



# Parcours Data Scientist

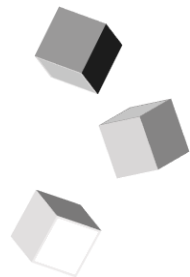
Projet N°3 :

Concevez une application au service de la santé publique

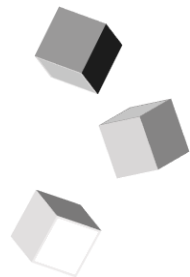
Daniel CHASTANET



## Sommaire



- ❑ Rappel de la problématique et idée d'application
- ❑ Nettoyage et analyse pré-exploratoire
- ❑ Analyse exploratoire
- ❑ Remplacement des données manquantes
- ❑ Mise en œuvre de l'application
- ❑ Conclusion



# Partie 1 : Présentation de l'idée d'application



## Projet de création d'application en lien avec la nutrition/santé

### BDD : Open Food Fact

- Automatisation du traitements des données
- Visualisations simples pour présenter les données
- Hypothèses et tests statistiques appropriés
- Elaborer une idée d'application



## Etude de ce qu'on trouve sur le « marché »

### Nutriscore :



Système d'étiquetage nutritionnel pour faciliter la compréhension des informations nutritionnelles (état : ministère de la Santé / volontariat)

### Eco score :



Indicateur représentant l'impact environnemental des produits alimentaires (état : ADEME / volontariat)

### Yuka :

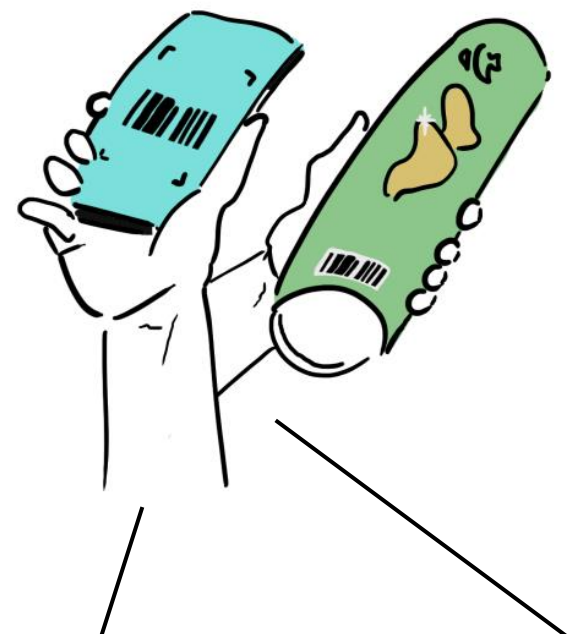


décrypte les étiquettes des produits alimentaires et cosmétiques et analyse leur impact sur la santé.

Une base de données complète : (1,5 million de produits alimentaires, 500 000 références cosmétiques, 800 nouveaux produits chaque jour)

Recommandation

Financement indépendant



Scan du code bar

**OC Nutrition**

**Nutri-Score** **Nova-Score** **Eco-Score** **k-score** **Addi-Score** **Méga-score**

**Similar products**

Chips éco	M	A	A	E	A	D
Chips chères	M	A	A	E	A	D
Chips colorées	M	A	A	E	A	D

**OC Nutrition**

**Nutri-Score** **Trans-Score**

**Nutri-Score détails**

Qualités pour 100g

- Fibres: Excellente quantité, 6,2 g
- Sucre: Peu de sucre, 6,3 g
- Calories: Peu calorique, 153 kcal
- Graisses saturées: Pas de graisse sat., 0 g

**Similar products**

Chips éco	M	A	A	E	A	D
Chips chères	M	A	A	E	A	D
Chips colorées	M	A	A	E	A	D

**OC Nutrition**

Les fibres sont des substances d'origine végétale indispensables au bon fonctionnement de l'intestin. On en distingue deux types : solubles et insolubles...

**Nutri-Score détails**

Qualités pour 100g

- Fibres: Excellente quantité, 6,2 g
- Sucre: Peu de sucre, 6,3 g
- Calories: Peu calorique, 153 kcal
- Graisses saturées: Pas de graisse sat., 0 g

**Similar products**

Chips éco	M	A	A	E	A	D
Chips chères	M	A	A	E	A	D
Chips colorées	M	A	A	E	A	D



# L'application



Nutri-Score



k-score



Addi-Score



Nova-Score

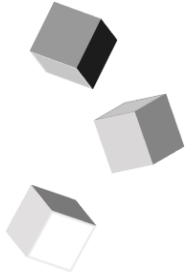


Eco-Score



Méga-score

Nom / Définition	Colonnes du dataframe
Nutri-Score : (Existant) : Barème calories, acides gras, sucres, fibres...	nutriscore_score / nutriscore_grade
k-score (non existant) : calories et mauvais gras !	energy_100g / fat_100g / saturated-fat_100g / carbohydrates_100g / sugars_100g / fiber_100g / proteins_100g / salt_100g
Addi-Score (mi - existant) : Nombre d'additifs (et dangerosité)	additives_n / Additives / additives_tags / additives_en
Nova-Score (existant) : Nombre de transformation et type	nova_group
Eco Score (Existant) : représente l'impact environnemental des produits alimentaires. Émissions de gaz à effet de serre (CO2) / Destruction de la couche d'ozone / Émissions de particules fines / Oxydation photochimique / Acidification / Radioactivité / etc...	ecoscore_score_fr / ecoscore_grade_fr
$\Sigma$ Somme des autres score	



# Partie 2 : Nettoyage de la base de données





# Prendre connaissances de la BDD

## Légende des données : dimension : (1918228, 186)

code : barcode of the product (can be EAN-13 or internal codes for some food stores)  
003. quantity : quantity and unit  
004. url : url of the product page on Open Food Facts  
005. creator : contributor who first added the product  
006. created\_t : date that the product was added (UNIX timestamp format)  
007. created\_datetime : date that the product was added (iso8601 format: yyyy-mm-ddThh:mn:ssZ)  
008. last\_modified\_t : date that the product page was last modified  
009. last\_modified\_datetime : date that the product was lastly modified (UNIX timestamp format)  
010. product\_name : name of the product  
011. abbreviated\_product\_name  
012. generic\_name : generic name of the product  
013. quantity : en g/l/capsules... ou no unit  
014. packaging : type --> shape, material : plastique / bocal etc...  
015. packaging\_tags : same que packaging avec des tirets à la place des espaces  
016. packaging\_text : '1 Pot en verre à recycler', ..., '1 PET packet'  
017. brands : 'courte paille', ..., 'miellerie de la natouze', 'Biocoop Bordeaux lac'  
018. brands\_tags : same stuff que 'brands' avec le tiret  
019. categories : 'Epicerie, Condiments, Sauces, Moutardes'  
...  
040. ingredients\_text : description des ingrédients avec pourcent ou pas  
041. allergens : 'en:mustard', 'en:eggs,en:mustard', ..., 'fr:Gs1:T4078:Al,fr:Gs1:T4078:BA'  
042. allergens\_en : ...  
043. traces : 'fr:Fruits à coques et/ou cacahuètes'  
044. traces\_tags : ...  
...

050. additives : empty...  
051. additives\_tags : 'en:e150,en:e160a,en:e202'  
052. additives\_en : ...  
053. ingredients\_from\_palm\_oil\_n : in number from nan, 0 to 5  
054. ingredients\_from\_palm\_oil : vide  
055. ingredients\_from\_palm\_oil\_tags : 'stearine-de-palme,oleine-de-palme,'e304-palmitate-d-ascorbyle'  
056. ingredients\_that\_may\_be\_from\_palm\_oil\_n : number  
057. ingredients\_that\_may\_be\_from\_palm\_oil  
058. ingredients\_that\_may\_be\_from\_palm\_oil\_tags  
059. nutriscore\_score : nan, 18., 1., 14., -2.  
060. nutriscore\_grade : nan, 'd', 'b', 'a', 'c', 'e' <https://fr.openfoodfacts.org/nutriscore>  
061. nova\_group : Une classification en 4 groupes pour mettre en évidence le degré de transformation des aliments nan, 1 à 4 <https://fr.openfoodfacts.org/nova>  
062. pnns\_groups\_1 : Programme national nutrition santé (PNNS) / 'Fat and sauces', 'Composite foods', 'Sugary snacks'  
063. pnns\_groups\_2 : 'Dressings and sauces', 'One-dish meals'  
...  
080. energy-kj\_100g : you know...  
energy-kcal\_100g  
energy\_100g  
...  
-linoleic-acid\_100g  
-arachidonic-acid\_100g  
-gamma-linolenic-acid\_100g  
-dihomo-gamma-linolenic-acid\_100g  
omega-9-fat\_100g  
-oleic-acid\_100g  
-elaidic-acid\_100g



## Nettoyage des données

Very large data frame (au moins pour mon pc) : `dimension : (1918228, 186)`

### Filtrage colonnes

➡ Removing columns with 80% and more missing values :

`(1918228, 57)`

➡ Keeping only interesting columns (for the app) :

`(1918228, 21)`

### Filtrage lignes

➡ Filtrage par « `ecoscore_score` » `!= isna()` :

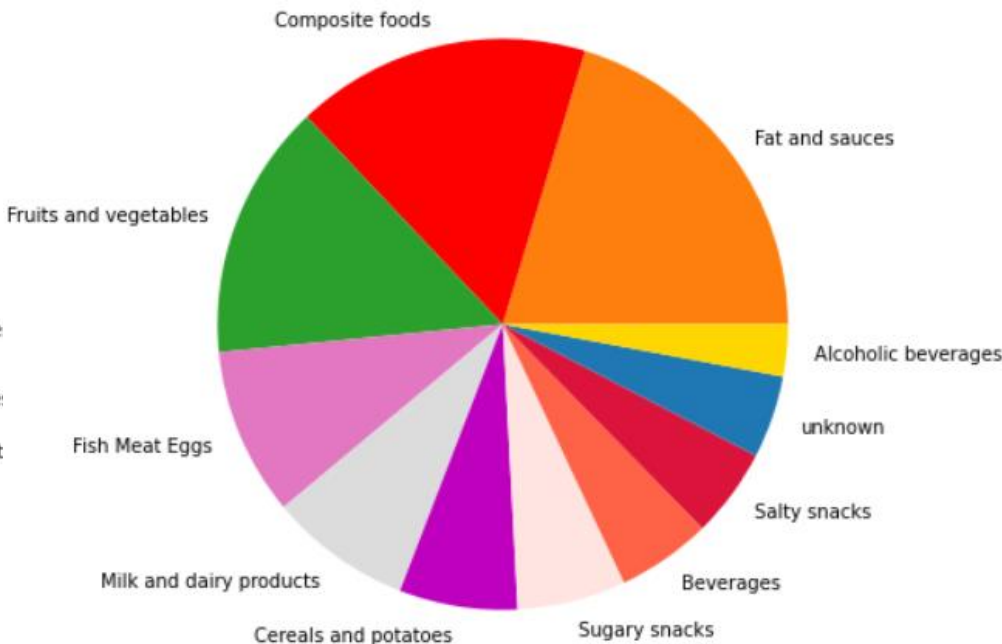
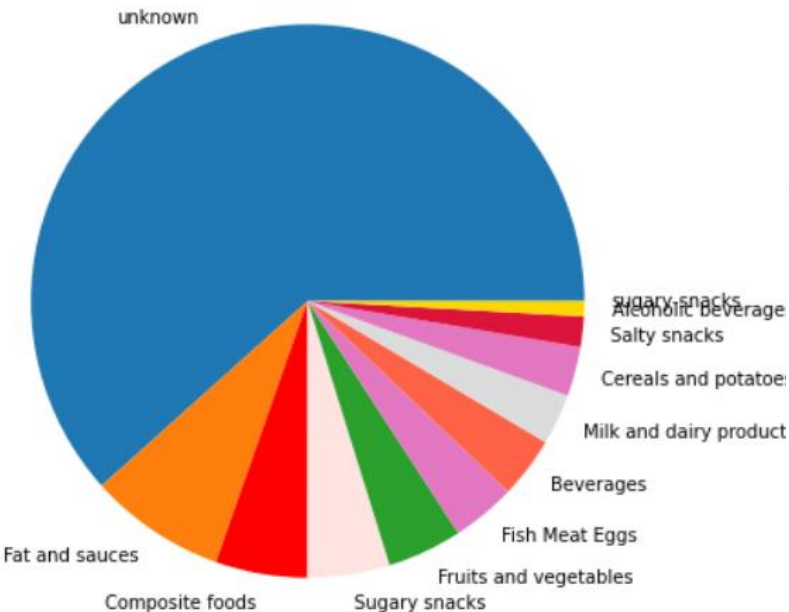
`(462847, 21)`

% de données manquantes

code	0.000000
product_name	0.725942
ingredients_text	45.071698
additives_n	45.071482
additives_tags	71.674009
additives_en	71.674009
nutriscore_score	22.651978
nutriscore_grade	22.651978
nova_group	49.852327
pnns_groups_1	0.000000
pnns_groups_2	0.000000
ecoscore_score_fr	0.000000
ecoscore_grade_fr	0.000000
energy_100g	15.015113
fat_100g	15.299440
saturated-fat_100g	16.818301
carbohydrates_100g	15.357559
sugars_100g	16.347519
fiber_100g	65.960026
proteins_100g	15.245211
salt_100g	16.701199



Do I cut too much ? Let's check the pnns\_groups (pnns\_1)

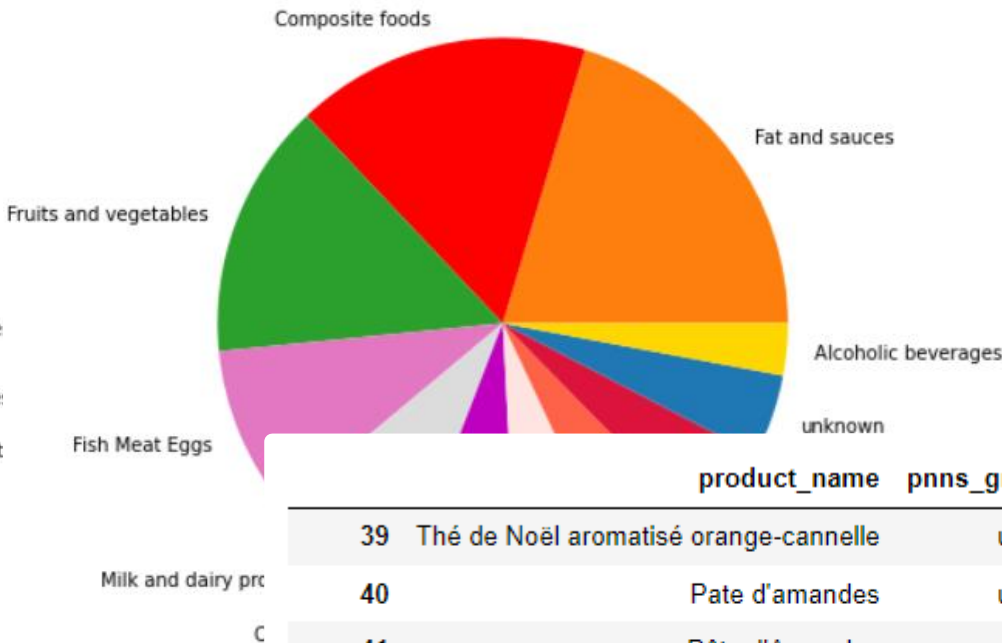
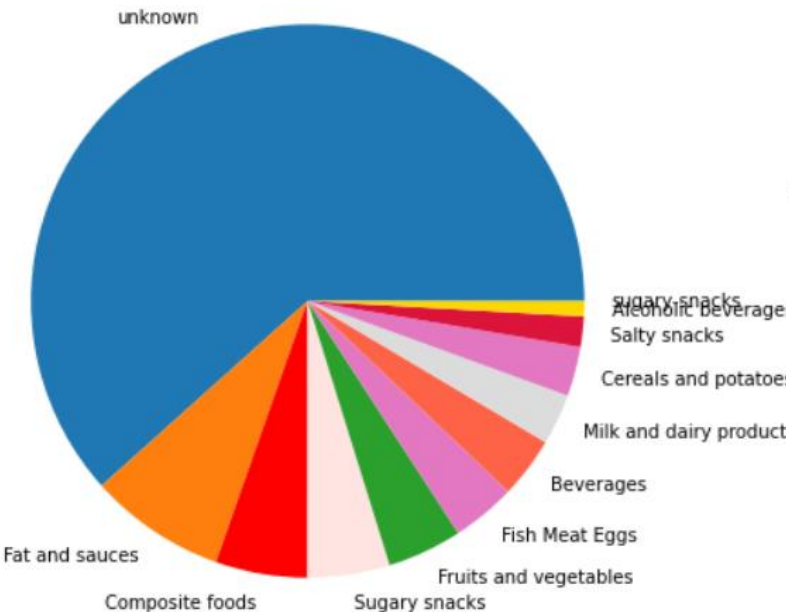


Waters and flavored waters		6863	
	pnns_groups_1	n	n_filtrés
0	Alcoholic beverages	17040	13591.0
1	Beverages	68447	23376.0
2	Cereals and potatoes	83804	43985.0
3	Composite foods	55513	21630.0
4	Fat and sauces	71737	37372.0
5	Fish Meat Eggs	102955	67140.0
6	Fruits and vegetables	56971	31006.0
7	Milk and dairy products	92067	77362.0
8	Salty snacks	33976	28531.0
9	Sugary snacks	153243	94136.0
10	sugary-snacks	2	NaN
11	unknown	1181969	24718.0



# Nettoyage des données

Do I cut too much ? Let's check the pnns\_groups (pnns\_1)



	product_name	pnns_groups_1	ecoscore_score_fr
39	Thé de Noël aromatisé orange-cannelle	unknown	92.0
40	Pate d'amandes	unknown	5.0
41	Pâte d'Amandes	unknown	10.0
57	Tisane nerf - sommeil	unknown	94.0
73	Pomme dauphine	unknown	38.0

Waters and flavored waters		6863		
	pnns_groups_1	n	n_filtrés	
0	Alcoholic beverages	17040	13591.0	
1	Beverages	68447	23376.0	
2	Cereals and potatoes	83804	43985.0	
3	Composite foods	55513	21630.0	
4	Fat and sauces	71737	37372.0	
5	Fish Meat Eggs	102955	67140.0	
		56971	31006.0	
		92067	77362.0	
		33976	28531.0	
		153243	94136.0	
		2	NaN	
		1181969	24718.0	

pnns cat : 12 / 42  
(1918228, 21)



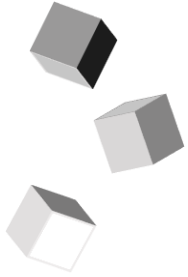
## Traitement valeurs abh rantes

(462847, 21)

### Filtrage valeurs aberrantes :

- ☐ Nombre de Kalories aberrantes (trop grande) 900kcal pour 100g au maximum
- ☐ 150 g de X dans 100 g de produit (X pouvant  tre du sel, sucre, fibre, etc)
- ☐ Valeurs n gatives en g pour 100 g
- ☐ Produit dont la sommes des composants d passe les 100g pour 100g
- ☐ Drop duplicate

(452724, 21)

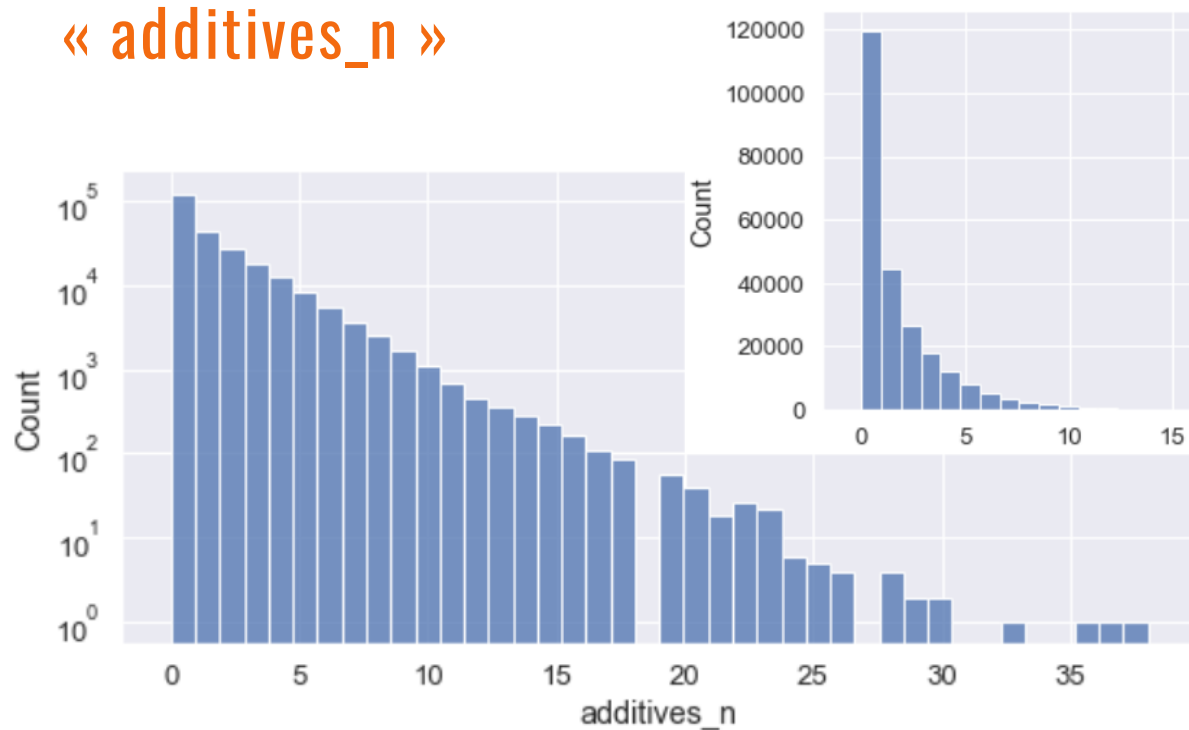


# Partie 3 : Analyse exploratoire



## Analyses univariées

« additives\_n »



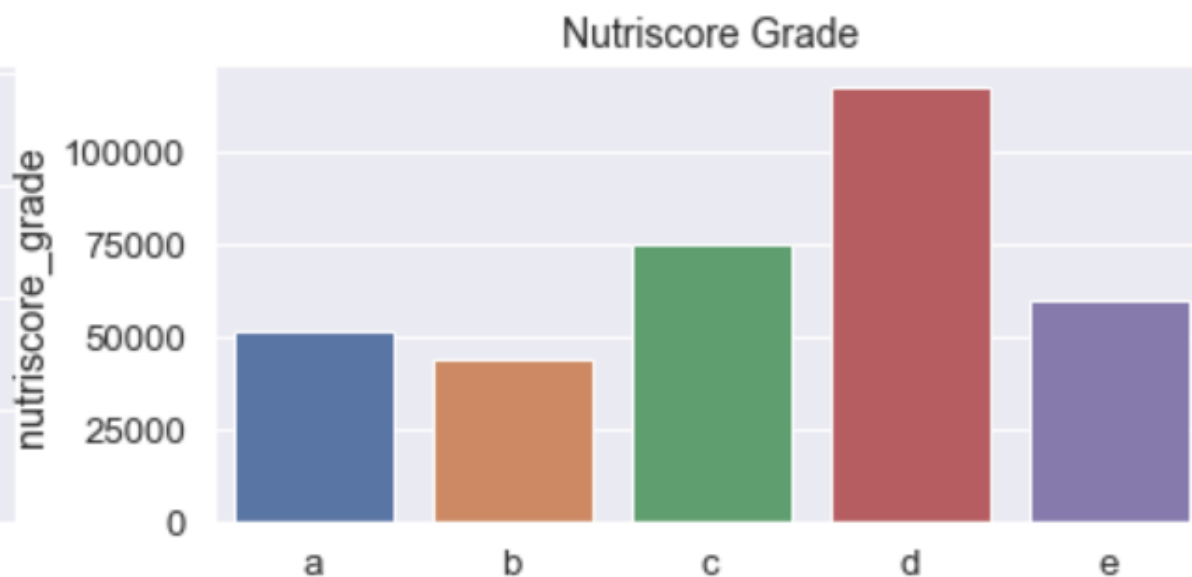
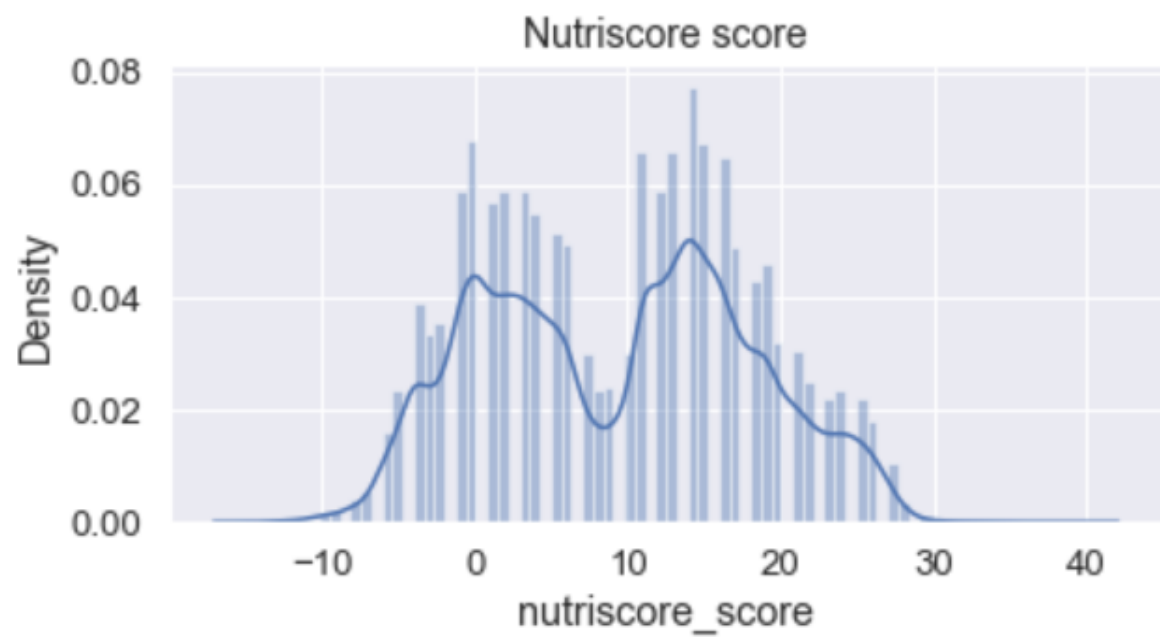
nombre d'additifs dans le produit

	additives_n	n	f	F
0	0.0	119387	0.263708	0.263708
1	1.0	44282	0.097812	0.361520
2	2.0	26805	0.059208	0.420729
3	3.0	17843	0.039413	0.460141
4	4.0	12395	0.027379	0.487520
5	5.0	8245	0.018212	0.505732
6	6.0	5507	0.012164	0.517896
7	7.0	3651	0.008065	0.525961
8	8.0	2466	0.005447	0.531408
9	9.0	1672	0.003693	0.535101
10	10.0	1103	0.002436	0.537537



## Analyses univariées

«nutriscore\_score », « nutriscore\_grade »

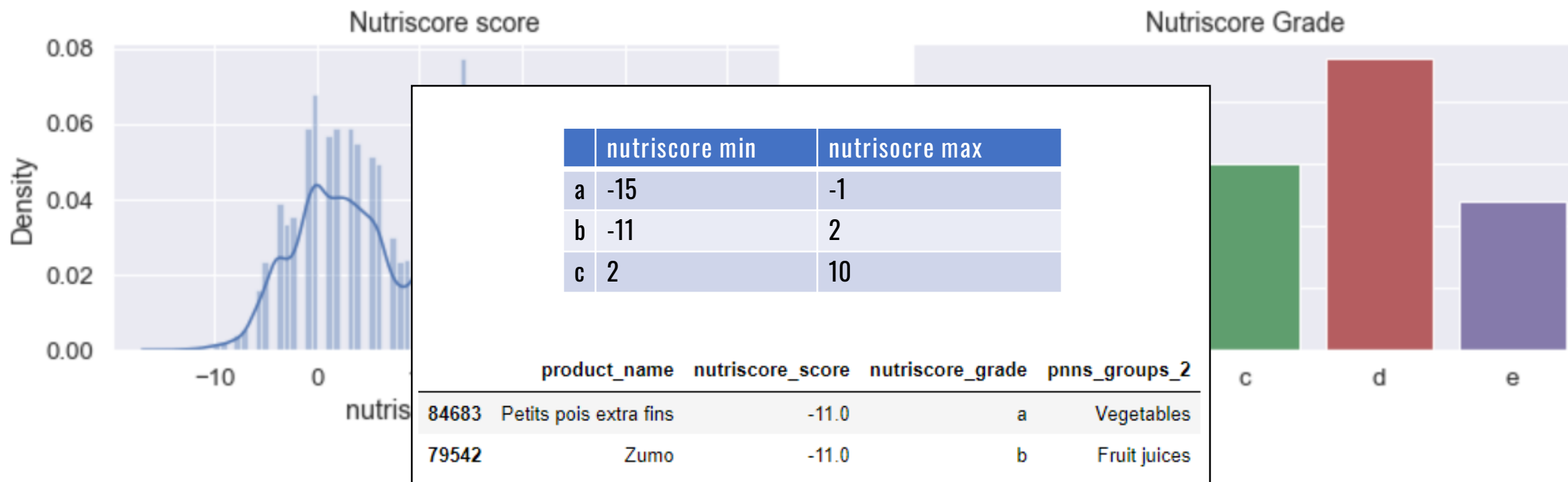






## Analyses univariées

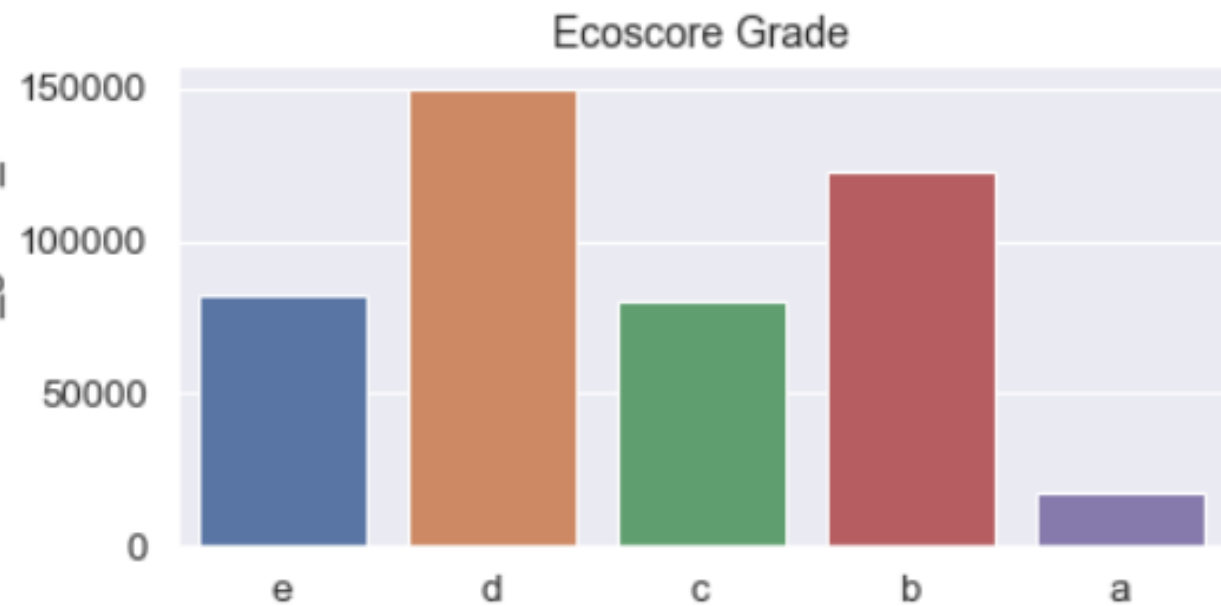
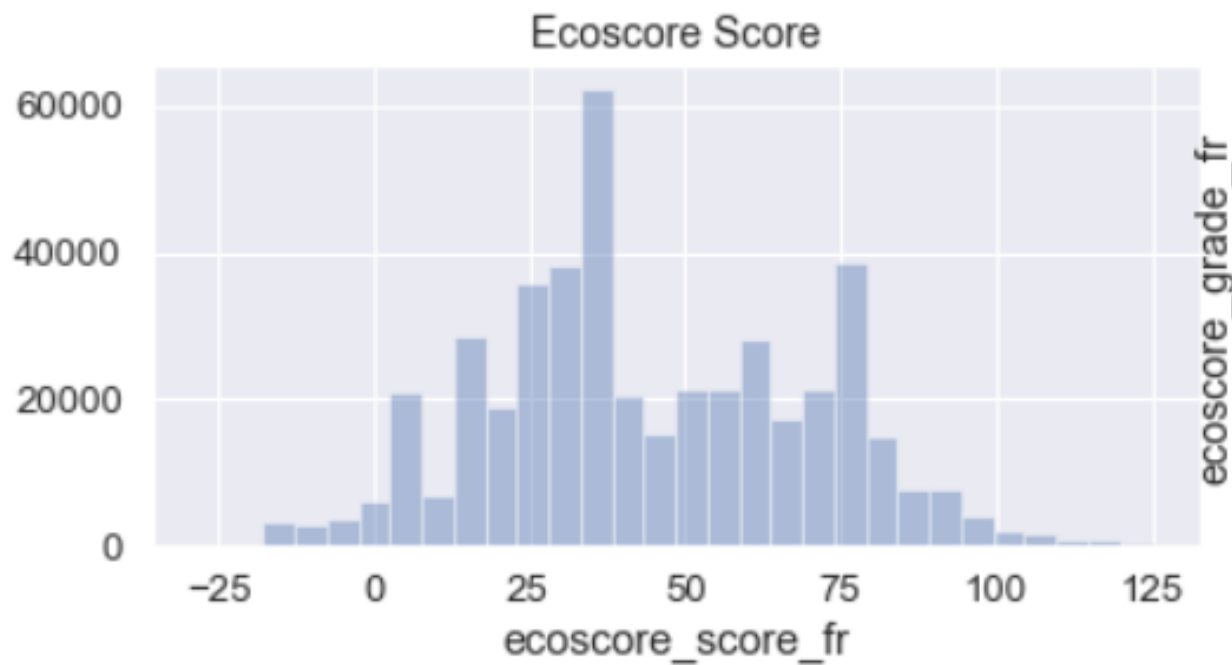
«nutriscore\_score », « nutriscore\_grade »





## Analyses univariées

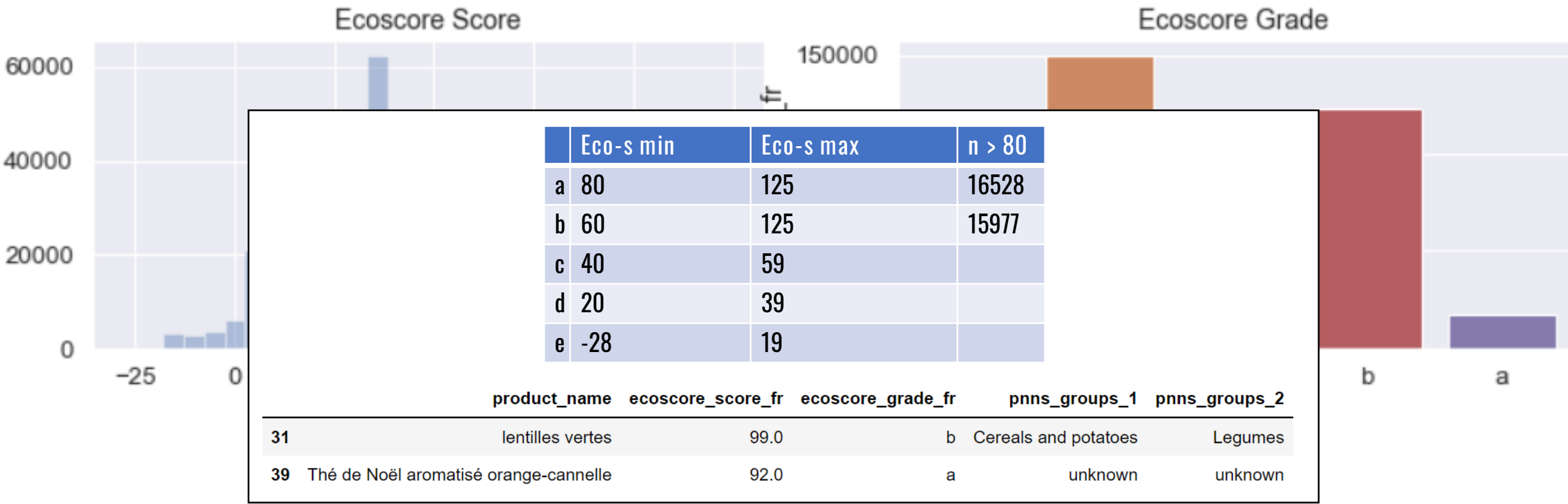
«ecoscore\_score », « ecoscore\_grade »





# Analyses univariées

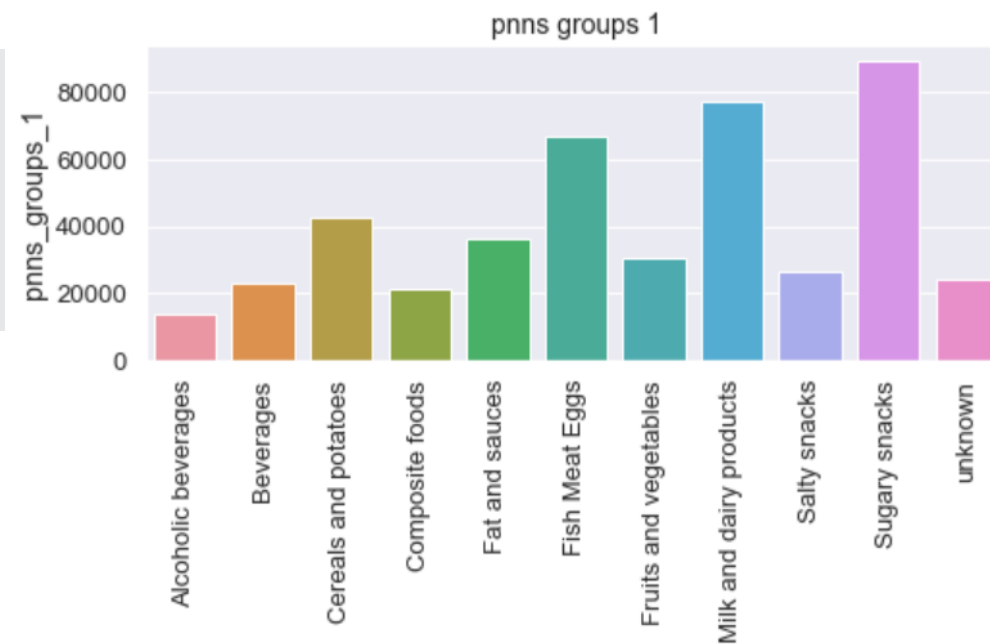
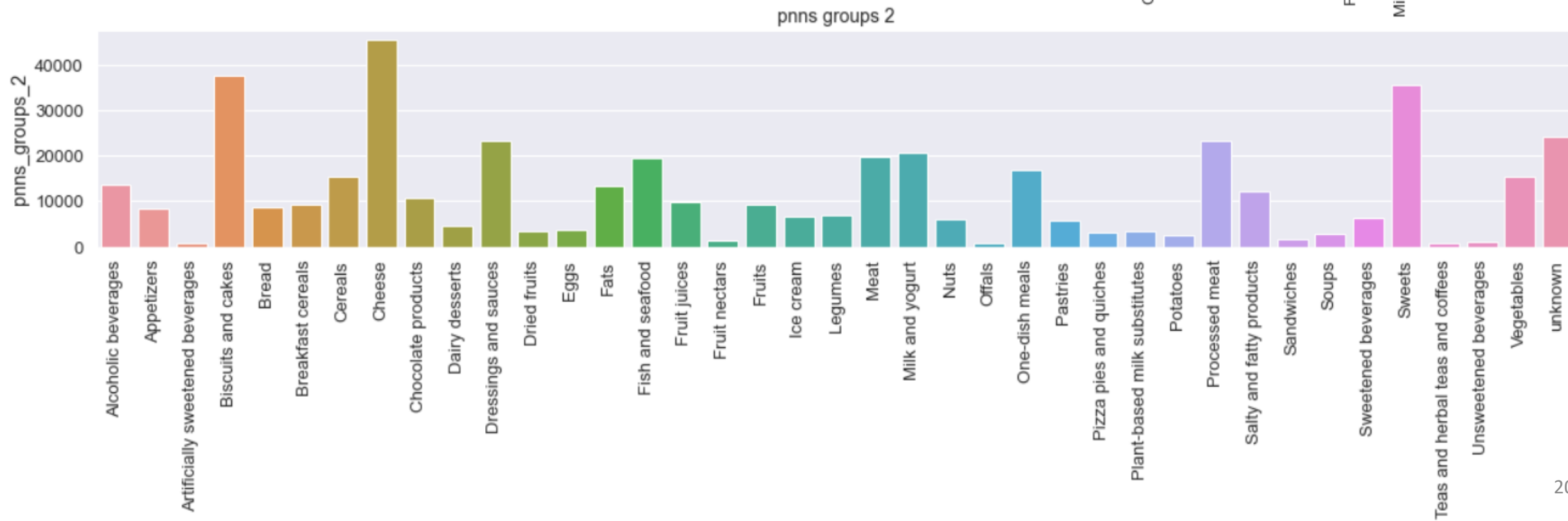
«ecoscore\_score », « ecoscore\_grade »





## Analyses univariées

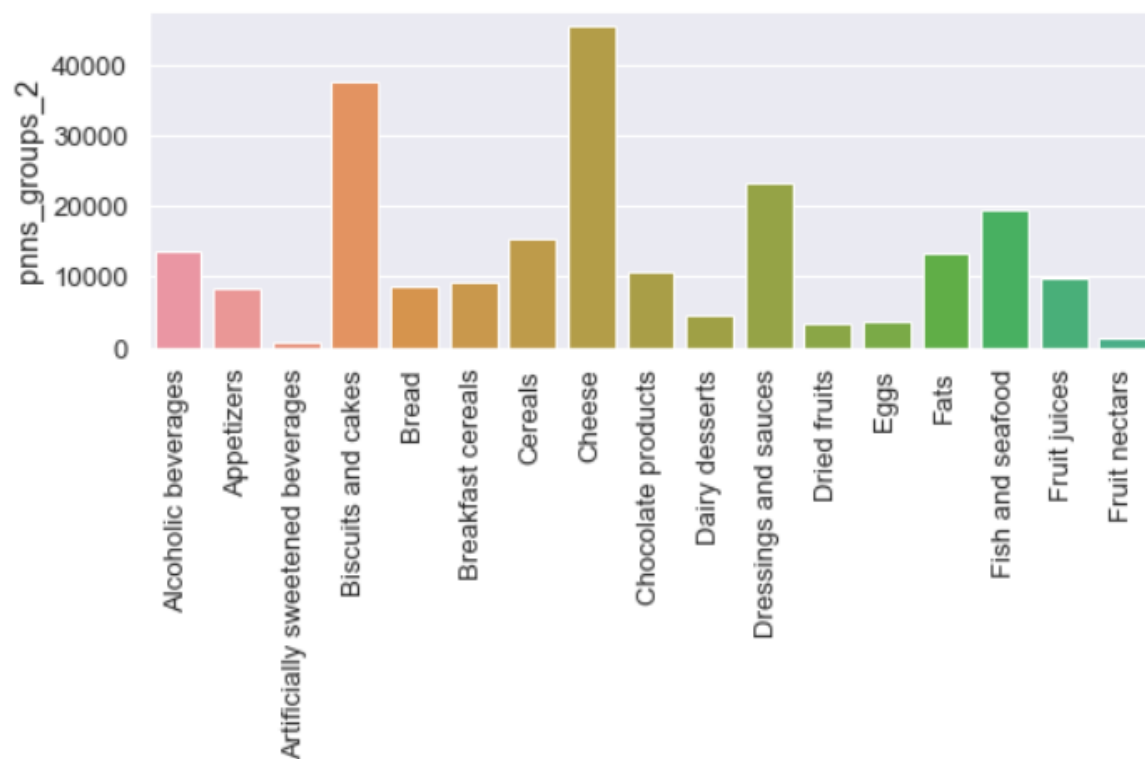
« pnns\_groups\_1 », « pnns\_groups\_2 »  
Nombre d'individus





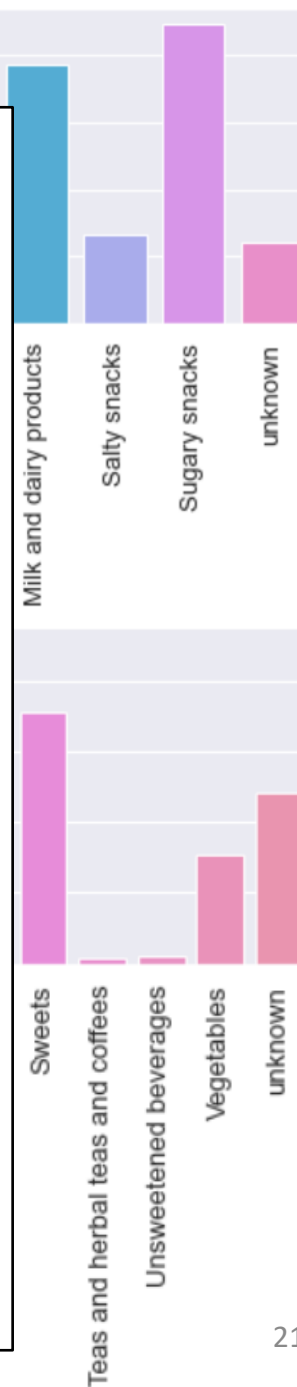
## Analyses univariées

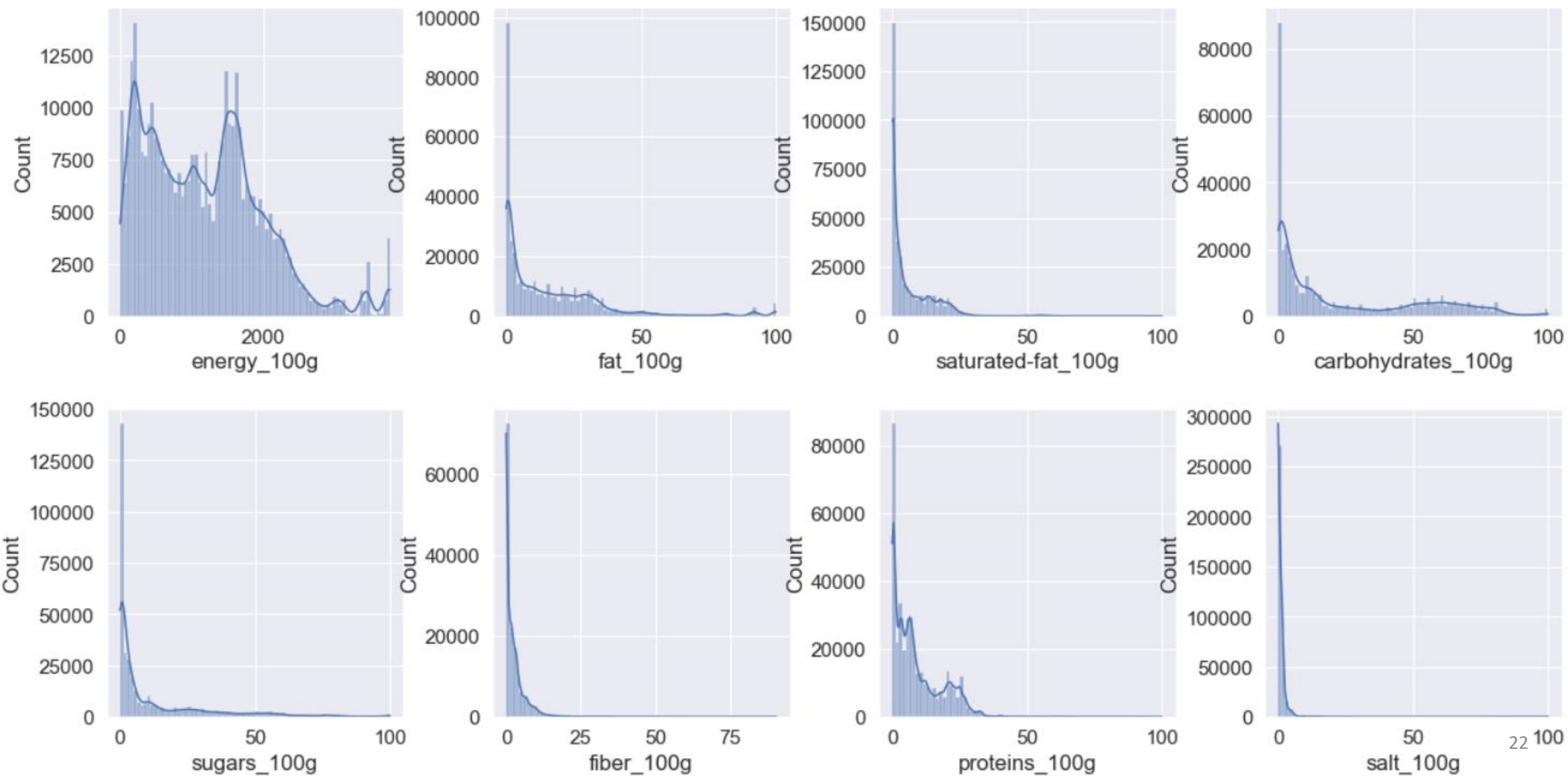
« pnns\_groups\_1 », « pnns\_groups\_2 »  
Nombre d'individus



pnns_groups_1	pnns_groups_2	
Alcoholic beverages Beverages	Alcoholic beverages	13580
	Artificially sweetened beverages	871
	Fruit juices	9722
	Fruit nectars	1392
	Plant-based milk substitutes	3379
	Sweetened beverages	6224
	Teas and herbal teas and coffees	654
	Unsweetened beverages	1066
Cereals and potatoes	Bread	8735
	Breakfast cereals	9205
	Cereals	15359
	Legumes	6898
Composite foods	Potatoes	2640
	One-dish meals	16896
	Pizza pies and quiches	3071
	Sandwiches	1555
Fat and sauces	Dressings and sauces	23193
Fish Meat Eggs	Fats	13339
	Eggs	3669
	Fish and seafood	19370
	Meat	19814
	Offals	884
	Processed meat	23125
	Dried fruits	3469
	Fruits	9259
	Soups	2823
Fruits and vegetables	Vegetables	15253
	Cheese	45302
	Dairy desserts	4602
	Ice cream	6704
	Milk and yogurt	20522
	Appetizers	8267
	Nuts	6074
	Salty and fatty products	12096
	Biscuits and cakes	37673
Salty snacks	Chocolate products	10745
	Pastries	5799
	Sweets	35414
	unknown	24081

pnns groups 1

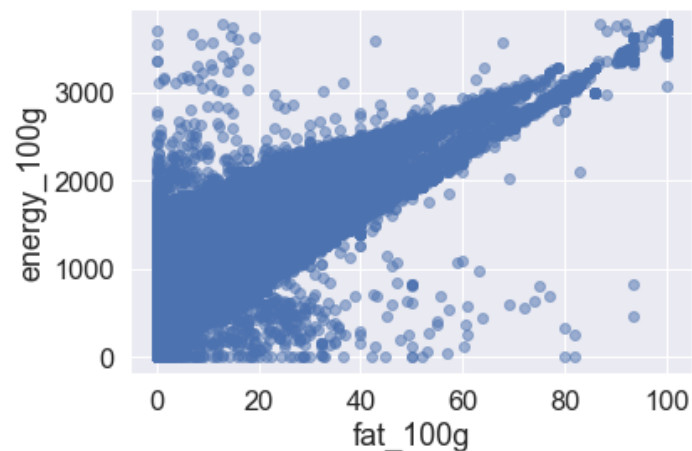




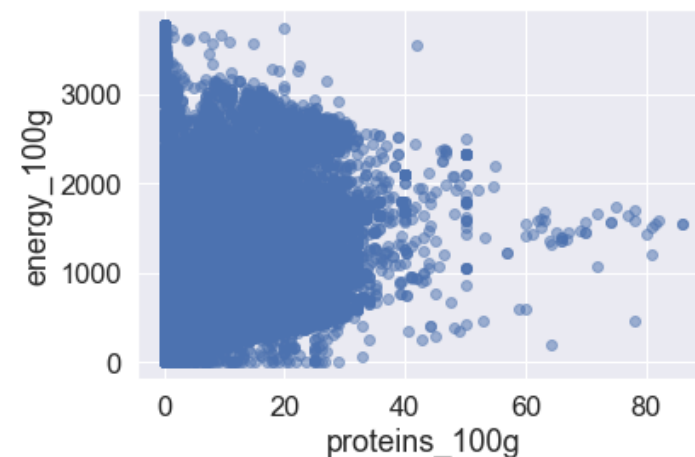


## Analyses multivariées

### Corrélation des catégorie « nutriscore et ...\_100g »



coefficient de Pearson : 0.7727228924095313  
covariance : 9814.94379179825



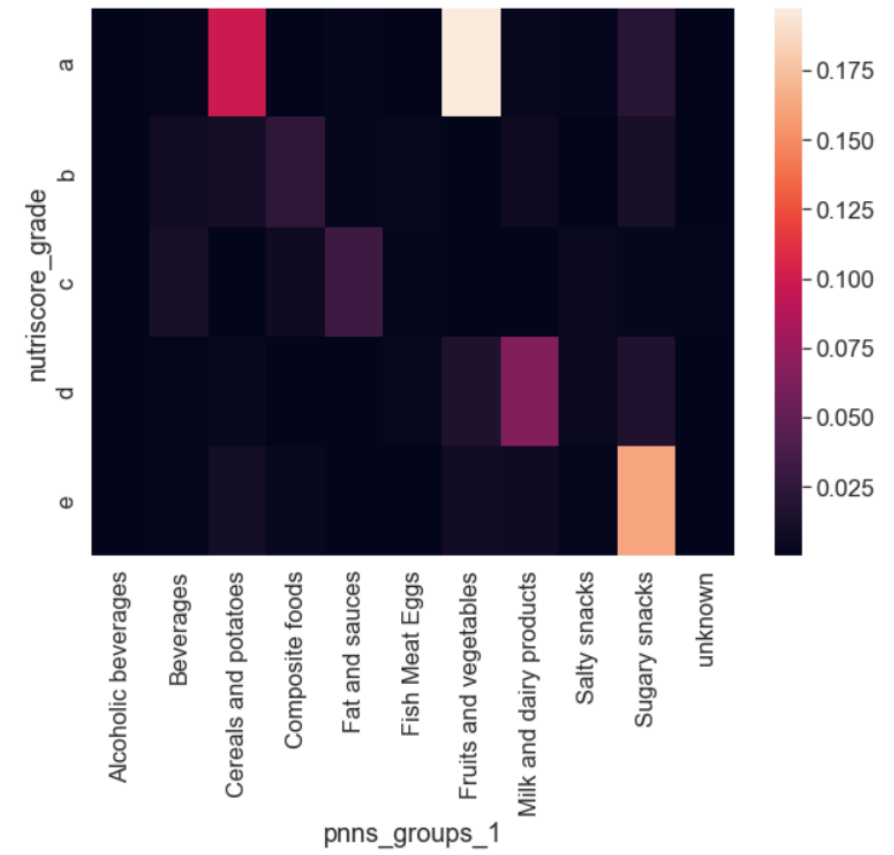
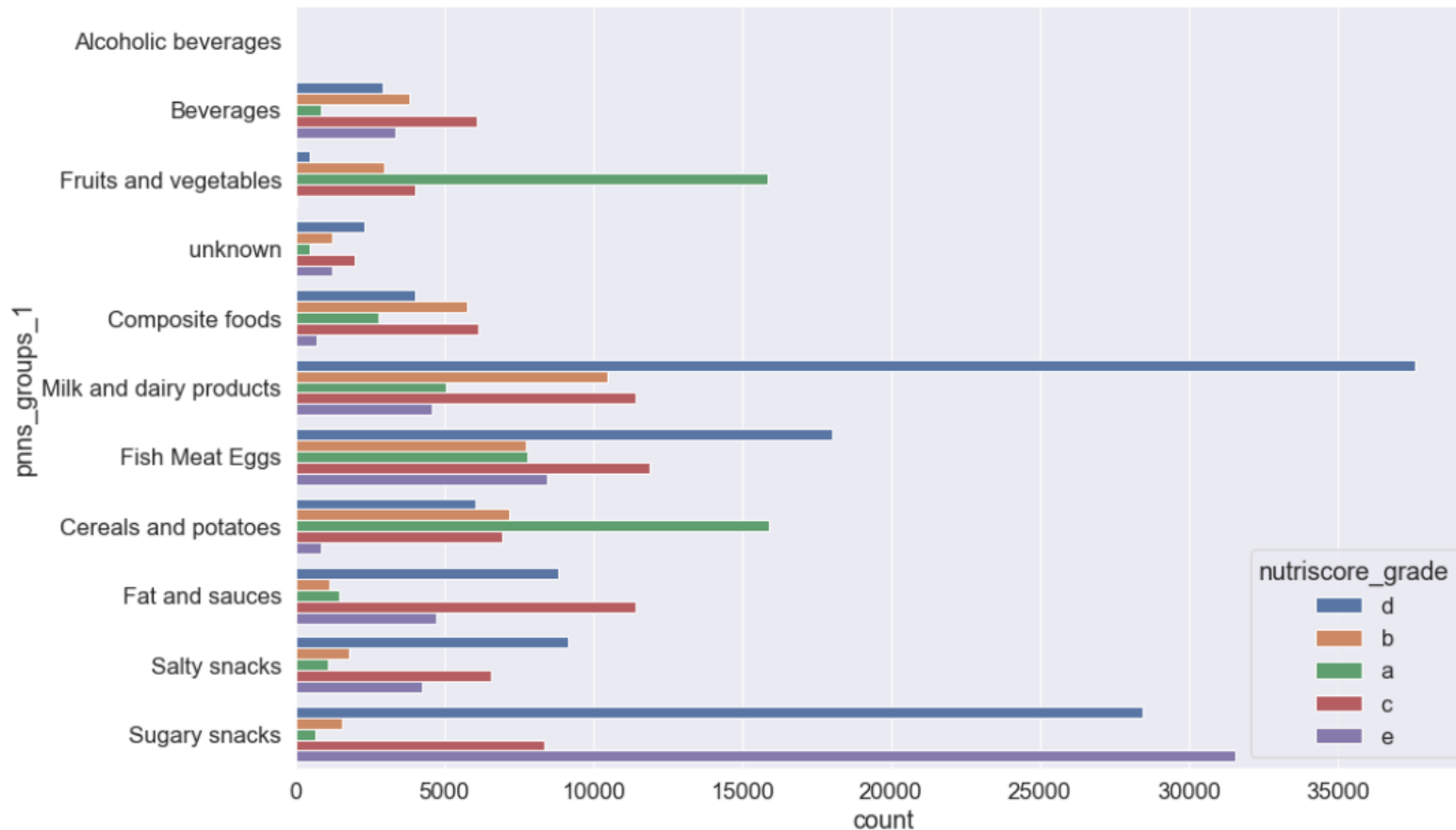
coefficient de Pearson : 0.25738134613637953  
covariance : 1516.176868031524

	product_name	energy_100g	fat_100g
3290	Pure Sesame Oil	7.14	100.000000
5645	Organic Extra Virgin Olive Oil	558.00	93.333333
19086	California Extra Virgin Olive Oil	335.00	93.333333
28228	Premium Extra Virgin Olive Oil	0.00	93.333333
37080	Extra Virgin Olive Oil	463.00	93.333333



## Analyses multivariées

### «nutriscore\_grade» vs «pnns\_1»

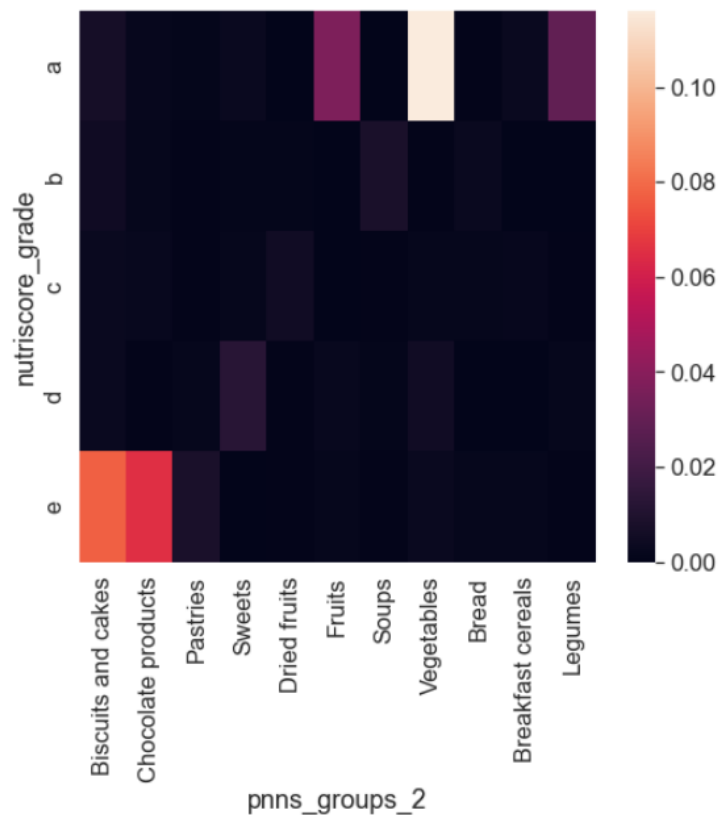




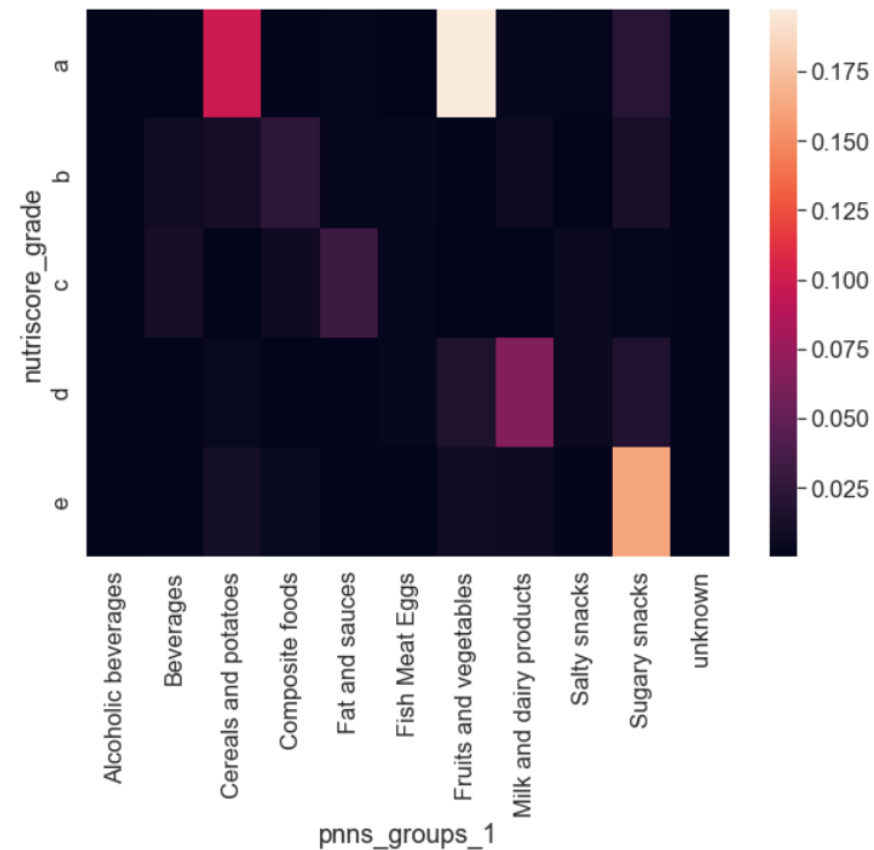


# Analyses multivariées

## «nutriscore\_grade » vs « pnns\_1 »



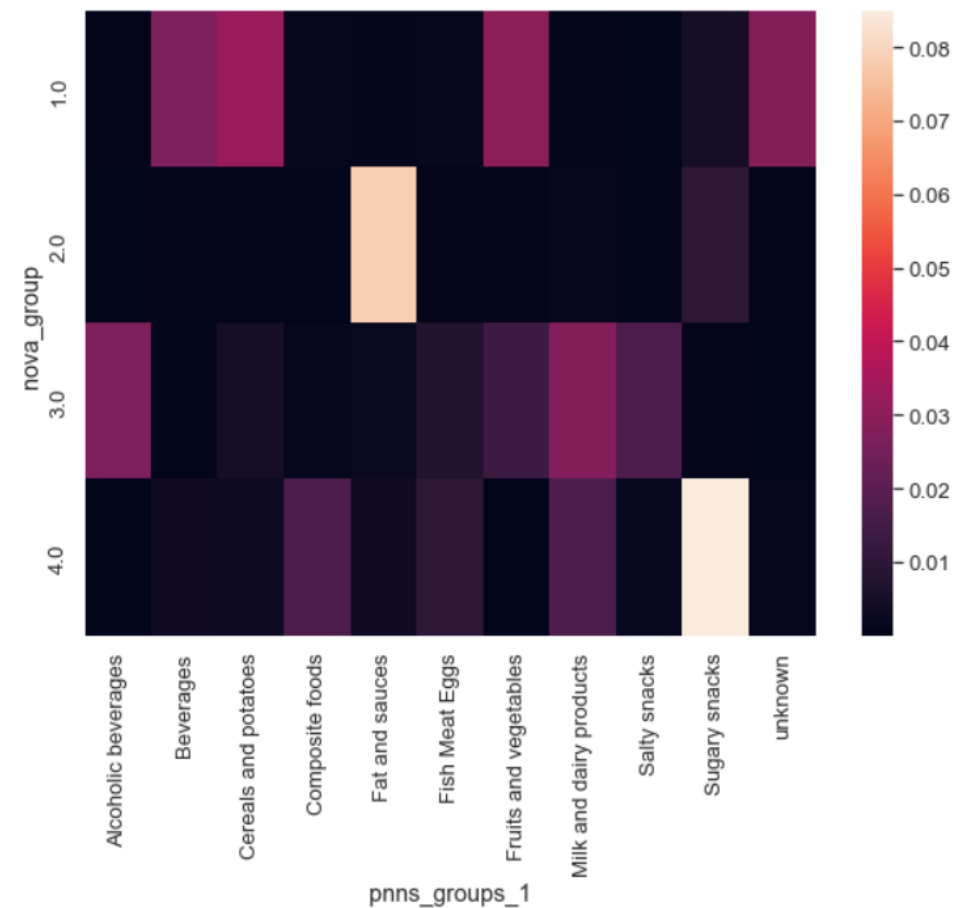
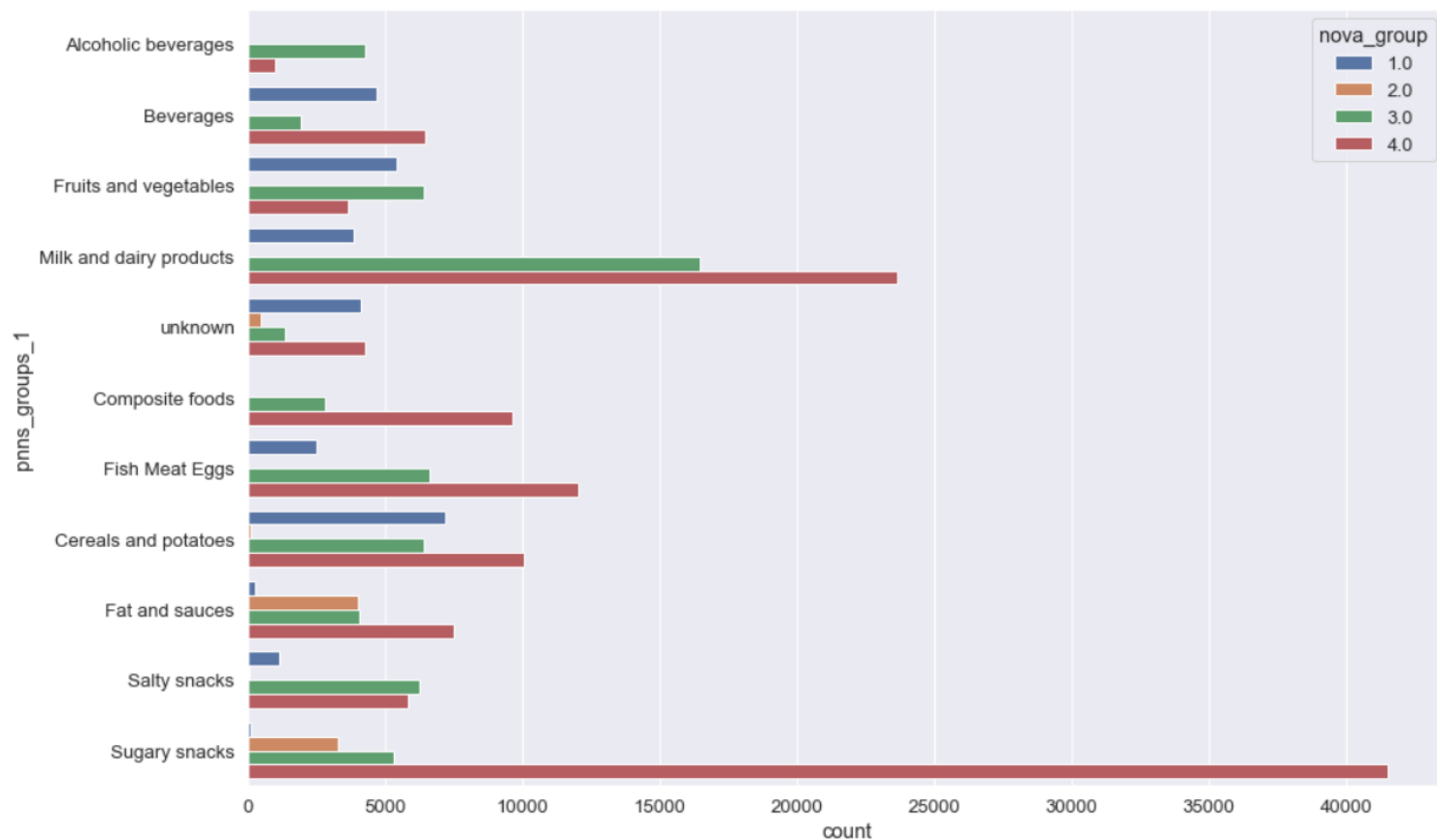
Sugary snacks	Biscuits and cakes	37673
	Chocolate products	10745
	Pastries	5799
	Sweets	35414
Fruits and vegetables	Dried fruits	3469
	Fruits	9259
	Soups	2823
	Vegetables	15253
Cereals and potatoes	Bread	8735
	Breakfast cereals	9205
	Cereals	15359
	Legumes	6898
	Potatoes	2640





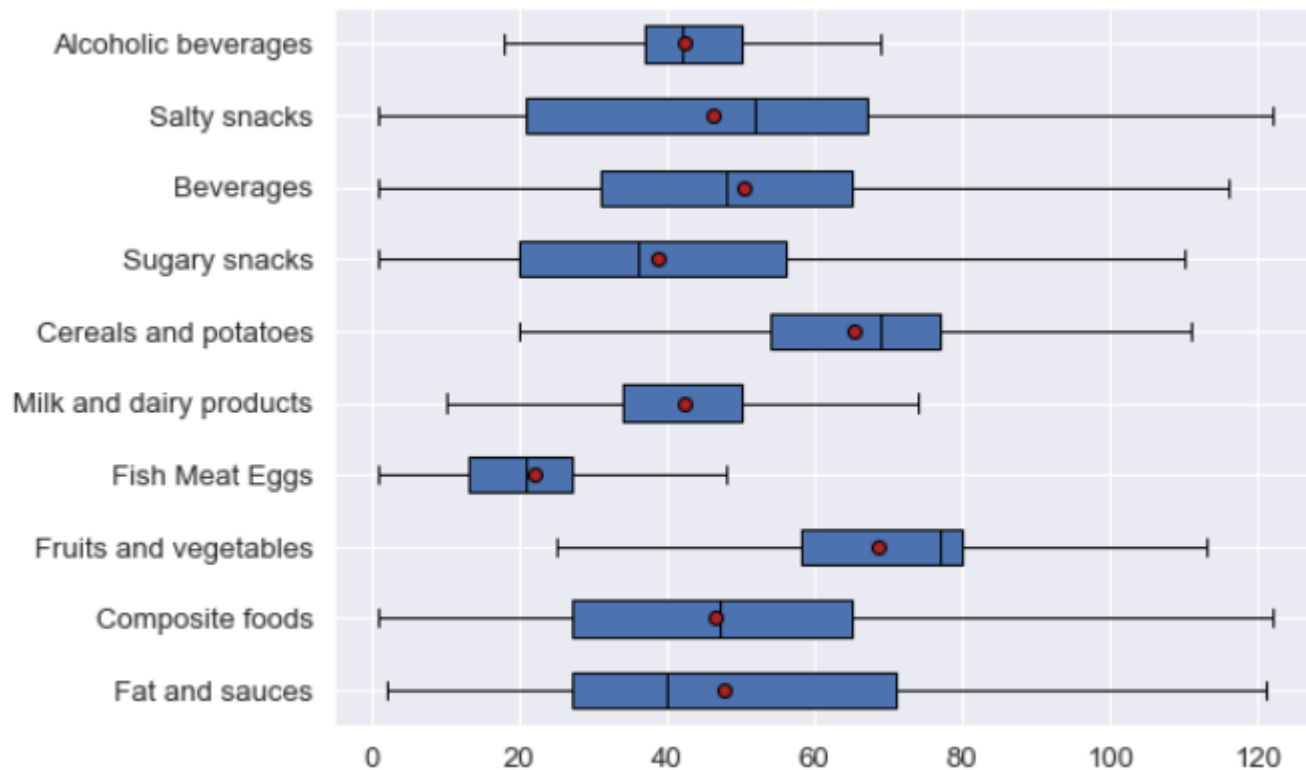
# Analyses multivariées

## «nova\_score» vs «pnns\_1»





## Analyses multivariées



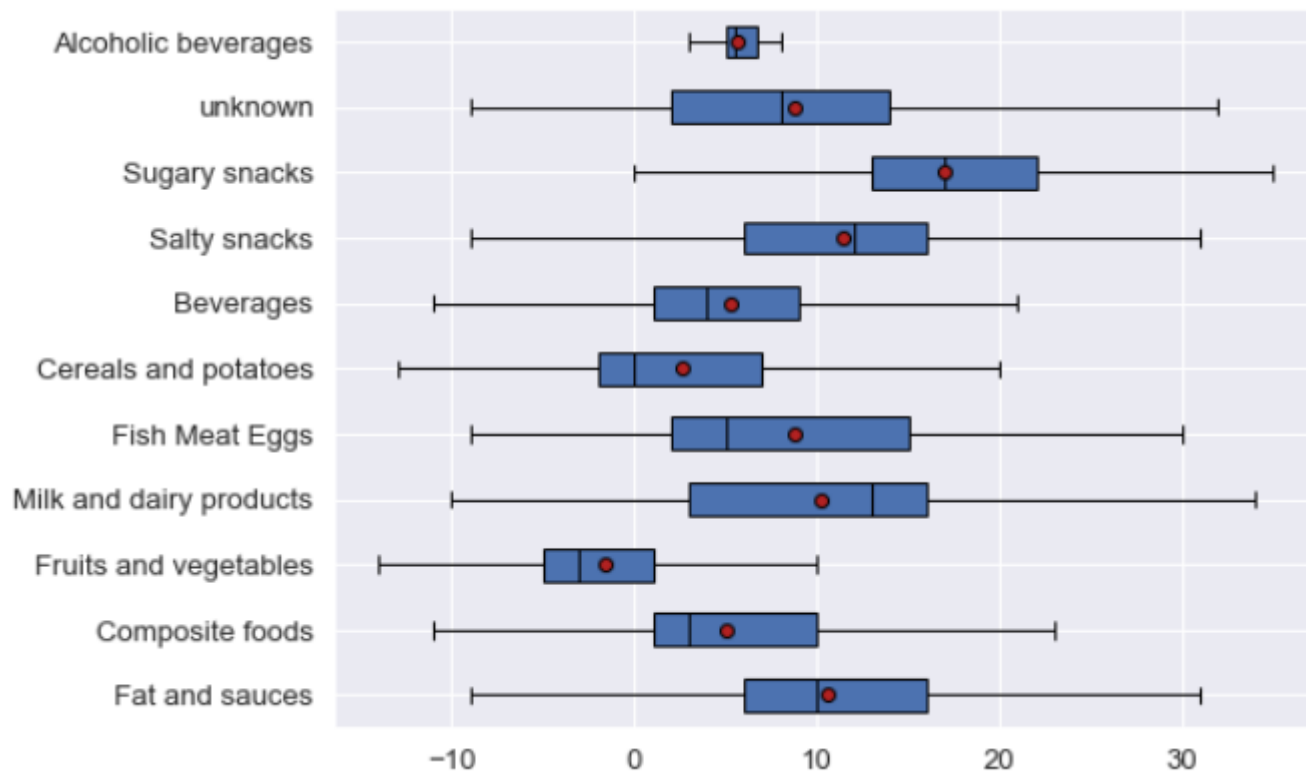
« ecoscore\_score », « pnns\_groups\_1 »

	sum_sq	df	F	PR(>F)
C(treatments)	9.029856e+07	10.0	18992.520639	0.0
Residual	2.152391e+08	452713.0	NaN	NaN

	group1	group2	Diff	Lower	Upper	q-value	p-value
39	Milk and dairy products	Alcoholic beverages	0.246681	-0.40645	0.899813	1.719239	0.9



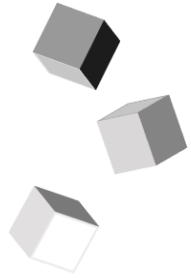
## Analyses multivariées



«nutriscore\_score », « pnns\_groups\_1 »

	sum_sq	df	F	PR(>F)
<b>C(treatments)</b>	9.482769e+06	10.0	19509.406494	0.0
<b>Residual</b>	1.694167e+07	348550.0	NaN	NaN

	group1	group2	p-value
9	Fat and sauces	Alcoholic beverages	0.788261
14	Composite foods	Beverages	0.264148
18	Composite foods	Alcoholic beverages	0.900000
26	Fruits and vegetables	Alcoholic beverages	0.281017
33	Milk and dairy products	Alcoholic beverages	0.868963
38	Fish Meat Eggs	unknown	0.900000
39	Fish Meat Eggs	Alcoholic beverages	0.900000
44	Cereals and potatoes	Alcoholic beverages	0.900000
48	Beverages	Alcoholic beverages	0.900000
51	Salty snacks	Alcoholic beverages	0.617893
54	unknown	Alcoholic beverages	0.900000



# Partie 4 : Données manquantes

k-score

Méga-score

Application



## Remplissage des données

### Méthode de remplacement des données

```
product_name      0.733560
ingredients_text   45.679708
additives_n       45.679487
nutriscore_score  23.008058
nutriscore_grade  23.008058
nova_group        50.429842
pnns_groups_1     0.000000
pnns_groups_2     0.000000
ecoscore_score_fr 0.000000
ecoscore_grade_fr 0.000000
energy_100g       15.339368
fat_100g          15.635354
saturated-fat_100g 17.146208
carbohydrates_100g 15.692342
sugars_100g       16.669759
fiber_100g        67.034661
proteins_100g     15.576378
salt_100g         17.040183
```



Meanvalue

`sklearn.impute.IterativeImputer`

`KNeighborsClassifier`

Estimation via règle de calcul

`sklearn.impute.IterativeImputer`

	mean	std
pnns_groups_2		
Alcoholic beverages	0.085493	0.467025
Appetizers	1.907947	1.864349
Artificially sweetened beverages	4.153846	1.965620
Biscuits and cakes	3.876277	2.886765
Bread	1.975615	2.003255
Breakfast cereals	0.867898	1.064692
Cereals	0.154698	0.783192
Cheese	1.060064	1.350985
Chocolate products	0.994323	1.052275
Dairy desserts	1.817036	1.935492
Dressings and sauces	1.725133	1.429991
Dried fruits	0.156529	0.543390
Eggs	0.145271	0.793507



## Remplissage des données

### Méthode de remplacement des données

product_name	0.733560
ingredients_text	45.679708
additives_n	45.679487
nutriscore_score	23.008058
nutriscore_grade	23.008058
nova_group	50.429842
pnns_groups_1	0.000000
pnns_groups_2	0.000000
ecoscore_score_fr	0.000000
ecoscore_grade_fr	0.000000
energy_100g	15.339368
fat_100g	15.635354
saturated-fat_100g	17.146208
carbohydrates_100g	15.692342
sugars_100g	16.669759
fiber_100g	67.034661
proteins_100g	15.576378
salt_100g	17.040183

Meanvalue

`sklearn.impute.IterativeImputer`

`KNeighborsClassifier`

Estimation via règle de calcul

`sklearn.impute.IterativeImputer`

Multivariate imputer that estimates each feature from all the others.

A strategy for imputing missing values by modeling each feature with missing values as a function of other features in a round-robin fashion.



## Remplissage des données

### Méthode de remplacement des données

Meanvalue

`sklearn.impute.IterativeImputer`

`KNeighborsClassifier`

Algorithme des K plus proche  
voisins

Estimation via règle de calcul

`sklearn.impute.IterativeImputer`

product_name	0.733560
ingredients_text	45.679708
additives_n	45.679487
nutriscore_score	23.008058
nutriscore_grade	23.008058
nova_group	50.429842
pnns_groups_1	0.000000
pnns_groups_2	0.000000
ecoscore_score_fr	0.000000
ecoscore_grade_fr	0.000000
energy_100g	15.339368
fat_100g	15.635354
saturated-fat_100g	17.146208
carbohydrates_100g	15.692342
sugars_100g	16.669759
fiber_100g	67.034661
proteins_100g	15.576378
salt_100g	17.040183





## Remplissage des données

### Méthode de remplacement des données

Meanvalue

`sklearn.impute.IterativeImputer`

KNeighborsClassifier

Estimation via règle de calcul

`sklearn.impute.IterativeImputer`

Estimation basée sur les valeurs  
nutritionnelles des aliments

product_name	0.733560
ingredients_text	45.679708
additives_n	45.679487
nutriscore_score	23.008058
nutriscore_grade	23.008058
nova_group	50.429842
pnns_groups_1	0.000000
pnns_groups_2	0.000000
ecoscore_score_fr	0.000000
ecoscore_grade_fr	0.000000
energy_100g	15.339368
fat_100g	15.635354
saturated-fat_100g	17.146208
carbohydrates_100g	15.692342
sugars_100g	16.669759
fiber_100g	67.034661
proteins_100g	15.576378
salt_100g	17.040183



## Remplissage des données

### Méthode de remplacement des données

product\_name 0.733560  
ingredients\_text 45.679708  
additives\_n 45.679487  
nutriscore\_score 23.008058  
nutriscore\_grade 23.008058  
nova\_group 50.429842  
pnns\_groups\_1 0.000000  
pnns\_groups\_2 0.000000  
ecoscore\_score\_fr 0.000000  
ecoscore\_grade\_fr 0.000000  
energy\_100g 15.339368  
fat\_100g 15.635354  
saturated-fat\_100g 17.146208  
carbohydrates\_100g 15.692342  
sugars\_100g 16.669759  
fiber\_100g 67.034661  
proteins\_100g 15.576378  
salt\_100g 17.040183



	energy_100g	fat_100g	carbohydrates_100g
count	452724.000000	452724.000000	452724.000000
mean	1178.044098	16.107308	24.705077
std	755.527756	18.428249	25.436517
min	-440.726660	-1.792836	-70.347807
25%	565.000000	1.900000	2.900000
50%	1169.654285	14.120000	17.500000
75%	1610.000000	21.700000	41.000000
max	4012.088556	100.000000	100.000000

product\_name 0.733560  
ingredients\_text 45.679708  
additives\_n 0.000000  
nutriscore\_score 0.000000  
nutriscore\_grade 23.008058  
nova\_group 0.000000  
pnns\_groups\_1 0.000000  
pnns\_groups\_2 0.000000  
ecoscore\_score\_fr 0.000000  
ecoscore\_grade\_fr 0.000000  
energy\_100g 0.000000  
fat\_100g 0.000000  
saturated-fat\_100g 0.000000  
carbohydrates\_100g 0.000000  
sugars\_100g 0.000000  
fiber\_100g 0.000000  
proteins\_100g 0.000000  
salt\_100g 0.000000



## k-score et Méga-score

**k\_score =**  $\Sigma$  sur les différentes catégories ‘...100g’ (energy, fat, etc... )  
par rapport aux apports journaliers recommandés (fat\_100g - 35)/10  
bonus / malus ( + saturated\_fat / - fiber)

**Méga\_score =**  $\Sigma$  sur les différents autres score (nutriscore\_score / nova\_group etc)

	k_score	mega_score
count	452724.000000	452724.000000
mean	-6.531733	5.936741
std	5.747572	32.093599
min	-24.400212	-118.960048
25%	-11.515675	-16.507641
50%	-6.585110	9.528842
75%	-2.793273	29.192871
max	14.847720	99.467778



## L'application

### Scan du code bare

#### Broccoli- Röshen

```
----
mega_score      -64.675226
nutriscore_score -7.000000
nova_group      1.000000
ecoscore_score_fr 83.000000
additives_n     0.000000
k_score         -15.324774
Name: 4421, dtype: float64
----
```

ingrédients :  
Brokkoliröschen

-----  
suggestion d'autres produit avec un meilleur score :

```
produit N° 1 :
----
Ajo morado
----
mega_score      -70.757388
nutriscore_score 7.678495
nova_group      1.000000
ecoscore_score_fr 95.000000
additives_n     0.000000
k_score         -6.564118
Name: 411503, dtype: float64
-----
```

```
produit N° 2 :
----
Ajo negro
----
mega_score      -84.181571
nutriscore_score -8.000000
nova_group      1.000000
ecoscore_score_fr 95.000000
additives_n     0.000000
k_score         -8.818429
Name: 429194, dtype: float64
-----
```

```
produit N° 3 :
----
Annalisa Canned Tomatoes
----
mega_score      -76.611394
nutriscore_score -5.000000
nova_group      1.000000
ecoscore_score_fr 97.000000
additives_n     0.000000
k_score         -15.388606
Name: 378609, dtype: float64
```

```
produit N° 1 :
----
Betteraves Rouges
----
mega_score      -83.441455
nutriscore_score -7.000000
nova_group      1.000000
ecoscore_score_fr 101.000000
additives_n     1.000000
k_score         -13.558545
Name: 94016, dtype: float64
-----
```

```
produit N° 2 :
----
Chopped Tomatoes
----
mega_score      -67.670297
nutriscore_score -5.000000
nova_group      1.000000
ecoscore_score_fr 89.000000
additives_n     1.000000
k_score         -15.329703
Name: 327748, dtype: float64
-----
```

```
produit N° 3 :
----
Bio Zwiebeln
----
mega_score      -92.757388
nutriscore_score 7.678495
nova_group      1.000000
ecoscore_score_fr 117.000000
additives_n     0.000000
k_score         -6.564118
Name: 301555, dtype: float64
```



## L'application

### 2 autres exemples

#### Lentilles cuisinées aux carottes et oignons

```
----
mega_score      -101.058968
nutriscore_score -4.000000
nova_group      1.000000
ecoscore_score_fr 123.000000
additives_n      2.000000
k_score         -13.941032
Name: 125260, dtype: float64
----
ingrédients :
ingrédients à renseigner
-----
suggestion d'autres produit avec un meilleur score :
produit N° 1 :
```

#### Couscous biologique de blé khorazan

```
----
mega_score      -114.430635
nutriscore_score -5.000000
nova_group      1.000000
ecoscore_score_fr 123.000000
additives_n      0.000000
k_score         -3.569365
Name: 183210, dtype: float64
-----
pas d'autre meilleur produit
pas d'autre meilleur produit
```

#### Vinagre balsámico de Módena

```
----
mega_score      -14.367607
nutriscore_score 8.229934
nova_group      4.000000
ecoscore_score_fr 71.000000
additives_n      2.000000
k_score         -6.402460
Name: 425260, dtype: float64
----
ingrédients :
ingrédients à renseigner
-----
suggestion d'autres produit avec un meilleur score :
produit N° 1 :
```

#### Peppery Moroccan Ketchup

```
----
mega_score      -24.90817
nutriscore_score 7.00000
nova_group      3.00000
ecoscore_score_fr 74.00000
additives_n      0.00000
k_score         -12.09183
Name: 43227, dtype: float64
```

#### produit N° 2 :

```
----
Knoblauch Sauce
----
mega_score      -27.831084
nutriscore_score 7.000000
nova_group      4.000000
ecoscore_score_fr 83.000000
additives_n      2.000000
k_score         -6.168916
Name: 314984, dtype: float64
-----
```

#### produit N° 3 :

```
----
Purée de tomates homogeneizado
----
mega_score      -53.423089
nutriscore_score 8.103689
nova_group      1.000000
ecoscore_score_fr 78.000000
additives_n      0.000000
k_score         -6.473222
Name: 371904, dtype: float64
```



## Conclusion

### Conclusion

Analyses univariées  
Analyses multivariées  
Prémises d'une application

### Limitation de l'étude

Beaucoup de valeurs manquantes  
Encore des valeurs aberrantes  
Limitations métier sur le k et Méga Score  
défauts techniques

### Perspectives

Une applications avec choix pour l'utilisateur  
De nouveaux scores



Fin de la présentation

Merci de m'avoir écouté !





## Remplissage des données

### Méthode de remplacement des données `sklearn.impute.IterativeImputer`

La supposition de base est que les données manquantes peuvent être expliquées par les autres colonnes et qu'on peut faire une estimation éclairée

Étape 1 : on remplit tout avec le mean de la colonne.

Étape 2 : on remplace les valeurs de la colonne de gauche et ainsi de suite jusqu'à la colonne de droite (3 étapes pour 3 colonnes)

Étape 3 : On soustrait les 2 « derniers » datasets

Etc etc... On s'arrête après 10 itérations

#### iteration 2

age	experience	salary
25	1.8538	50
27	3	72.7748
29	5	110
31	7	140
33	9	170
36.2532	11	200

First

#### iteration 3

age	experience	salary
25	0.9172	50
27	3	80.7385
29	5	110
31	7	140
33	9	170
34.8732	11	200

Second

#### iteration 4

age	experience	salary
25	1.0015	50
27	3	79.9876
29	5	110
31	7	140
33	9	170
35.0019	11	200

Third

After all imputations

age	experience	salary
25	0.9172	50
27	3	80.7385
29	5	110
31	7	140
33	9	170
34.8732	11	200

Second

After all imputations

age	experience	salary
25	1.0015	50
27	3	79.9876
29	5	110
31	7	140
33	9	170
35.0019	11	200

Third

After all imputations

age	experience	salary
25	0.9999	50
27	3	80.0007
29	5	110
31	7	140
33	9	170
34.9998	11	200

Fourth

After Second - First

age	experience	salary
0	0.9366	0
0	0	7.9637
0	0	0
0	0	0
0	0	0
1.38	0	0

Difference matrix

After Third - Second

age	experience	salary
0	0.0842	0
0	0	0.751
0	0	0
0	0	0
0	0	0
0.1287	0	0

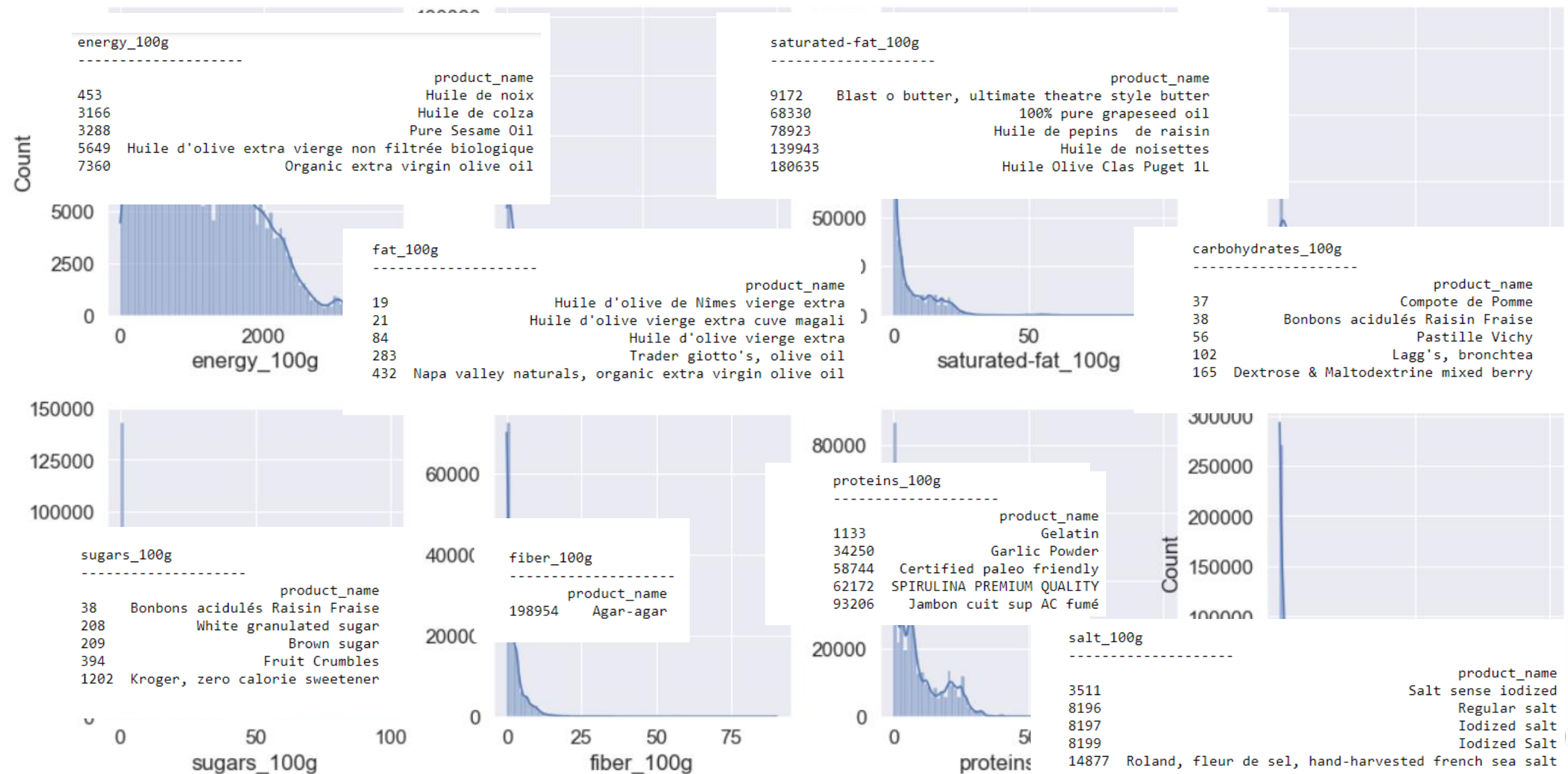
Difference matrix

After Fourth - Third

age	experience	salary
0	0.0016	0
0	0	0.0131
0	0	0
0	0	0
0	0	0
0.002	0	0

Difference matrix







## k-score et Méga-score

$$\text{k\_score} = \frac{(\text{energy\_100} - 2403.5)}{378} + \frac{(\text{'fat\_100g'} - 35)}{10} + \frac{(\text{'saturated-fat\_100g'})}{10} + \frac{(\text{carbohydrates\_100g} - 53)}{10} + \frac{(\text{sugars\_100g})}{10} + \frac{\text{salt\_100g}}{10} + \frac{(\text{'proteins\_100g'} - 12)}{10} - \frac{\text{fiber\_100g}}{10}$$

$$\text{Méga\_score} = \text{nutriscore\_score} + \text{nova\_group} * 10 + \text{'additives\_n'} - \text{écoscore\_score\_fr} - \text{k\_score}$$



## k-score et Méga-score

		k_score		mega_score	
		mean	std	mean	std
pnnns_groups_1	pnnns_groups_2				
Alcoholic beverages	Alcoholic beverages	-7.471919	2.473715	6.695231	13.553697
Beverages	Artificially sweetened beverages	-14.224208	3.300956	11.595462	19.434894
	Fruit juices	-13.034468	2.379500	-12.588925	14.911204
	Fruit nectars	-13.148338	2.351942	15.309572	17.265314
	Plant-based milk substitutes	-13.331134	2.797635	-0.713888	31.139878
	Sweetened beverages	-9.016513	7.004630	-1.668555	25.950041
	Teas and herbal teas and coffees	-13.930828	3.656838	-23.710676	21.383530
	Unsweetened beverages	-12.262877	4.024396	-27.722421	28.418722
Cereals and potatoes	Bread	-5.460284	2.736713	-20.101231	20.543182
	Breakfast cereals	-2.694466	2.551121	-13.707322	25.512619
	Cereals	-4.434512	2.972469	-40.348471	27.453369
	Legumes	-9.633180	5.043828	-30.541920	26.920877
	Potatoes	-10.293894	3.146938	-30.952647	19.316018
Composite foods	One-dish meals	-10.664973	3.378885	10.472520	29.435423
	Pizza pies and quiches	-7.326599	3.126628	13.026982	18.760789
	Sandwiches	-8.303774	1.597796	26.665906	22.498428
Fat and sauces	Dressings and sauces	-8.965931	4.700407	-2.467008	26.640280
	Fats	2.537131	4.217577	0.489337	17.223873

Fish Meat Eggs	Eggs	-10.375892	3.036371	-15.792344	17.155765
	Fish and seafood	-10.267857	2.267553	28.319594	21.621744
	Meat	-9.709818	2.674299	28.852349	19.030149
	Offals	-9.197778	1.815217	39.810246	17.898896
Fruits and vegetables	Processed meat	-8.770102	2.970108	52.239887	15.285763
	Dried fruits	-2.338639	3.323005	3.927407	23.729642
	Fruits	-9.221205	5.048805	-34.851242	27.641035
	Soups	-13.980892	2.791174	-31.092377	13.157553
Milk and dairy products	Vegetables	-12.978325	3.791422	-34.019588	27.392035
	Cheese	-5.963484	2.454805	20.795047	12.172390
	Dairy desserts	-10.470221	3.069992	-0.716787	24.612471
	Ice cream	-7.041939	2.651016	2.760947	11.875592
Salty snacks	Milk and yogurt	-12.056870	3.023035	-1.919599	22.294138
	Appetizers	-2.878896	3.280413	-5.929949	18.431676
	Nuts	-1.097395	3.082015	13.828993	20.384899
	Salty and fatty products	-7.289247	4.102381	8.693991	30.565605
Sugary snacks	Biscuits and cakes	-0.070130	3.395985	20.943153	17.956378
	Chocolate products	3.549001	3.178363	46.828448	13.245075
	Pastries	-3.283363	2.439530	20.143344	12.532068
	Sweets	-1.779095	5.277827	7.743143	30.659600
unknown	unknown	-8.368471	5.427680	-15.748180	32.993252



## Backup slides

# Légende des données (bdd complète):

stores	87.459259	ingredients_from_palm_oil	100.000000
countries	0.307002	ingredients_from_palm_oil_tags	99.252748
countries_tags	0.307211	ingredients_that_may_be_from_palm_oil_n	63.014094
countries_en	0.307211	ingredients_that_may_be_from_palm_oil	100.000000
ingredients_text	63.014146	ingredients_that_may_be_from_palm_oil_tags	97.796821
allergens	90.224103	nutriscore_score	63.431667
allergens_en	100.000000	nutriscore_grade	63.431667
traces	94.810210	nova_group	67.690389
traces_tags	93.564425	pnns_groups_1	0.026274
traces_en	93.564425	pnns_groups_2	0.026170
serving_size	74.368271	states	0.000000
serving_quantity	74.603436	states_tags	0.000000
no_nutriments	100.000000	states_en	0.000000
additives_n	63.014094	brand_owner	84.890117
additives	100.000000	ecoscore_score_fr	75.871116
additives_tags	78.625534	ecoscore_grade_fr	75.871116
additives_en	78.625534	main_category	52.830998
ingredients_from_palm_oil_n	63.014094	main_category_en	52.830998
		image_url	23.825948

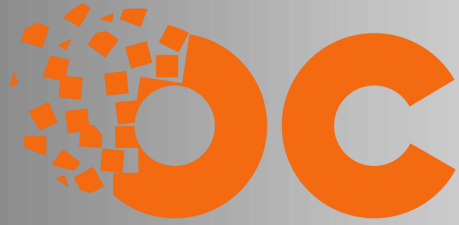
% de données manquantes	
code	0.000000
url	0.000000
creator	0.000209
created_t	0.000000
created_datetime	0.000000
last_modified_t	0.000000
last_modified_datetime	0.000000
product_name	4.169265
abbreviated_product_name	99.633360
generic_name	93.772638
quantity	75.140807
packaging	84.390594
packaging_tags	84.392262
packaging_text	99.704884
brands	48.403787

brands_tags	48.407124
categories	52.830685
categories_tags	52.830998
categories_en	52.830998
origins	95.726942
origins_tags	95.733719
origins_en	95.733719
manufacturing_places	93.698664
manufacturing_places_tags	93.701531
labels	78.623344
labels_tags	78.623761
labels_en	78.623761
emb_codes	93.806784
emb_codes_tags	93.808974
first_packaging_code_geo	96.243356
cities	100.000000
cities_tags	95.952097
purchase_places	91.864106

image_small_url	23.825948	-arachidic-acid_100g	99.994370
image_ingredients_url	63.753944	-behenic-acid_100g	99.997863
image_ingredients_small_url	63.753944	-lignoceric-acid_100g	99.999270
image_nutrition_url	48.846644	-cerotic-acid_100g	99.999166
image_nutrition_small_url	48.846644	-montanic-acid_100g	99.998853
energy-kj_100g	92.815609	-melissic-acid_100g	99.999114
energy-kcal_100g	23.275492	monounsaturated-fat_100g	97.503842
energy_100g	20.563353	polyunsaturated-fat_100g	97.505771
energy-from-fat_100g	99.949328	omega-3-fat_100g	99.893235
fat_100g	20.976703	-alpha-linolenic-acid_100g	99.960432
saturated-fat_100g	23.210275	-eicosapentaenoic-acid_100g	99.993588
-butyric-acid_100g	99.998332	-docosahexaenoic-acid_100g	99.988479
-caproic-acid_100g	99.999687	omega-6-fat_100g	99.971745
-caprylic-acid_100g	99.999635	-linoleic-acid_100g	99.974925
-capric-acid_100g	99.999427	-arachidonic-acid_100g	99.993379
-lauric-acid_100g	99.998853	-gamma-linolenic-acid_100g	99.999114
-myristic-acid_100g	99.999166	-dihomo-gamma-linolenic-acid_100g	99.999374
-palmitic-acid_100g	99.998280	omega-9-fat_100g	99.994318
-stearic-acid_100g	99.999166	-oleic-acid_100g	99.996507

-elaidic-acid_100g	99.999896	proteins_100g	20.913729
-gondoic-acid_100g	99.997654	casein_100g	99.997081
-mead-acid_100g	99.999114	serum-proteins_100g	99.996559
-erucic-acid_100g	99.999531	nucleotides_100g	99.997706
-nervonic-acid_100g	99.999218	salt_100g	24.773228
trans-fat_100g	86.194133	sodium_100g	24.773385
cholesterol_100g	85.987849	alcohol_100g	98.858686
carbohydrates_100g	20.998286	vitamin-a_100g	88.919044
sugars_100g	21.824621	beta-carotene_100g	99.995725
-sucrose_100g	99.991398	vitamin-d_100g	99.483325
-glucose_100g	99.994891	vitamin-e_100g	99.813057
-fructose_100g	99.995829	vitamin-k_100g	99.935826
-lactose_100g	99.954437	vitamin-c_100g	88.516537
-maltose_100g	99.996403	vitamin-b1_100g	98.762869
-maltodextrins_100g	99.991033	vitamin-b2_100g	98.811924
starch_100g	99.972683	vitamin-b2_100g	98.756665
polyols_100g	99.808678	vitamin-b6_100g	99.160736
fiber_100g	74.775887	vitamin-b9_100g	99.460700
-soluble-fiber_100g	99.811805	folates_100g	99.558134
-insoluble-fiber_100g	99.825933		

caffeine_100g	99.975029	vitamin-b12_100g	99.327661
taurine_100g	99.989574	biotin_100g	99.937703
ph_100g	99.989522	pantothenic-acid_100g	99.677828
fruits-vegetables-nuts_100g	99.541608	silica_100g	99.991920
fruits-vegetables-nuts-dried_100g	99.980503	bicarbonate_100g	99.977792
fruits-vegetables-nuts-estimate_100g	99.383285	potassium_100g	95.195097
collagen-meat-protein-ratio_100g	99.983631	chloride_100g	99.958347
cocoa_100g	99.680121	calcium_100g	85.893908
chlorophyl_100g	99.999739	phosphorus_100g	99.274070
carbon-footprint_100g	99.975185	iron_100g	86.184072
carbon-footprint-from-meat-or-fish_100g	99.387925	magnesium_100g	99.227829
nutrition-score-fr_100g	63.431459	zinc_100g	99.464245
nutrition-score-uk_100g	99.999583	copper_100g	99.774219
glycemic-index_100g	99.999791	manganese_100g	99.786574
water-hardness_100g	99.999948	fluoride_100g	99.979617
choline_100g	99.996768	selenium_100g	99.867065
phyloquinone_100g	99.910751	chromium_100g	99.989052
beta-glucan_100g	99.998071	molybdenum_100g	99.983735
inositol_100g	99.996455	iodine_100g	99.880515
carnitine_100g	99.997967		



## Backup slides

Quelles sont les modalités de calcul du nutri score ?

Le Nutri Score est calculé en fonction de plusieurs composantes :

- le nombre de calories contenues dans le produit fini,
- la quantité d'acides gras,
- la dose de sucres,
- la proportion de fibres...

Pour l'attribution des points d'un aliment, il faut évaluer le rapport entre ses composants nocifs (les calories, le sel, les sucres...) et ses véritables apports nutritionnels (fruits, légumes, fibres, protéines).

### CALCUL DU NUTRI-SCORE

01

Attribution des points à différents facteurs nutritionnels



#### Nutriments à limiter

Points N	Sécher pour les boissons		Sécher pour les matières grasses		Sécher pour les matières grasses	
Points	Energie (kJ)	Sucres (g)	Energie (kJ)	Sucres (g)	Grasses saturées (g)	Sodium (mg)
0	≤ 335	≤ 4,5	≤ 0	≤ 0	≤ 1	≤ 10
1	> 335	> 4,5	≤ 30	≤ 1,5	> 1	< 16
2	> 670	> 9	≤ 60	≤ 3	> 2	< 22
3	> 1005	> 13,5	≤ 90	≤ 4,5	> 3	< 28
4	> 1340	> 18	≤ 120	≤ 6	> 4	< 34
5	> 1675	> 22,5	≤ 150	≤ 7,5	> 5	< 40
6	> 2010	> 27	≤ 180	≤ 9	> 6	< 46
7	> 2345	> 31	≤ 210	≤ 10,5	> 7	< 52
8	> 2680	> 36	≤ 240	≤ 12	> 8	< 58
9	> 3015	> 40	≤ 270	≤ 13,5	> 9	< 64
10	> 3350	> 45	> 270	> 13,5	> 10	> 64
Gamme (points)	0 à 10	0 à 10	0 à 10	0 à 10	0 à 10	0 à 10
Total	Somme des points pour l'énergie, les sucres, les graisses saturées et le sodium					

#### Nutriments, aliments à encourager

Points P	Sécher pour les boissons		Sécher pour les matières grasses	
Points	Fruits, légumes (%)	Fibres (g)	Protéines (g)	
0	≤ 40	≤ 0,7	≤ 1,6	
1	> 40	> 0,7	> 1,6	
2	> 60	> 1,4	> 3,2	
3	-	> 2,1	> 4,8	
4	-	> 2,8	> 6,4	
5	> 80	> 3,5	> 8,0	
6	-	-	-	
7	-	-	-	
8	-	-	-	
9	-	-	-	
10	-	> 80	-	
Gamme (points)	0 à 5	0 à 10	0 à 5	0 à 5
Total	Somme des points pour les consommations de fruits et légumes, les fibres et les protéines			

02

Choix de la méthode de calcul du score final

Points Négatifs N ≥ 11

Points Négatifs N < 11 ou pour les fromages

Score des fruits et légumes = 5

Score des fruits et légumes < 5

Points N - Points P

Points N - (points fibres + points fruits et légumes)

Points N - Points P

03

Attribution d'une couleur et d'une lettre

Quels produits concernés ?

Score final variant de -15 (qualité nutritionnelle élevée) à 40 (faible qualité nutritionnelle)

Aliments solides	Boissons	Logo
Min à -1	Eaux toujours en A	
0 à 2	Min à 1	
3 à 10	2 à 5	
11 à 18	6 à 9	
19 à max	10 à max	

- Tous les aliments transformés, excepté les herbes aromatiques, thés, cafés, levures...
- Toutes les boissons, excepté les boissons alcoolisées
- Excepté les produits dont la face la plus grande a une surface inférieure à 25 cm<sup>2</sup>