



OPENCLASSROOMS

## Projet 3

-

Préparez des données pour un  
organisme de santé publique



CentraleSupélec

Pierrick BERTHE

Formation Expert en Data Science  
*Openclassrooms – CentraleSupélec*

*août 2023 → avril 2024*



# Sommaire



I – Problématique

II – Présentation du jeu de données

III - Nettoyage des données

IV – Analyses univariées

V – Analyse multivariée

VI – Faisabilité du système d’auto-complétion

VII – RGPD

VIII - Conclusion



# Problématique



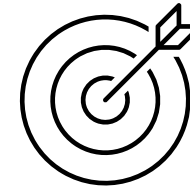
**Santé publique France** est un établissement public à caractère administratif français, placé sous la tutelle du ministère chargé de la santé, créé le 1er mai 2016.



Cet établissement souhaite améliorer sa base de données **Open Food Facts** qui est mise à disposition de particuliers et d'organisations afin de leur permettre de connaître la qualité nutritionnelle de produits. Mise en place d'un système **d'auto-complétion** pour aider les usagers à remplir plus efficacement la base de données.

## Missions :

- Nettoyage des données
- Analyses univariées
- Analyse multivariée
- Etudier la faisabilité du système d'auto-complétion
- Rappel RGPD





# Sommaire



I – Problématique

**II – Présentation du jeu de données**

III - Nettoyage des données

IV – Analyses univariées

V – Analyse multivariée

VI – Faisabilité du système d'auto-complétion

VII – RGPD

VIII - Conclusion



# Présentation du jeu de données



fr.openfoodfacts.org.products.csv

320\_749 lignes



162 colonnes

## ➤ Descriptif des produits alimentaires :

- Les **informations générales** des produits : code barre, nom, date de modification, etc.
- Un ensemble de **tags** : catégorie du produit, localisation, origine, etc.
- Les **ingrédients** composant les produits et leurs additifs éventuels.
- Des **informations nutritionnelles** : quantité en grammes d'un nutriment pour 100 grammes du produit.

## ➤ Valeurs manquantes :

- 39\_604\_863 de NaN pour 51\_961\_338 observations (76.22 %)
- 154 / 162 colonnes concernées

## ➤ Doublons

Pas de doublons sur la colonne du code barre.



# Sommaire



I – Problématique

II – Présentation du jeu de données

**III - Nettoyage des données**

IV – Analyses univariées

V – Analyse multivariée

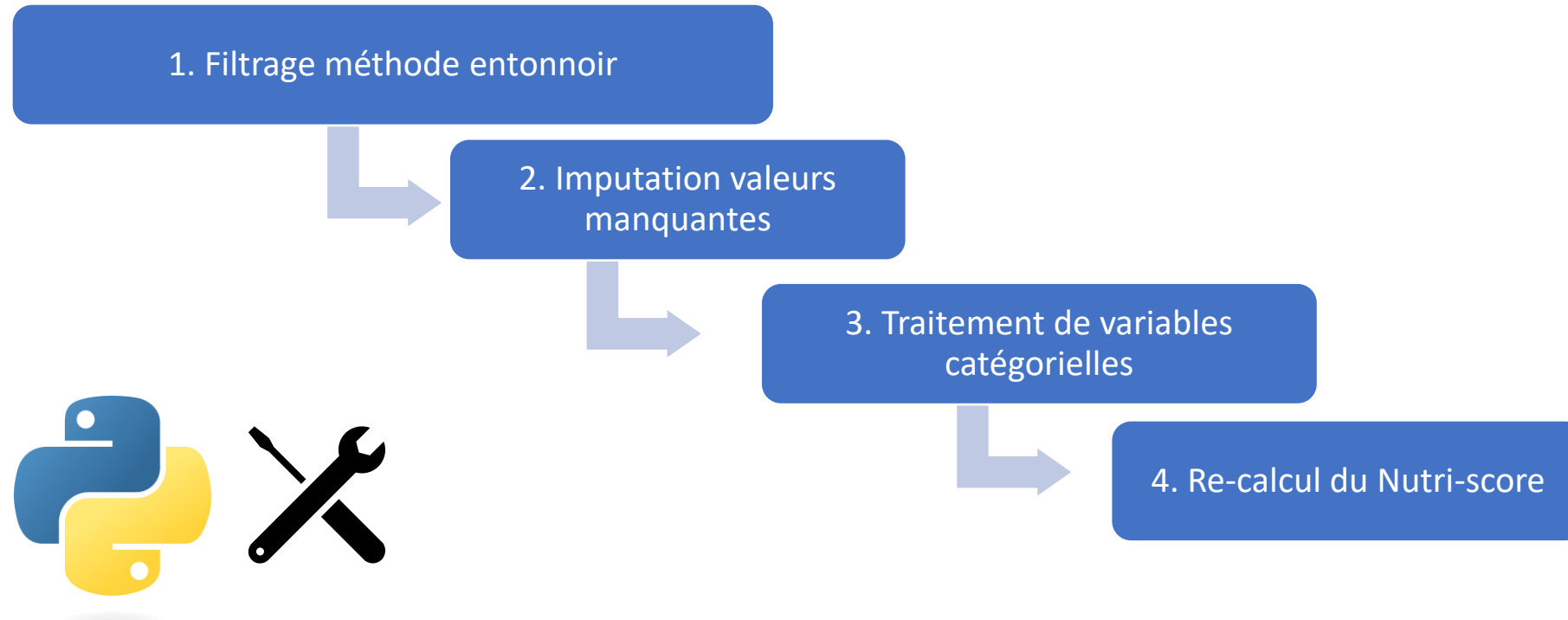
VI – Faisabilité du système d'auto-complétion

VII – RGPD

VIII - Conclusion



# Nettoyage des données





# Nettoyage des données



Nom	Utilisation	Fonctions spécifiques
Anaconda	Gestion de package Gestion d'environnement virtuel	Conda : installation de package via le terminal
Visual Studio Code 1.83.0	Structurer la démarche Exécuter code par étape Expliquer la démarche (markdown)	
Python 3.11.6	Appel aux librairies Boucles for pour générer plusieurs calculs et graphiques	Boucles, listes, dictionnaires, librairies, méthodes
Pandas 2.1.0	Manipulation de données Représentation des données	Manipulation de Dataframe : création, copie, filtres, tris, description, concaténation
Matplotlib 3.7.2 Seaborn 0.12.2	Génération de graphiques de visualisation	Barplot, scatterplot, lineplot, distplot, heatmap
Numpy 1.25.2	Manipulation de matrices et fonctions mathématiques	Histogram, argmax, arange, object, number
Missingno 0.5.2	Représentation graphique pour valeurs manquantes	Matrice de NaN
Sklearn 1.3.0	Apprentissage automatique et modélisation statistique	SimpleImputer, KNNImputer, StandardScaler, PCA
Scipy 1.11.2	Calculs de mathématiques complexes ou de problèmes scientifiques	Stats, chi2_contingency, shapiro, kruskal





# Nettoyage des données



1/ Filtrage méthode « entonnoir »

a) Filtrage colonne par % remplissage

=> Suppression des colonnes < 25% de remplissage



# Nettoyage des données



## 1/ Filtrage méthode « entonnoir »

### a) Filtrage colonne par % remplissage

=> Suppression des colonnes < 25% de remplissage

### b) Filtrage colonnes redondantes

=> Suppression des colonnes identiques:

- Tags (par ex : « *countries* » / « *countries\_tags* » / « *countries\_fr* »)
- Temps (par ex : « *created\_t* » / « *created\_datetime* »)



# Nettoyage des données



## 1/ Filtrage méthode « entonnoir »

a) Filtrage colonne par % remplissage

=> Suppression des colonnes < 25% de remplissage

b) Filtrage colonnes redondantes

=> Suppression des colonnes identiques:

- Tags (par ex : « countries » / « countries\_tags » / « countries\_fr »)
- Temps (par ex : « created\_t » / « created\_datetime »)

c) Filtrage des colonnes inutiles

=> Suppression colonne : « url » / « creator » / « Nutri-score\_UK »



# Nettoyage des données



## 1/ Filtrage méthode « entonnoir »

a) Filtrage colonne par % remplissage

=> Suppression des colonnes < 25% de remplissage

b) Filtrage colonnes redondantes

=> Suppression des colonnes identiques:

- Tags (par ex : « countries » / « countries\_tags » / « countries\_fr »)
- Temps (par ex : « created\_t » / « created\_datetime »)

c) Filtrage des colonnes inutiles

=> Suppression colonne : « url » / « creator » / « Nutri-score\_UK »

d) Filtrage par pays

=> Conservation des produits disponibles en France



# Nettoyage des données



## 1/ Filtrage méthode « entonnoir »

a) Filtrage colonne par % remplissage

=> Suppression des colonnes < 25% de remplissage

b) Filtrage colonnes redondantes

=> Suppression des colonnes identiques:

- Tags (par ex : « countries » / « countries\_tags » / « countries\_fr »)
- Temps (par ex : « created\_t » / « created\_datetime »)

c) Filtrage des colonnes inutiles

=> Suppression colonne : « url » / « creator » / « Nutri-score\_UK »

d) Filtrage par pays

=> Conservation des produits disponibles en France

e) Filtrage valeurs aberrantes

=> Application de **règles métiers** :

- $0 \leq \text{Energie sur 100g} \leq 3_762 \text{ kJ}$
- $0 \leq \text{Quantité sur 100g} \leq 100\text{g}$
- Acide gras trans / saturé  $\leq$  Lipide
- Sucres  $\leq$  Glucide



# Nettoyage des données



## 2. Imputation valeurs manquantes

### a) Colonnes essentielles

=> **Suppression des produits** où il manque le code-barre, le nom du produit ou le pays



# Nettoyage des données



## 2. Imputation valeurs manquantes

### a) Colonnes essentielles

=> **Suppression des produits** où il manque le code-barre, le nom du produit ou le pays

### b) Ingrédients mineurs

=> **Imputation par la valeur « 0 »** des valeurs manquantes des ingrédients dont l'affichage sur l'emballage n'est pas obligatoire  
*(fibres, sodium, calcium, vitamines, etc...)*

### c) Ingrédients majeurs

=> **Imputation par l'algorithme KNN** des valeurs manquantes des ingrédients dont l'affichage sur l'emballage est obligatoire  
*(énergie, lipide, acides gras saturés, glucide, sucre, protéine, sel)*



# Nettoyage des données



## 3. Traitement de variables catégorielles

a) Labélisation valeurs  
manquantes

=> Attribution valeur « inconnu »





# Nettoyage des données



## 3. Traitement de variables catégorielles

a) Labélisation valeurs manquantes

=> Attribution valeur « inconnu »



b) Séparation des chaînes de caractère

=> séparation des valeurs par la virgule si nécessaire

`categories_fr`  
Plats préparés,Plats à base de viande,Plats à ...



# Nettoyage des données



## 3. Traitement de variables catégorielles

a) Labélisation valeurs manquantes

=> Attribution valeur « inconnu »

b) Séparation des chaînes de caractère

=> séparation des valeurs par la virgule si nécessaire

`categories_fr`  
Plats préparés,Plats à base de viande,Plats à ...

c) Suppression sources d'erreur

=> Uniformisation bas de casse, suppression des accents, des caractères spéciaux, des espaces vides, des préfixes relatifs aux pays



# Nettoyage des données



## 3. Traitement de variables catégorielles

a) Labélisation valeurs manquantes

=> Attribution valeur « inconnu »

b) Séparation des chaînes de caractère

=> séparation des valeurs par la virgule si nécessaire

c) Suppression sources d'erreur

=> Uniformisation bas de casse, suppression des accents, des caractères spéciaux, des espaces vides, des préfixes relatifs aux pays

Comptage des occurrences

=> Pour chaque valeur unique



# Nettoyage des données



Comptage des occurrences

« categories\_fr » => 6\_518 catégories



« main\_categories\_fr » => 956 catégories





# Nettoyage des données



## Comptage des occurrences

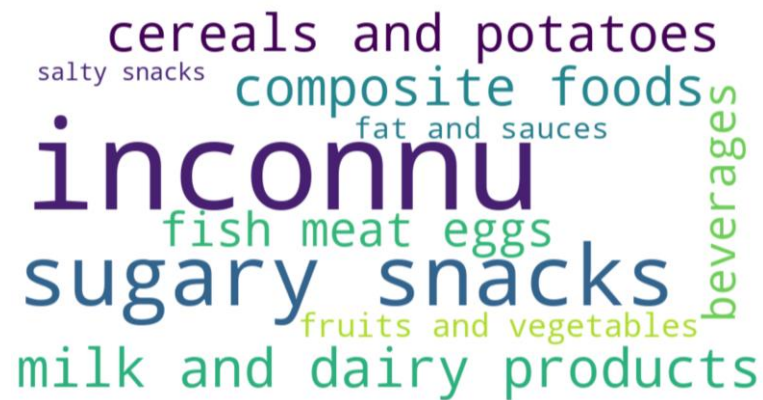
« categories\_fr » => 6\_518 catégories



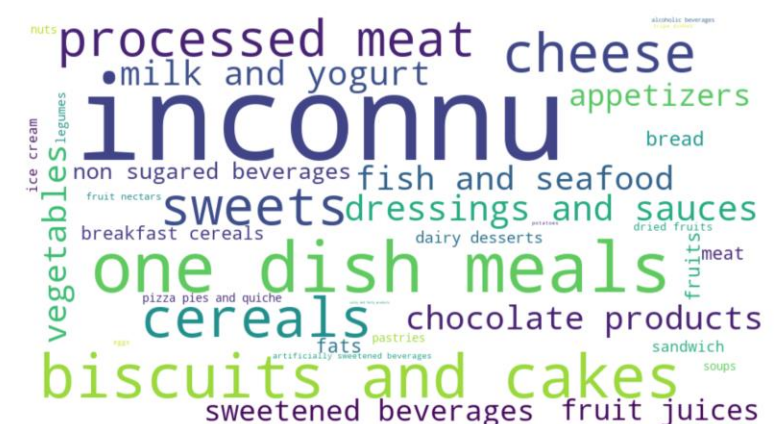
« main\_categories\_fr » => 956 catégories



« pnns\_group\_1 » => 10 catégories ✓



« pnns\_group\_2 » => 37 catégories ✓





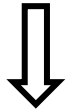
# Nettoyage des données



## 4. Re-calcul du Nutri-score



Méthode de calcul de 2017



Méthode de calcul de 2023





# Nettoyage des données



## 4. Re-calcul du Nutri-score



Méthode ~~de calcul~~ de 2017



Méthode de calcul de 2023



Élément /100g	Points
Energie (KJ)	0 à 10
Sucres (g)	0 à 15*
Acides gras saturés (g)	0 à 10
Sel (g)	0 à 20



Composante N  
0 à 55 points



Règles de calcul

Score nutritionnel



Élément /100g	Points
Protéines (g)**	0 à 7***
Fibres (g)	0 à 5
Fruits, légumes et légumineuses (%)	0 à 5****



Composante P  
0 à 17 points

Moins bonne qualité  
nutritionnelle

NUTRI-SCORE  
A B C D E

NUTRI-SCORE  
A B C D E

NUTRI-SCORE  
A B C D E

NUTRI-SCORE  
A B C D E

NUTRI-SCORE  
A B C D E

Meilleure qualité  
nutritionnelle

\*10 points maximum pour les boissons

\*\*En fonction des points "défavorables" N et du nombre de points accordés pour la composante "fruits, légumes et légumineuses", les protéines sont prises en compte ou non.

\*\*\*2 points maximum pour les viandes rouges

\*\*\*\*6 points maximum pour les boissons

La présence d'édulcorants non nutritifs pour les boissons ajoute 4 points au total N

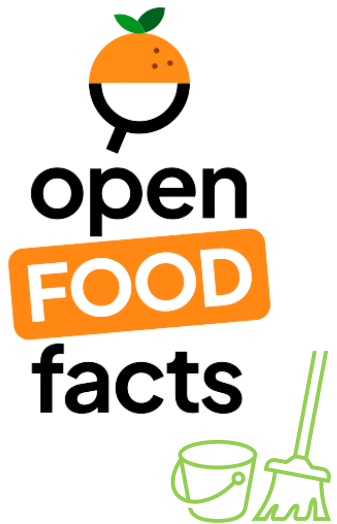
=> Imputation des valeurs manquantes par les valeurs recalculées :

- Score nutritionnel

- Nutri-score



# Nettoyage des données



fr.openfoodfacts.org.products.csv  
dataset\_openfoodfacts\_cleaned.csv



## ➤ Descriptif des produits alimentaires :

- Les **informations générales** des produits : code barre, nom, date de modification, etc.
- Un ensemble de **tags** : catégorie du produit, pays de vente.
- Des **informations nutritionnelles** : quantité en grammes d'un nutriment pour 100 grammes du produit.

## ➤ Valeurs manquantes :

- ~~— 39\_604\_863 de NaN pour 51\_961\_338 observations (76.22 %)~~
- 110\_059 de NaN pour 2\_222\_496 observations (4.95 %)
- ~~— 154 / 162 colonnes concernées~~
- 5 / 36 colonnes concernées





# Sommaire



I – Problématique

II – Présentation du jeu de données

III - Nettoyage des données

**IV – Analyses univariées**

V – Analyse multivariée

VI – Faisabilité du système d’auto-complétion

VII – RGPD

VIII - Conclusion



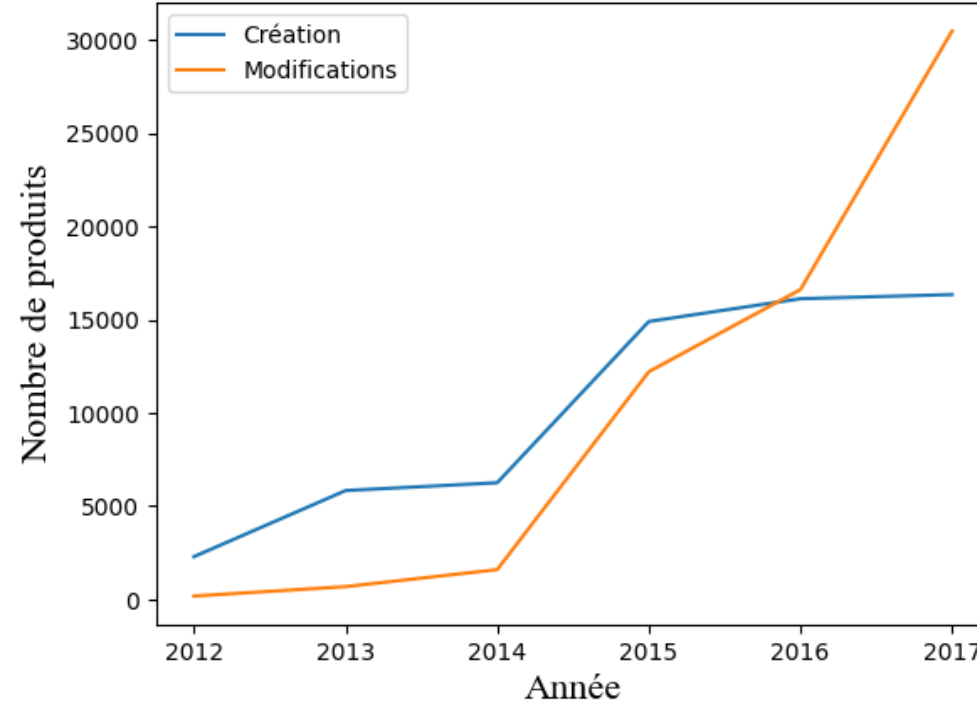
# Analyses univariées



Nombre de création de nouvelle ligne  
de produit et de modification sur le  
jeu de données



Evolution des créations et modifications de produits par année



**=> Le jeu de données n'est plus mis à jour depuis 6 ans**



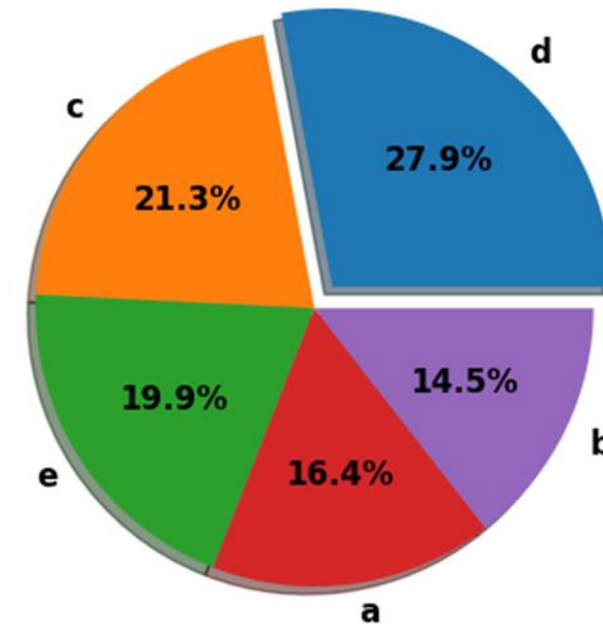
# Analyses univariées



Nombre de produits associés à chaque  
Nutri-score



Répartition des Nutri-scores



=> Les lettres C et D représentent la moitié des produits (49,2%)



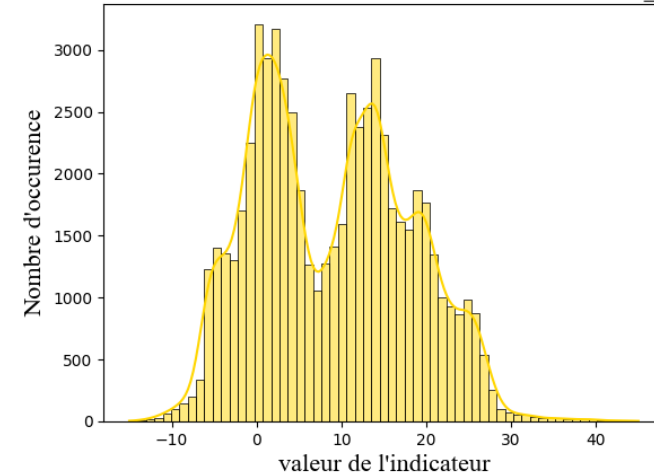
# Analyses univariées



Distribution du score nutritionnel



Distribution de la colonne 'nutrition-score-fr\_100g'



moyenne	médiane	classe modale	ecart-type	CV	skewness	kurtosis
0	8.9	9.0 (-0.2, 0.7): 3210 élément(s)	9.2	1.03	0.22	-0.8

- **Aspect** : distribution bimodale => 2 zones de fortes concentrations (autour de 0 et 15).
- **Tendance centrale** : La moyenne et la médiane sont proches autour du score 9.
- **Dispersion** : les données forment un groupe très hétérogène (Coeff. variation > 100%).
- **Forme** : données étalées à droite (skewness>0) et plus aplaties que la loi normale (kurtosis<0).

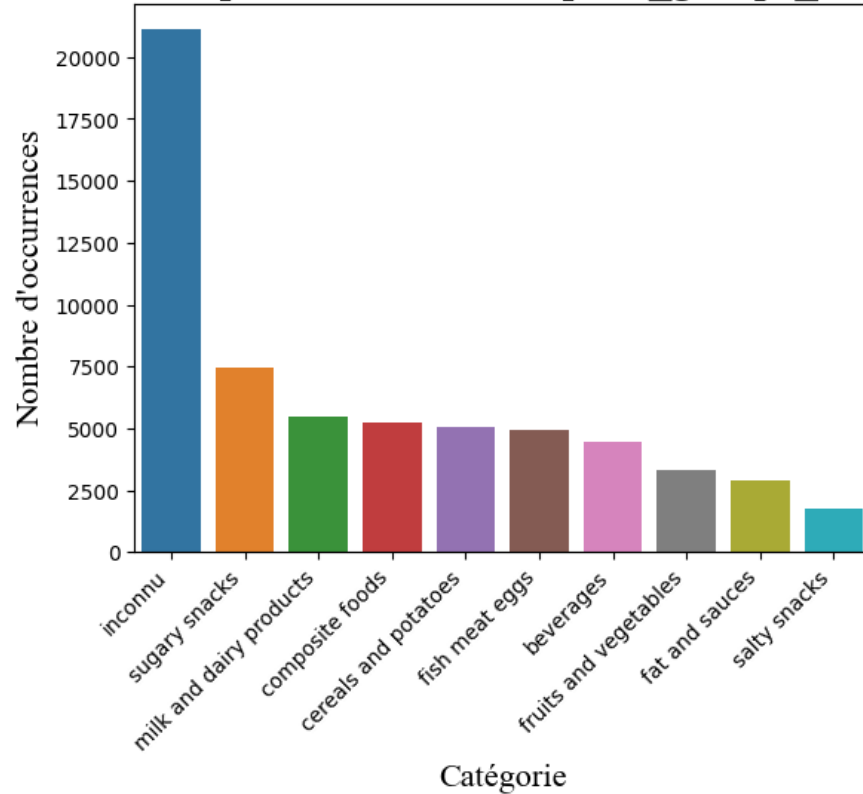


# Analyses univariées

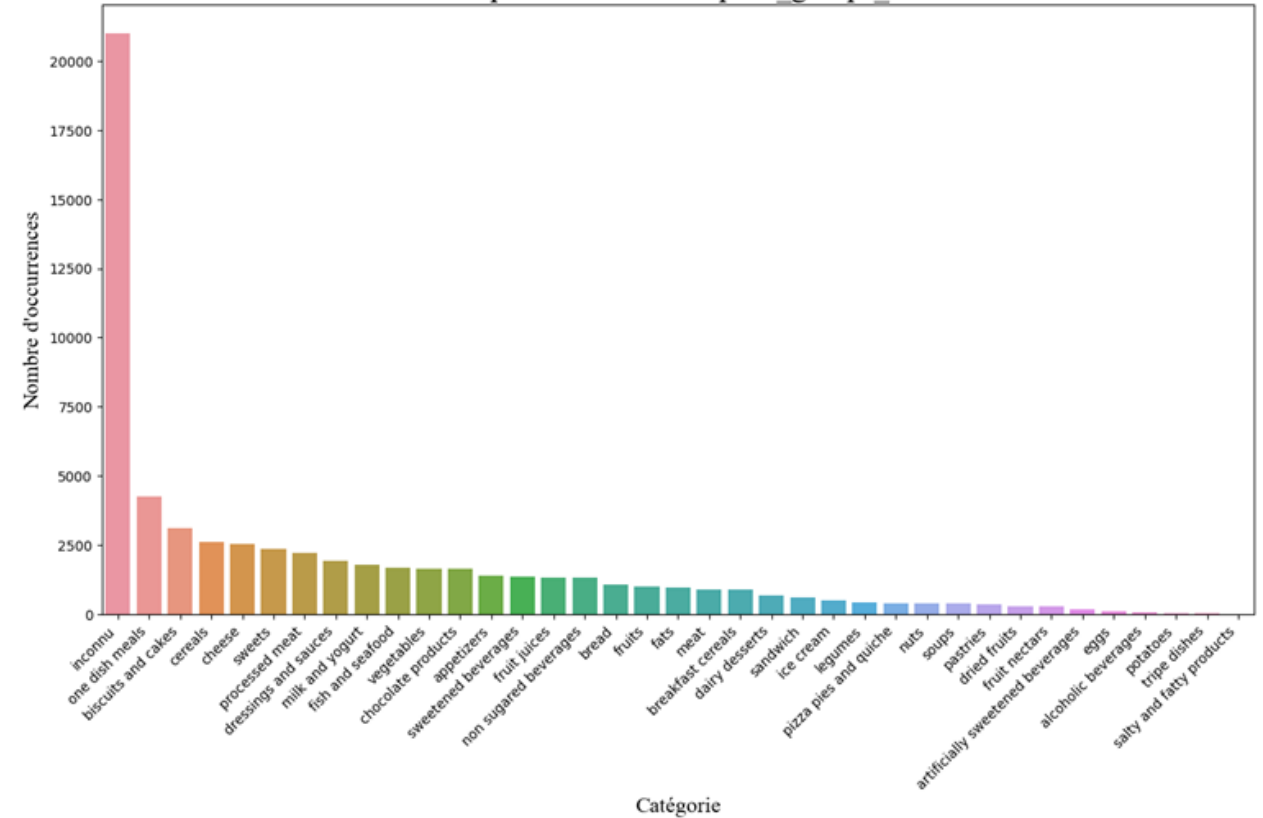


Récurrance des catégories

Barplot de la colonne 'pnns\_groups\_1'



Barplot de la colonne 'pnns\_groups\_2'



=> env. 30% des produits sont de catégorie « inconnue »



# Sommaire



I – Problématique

II – Présentation du jeu de données

III - Nettoyage des données

IV – Analyses univariées

**V – Analyse multivariée**

VI – Faisabilité du système d’auto-complétion

VII – RGPD

VIII - Conclusion



# Analyse multivariée



Analyse en Composantes Principales  
(ACP)



21 composantes expliquent 100% des données

9 composantes > seuil de Kaiser de 4,76



# Analyse multivariée

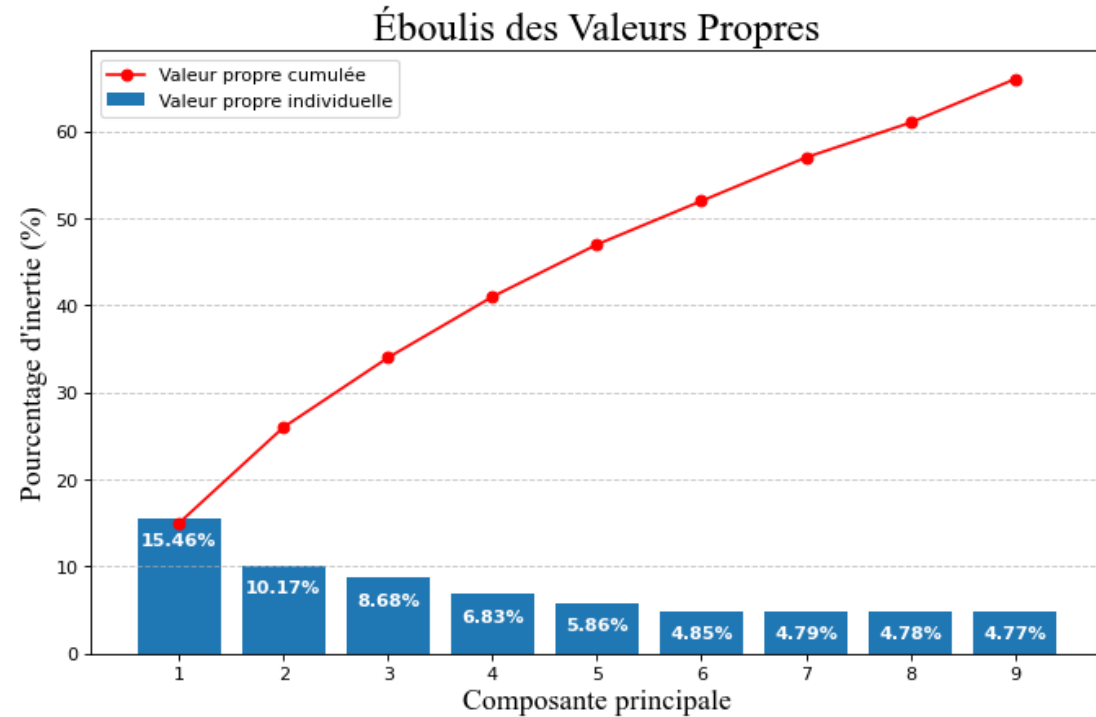


## Analyse en Composantes Principales (ACP)



21 composantes expliquent 100% des données

9 composantes > seuil de Kaiser de 4,76

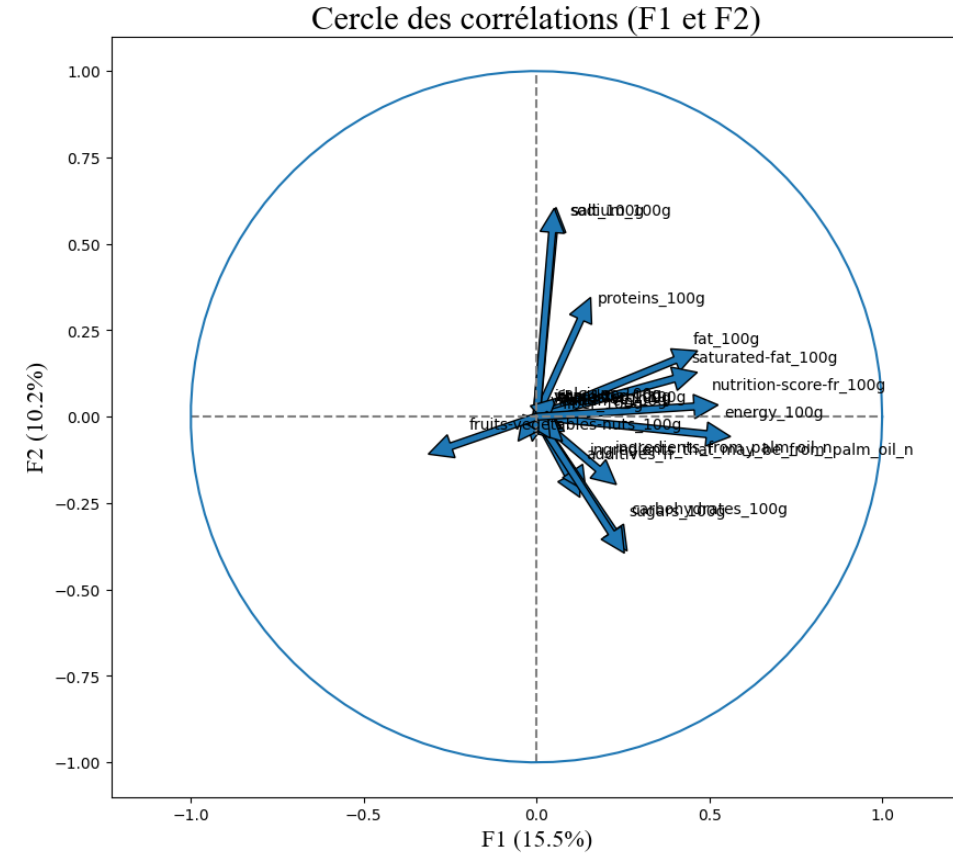
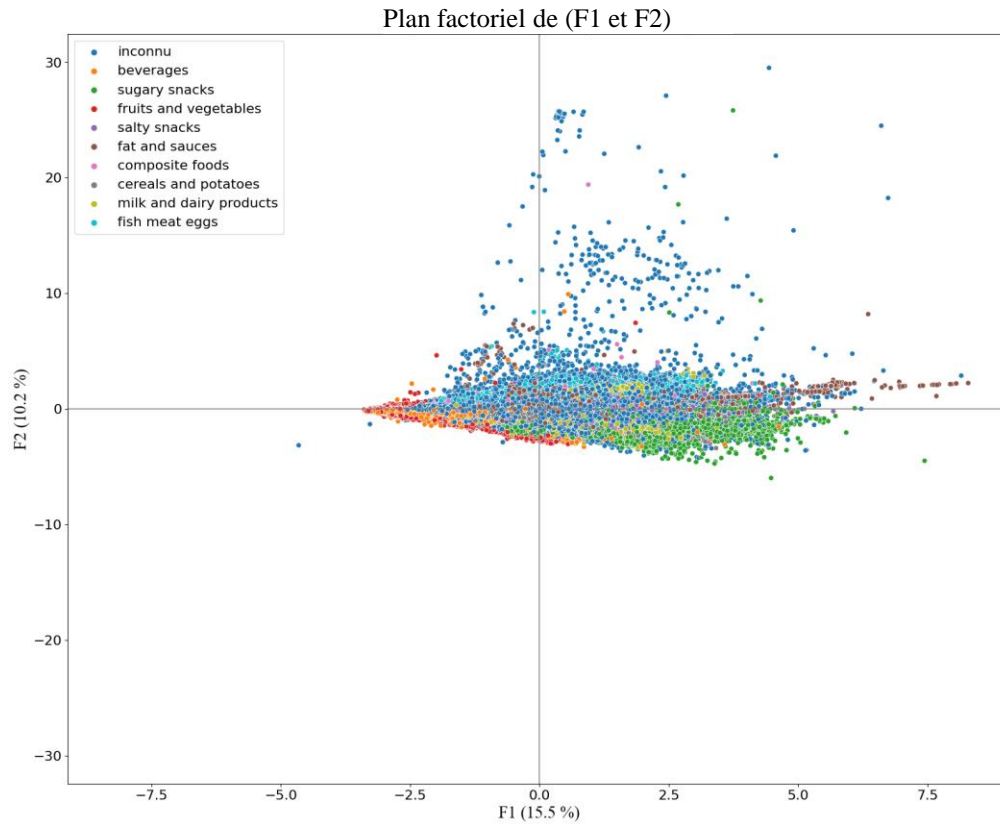


**=> Les 5 premières composantes expliquent env. 50% de l'inertie totale**





# Analyse multivariée

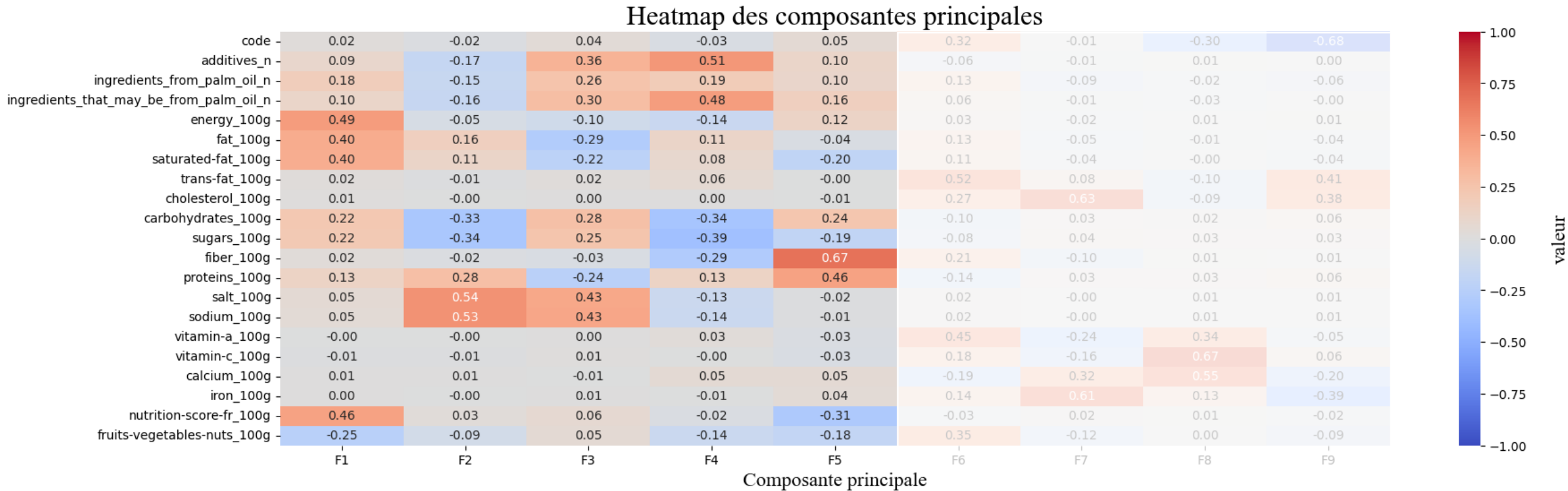




# Analyse multivariée



Nom de la variable

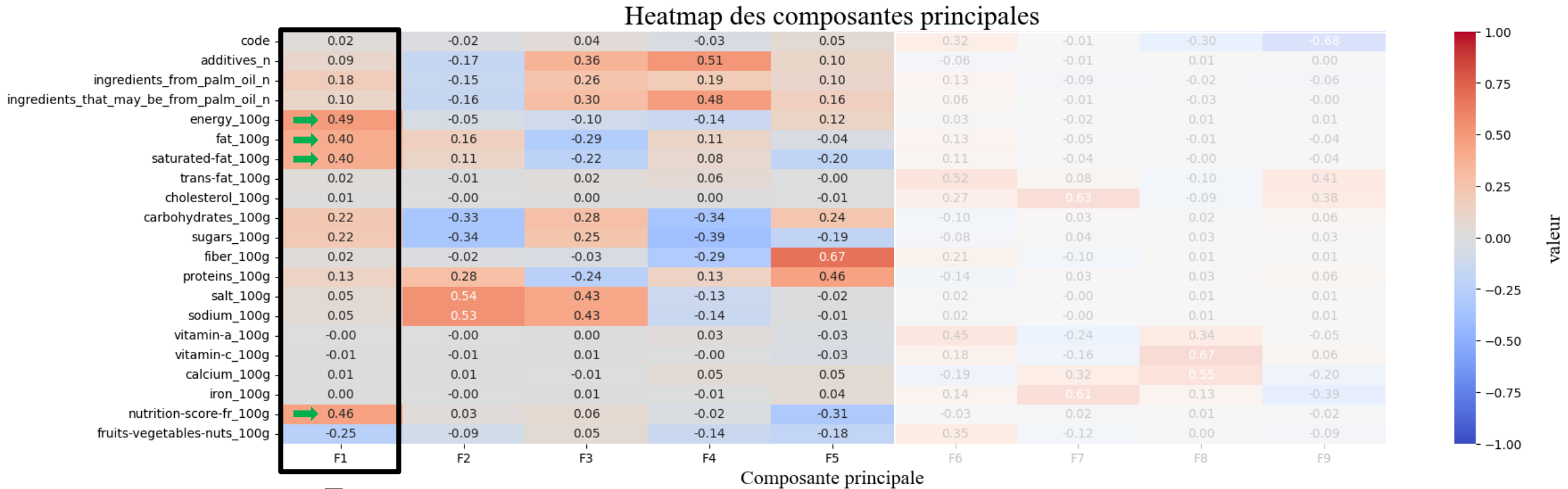




# Analyse multivariée



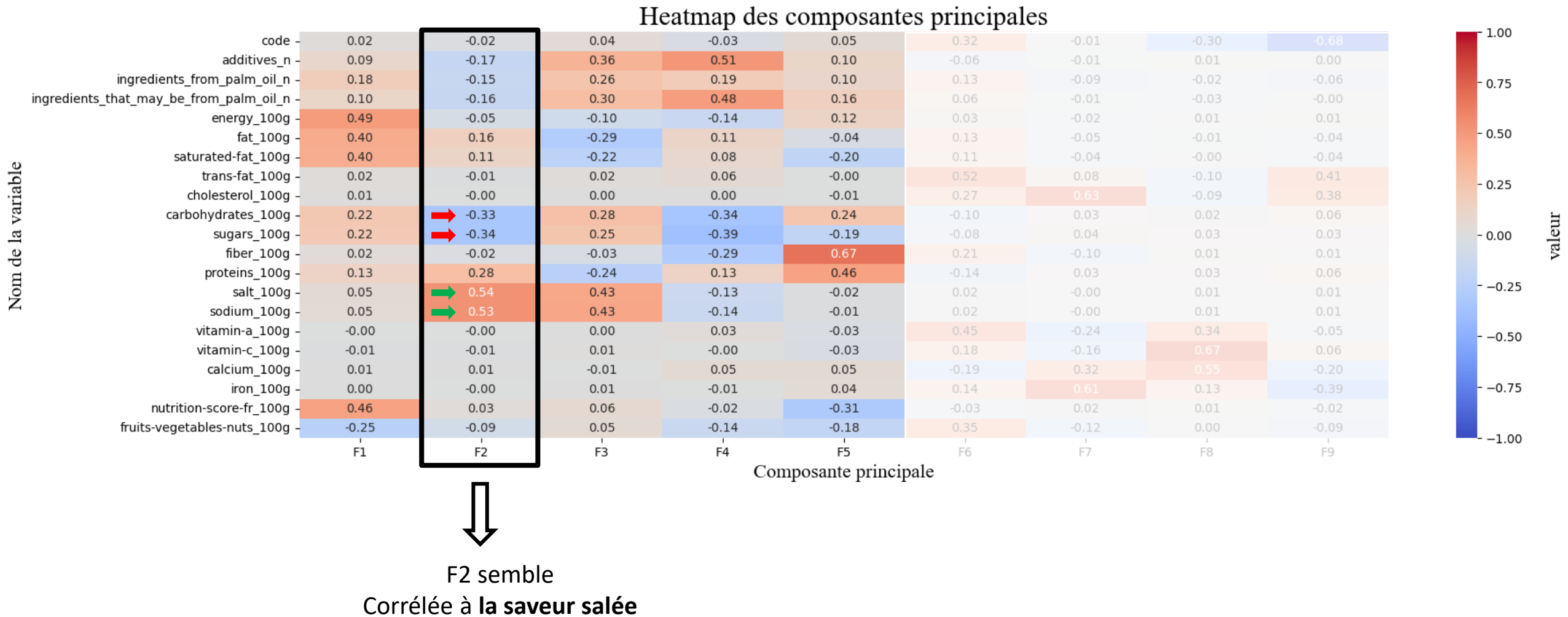
Nom de la variable



F1 semble  
Corrélée à l'apport calorique

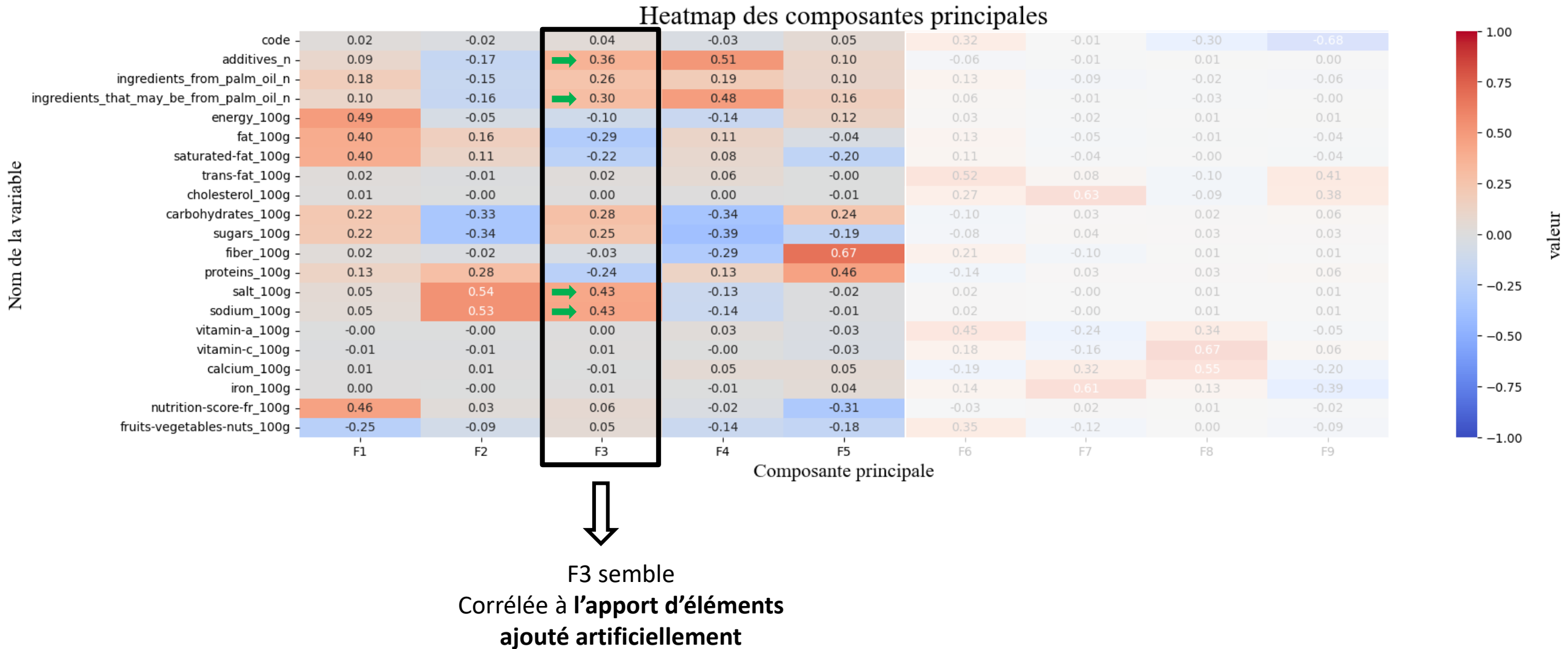


# Analyse multivariée





# Analyse multivariée

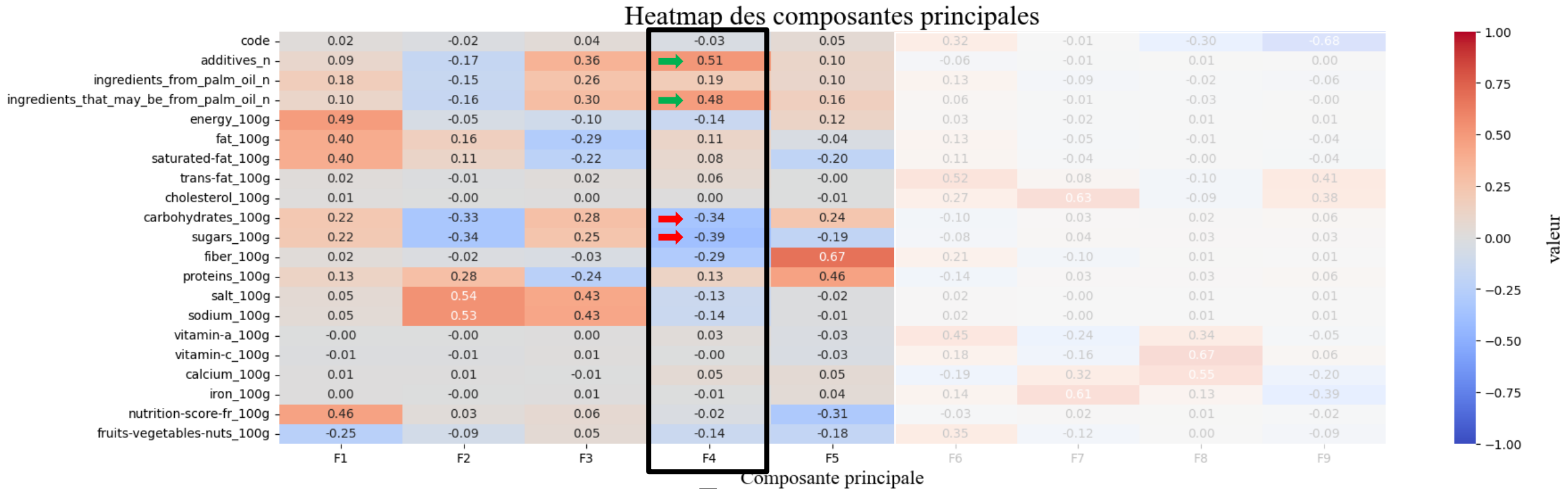




# Analyse multivariée



Nom de la variable



valeur

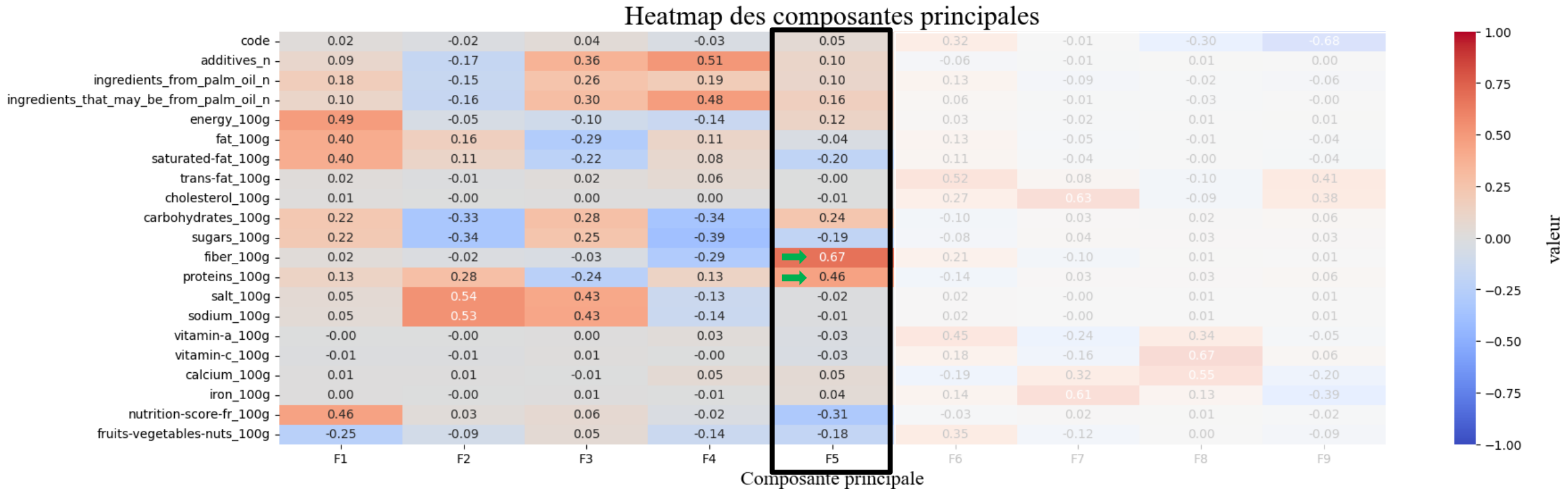
F4 semble  
Corrélée à l'apport d'édulcorants



# Analyse multivariée



Nom de la variable



valeur

F5 semble

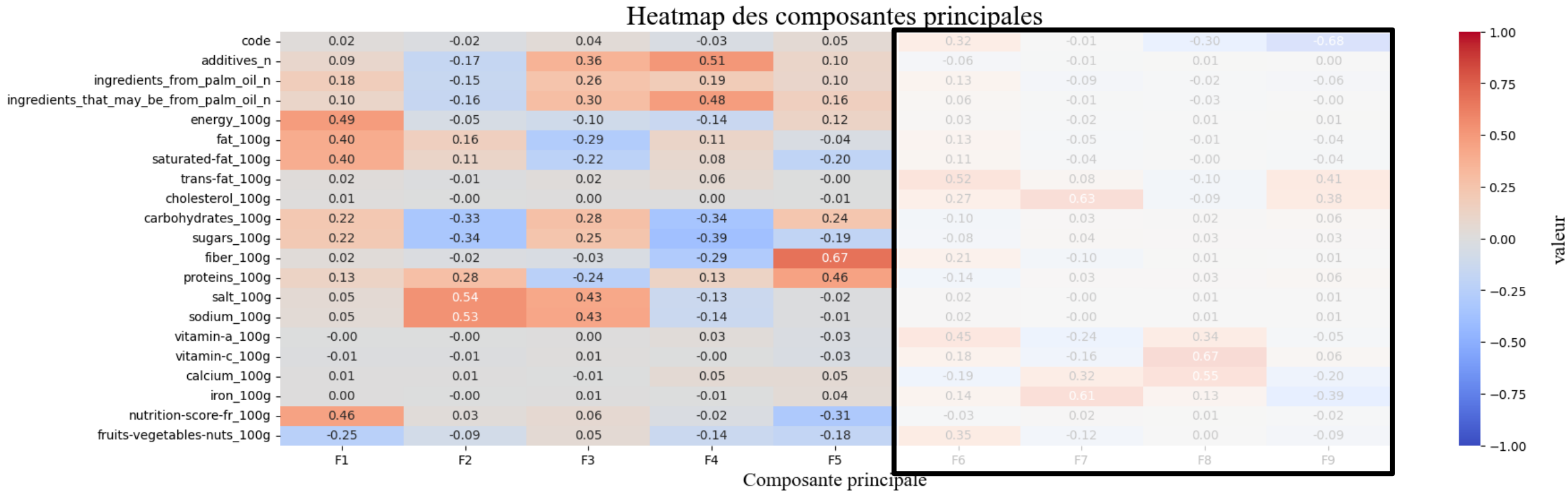
Corrélée à la quantité des nutriments de la composante P (+) du score nutritionnel



# Analyse multivariée



Nom de la variable







# Sommaire



I – Problématique

II – Présentation du jeu de données

III - Nettoyage des données

IV – Analyses univariées

V – Analyse multivariée

**VI – Faisabilité du système d’auto-complétion**

VII – RGPD

VIII - Conclusion



Observations évaluant la pertinence / faisabilité de l'idée  
application du client

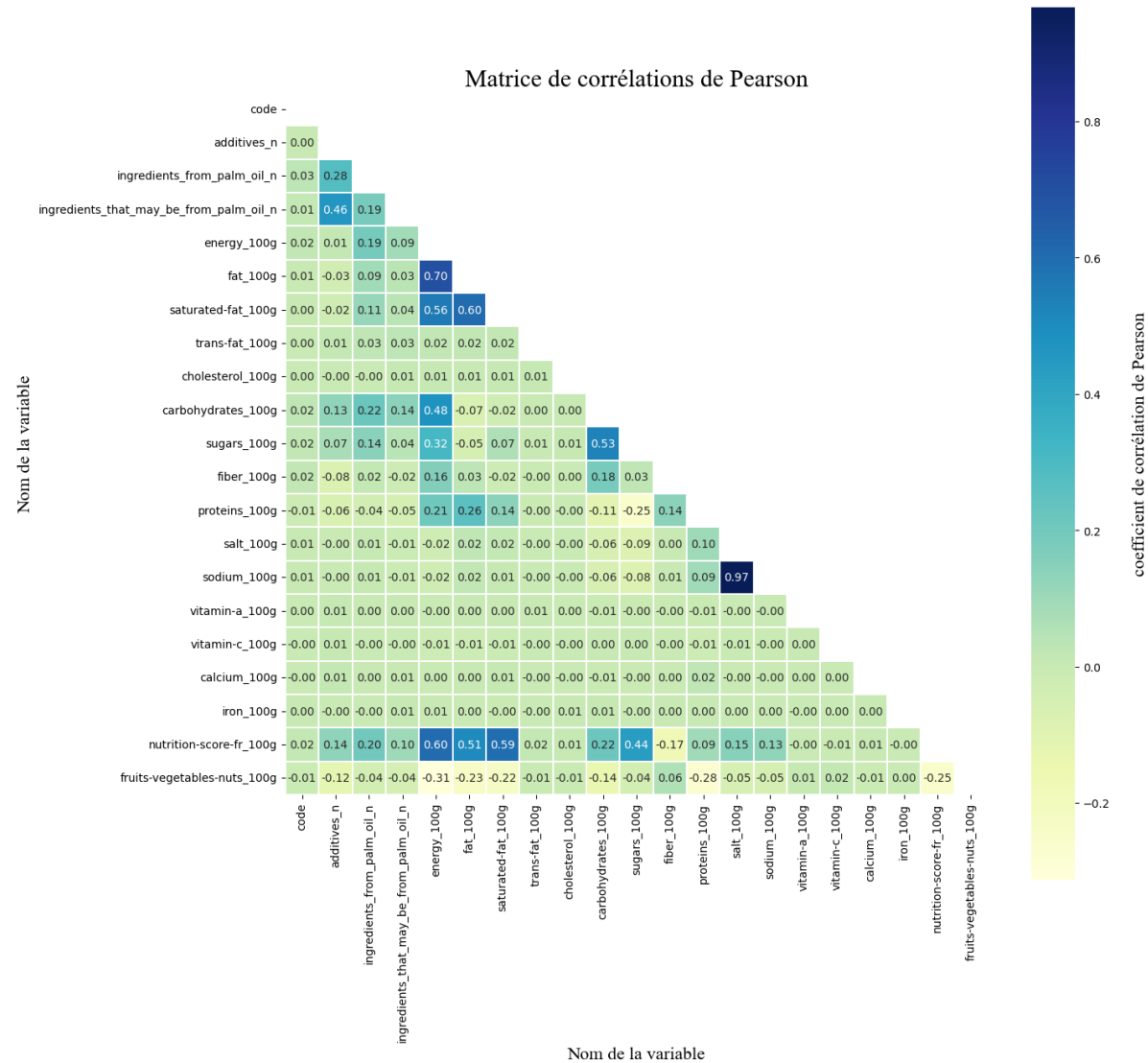
12345  
67890  
+ - × ÷

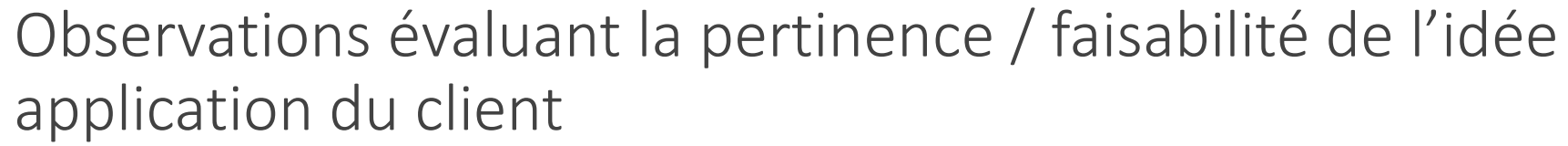




# Observations évaluant la pertinence / faisabilité de l'idée application du client

12345  
67890  
+ - x ÷

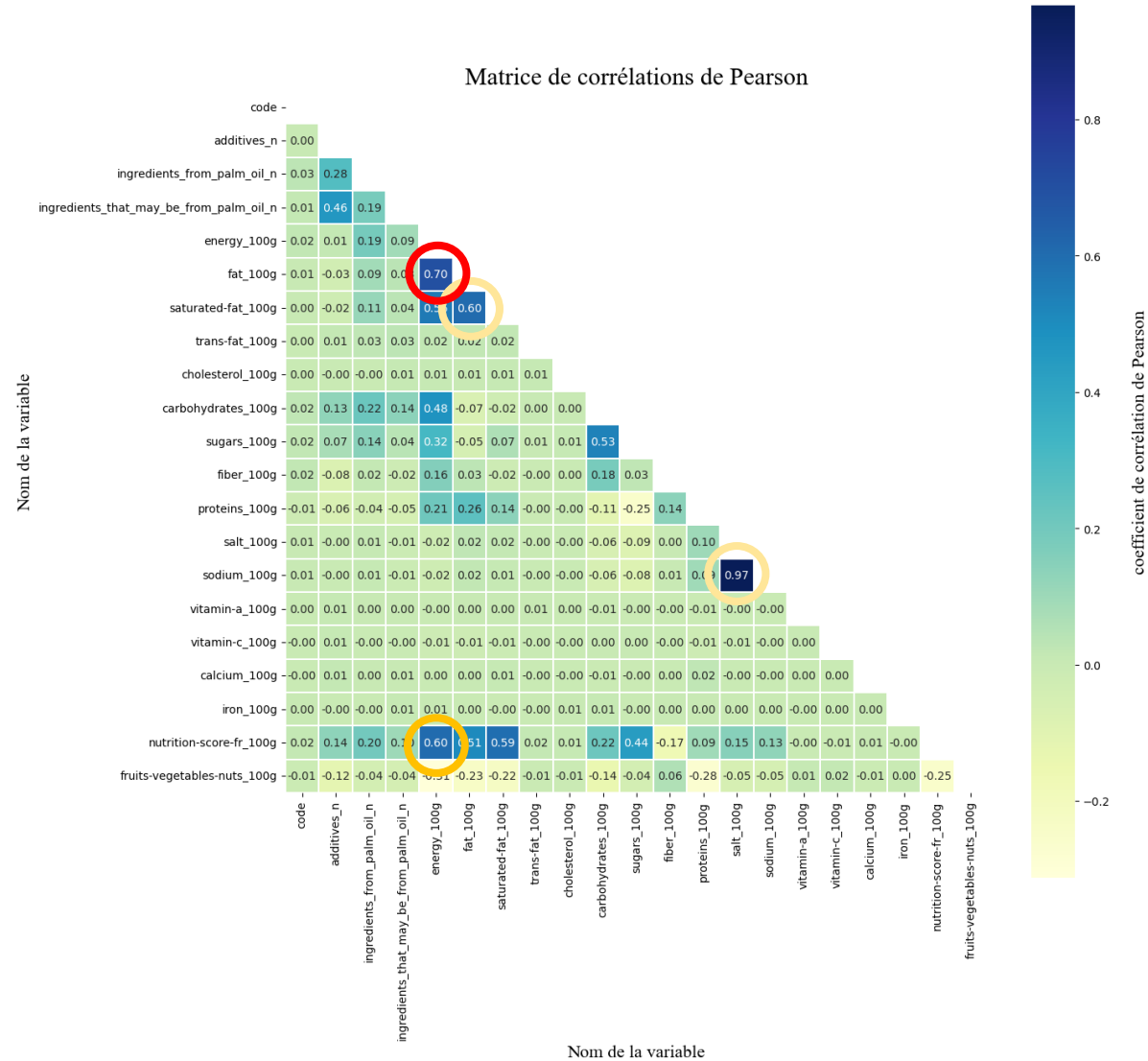






# Observations évaluant la pertinence / faisabilité de l'idée application du client

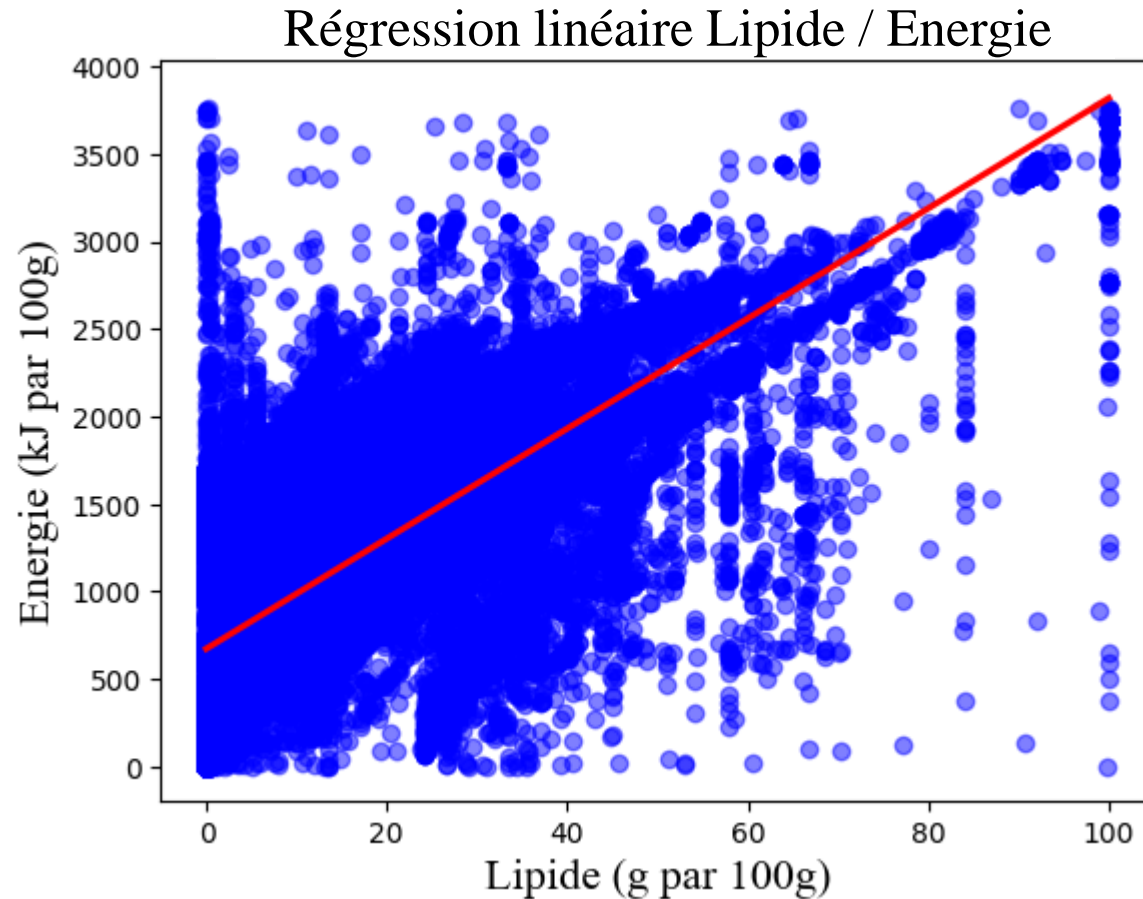
1 2 3 4 5  
6 7 8 9 0  
+ - \* /





# Observations évaluant la pertinence / faisabilité de l'idée application du client

12345  
67890  
+ - × ÷



Hypothèse nulle ( $H_0$ ) : Pas de corrélation entre les 2 variables.  
Hypothèse alternative ( $H_A$ ) : Corrélation entre les 2 variables.

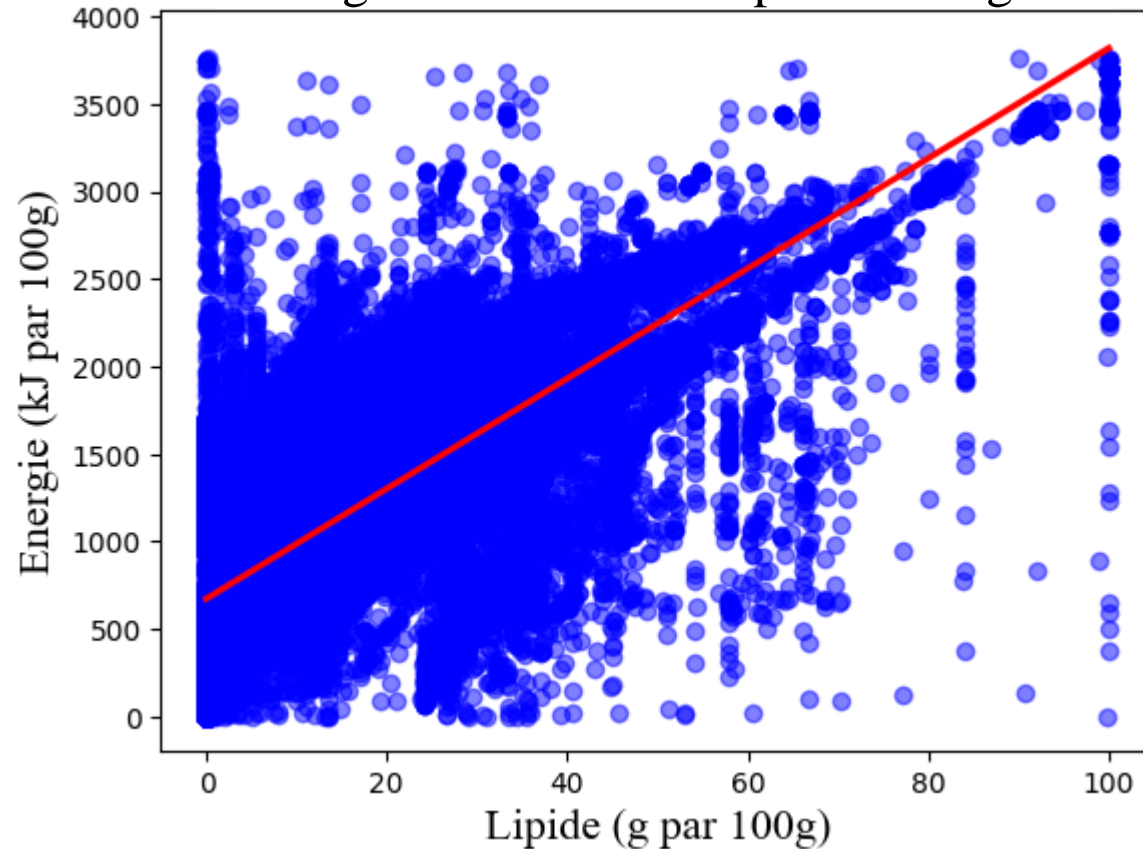


# Observations évaluant la pertinence / faisabilité de l'idée application du client

12345  
67890  
+ - × ÷



Régression linéaire Lipide / Energie



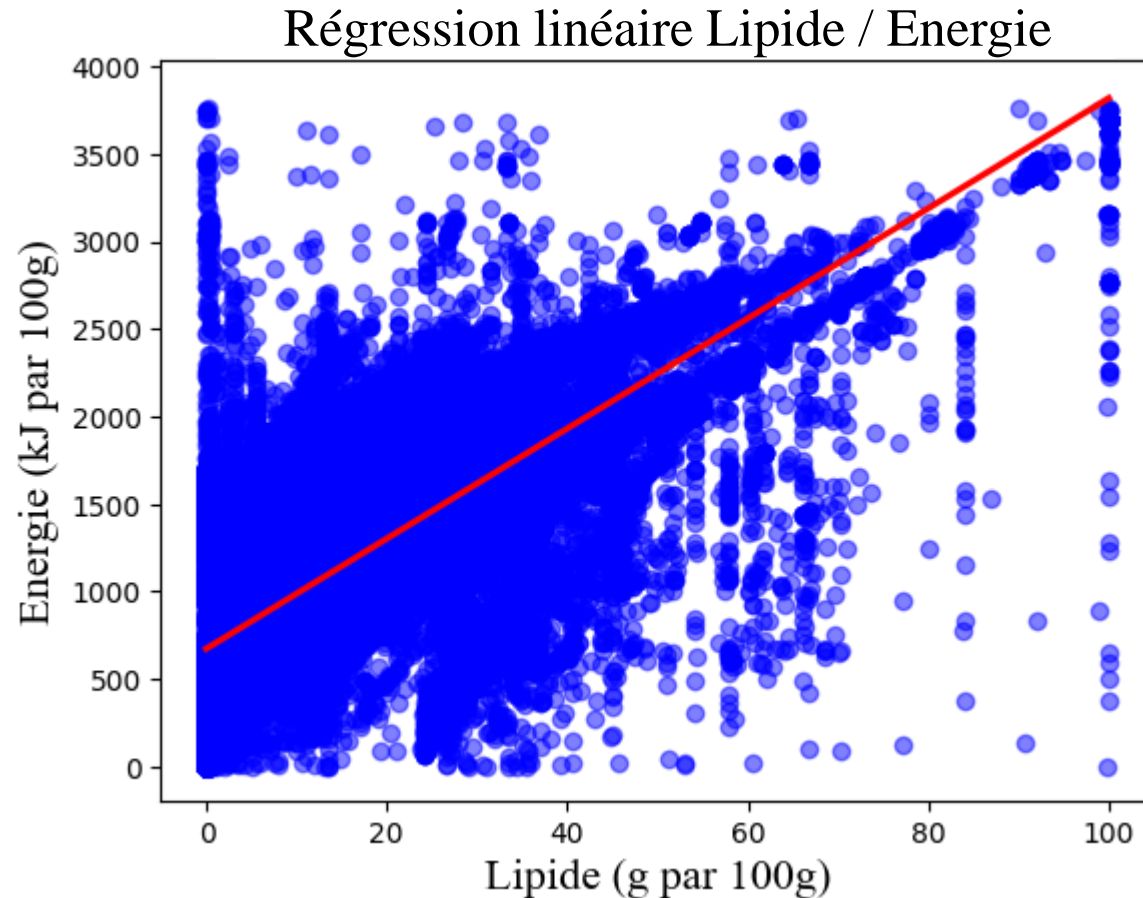
~~Hypothèse nulle (H0) : Pas de corrélation entre les 2 variables.~~  
**Hypothèse alternative (HA) : Corrélation entre les 2 variables.**

**p\_value = 0.00**



# Observations évaluant la pertinence / faisabilité de l'idée application du client

12345  
67890  
+ - \* ÷



Hypothèse nulle ( $H_0$ ) : Pas de corrélation entre les 2 variables.  
Hypothèse alternative ( $H_A$ ) : Corrélation entre les 2 variables.

$p\_value = 0.00$

Energie =  $31.5 * \text{lipide} + 672.9$   
Corrélation Pearson ( $r$ ) = 0.70

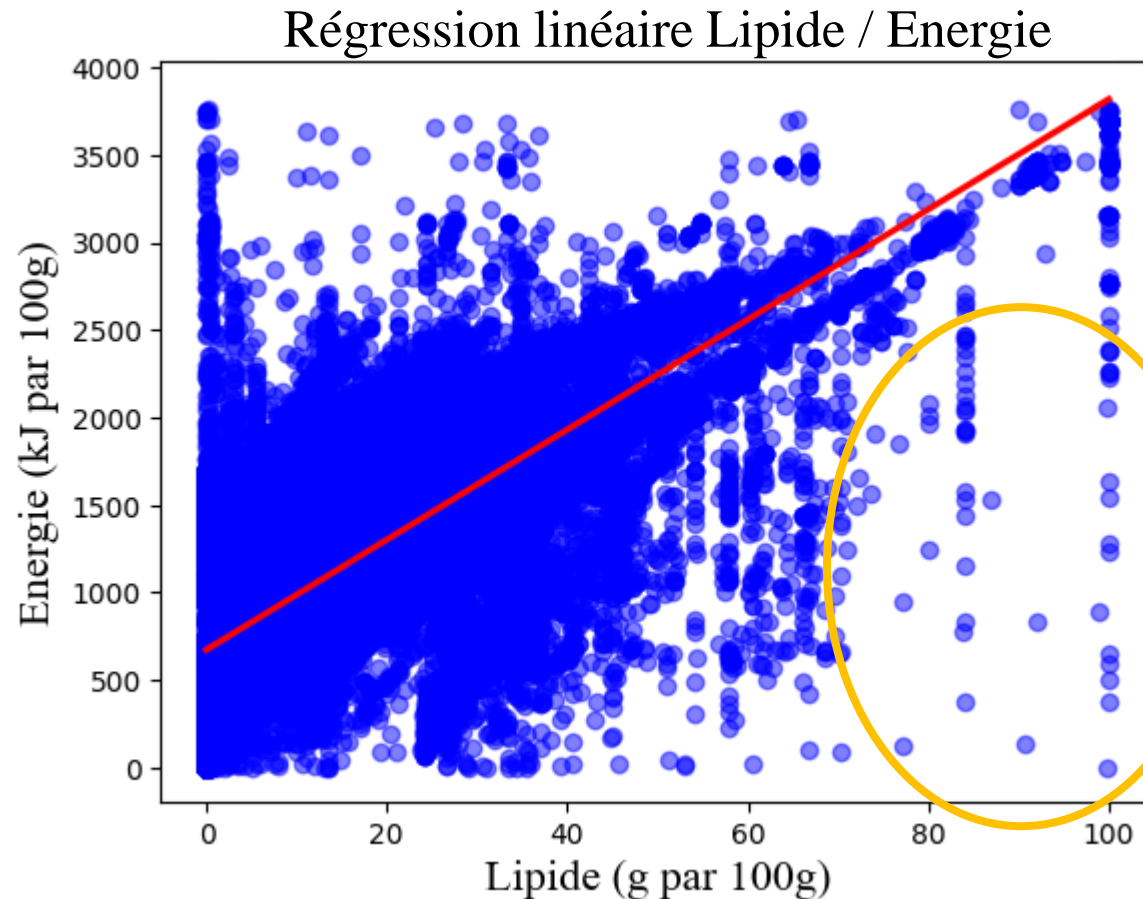
=> Il existe une corrélation positive entre la quantité de lipide et d'énergie ( $p\text{-value} < 0,05$  ;  $r=0,70$ )





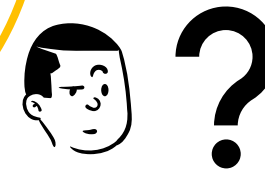
# Observations évaluant la pertinence / faisabilité de l'idée application du client

12345  
67890  
+ - \* ÷



Hypothèse nulle ( $H_0$ ) : Pas de corrélation entre les 2 variables.  
Hypothèse alternative ( $H_A$ ) : Corrélation entre les 2 variables.

Energie =  $31.5 * \text{lipide} + 672.9$   
Corrélation Pearson ( $r$ ) = 0.70  
 $p\_value = 0.00$



=> Il existe une corrélation positive entre la quantité de lipide et d'énergie ( $p\text{-value} < 0,05$  ;  $r=0,70$ )



# Observations évaluant la pertinence / faisabilité de l'idée application du client

ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ





# Observations évaluant la pertinence / faisabilité de l'idée application du client

ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ



$\chi^2$

Hypothèse nulle (H0) : Pas de corrélation entre les 2 variables.

Hypothèse alternative (HA) : Corrélation entre les 2 variables.



# Observations évaluant la pertinence / faisabilité de l'idée application du client

ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ



$\chi^2$

Hypothèse nulle (H0) : Pas de corrélation entre les 2 variables.

Hypothèse alternative (HA) : Corrélation entre les 2 variables.

Nutri-score / catégorie pnns 2

Test du chi-deux :

$\text{Chi}^2 = 40894.46$

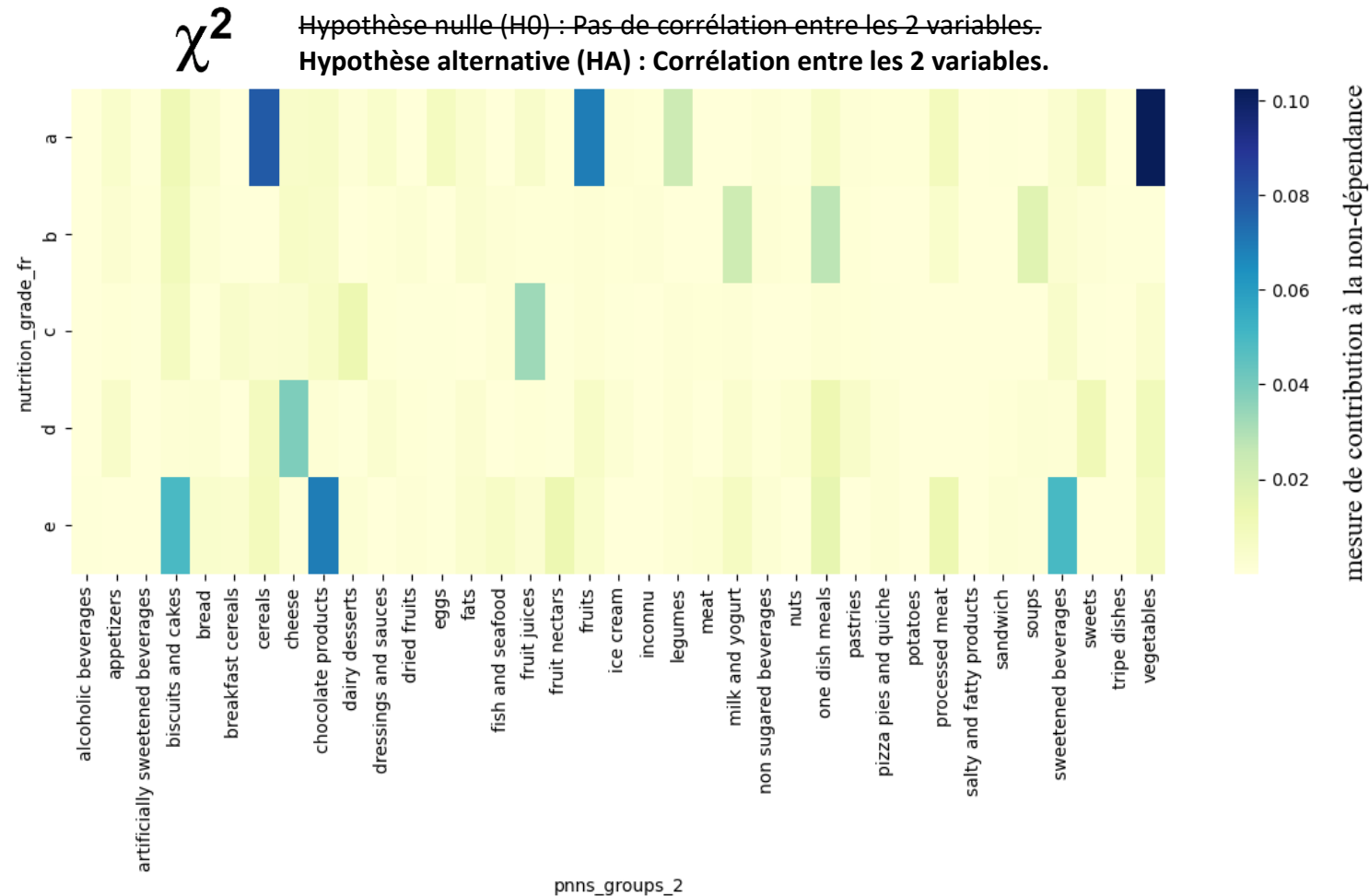
$p\text{-value} = 0.00$

**=> Il existe une corrélation entre le Nutri-score et la catégorie pnns 2 ( $p\text{-value} < 0,05$ )**



# Observations évaluant la pertinence / faisabilité de l'idée application du client

ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ

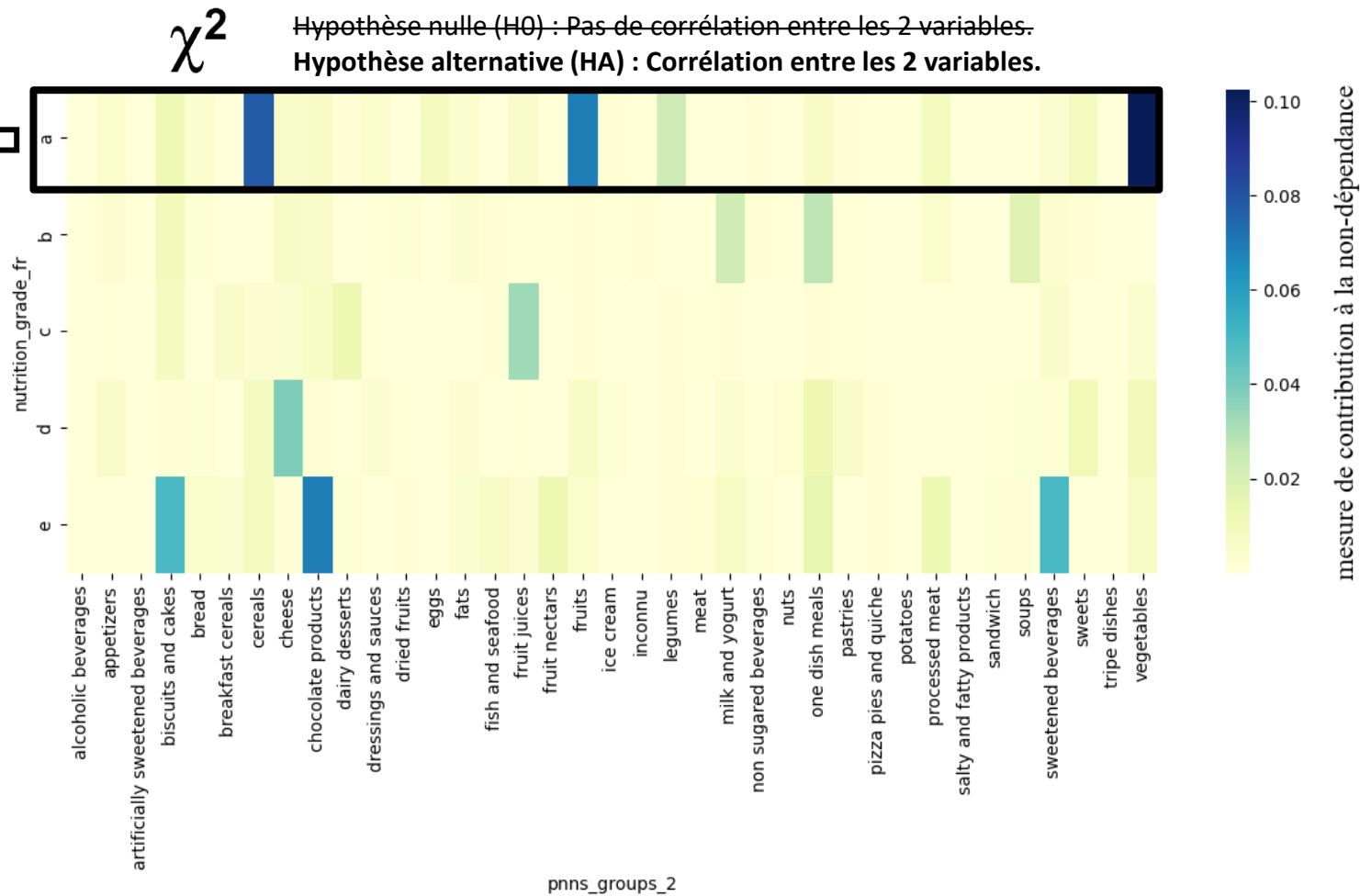


=> Il existe une corrélation entre le Nutri-score et la catégorie pnn2 ( $p\text{-value} < 0,05$ )



# Observations évaluant la pertinence / faisabilité de l'idée application du client

ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ



=> Il existe une corrélation entre le Nutri-score et la catégorie pnn2 ( $p\text{-value} < 0,05$ )



# Observations évaluant la pertinence / faisabilité de l'idée application du client

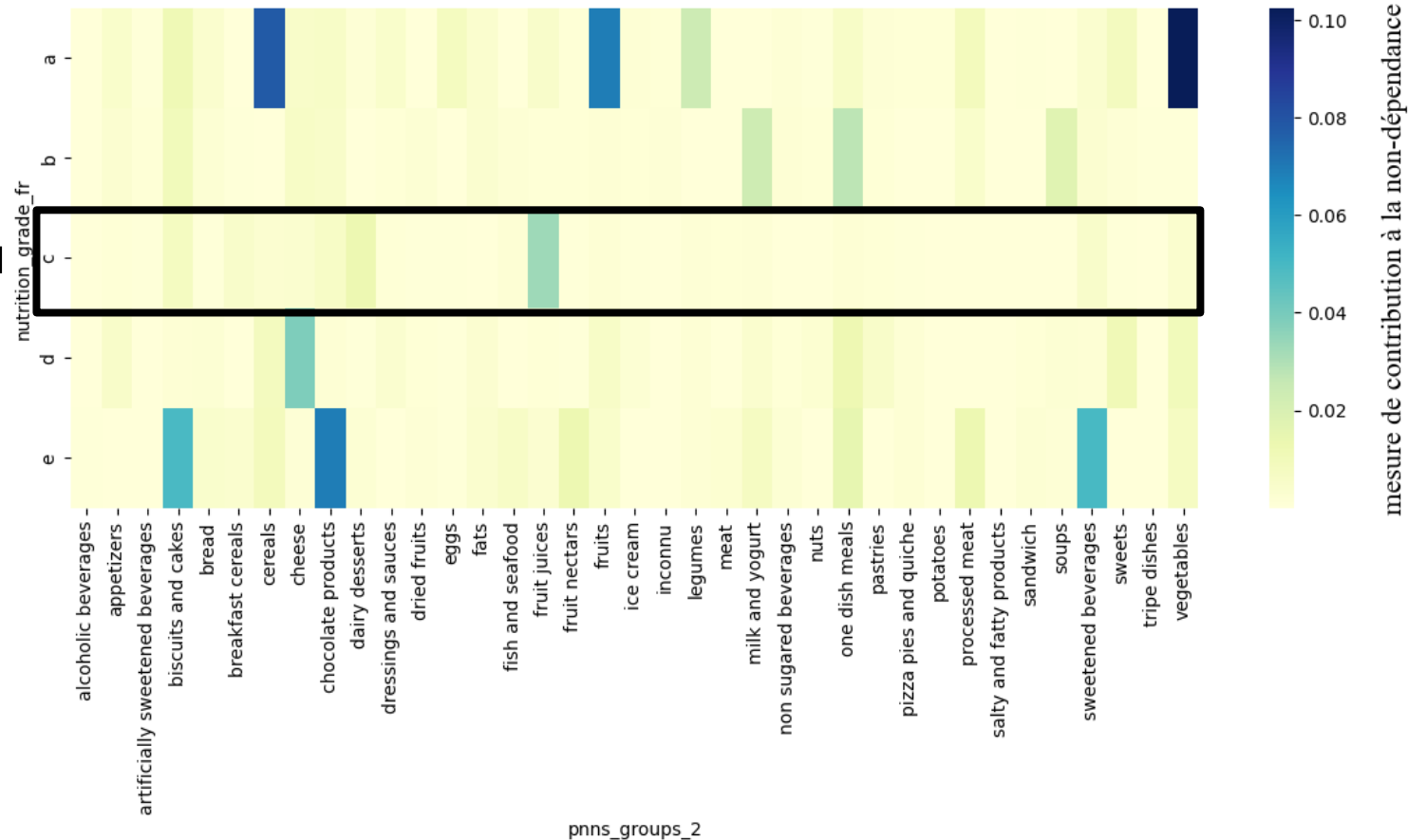
ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ



$\chi^2$

Hypothèse nulle (H0) : Pas de corrélation entre les 2 variables.

Hypothèse alternative (HA) : Corrélation entre les 2 variables.



=> Il existe une corrélation entre le Nutri-score et la catégorie pnnns 2 ( $p\text{-value} < 0,05$ )



# Observations évaluant la pertinence / faisabilité de l'idée application du client

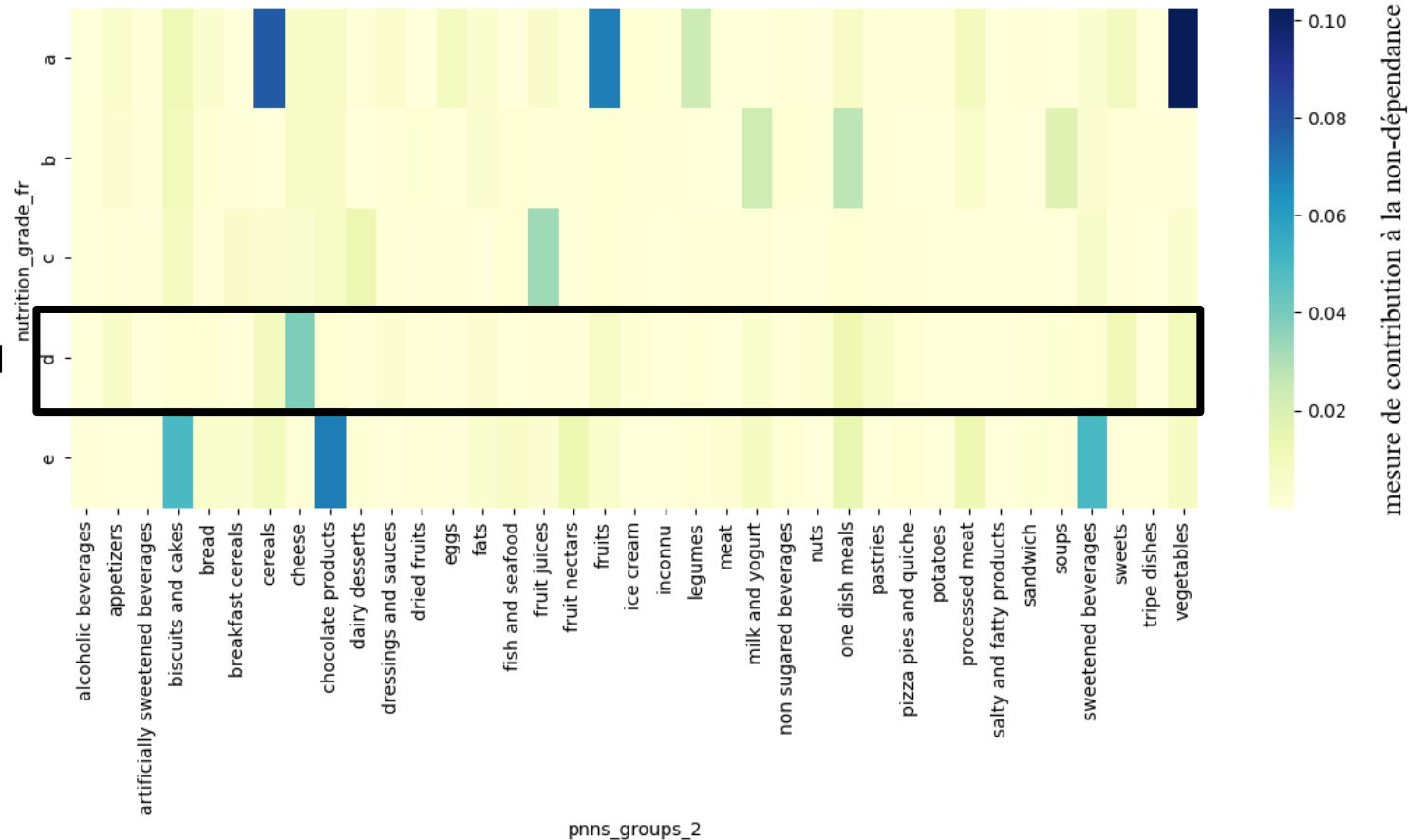
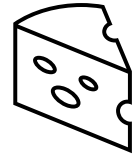
ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ



$\chi^2$

Hypothèse nulle (H0) : Pas de corrélation entre les 2 variables.

Hypothèse alternative (HA) : Corrélation entre les 2 variables.



=> Il existe une corrélation entre le Nutri-score et la catégorie pnnns 2 ( $p\text{-value} < 0,05$ )





# Observations évaluant la pertinence / faisabilité de l'idée application du client

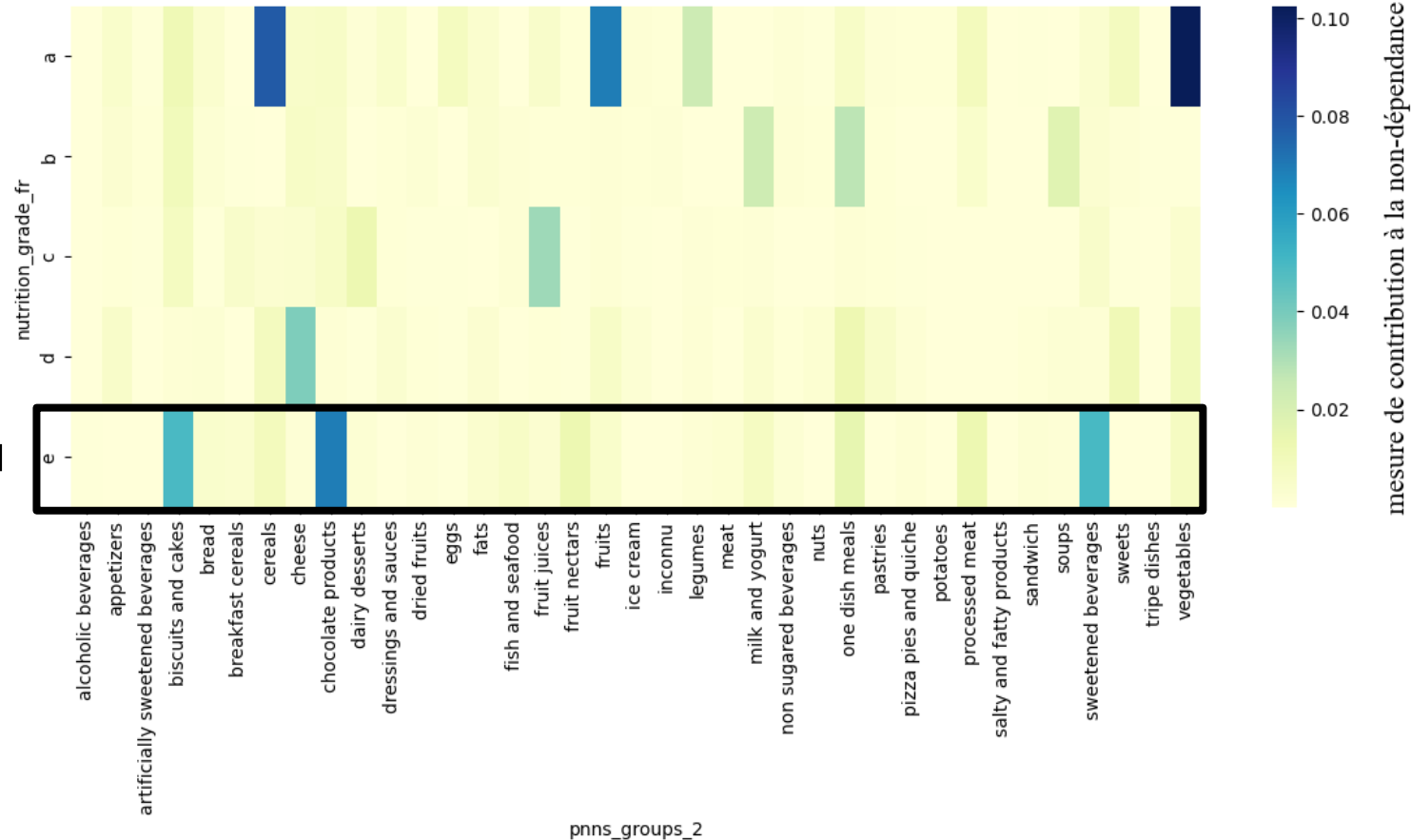
ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ



$\chi^2$

Hypothèse nulle (H0) : Pas de corrélation entre les 2 variables.

Hypothèse alternative (HA) : Corrélation entre les 2 variables.



=> Il existe une corrélation entre le Nutri-score et la catégorie pnns 2 ( $p\text{-value} < 0,05$ )



# Observations évaluant la pertinence / faisabilité de l'idée application du client

ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ



12345  
67890  
+ - \* ÷





# Observations évaluant la pertinence / faisabilité de l'idée application du client

ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ



12345  
67890  
+ - \* ÷



12345  
67890  
+ - \* ÷

Test de normalité des données => **Test de Shapiro-Wilk**



Hypothèse nulle (H0) : la population est distribuée normalement selon la loi gaussienne.

Hypothèse alternative (HA) : la population n'est pas distribuée normalement selon la loi gaussienne.



# Observations évaluant la pertinence / faisabilité de l'idée application du client



**12345 67890 +-\*/**    Test de normalité des données => **Test de Shapiro-Wilk**

↓

~~Hypothèse nulle (H0) : la population est distribuée normalement selon la loi gaussienne.~~  
**Hypothèse alternative (HA) : la population n'est pas distribuée normalement selon la loi gaussienne.**

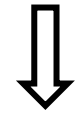
	Column	SW-statistic	P-value
0	code	0.12436	0.0
1	additives_n	0.64575	0.0
2	ingredients_from_palm_oil_n	0.22780	0.0
3	ingredients_that_may_be_from_palm_oil_n	0.30616	0.0
4	energy_100g	0.94936	0.0
5	fat_100g	0.77444	0.0
6	saturated-fat_100g	0.64348	0.0
7	trans-fat_100g	0.00292	0.0
8	cholesterol_100g	0.00055	0.0
9	carbohydrates_100g	0.86016	0.0
10	sugars_100g	0.71868	0.0
11	fiber_100g	0.44373	0.0
12	proteins_100g	0.81852	0.0
13	salt_100g	0.18533	0.0
14	sodium_100g	0.17332	0.0
15	vitamin-a_100g	0.00113	0.0
16	vitamin-c_100g	0.00447	0.0
17	calcium_100g	0.00815	0.0
18	iron_100g	0.00087	0.0
19	nutrition-score-fr_100g	0.97473	0.0
20	fruits-vegetables-nuts_100g	0.73058	0.0



# Observations évaluant la pertinence / faisabilité de l'idée application du client



**1 2 3 4 5**  
**6 7 8 9 0**  
**+ - x ÷**    Test de normalité des données => **Test de Shapiro-Wilk**



~~Hypothèse nulle (H0) : la population est distribuée normalement selon la loi gaussienne.~~  
**Hypothèse alternative (HA) : la population n'est pas distribuée normalement selon la loi gaussienne.**



**Test non-paramétrique de Kruskal-Wallis**  
**kw**

	Column	SW-statistic	P-value
0	code	0.12436	0.0
1	additives_n	0.64575	0.0
2	ingredients_from_palm_oil_n	0.22780	0.0
3	ingredients_that_may_be_from_palm_oil_n	0.30616	0.0
4	energy_100g	0.94936	0.0
5	fat_100g	0.77444	0.0
6	saturated-fat_100g	0.64348	0.0
7	trans-fat_100g	0.00292	0.0
8	cholesterol_100g	0.00055	0.0
9	carbohydrates_100g	0.86016	0.0
10	sugars_100g	0.71868	0.0
11	fiber_100g	0.44373	0.0
12	proteins_100g	0.81852	0.0
13	salt_100g	0.18533	0.0
14	sodium_100g	0.17332	0.0
15	vitamin-a_100g	0.00113	0.0
16	vitamin-c_100g	0.00447	0.0
17	calcium_100g	0.00815	0.0
18	iron_100g	0.00087	0.0
19	nutrition-score-fr_100g	0.97473	0.0
20	fruits-vegetables-nuts_100g	0.73058	0.0



# Observations évaluant la pertinence / faisabilité de l'idée application du client



**kw**

Hypothèse nulle (H0) : la valeur de la variable qualitative n'influence pas les valeurs de la variable quantitative.  
Hypothèse alternative (HA) : la valeur de la variable qualitative influence les valeurs de la variable quantitative.

ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ



12345  
67890  
+ - \* ÷



# Observations évaluant la pertinence / faisabilité de l'idée application du client



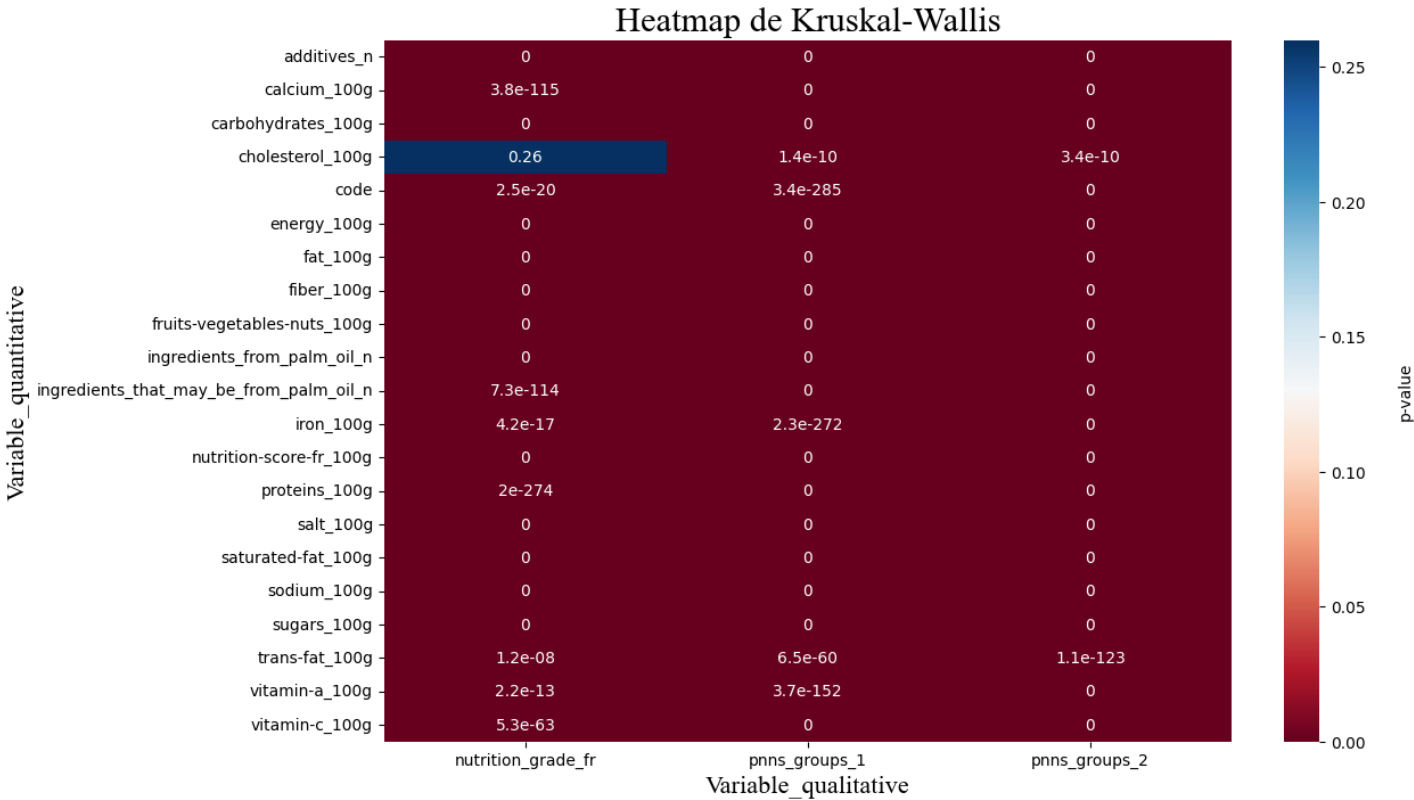
ABCDEF  
GHIKLM  
NOPQRS  
TVXYZ

VS

12345  
67890  
+-x÷

kw

Hypothèse nulle (H0) : la valeur de la variable qualitative n'influence pas les valeurs de la variable quantitative.  
Hypothèse alternative (HA) : la valeur de la variable qualitative influence les valeurs de la variable quantitative.



=> Les valeurs des variables quantitatives sont influencées par les variables qualitatives relatives au Nutri-score et aux pnns des catégories d'aliment, à l'exception de la quantité de cholestérol qui n'est pas influencé par le Nutri-score.



# Sommaire



I – Problématique

II – Présentation du jeu de données

III - Nettoyage des données

IV – Analyses univariées

V – Analyse multivariée

VI – Faisabilité du système d’auto-complétion

**VII – RGPD**

VIII - Conclusion





Voici les 5 grand principe RGPD (Règlement sur la Protection des Données Personnelles) :

1. Finalité
2. Pertinence
3. Durée limitée de conservation
4. Sécurité
5. Droits des personnes





# Sommaire



I – Problématique

II – Présentation du jeu de données

III - Nettoyage des données

IV – Analyses univariées

V – Analyse multivariée

VI – Faisabilité du système d’auto-complétion

VII – RGPD

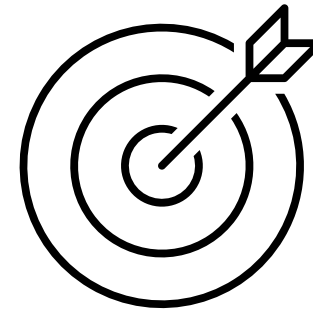
**VIII - Conclusion**



# Conclusion



- Nettoyer le jeu de données ✓
- Analyse univariée des variables importantes ✓
- Analyse multivariée et résultats statistiques ✓
- Conclusion sur la faisabilité du système d'auto-complétion ✓
- Rappel RGPD ✓





# Conclusion



## **Faisabilité système d'auto-complétion:**

- Corrélation données quantitatives
- Corrélation données qualitatives
- Influence des données qualitatives sur les données quantitatives



## **Limites :**

- Données obsolètes (pas mis à jour depuis 2017)
- 30% des produits sont de catégorie « inconnu »
- Application de règle métier pour détecter des variables aberrantes (lipide / énergie)
- Mise à jour des Nutri-score (méthode de calcul 2023)



**La création d'un système d'auto-complétion fonctionnel pour notre client "Santé Publique France" n'est pas encore possible avec ce jeu de données. Plusieurs travaux d'amélioration de sa qualité sont à prévoir encore en amont.**

OPENCLASSROOMS

Merci pour votre attention



CentraleSupélec

Pierrick BERTHE

Formation Expert en Data Science  
*Openclassrooms – CentraleSupélec*

*août 2023 → avril 2024*