

# Projet de conception d'une application au service de la santé publique

Cyril REGAN

Base de données :

<https://static.openfoodfacts.org/data/en.openfoodfacts.org.products.csv>

10 mars 2020



# Table of Contents

## Problématique et jeu de données

### Sélection des données

Nettoyage des données

Train Validation et Test

### Analyse des données

Tests

Tests de Kolmogorov

Classification

Test de Khi2

Test d'ANOVA

Test de Student

ACP

Eboulis des valeurs propres

Cercles et projection sur (F1, F2)

Cercles et projection sur (F3, F4)

### Construction du modèle

Imputation des variables

Normalisation et réduction

Combinaisons de knn

Construction du nutriscore

Outliers multivariées :

Identification des modèles :

Evaluation du modèle retenu :

Construction de la fonction :

### Conclusion



## Problématique



Appel à projets pour trouver des idées innovantes d'applications en lien avec l'alimentation.



Construire une fonction qui prédit le nutriscore en fonction de caractéristiques partielles



## Présentation du jeu de données



Données brutes : Tableau de 1113279 lignes × 178 colonnes (2.3 Go).



le nutriscore est calculé suivant 10 variables :

Augmente le nutriscore :	Valeur éner- gétique	Acides gras saturés	Sucres	Sodium
Baisse le nutriscore :		Fibres	Fruits Légumes / Fruits à coque	Protéines



# Table of Contents

## Problématique et jeu de données

## Sélection des données

Nettoyage des données

Train Validation et Test

## Analyse des données

Tests

Tests de Kolmogorov

Classification

Test de Khi2

Test d'ANOVA

Test de Student

ACP

Eboulis des valeurs propres

Cercles et projection sur (F1,  
F2)

Cercles et projection sur (F3,  
F4)

## Construction du modèle

Imputation des variables

Normalisation et réduction

Combinaisons de knn

Construction du nutriscore

Outliers multivariées :

Identification des modèles :

Evaluation du modèle re-  
tenu :

Construction de la fonction :

## Conclusion



## Nettoyage des données

Les variables retenues sont :

product_name	nutriscore_score	energy_100g
saturated_fat_100g	sugars_100g	fiber_100g
proteins_100g	sodium_100g	Fruits(+ de 99% NaN)

$X = \text{NaN}$  pour  $X \in [\text{saturated\_fat\_100g}, \text{sugars\_100g}, \text{fiber\_100g}, \text{proteins\_100g}, \text{sodium\_100g}]$  tel que :

$$X > 100 \text{ ou } < 0$$

et  $\text{energy\_100g} = \text{NaN}$  pour le 1<sup>er</sup> et le dernier centile.



## Entrainement Validation et Test

Nettoyage et outliers	DATA	
	Data with NaN	Data complete
Analyse Statistique	Data complete	
Construction modèles =>	Imputation <i>non</i> nutri Prediction nutri	<div> <math>NaN_{\notin \text{nutri}}</math> </div> <div>Train + validation</div>
Test final =>	Train + validation	
		Test



# Table of Contents

## Problématique et jeu de données

## Sélection des données

Nettoyage des données

Train Validation et Test

## Analyse des données

### Tests

Tests de Kolmogorov

Classification

Test de Khi2

Test d'ANOVA

Test de Student

### ACP

Eboulis des valeurs propres

Cercles et projection sur (F1, F2)

Cercles et projection sur (F3, F4)

## Construction du modèle

Imputation des variables

Normalisation et réduction

Combinaisons de knn

Construction du nutriscore

Outliers multivariées :

Identification des modèles :

Evaluation du modèle retenu :

Construction de la fonction :

## Conclusion



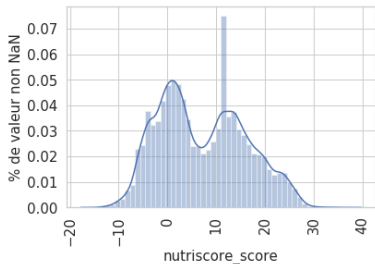


## Tests de Kolmogorov : normalité

**p-valeur** : *plus petite valeur du niveau de test  $\alpha$  conduisant au rejet de l'hypothèse nulle  $H_0$*  :

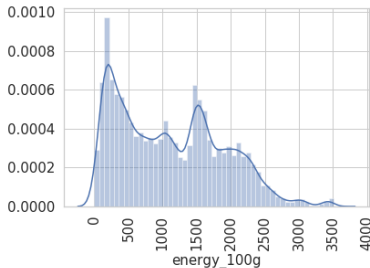
$$\text{Rejet de } H_0 \text{ au niveau de test } \alpha \iff \text{p-valeur} < \alpha$$

**=>  $H_0$**  : la distribution est normale



nutriscore\_100g

p-valeur = 0 :  $H_0$  est rejetée.



energy\_100g

p-valeur = 0 :  $H_0$  est rejetée.



## Classification

Classification des variables qui augmentent le nutriscore selon



energy_100g	saturated_fat_100g	sugars_100g	sodium_100g
Classe 0			
1005	3	13.5	0.27
Classe 3			
2345	7	31	0.63
Classe 7			



## Classification des variables qui diminuent le nutriscore selon



fiber_100g	proteins_100g
Classe 0	
1.9	3.0
Classe 2	
3.7	6.4
Classe 4	



## Test de Khi2 : indépendance des classes proteins\_100g et energy\_100g

=> **H0** : les deux variables proteins\_100g et energy\_100g sont indépendantes.

proteins_100g	0	2	4	Total
energy_100g				
0	19398	16134	15039	50571
3	3922	11968	31677	47567
7	1492	619	4267	6378
Total	24812	28721	50983	104516

Tableau de contingence

La p-valeur = 0.0 : l'hypothèse nulle est rejetée



## Test d'analyse de la variance (ANOVA)

=> **H0** (multifactoriel) : les distributions du **nutri\_score** suivent la même loi normale suivant les modalités des facteurs **sugars\_100g** et **fiber\_100g** ou leur combinaison

Ensembles des groupes ci dessous :

sugars_100g	fiber_100g	N > 10	pval Kolmogorov
0	0	26580	p = 0
	2	7328	p = 6.63272e-91
	4	4488	p = 4.82879e-85
3	0	15530	p = 0
	2	9731	p = 0
	4	14961	p = 0
7	0	14480	p = 0
	2	5016	p = 6.63272e-91
	4	6402	p = 4.82879e-85



## Résultats de l'ANOVA :

	sum_sq	df	F	PR(>F)
C(sugars_100g)	3.525e+06	2.0	42159.100	<b>0.0</b>
C(fiber_100g)	5.348e+05	2.0	6395.308	<b>0.0</b>
C(sugars_100g) : C(fiber_100g)	8.343e+04	4.0	498.867	<b>0.0</b>
Residual	4.369e+06	104507.0	NaN	NaN

Sous l'hypothèse de normalité des groupes (non satisfaite) :

**L'hypothèse nulle rejetée** : le nutri\_score n'a pas la même loi normale suivant les modalités des facteurs sugars\_100g ou fiber\_100g ou leur combinaison.



## Test de Student : Comparaison de 2 moyennes

=> **HO** : les moyennes du nutri\_score sont égales dans les deux groupes('0') et ('3') de energy\_100g

Hypothèse pour le test de Student : Distribution normale des 2 groupes.

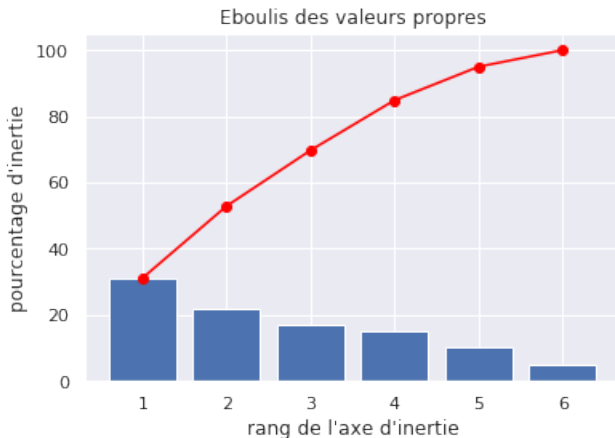
('0')	p = 0	HO Kolmogorov est rejetée
('3')	p = 0	HO Kolmogorov est rejetée
Test de Kolmogorov		

Sous l'hypothèse de normalité des groupes (non satisfaite) :

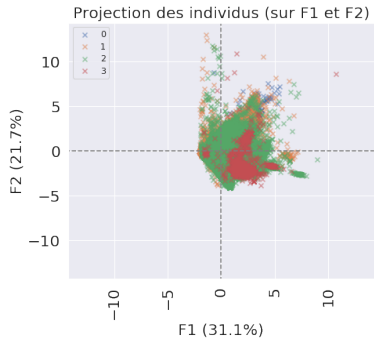
TTest	p = 0	H0 Student est rejetée.
-------	-------	-------------------------



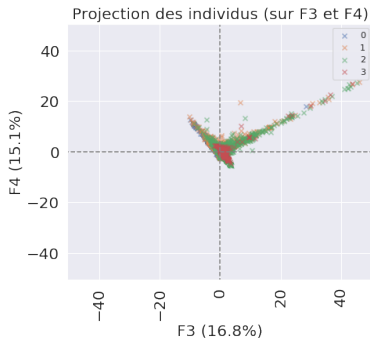
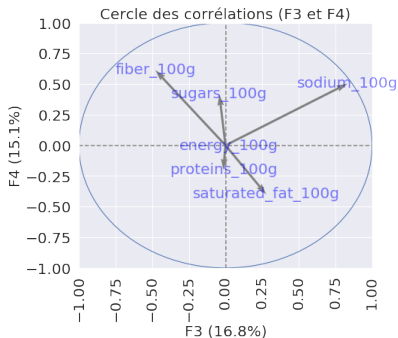
## Analyse par Composantes Principales (ACP) : Eboulis des valeurs propres







## Cercles des corrélations et projection des individus sur (F3, F4)



# Table of Contents

## Problématique et jeu de données

## Sélection des données

Nettoyage des données

Train Validation et Test

## Analyse des données

Tests

Tests de Kolmogorov

Classification

Test de Khi2

Test d'ANOVA

Test de Student

ACP

Eboulis des valeurs propres

Cercles et projection sur (F1,  
F2)

Cercles et projection sur (F3,  
F4)

## Construction du modèle

Imputation des variables

Normalisation et réduction

Combinaisons de knn

Construction du nutriscore

Outliers multivariées :

Identification des modèles :

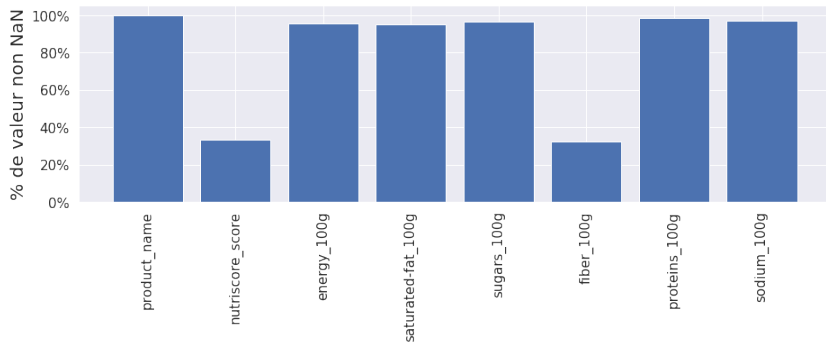
Evaluation du modèle re-  
tenu :

Construction de la fonction :

## Conclusion



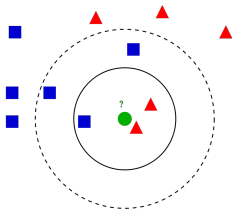
Imputation des variables : Normalisation et réduction des variables [energy\_100g : sodium\_100g] par une moyenne = 0 et une variance = 1.



## Combinaisons de knn

- L'imputation automatique par KNNImputer de sklearn impossible (Pb memoire).
- ⇒ Choix d'algorithme des k plus proches voisins (k-nearest neighbors algorithm : knn) successifs.

Pour chaque variable : imputation par knn (avec optimisation de l'hyperparamètre k) en fonction des données complètes des combinaisons des autres variables.





## Construction du modèle nutriscore : détection des valeurs abhérantes multivariées :

Si la masse de :

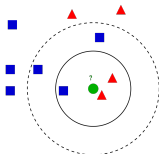
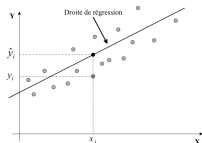
$$\sum (\text{saturated\_fat\_100g} + \text{sugars\_100g} + \text{fiber\_100g} + \text{proteins\_100g} + \text{sodium\_100g}) > \mathbf{100} \text{ ou } < \mathbf{0}$$

=> on supprime l'individu.



## Identification des modèles d'imputation du nutriscore :

Modèle d'apprentissage **supervisé** de **regression**



Sans ACP :  
 [energy\_100g :  
 sodium\_100g]

ou

avec ACP :  
 [F1 :F6]

Coéfficient de détermination :  $R^2 =$

Root Mean Square Error RMSE =

Mean Absolute Error :  
 MAE =

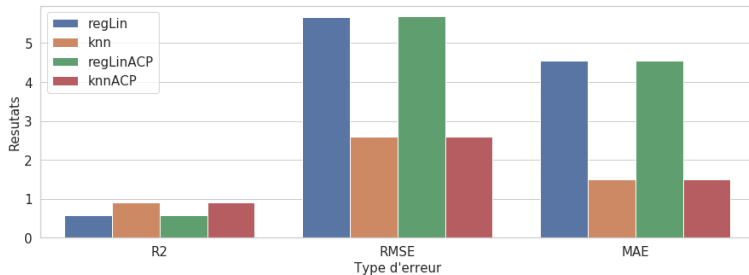
$$1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$



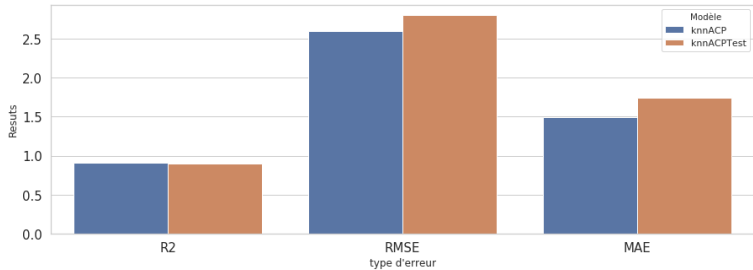




Modèle supervisé de regression retenu : **knn avec ACP**



## Evaluation du modèle retenu avec l'échantillon test [Schéma =>](#) :



## Construction de la fonction nutriscore :

Utilisation :

- Entrée : [energy\_100g :sodium\_100g]
- Sortie : le nutriscore

Fonctionnement :

- Imputation des variables manquantes avec KNNImputer de sklearn
- Réduction et normalisation des variables
- Projection sur les composantes principales
- Prediction du nutriscore avec le modèle retenu (knn avec ACP)



Test de la fonction sur des aliments de nutriscore élevé et faible.  
Données issues de l'échantillon de test :

product_name	nutriscore_score	energy_100g	saturated_fat_100g	sugars_100g	fiber_100g	proteins_100g	sodium_100g
White Chocolate	<b>28</b>	2406.	25.	50.	0.0	7.5	0.1016
Fat Free Skim Milk	<b>-1</b>	138.0	0.0	4.58	0.0	3.33	0.0508



White Chocolate de nutriscore = 28 :

calculNutriACP( 2406.0, 25.0, 50.0, 0.0, 7.5, 0.1016)	=	28
calculNutriACP( 2406.0, 25.0, 50.0, 0.0, 7.5, <b>np.nan</b> )	=	28
calculNutriACP( 2406.0, 25.0, 50.0, 0.0, <b>np.nan</b> , <b>np.nan</b> )	=	28
calculNutriACP( 2406.0, 25.0, 50.0, <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> )	=	27
calculNutriACP( 2406.0, 25.0, <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> )	=	23
calculNutriACP( 2406.0, <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> )	=	20

Fat Free Skim Milk = -1 :

calculNutriACP(138.0, 0.0, 4.58, 0.0, 3.33, 0.0508)	=	-1
calculNutriACP(138.0, 0.0, 4.58, 0.0, 3.33, <b>np.nan</b> )	=	-1
calculNutriACP(138.0, 0.0, 4.58, 0.0, <b>np.nan</b> , <b>np.nan</b> )	=	-1
calculNutriACP(138.0, 0.0, 4.58, <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> )	=	-1
calculNutriACP(138.0, 0.0, <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> )	=	0
calculNutriACP(138.0, <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> , <b>np.nan</b> )	=	4



# Table of Contents

## Problématique et jeu de données

## Sélection des données

Nettoyage des données

Train Validation et Test

## Analyse des données

Tests

Tests de Kolmogorov

Classification

Test de Khi2

Test d'ANOVA

Test de Student

ACP

Eboulis des valeurs propres

Cercles et projection sur (F1, F2)

Cercles et projection sur (F3, F4)

## Construction du modèle

Imputation des variables

Normalisation et réduction

Combinaisons de knn

Construction du nutriscore

Outliers multivariées :

Identification des modèles :

Evaluation du modèle retenu :

Construction de la fonction :

## Conclusion



## Conclusion

- **Analyse des données** : échantillons  $\neq$  loi normale  
 ⇒ tous les tests (Khi2/ANOVA/Student) rejettent l'hypothèse nulle d'indépendance des variables
- **Projection sur les composantes principales (ACP)**  
 ⇒ Corrélations entre (energy\_100g et saturated\_fat\_100g),  
 ⇒ Anticorrélation entre (proteins\_100g et sugars\_100g) et (entre fiber\_100g, sodium\_100g).
- **Construction du modèle du nutriscore**  
 ⇒ knn avec ACP est retenu. L'ACP améliore peu les résultats (dû au fait du faible nombre de variables).
- **Fonction calcul du nutriscore avec données partielles**  
 ⇒ Résultats cohérents pour produits de nutriscore élevé et faible



## Perspectives

- Contacter les développeurs d'appli comme Yuka pour améliorer notre outil et potentiellement leur service.





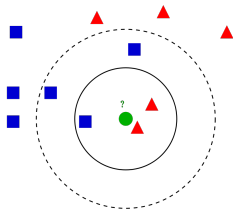
Merci de votre attention





## Combinaisons de knn

- L'imputation automatique par KNNImputer de sklearn impossible (Pb memoire).
- ⇒ Choix d'algorithme des k plus proches voisins (k-nearest neighbors algorithm : knn) successifs :



ooo

o

o

o

o

o

o

ooooooo

ooo

o

o

o

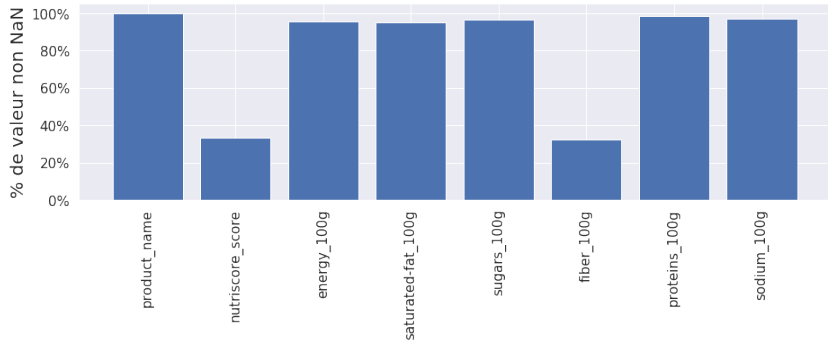
ooooooo

o

o

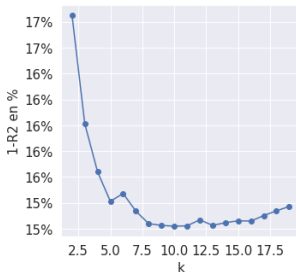
o

ooooo●oooo





- si `[Prediction] ≠ []` : On choisi k tel que `knn.score([fit])` soit le plus petit possible :



- On **impute** `[Prediction].energy_100g` avec `knn.predict()`
- ⇒ Le nombre de NaN de la variable **energy\_100g** restant = **34802** : L'imputation est incomplète



2. Imputation d'**energy\_100g** avec les variables d'entrainement et de test (1) (2) (3) (4) :

(1) saturated\_fat\_100g

(2) sugars\_100g

(3) fiber\_100g

(4) proteins\_100g

⇒ NaN d'**energy\_100g** restant = 34734

3. Imputation avec [(1) (2) (3) (5)], [(1) (2) (4) (5)], ...

4. Imputation avec [(1) (2) (3)], [(1) (2) (4)], ..., [(1) (2)], ..., [(1)], ...



30 knn plus tard... **energy\_100g** est complètement imputée.  
Le même algo est appliquée pour les autres variables. Infine :

