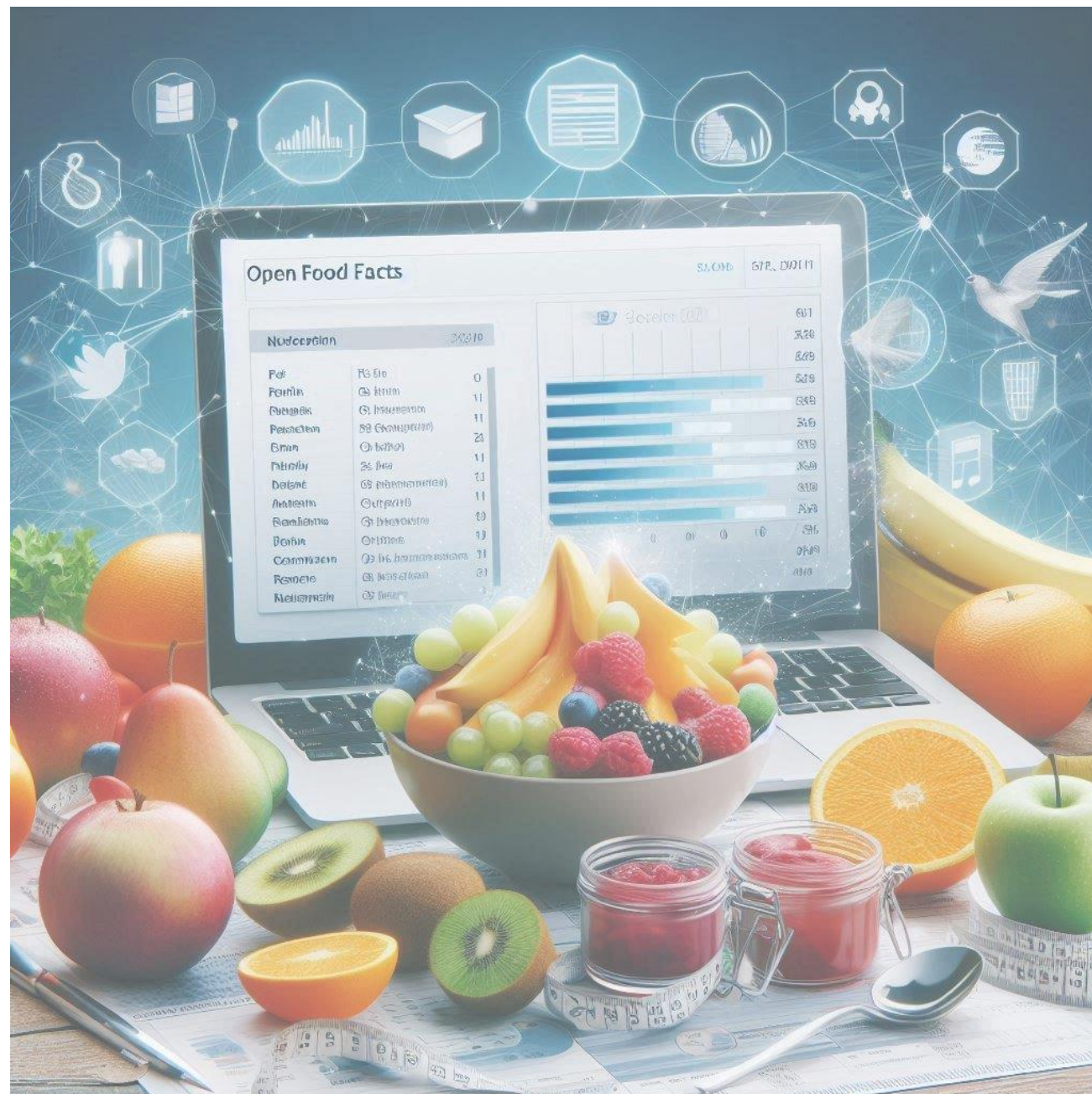


AI Engineer – P03
OpenClassrooms

**Préparez des
données pour un
organisme de santé
publique**



Introduction

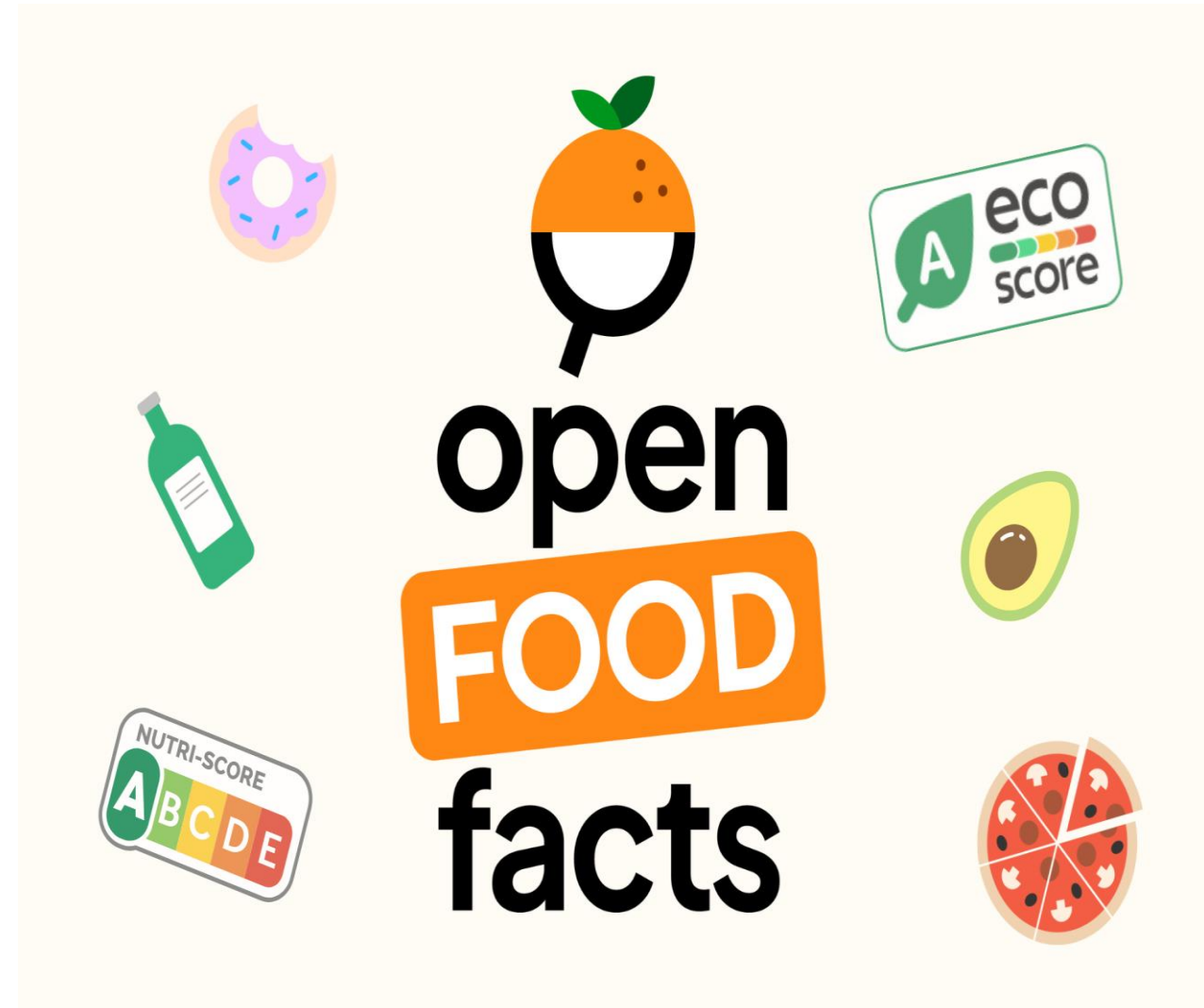
CONTEXTE DU PROJET

Santé publique France souhaite améliorer sa base de données Open Food Facts. Cette base de données open source est mise à la disposition de particuliers et d'organisations afin de leur permettre de connaître la qualité nutritionnelle de produits.

Aujourd'hui, pour ajouter un produit à la base de données d'Open Food Facts, il est nécessaire de remplir de nombreux champs textuels et numériques, ce qui peut conduire à des erreurs de saisie et à des valeurs manquantes dans la base.

L'agence Santé publique France nous confie la création d'un système de suggestion ou d'auto-complétion pour aider les usagers à remplir plus efficacement la base de données.

Nous sommes missionnés sur le projet de nettoyage et exploration des données, afin de déterminer la faisabilité de cette idée d'application.





Aujourd'hui



Demain

Sommaire

1^{er} Notebook => Nettoyage

Analyse exploratoire des données

Réduction du dataset & sélection de la cible

Gestion des valeurs :

- Les Outliers (Avant traitement)
- Les Outliers (Après traitement)
- Imputation avec KNN

2^{ème} Notebook => Analyse

Analyses univariées

Analyses bivariées

Analyses multivariées

Test Statistiques

ACP

Cercles de corrélation

Conclusion



1^{er} Notebook

Le Nettoyage

```
nb_lignes, nb_colonnes = data.shape
print(f"Nombre de lignes : {nb_lignes}")
print(f"Nombre de colonnes : {nb_colonnes}")
```

```
Nombre de lignes : 320772
Nombre de colonnes : 162
```

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 320772 entries, 0 to 320771
Columns: 162 entries, code to water-hardness_100g
dtypes: float64(106), object(56)
memory usage: 396.5+ MB
```

```
data.head()
```

Python

	code	url	creator	created_t	created_datetime	last_modified_t	last_modified_datetime	product_name	generic_name	quantity	...	ph_100g	fruits-vegetables-nuts_100g	collagen-meat-protein-ratio_100g	cocoa_100g	chlorophyl_100g	carbon-footprint_100g	nutrition-score-fr_100g	nutrition-score-uk_100g	glycemic-index_100g	hardne
0	0000000003087	http://world-fr.openfoodfacts.org/produit/0000...	openfoodfacts-contributors	1474103866	2016-09-17T09:17:46Z	1474103893	2016-09-17T09:18:13Z	Farine de blé noir	NaN	1kg	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	
1	0000000004530	http://world-fr.openfoodfacts.org/produit/0000...	usda-ndb-import	1489069957	2017-03-09T14:32:37Z	1489069957	2017-03-09T14:32:37Z	Banana Chips Sweetened (Whole)	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	14.0	14.0	NaN	
2	0000000004559	http://world-fr.openfoodfacts.org/produit/0000...	usda-ndb-import	1489069957	2017-03-09T14:32:37Z	1489069957	2017-03-09T14:32:37Z	Peanuts	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	0.0	0.0	NaN	
3	00000000016087	http://world-fr.openfoodfacts.org/produit/0000...	usda-ndb-import	1489055731	2017-03-09T10:35:31Z	1489055731	2017-03-09T10:35:31Z	Organic Salted Nut Mix	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	12.0	12.0	NaN	
4	00000000016094	http://world-fr.openfoodfacts.org/produit/0000...	usda-ndb-import	1489055653	2017-03-09T10:34:13Z	1489055653	2017-03-09T10:34:13Z	Organic Polenta	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	

5 rows × 162 columns

Analyse
exploratoire des
données

Réduction du
dataset et
sélection de la cible

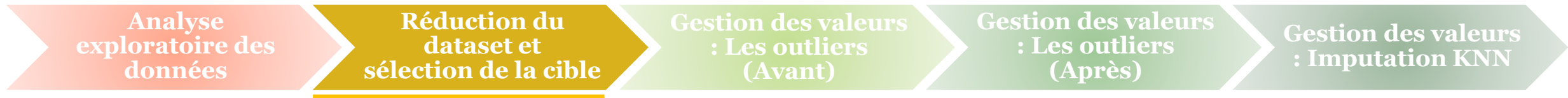
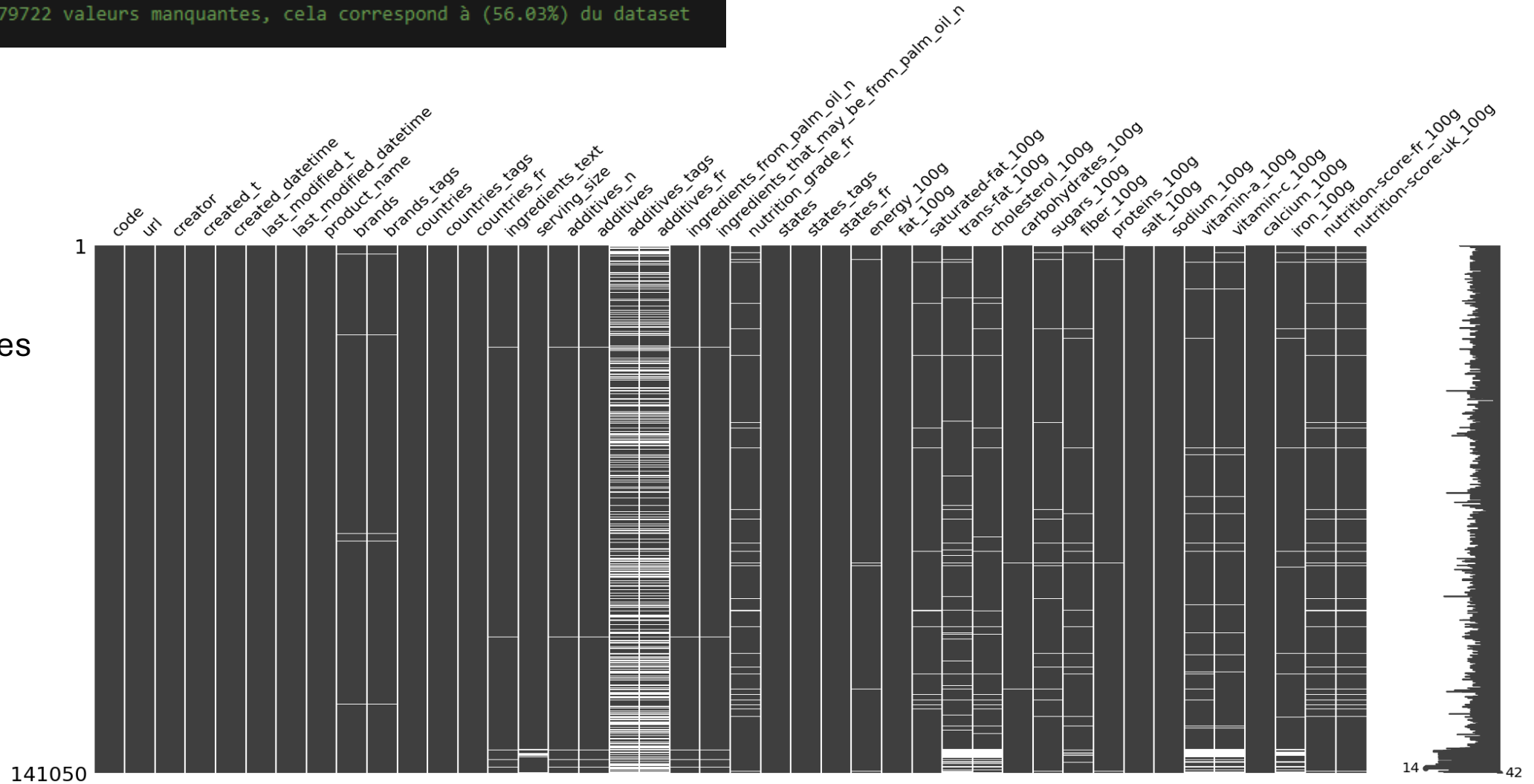
Gestion des valeurs
: Les outliers
(Avant)

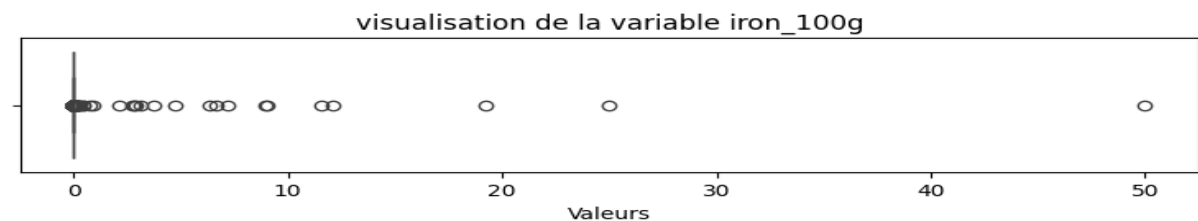
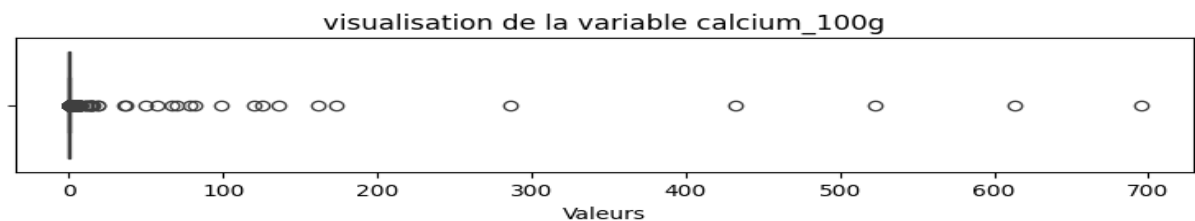
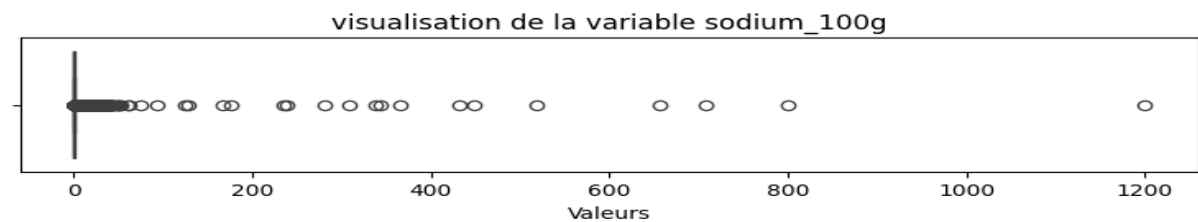
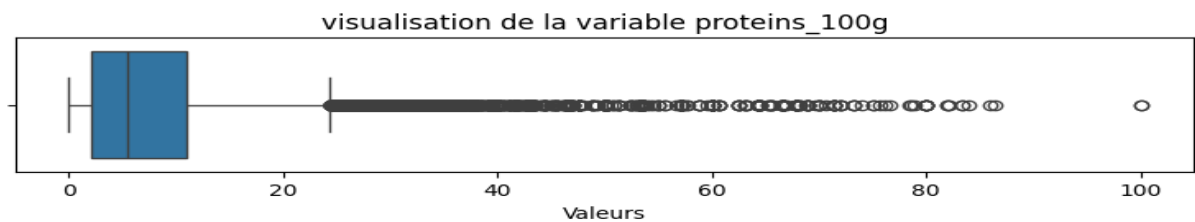
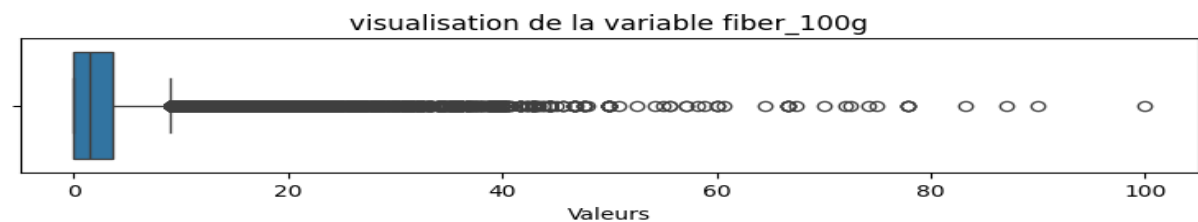
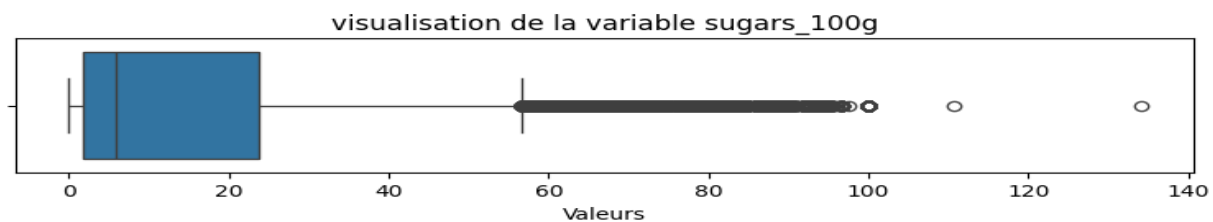
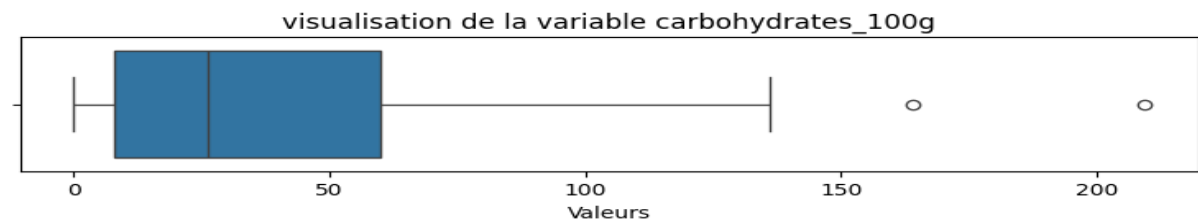
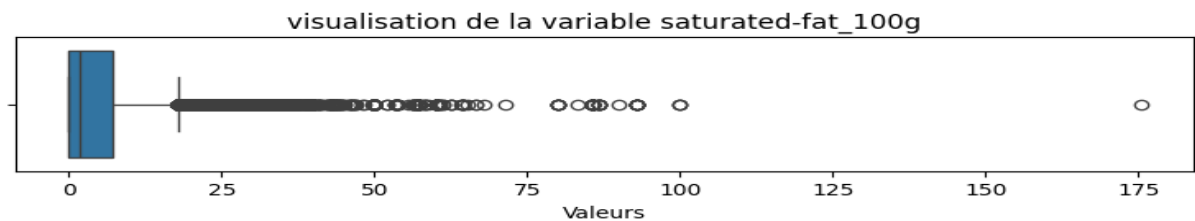
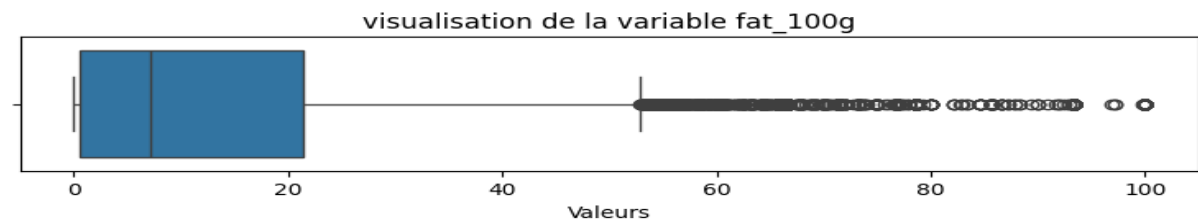
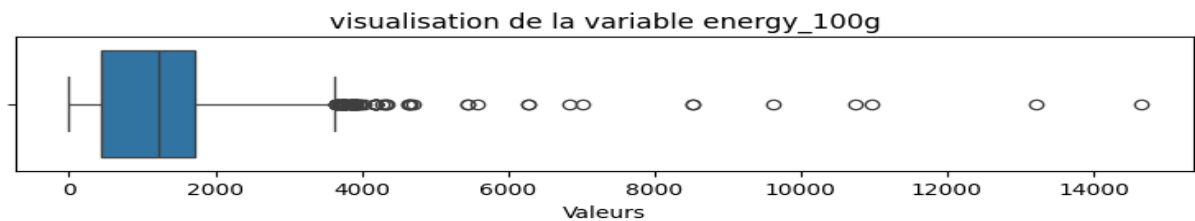
Gestion des valeurs
: Les outliers
(Après)

Gestion des valeurs
: Imputation KNN

```
#####  
# 1.2 : Choix de la cible  
#####  
  
# (quant) calcium_100g - 179722 valeurs manquantes, cela correspond à (56.03%) du dataset
```

Les colonnes avec plus de 60% de valeurs manquantes ont été supprimé





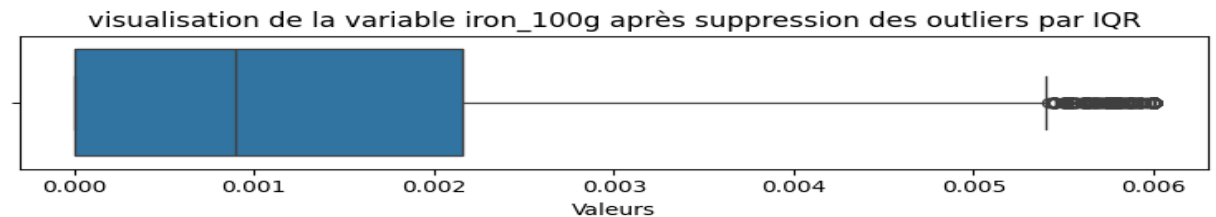
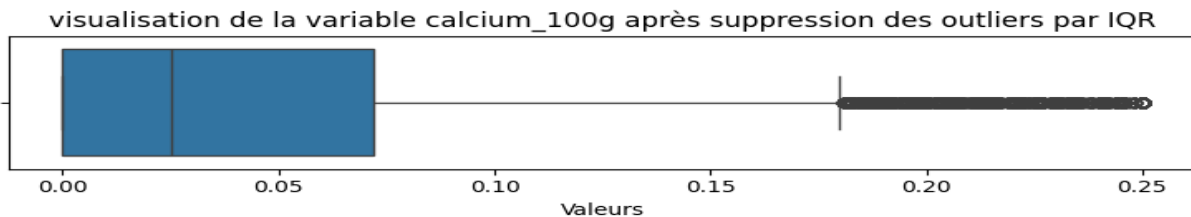
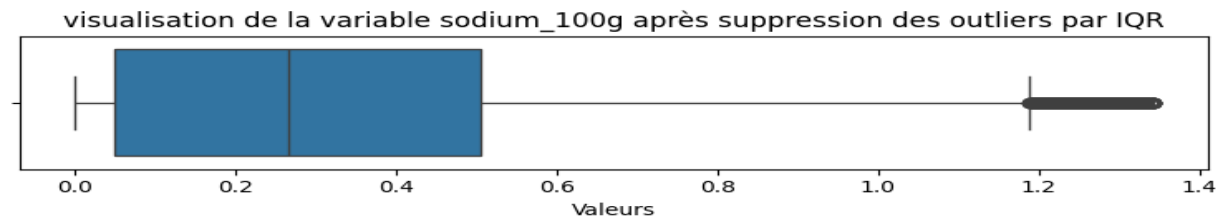
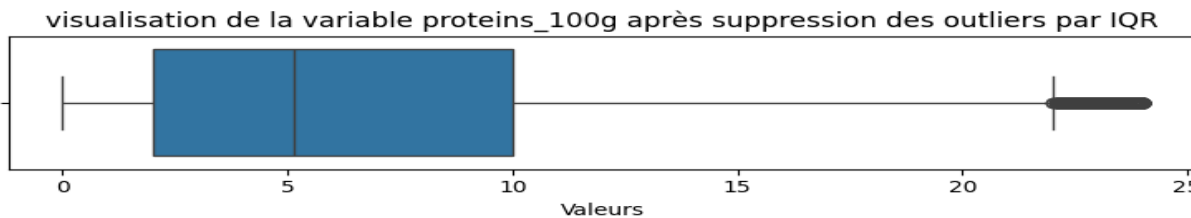
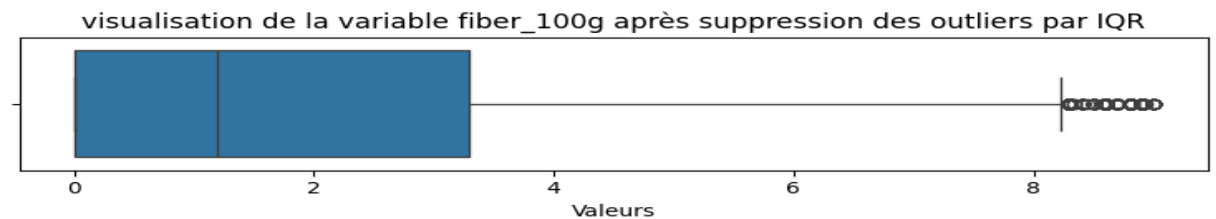
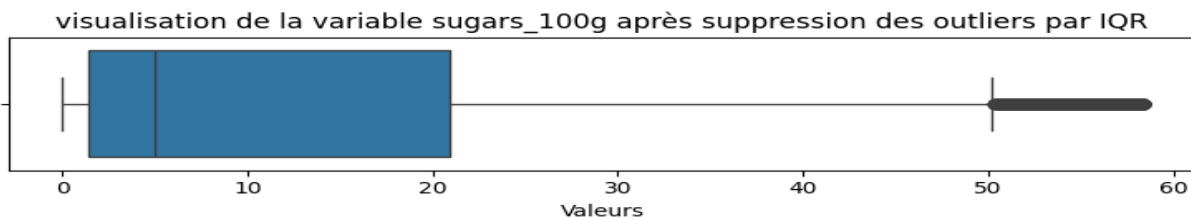
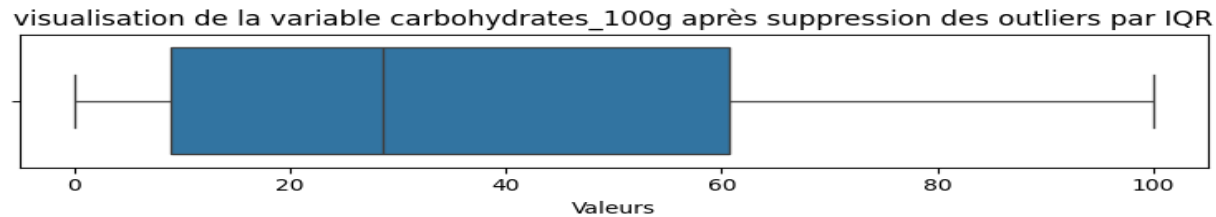
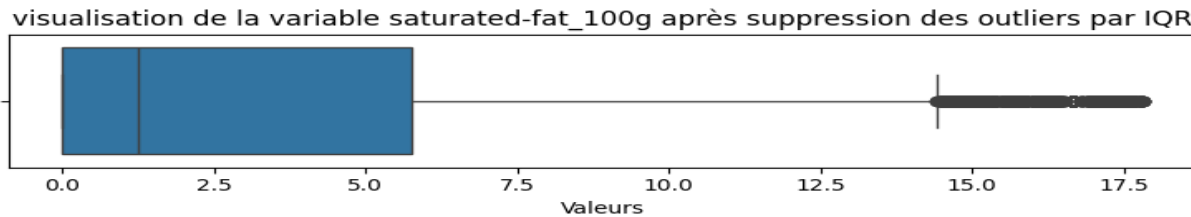
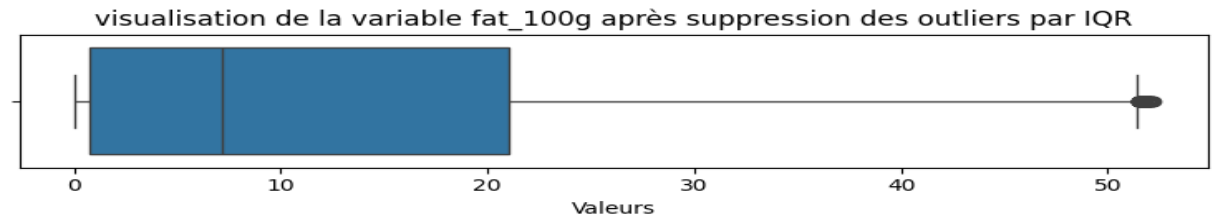
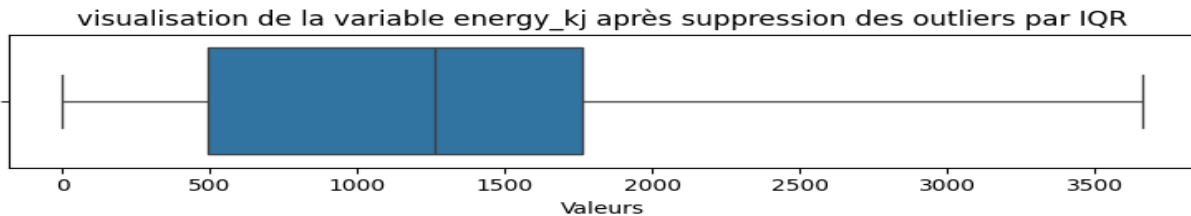
Analyse
exploratoire des
données

Réduction du
dataset et
sélection de la cible

Gestion des valeurs
: Les outliers
(Avant)

Gestion des valeurs
: Les outliers
(Après)

Gestion des valeurs
: Imputation KNN



Analyse
exploratoire des
données

Réduction du
dataset et
sélection de la cible

Gestion des valeurs
: Les outliers
(Avant)

Gestion des valeurs
: Les outliers
(Après)

Gestion des valeurs
: Imputation KNN

```

num_features = data.select_dtypes(include = ["number"]).columns
num_fill_rate = data[num_features].notnull().mean().sort_values(ascending = False)

print("Taux de remplissage des variables numériques (en %) : ")
print(round(num_fill_rate * 100, 2))

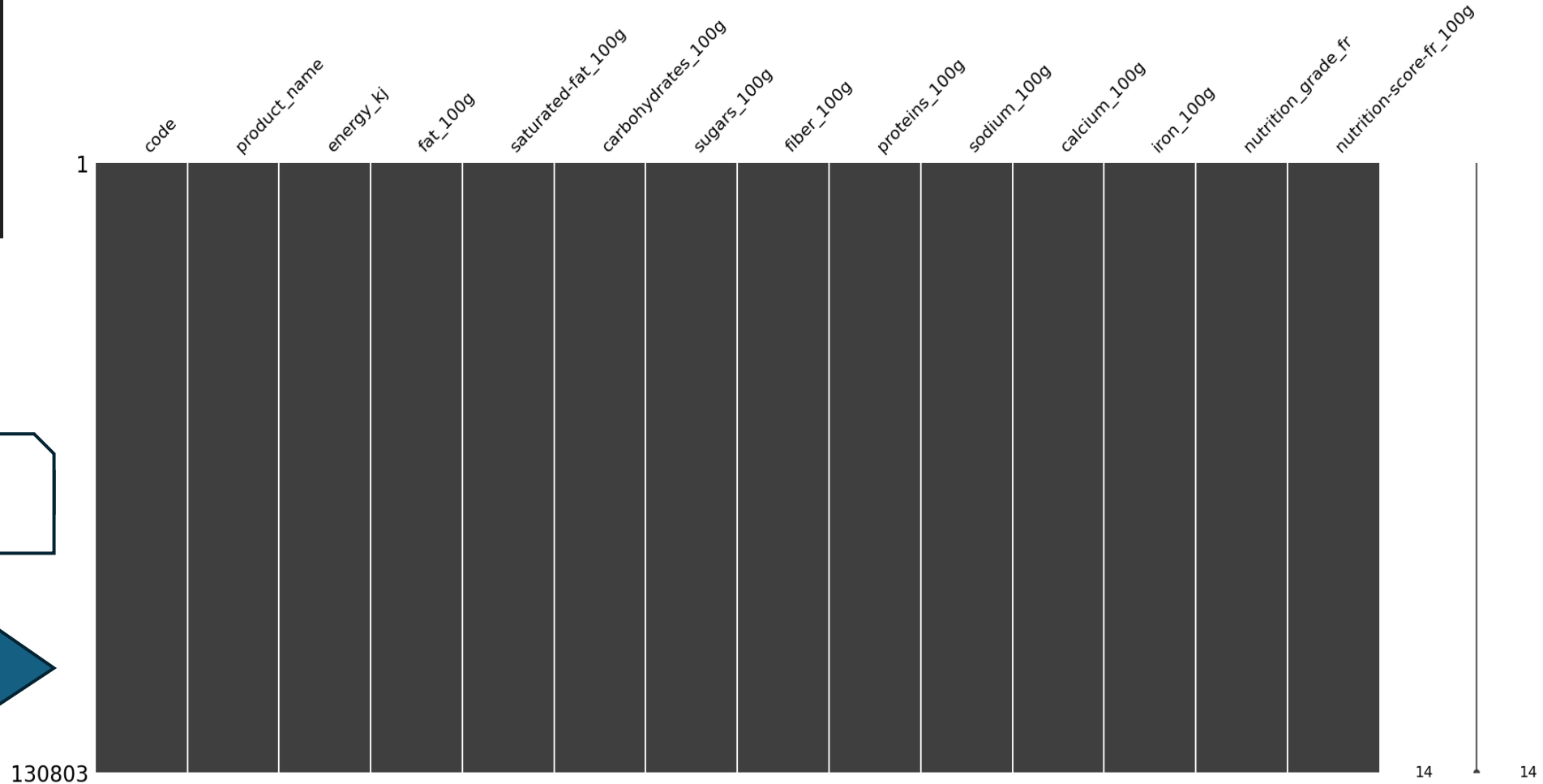
```

✓ 0.0s

Taux de remplissage des variables numériques (en %) :

carbohydrates_100g	100.00
energy_kj	100.00
fat_100g	97.75
iron_100g	95.73
sugars_100g	95.33
sodium_100g	94.88
proteins_100g	94.32
saturated-fat_100g	93.36
fiber_100g	92.36
calcium_100g	91.68
nutrition-score-fr_100g	83.78
dtype: float64	

Imputation
avec KNN



Analyse
exploratoire des
données

Réduction du
dataset et
sélection de la cible

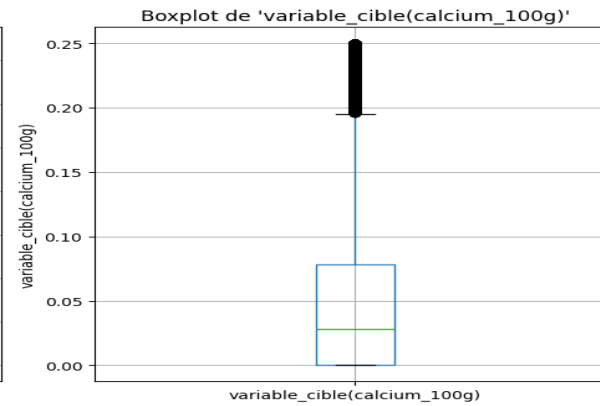
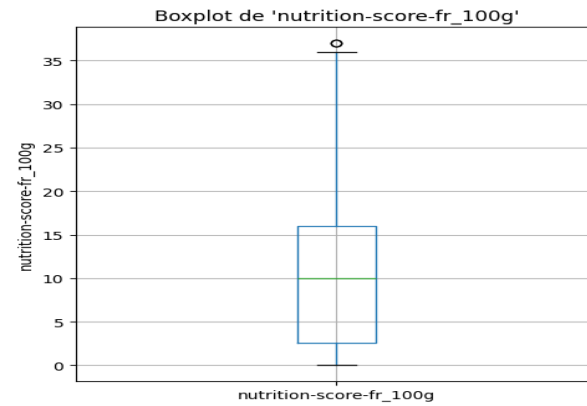
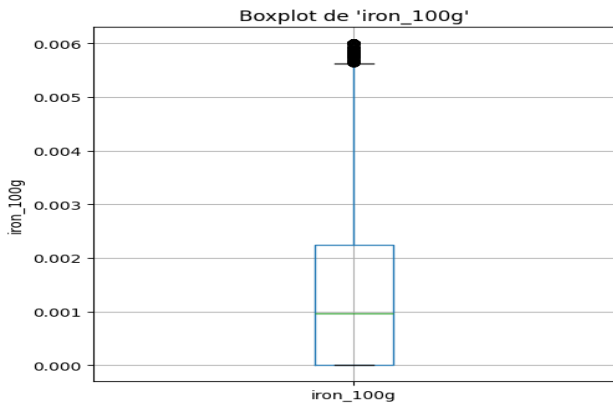
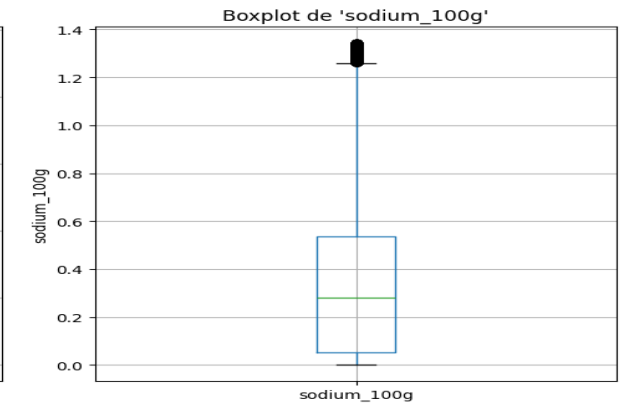
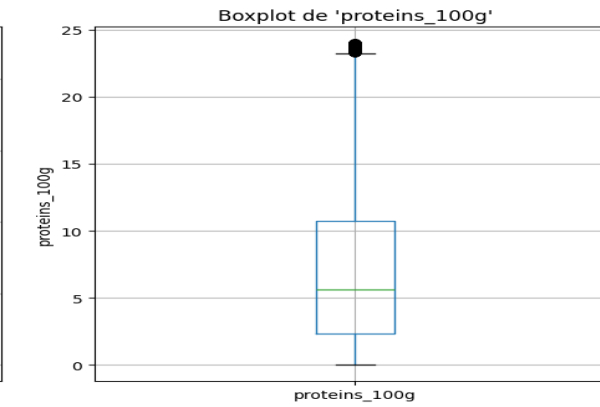
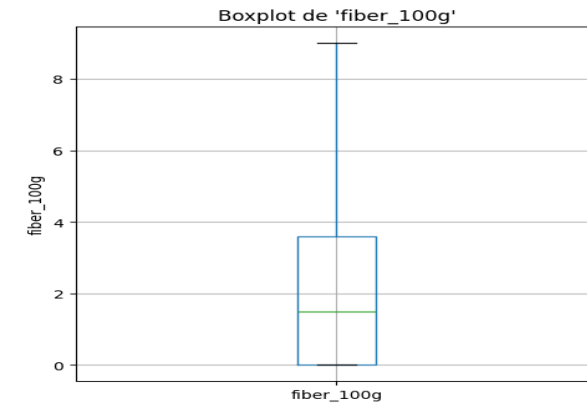
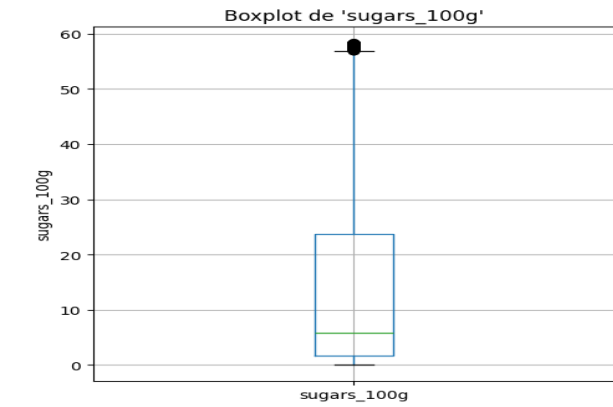
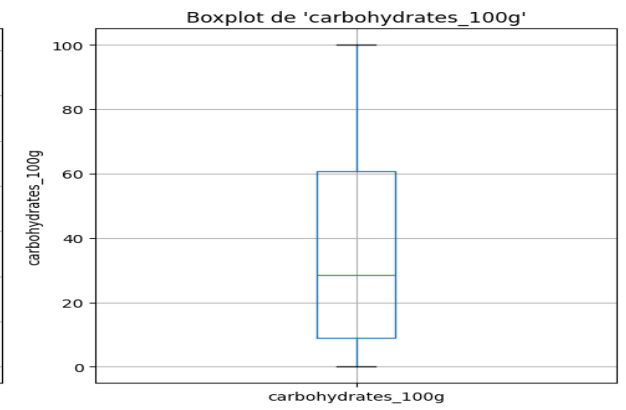
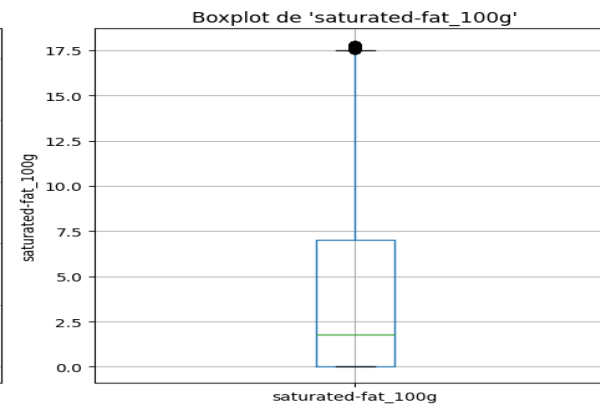
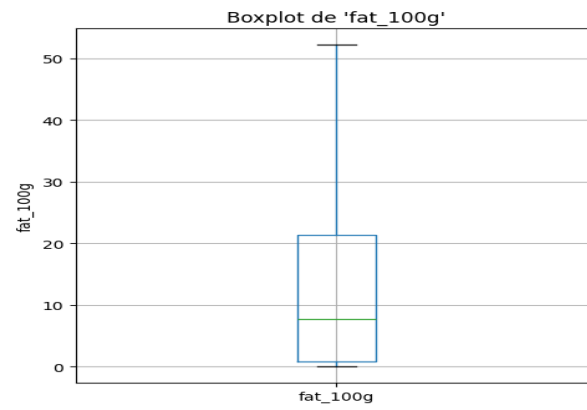
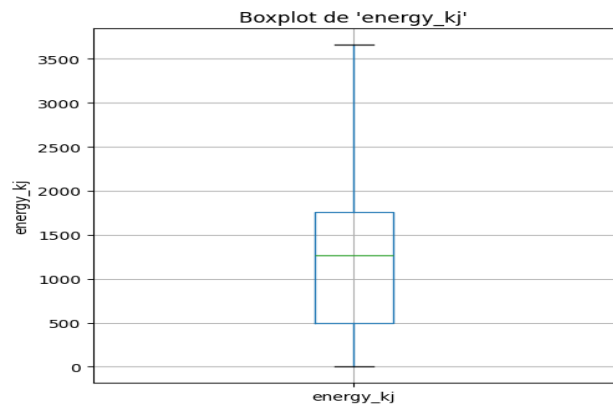
Gestion des valeurs
: Les outliers
(Avant)

Gestion des valeurs
: Les outliers
(Après)

Gestion des valeurs
: Imputation KNN

2^{ème} Notebook

L'analyse



Analyses
univariées

Analyses
bivariées

Analyses
multivariées

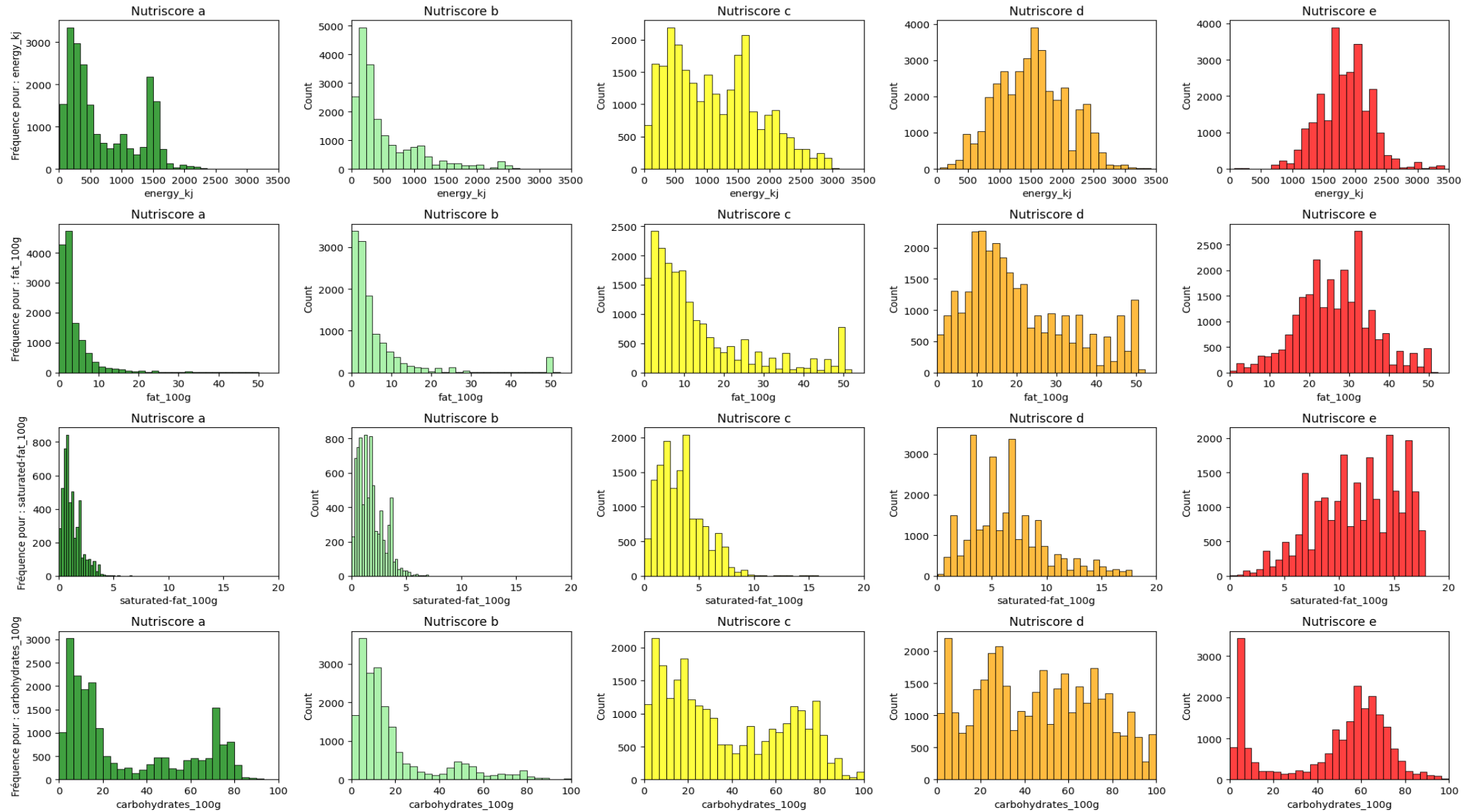
Test Statistiques

Analyse des
Composantes
Principales
(ACP)

Cercles de
corrélation

Conclusion

Distribution des nutriments par nutriscore



Analyses
univariées

Analyses
bivariées

Analyses
multivariées

Test Statistiques

Analyse des
Composantes
Principales
(ACP)

Cercles de
corrélation

Conclusion

Analyses bivariées

Pour la variable catégorielle

Test Shapiro-Wilk

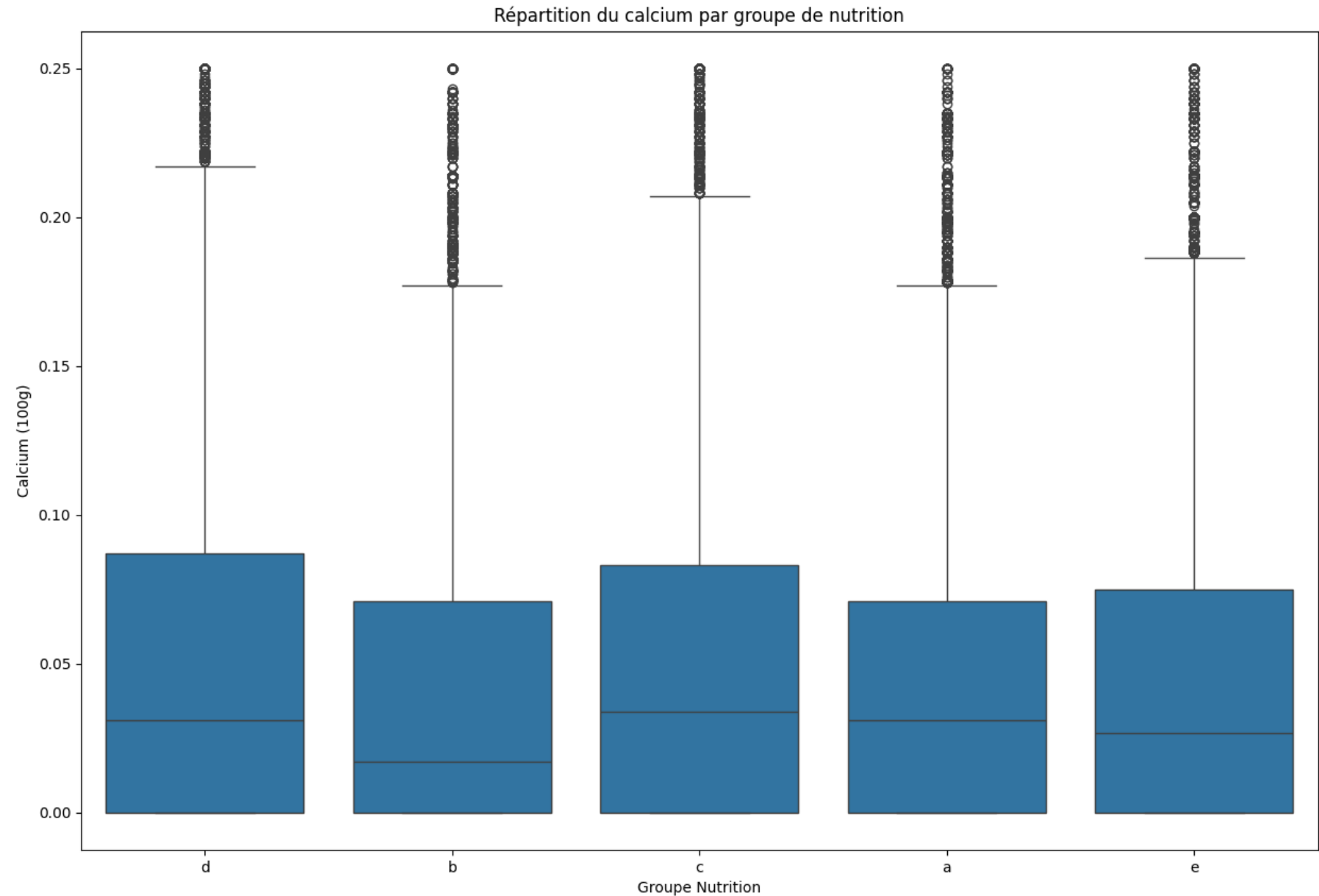
W-statistic = 0.00039062091216846007

p-value = 9.436948635914407e-211

Test Kruskal-Wallis

Statistique = 175.8750825499985

p-value = 5.731868141181024e-37



Analyses
univariées

Analyses
bivariées

Analyses
multivariées

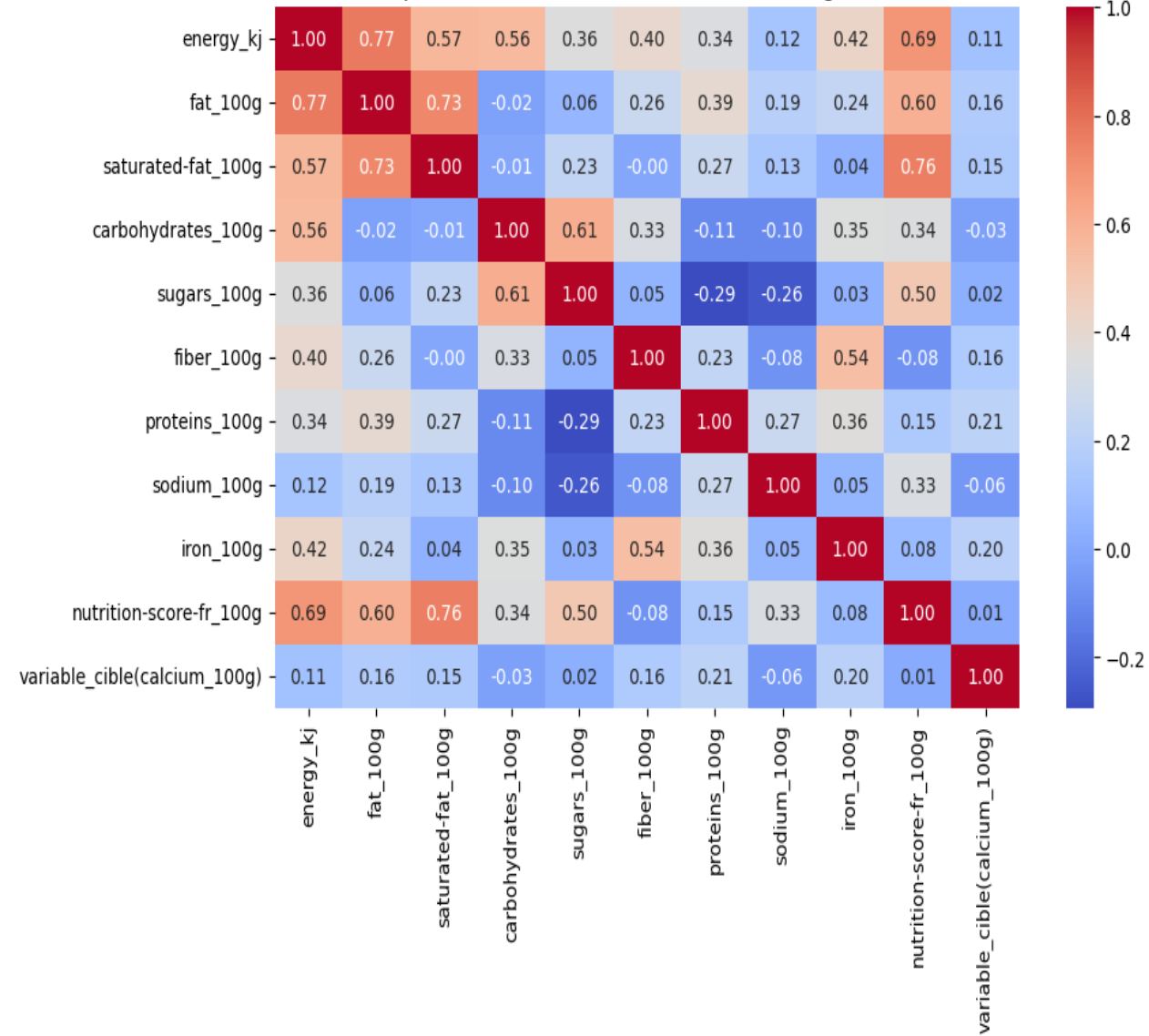
Test Statistiques

Analyse des
Composantes
Principales
(ACP)

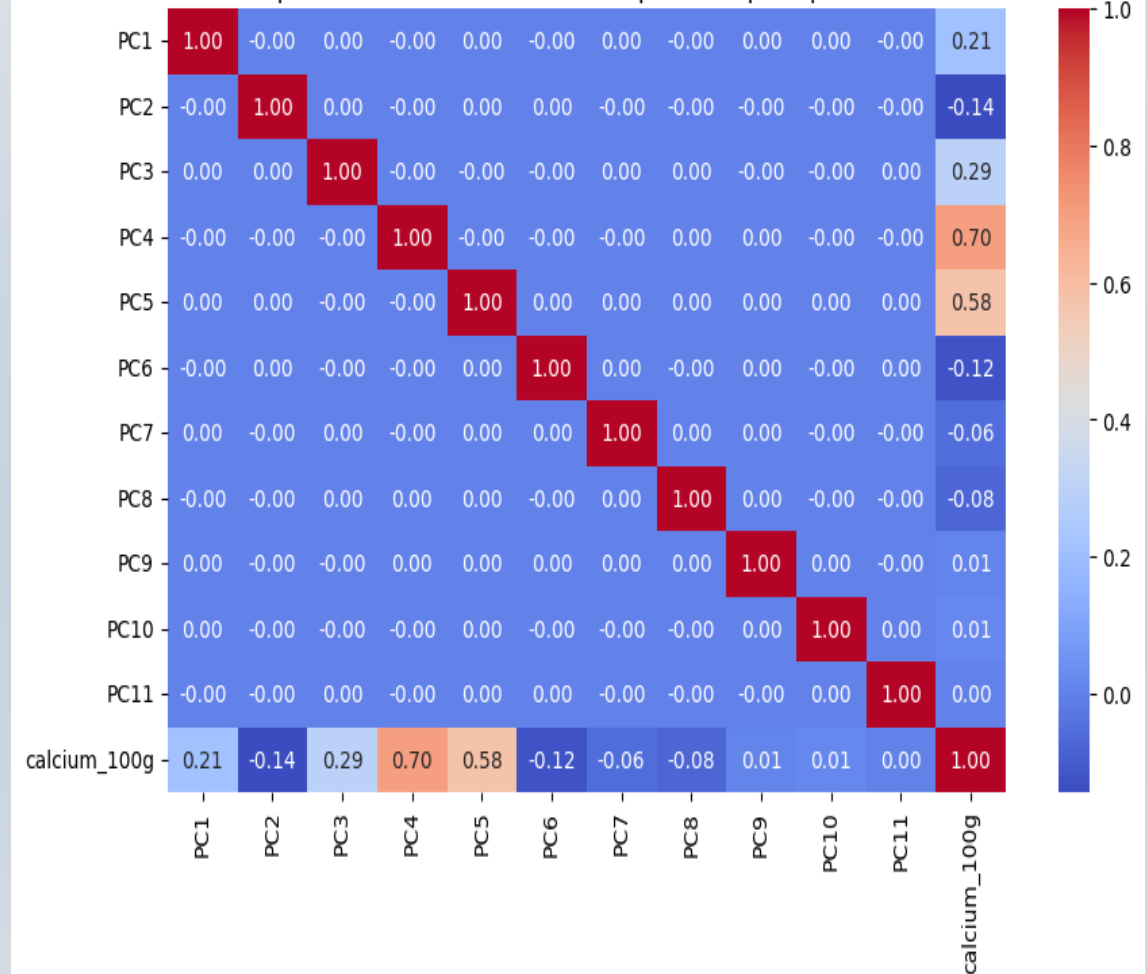
Cercles de
corrélation

Conclusion

Heatmap des corrélations entre les variables d'origine et la cible



Heatmap des corrélations entre les composantes principales et la cible

Analyses
univariéesAnalyses
bivariéesAnalyses
multivariées

Test Statistiques

Analyse des
Composantes
Principales
(ACP)Cercles de
corrélation

Conclusion

Tests statistiques

Pour la variable catégorielle

Test Kolmorov-Smirnov

PC1 : Statistique = 0.2165941598036073,
p-value = 0.0

PC3 : Statistique = 0.09339074217058335,
p-value = 0.0

PC4 : Statistique = 0.039900171733098566,
p-value = 2.236071262374604e-181

PC5 : Statistique = 0.06908826726005846,
p-value = 0.0

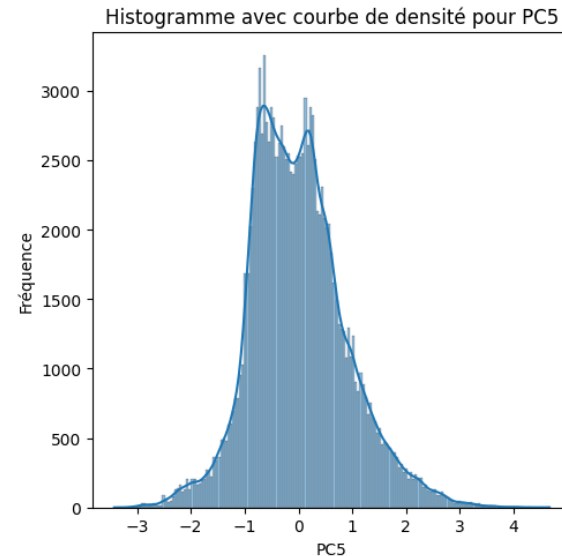
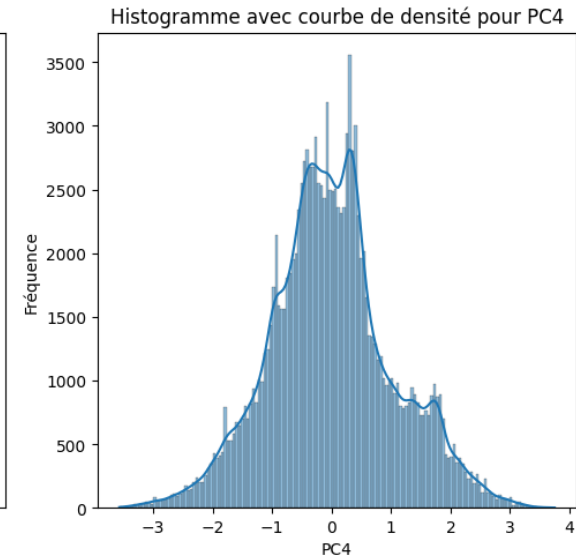
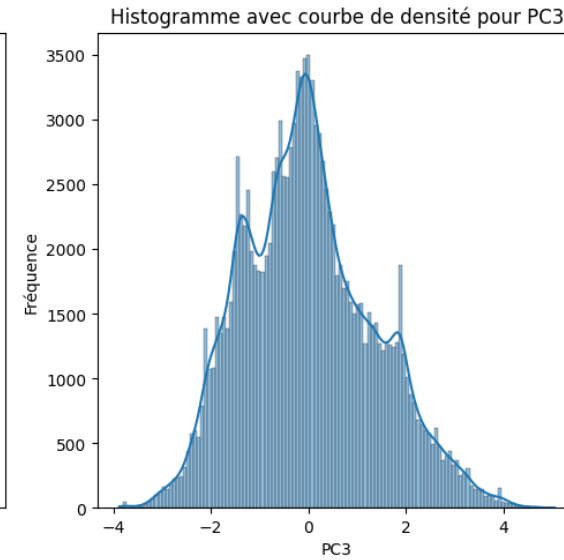
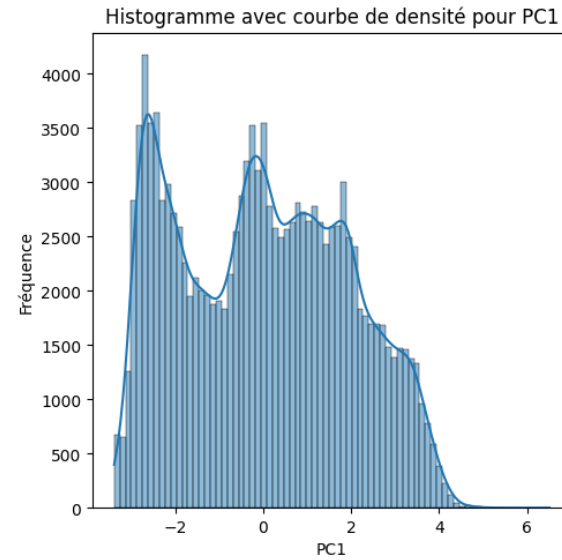
Test Spearman

PC1 et calcium 100g : 0.17159376490111275,
p-value : 0.0

PC3 et calcium 100g : 0.29003305611411057,
p-value : 0.0

PC4 et calcium 100g : 0.6703412004890147,
p-value : 0.0

PC5 et calcium 100g : 0.5003067599770551,
p-value : 0.0



Analyses
univariées

Analyses
bivariées

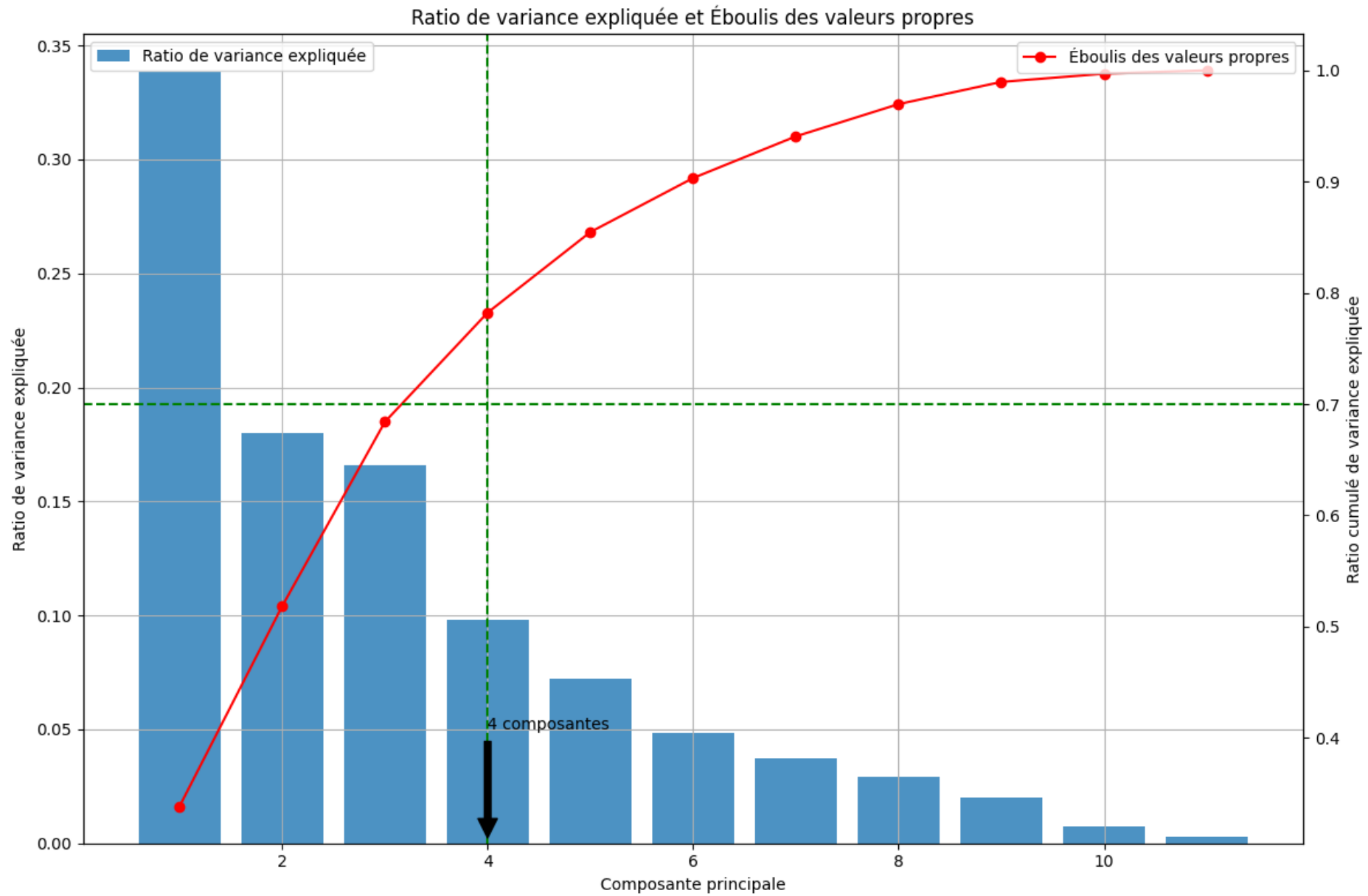
Analyses
multivariées

Test Statistiques

Analyse des
Composantes
Principales
(ACP)

Cercles de
corrélation

Conclusion



Analyses
univariées

Analyses
bivariées

Analyses
multivariées

Test Statistiques

Analyse des
Composantes
Principales
(ACP)

Cercles de
corrélation

Conclusion

Cercle de corrélation

Pour la variable catégorielle

PC1

Explique 33,8% de la variance

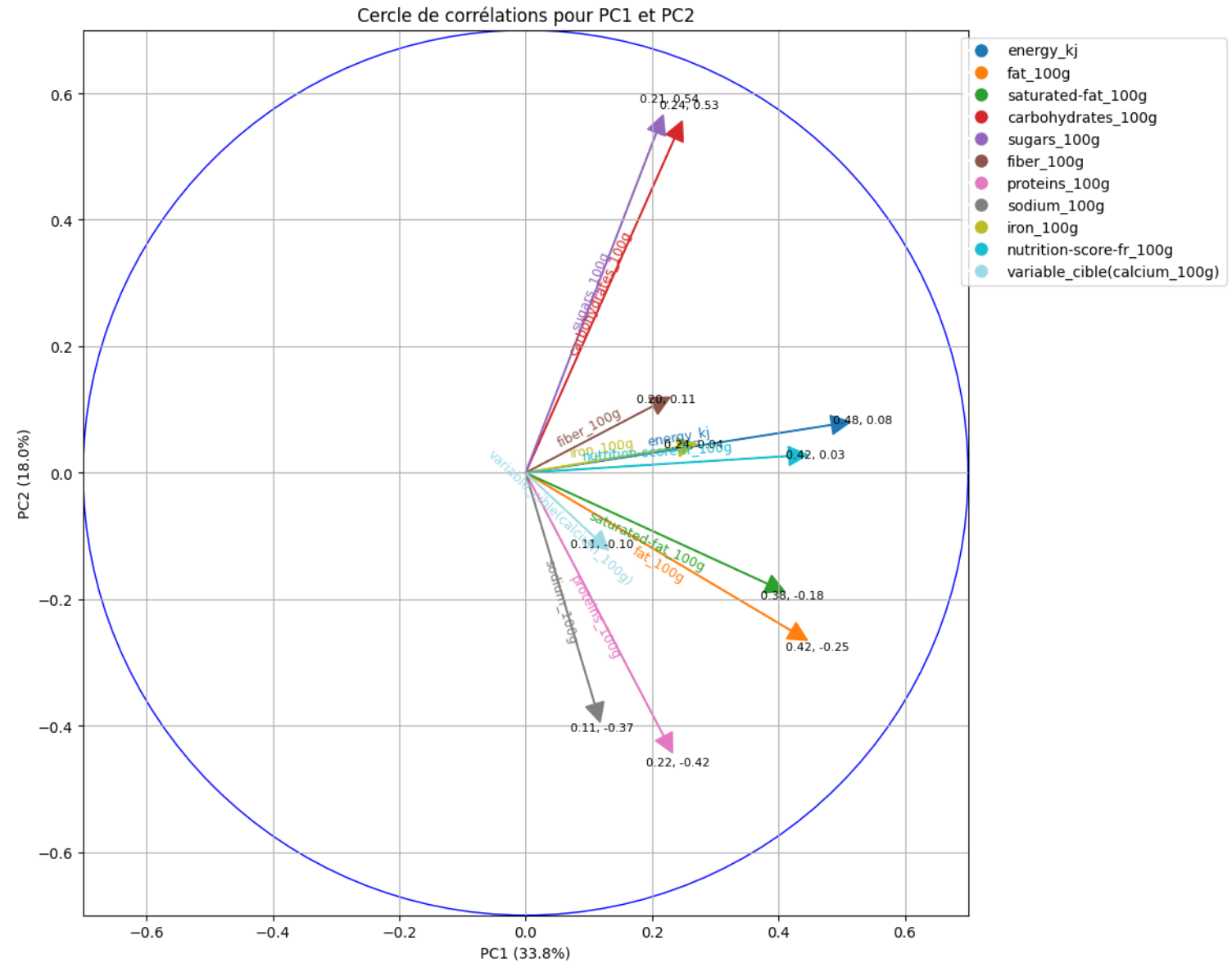
PC2

Explique 18,0% de la variance

Les variables qui doivent impérativement être rempli dans la future application sont :

- Energy kj
- Carbohydrates
- sugars
- saturated-fat
- fat
- proteins
- sodium

Ces informations permettront de définir le nutriscore (A, B, C, D et E), l'ajout de nouvelles variables comme les fibres, le calcium permettrait de compléter l'information nutritionnelles.



Analyses
univariées

Analyses
bivariées

Analyses
multivariées

Test Statistiques

Analyse des
Composantes
Principales
(ACP)

Cercles de
corrélation

Conclusion

Résumé du projet :

- * Améliorer la base de données d'Open Food Facts pour santé publique France
- * Mise en place d'un système de suggestion et d'auto-complétion pour réduire les erreurs de saisie

Réalisation :



Impact :



Réduction des champs à saisir, cela signifie moins de possibilité d'erreur de saisie donc une amélioration de la qualité des données

Nous avons maintenant une meilleure compréhension des facteurs influençant la qualité nutritionnelle des produits



Proposition d'amélioration :

- * Dans l'application final limité la saisie à une valeurs maximal (si 100g impossibilité de saisir 101)
- * Sensibiliser les utilisateurs sur l'importance de la qualité des données
- * Ajouter d'autre informations nutritionnelles comme les différentes vitamines

Analyses univariées

Analyses bivariées

Analyses multivariées

Test Statistiques

Analyse des Composantes Principales (ACP)

Cercles de corrélation

Conclusion

MERCI