Il y a une différence significative de grade Nutri-score entre les différents groupes ( $\chi^2 = 9166$ , p-value=0)
- Pour chaque catégorie, il existe des alternatives de produits à Nutri-score plus faible

# Annexe 6 : Analyse des produits d'origine France et du Nutri-score

29

- Quelle est la relation entre les produits Français et le Grade Nutri-score ?



- Il y a une différence significative de grade Nutri-score entre les différents groupes ( $\chi^2 = 669.4$ ,  $p\text{-value} < 10^{-142}$ )
- Pour chaque catégorie, il existe des alternatives de produits à Nutri-score plus faible