

#### Formation Ingénieur Machine Learning

Projet:

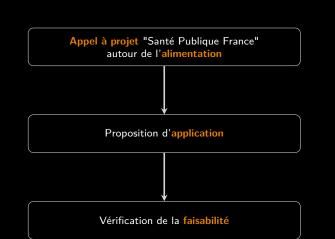
Concevez une application au service de la santé publique

12 Février 2023

### CONTEXTE

Appel à projet "Santé Publique France" autour de l'alimentation





#### CONTEXTE

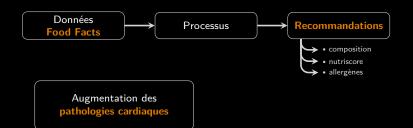
Application proposée

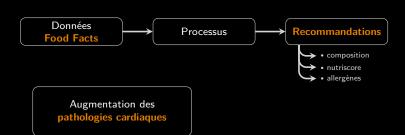
# Données Food Facts



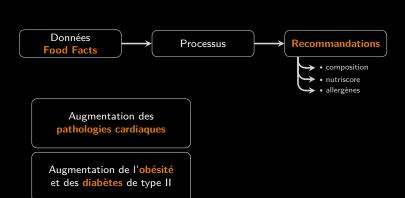




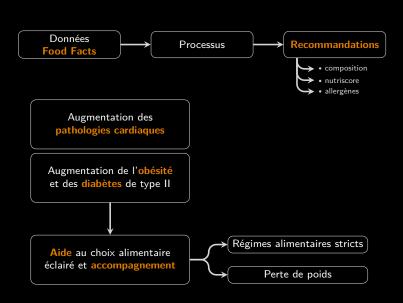


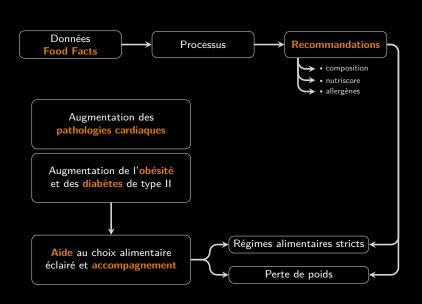


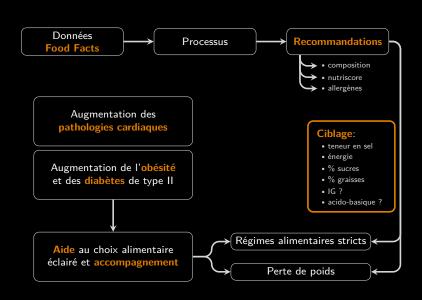
Augmentation de l'obésité et des diabètes de type II



Aide au choix alimentaire éclairé et accompagnement







Nom / pseudo: Toto Pathologies: type liste Allergènes: gluten, oeuf Alim. interdits: beurre
Allergènes: gluten, oeuf
Alim. interdits: beurre
Taille: 1 m 70
Poids: 80 kg
Poids cible: 70kg / 6 mois

# Profil utilisateur: Nom / pseudo: Toto type liste Allergènes: gluten, oeuf Alim. interdits: beurre 1 m 70 80 kg 70kg / 6 mois

#### Affichage produit:

Teneur en sel:

Labels:

Nutriscore:

Énergie pour 100g (kJ): Sucres pour 100g:

Graisses pour 100g:

Présence d'allergènes:

bio. non-OGM

Produits similaires recommandés





- Courbes IMC/poids fonction du temps
- Recommandations de produits en fonction de l'IMC, de la charge calorique et d'un objectif et d'un temps choisi

#### **CONTEXTE**

Data set

produits référencés pour les données étudiées

produits référencés ( pour les données étudiées

#### 191

produits référencés ( pour les données étudiées

### 191

variables

produits référencés

pour les données étudiées

191

variables

## Informations générales

- code
- nom
- ...

pour les produits référencés données étudiées

191

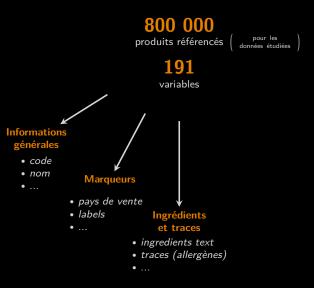
variables

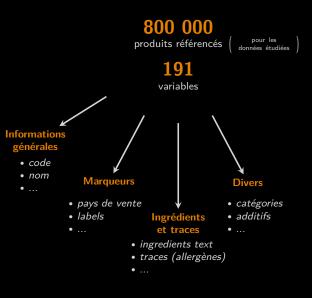
Informations générales • code

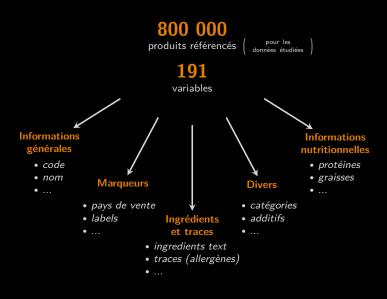
- nom

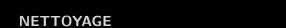
Marqueurs

- pays de vente
- labels

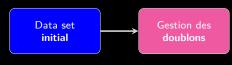




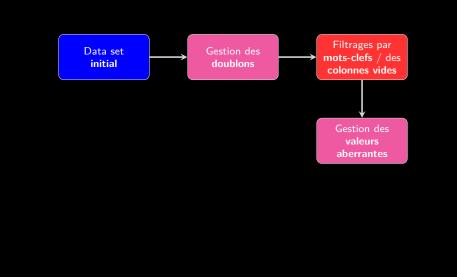


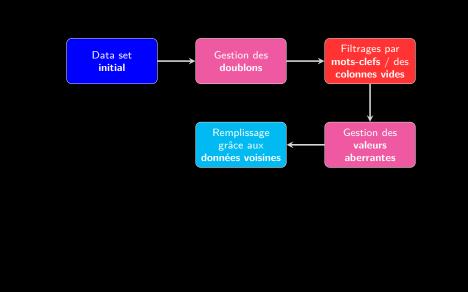


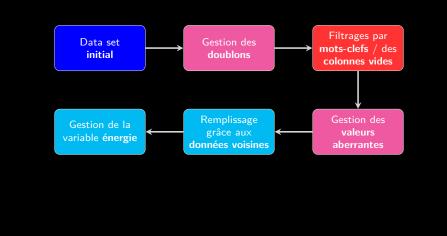
## Data set initial

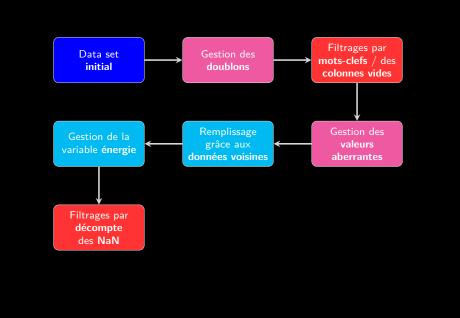


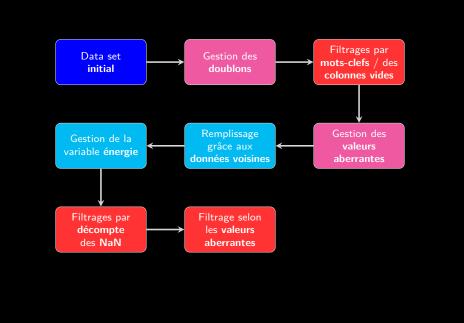


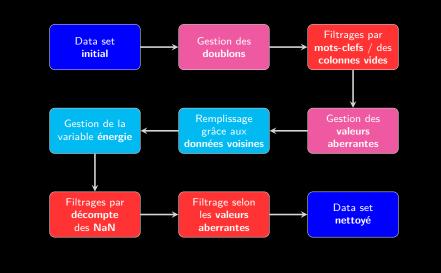


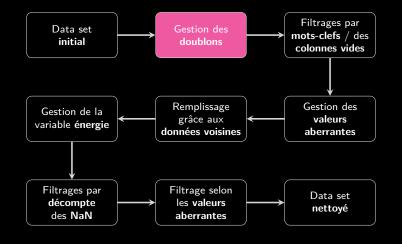












### Détection et tri ascendant par valeurs:

(création d'un DataFrame temporaire)

# **Détection** et tri ascendant par valeurs: (création d'un DataFrame temporaire)

index	code	last modified datetime	n filled
421527	31843340000818	2021-08-17t06:35:03z	28
349035	31843340000818	2022-02-11t08:47:36z	30
61995	3560070278831	2021-04-17t07:44:17z	41
188851	3560070278831	2022-02-10t18:03:06z	47
270028	3700320230572	2021-08-24t12:58:09z	16
749882	3700320230572	2021-08-24t12:58:58z	33
480000	7071688002962	2021-07-13t14:26:35z	40
477267	7071688002962	2021-07-13t14:26:35z	45

## **Détection** et tri ascendant par valeurs:

### $({\it cr\'eation}\ d'un\ {\it DataFrame}\ temporaire)$

#### même code

	meme code		
index	code	last modified datetime	n filled
421527	31843340000818	2021-08-17t06:35:03z	28
349035	31843340000818	2022-02-11t08:47:36z	30
61995	3560070278831	2021-04-17t07:44:17z	41
188851	3560070278831	2022-02-10t18:03:06z	47
270028	3700320230572	2021-08-24t12:58:09z	16
749882	3700320230572	2021-08-24t12:58:58z	33
480000	7071688002962	2021-07-13t14:26:35z	40
477267	7071688002962	2021-07-13t14:26:35z	45

# **Détection** et tri ascendant par valeurs: (création d'un DataFrame temporaire)

		+ ancien / + récent	
index		last modified datetime	n filled
421527	31843340000818	2021-08-17t06:35:03z	28
349035	31843340000818	2022-02-11t08:47:36z	30
61995	3560070278831	2021-04-17t07:44:17z	41
188851	3560070278831	2022-02-10t18:03:06z	47
270028	3700320230572	2021-08-24t12:58:09z	16
749882	3700320230572	2021-08-24t12:58:58z	33
480000	7071688002962	2021-07-13t14:26:35z	40
477267	7071688002962	2021-07-13t14:26:35z	45

# **Détection** et tri ascendant par valeurs:

(création d'un DataFrame temporaire)

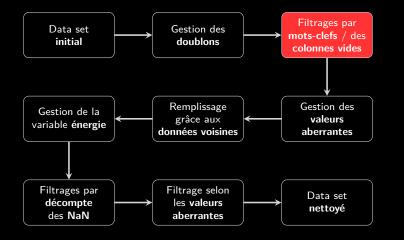
#### - rempli / + rempli

		Tell	ipii /   Telli
index		last modified datetime	n filled
421527	31843340000818	2021-08-17t06:35:03z	28
349035	31843340000818	2022-02-11t08:47:36z	30
61995	3560070278831	2021-04-17t07:44:17z	41
188851	3560070278831	2022-02-10t18:03:06z	47
270028	3700320230572	2021-08-24t12:58:09z	16
749882	3700320230572	2021-08-24t12:58:58z	33
480000	7071688002962	2021-07-13t14:26:35z	40
477267	7071688002962	2021-07-13t14:26:35z	45

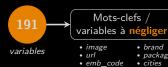
# **Détection** et tri ascendant par valeurs: (création d'un DataFrame temporaire)

index	code	last modified datetime	n filled
421527	31843340000818	2021-08-17t06:35:03z	28
349035	31843340000818	2022-02-11t08:47:36z	30
61995	3560070278831	2021-04-17t07:44:17z	41
188851	3560070278831	2022-02-10t18:03:06z	47
270028	3700320230572	2021-08-24t12:58:09z	16
749882	3700320230572	2021-08-24t12:58:58z	33
480000	7071688002962	2021-07-13t14:26:35z	40
477267	7071688002962	2021-07-13t14:26:35z	45









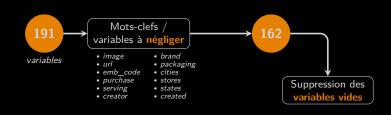
· packaging

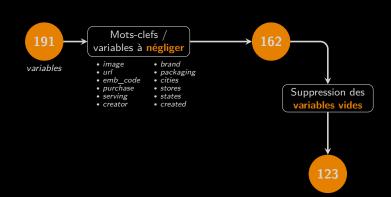
emb\_codepurchase stores

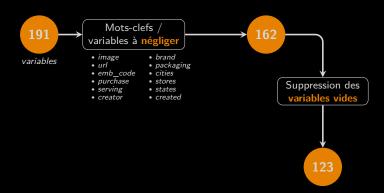
 serving states created

• creator



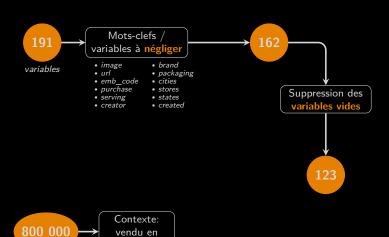






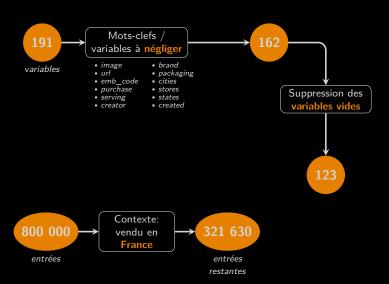
800 000

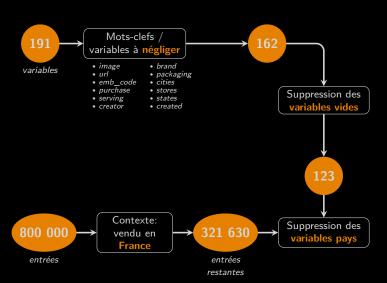
entrées

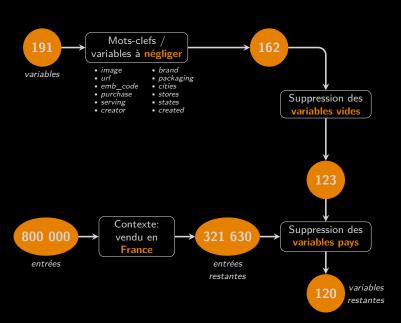


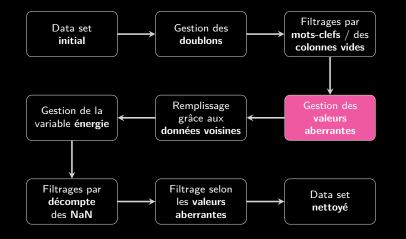
**France** 

entrées





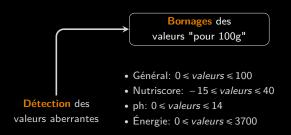


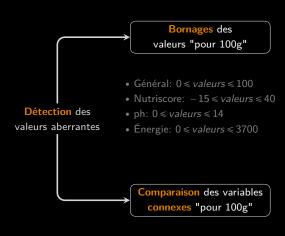


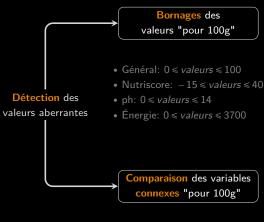
### Détection des

valeurs aberrantes

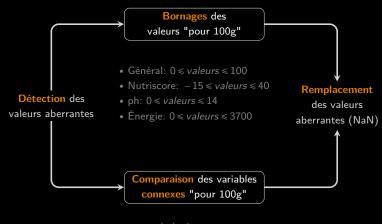




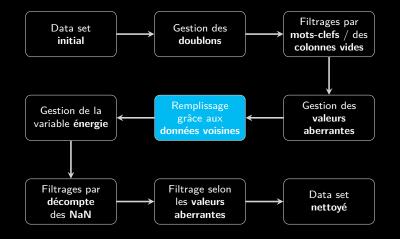




- carbohydrates ≥ sugars
- salt ≥ sodium
- fat > other fats



- carbohydrates ≥ sugars
- salt ≥ sodium
- fat > other fats



# **Détection** des NaN pour les variables:

pour les variables:

product name nutriscore

8576 116 701





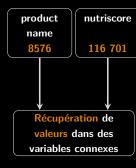
~isna\isna	product name	abbreviated product name	generic name
product name		310337	282502
abbreviated product name	141		344
generic name	57	28095	



~isna\isna	product name	abbreviated product name	generic name
product name		310337	282502
abbreviated product name	141		344
generic name	57	28095	



~isna\isna	nutriscore score	nutriscore grade	nutrition- score-fr 100g	nutrition- score-uk 100g
nutriscore score			91	117088
nutriscore grade			91	117088
nutrition- score-fr 100g				116997
nutrition- score-uk 100g				



~isna\isna	nutriscore score	nutriscore grade	nutrition- score-fr 100g	nutrition- score-uk 100g
nutriscore score			91	117088
nutriscore grade			91	117088
nutrition- score-fr 100g	1			116997
nutrition- score-uk 100g				

### Détection des NaN

pour les variables:

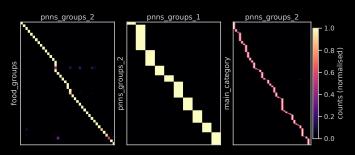


~isna\isna	nutriscore score	nutriscore grade	nutrition- score-fr 100g	nutrition- score-uk 100g
nutriscore score			91	117088
nutriscore grade			91	117088
nutrition- score-fr 100g	1			116997
nutrition- score-uk 100g				

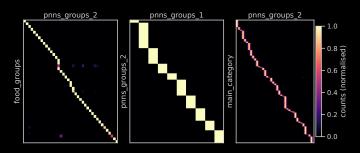
valeurs remplies

179 241 NaN pour la variable pnns groups 1

179 241 NaN pour la variable pnns groups 1



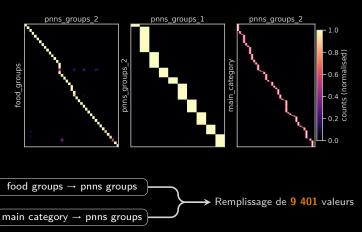
179 241 NaN pour la variable pnns groups 1

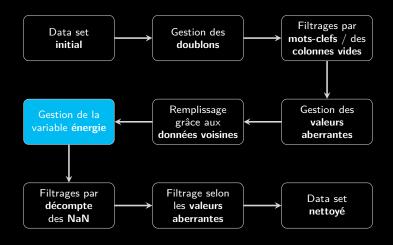


food groups → pnns groups

main category → pnns groups

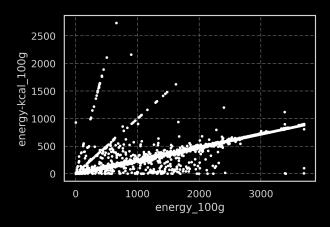
179 241 NaN pour la variable pnns groups 1





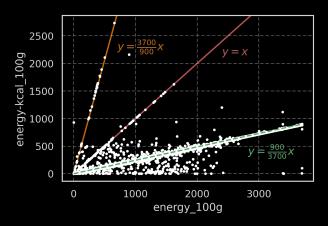
### Présence de valeurs incohérentes:

idéalement x en kJ, y en kcal et  $y \approx \frac{900}{3700}x$ 



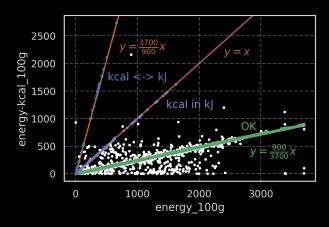
### Présence de valeurs incohérentes:

idéalement x en kJ, y en kcal et  $y \approx \frac{900}{3700}x$ 

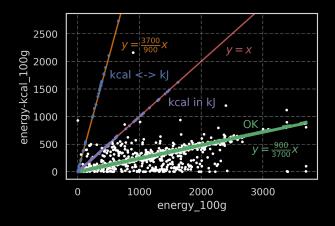


### Présence de valeurs incohérentes:

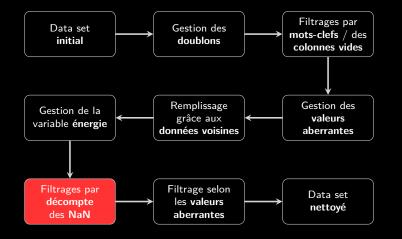
idéalement x en kJ, y en kcal et  $y \approx \frac{900}{3700}x$ 



Présence de valeurs incohérentes: idéalement x en kJ, y en kcal et  $y \approx \frac{900}{3700} x$ 



Solution: garder la valeur maximale (en kJ)



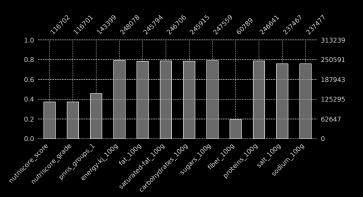








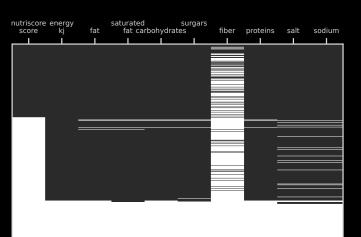
### Taux de remplissage des variables sélectionnées:

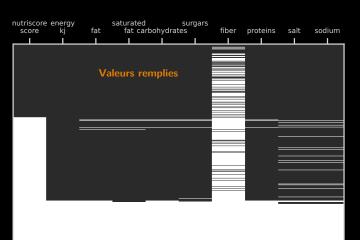


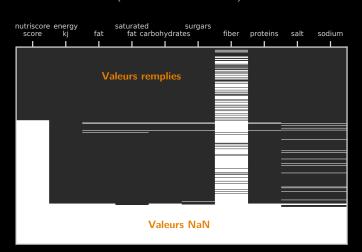


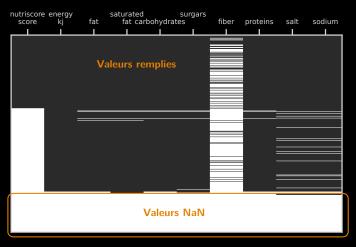
### Taux de remplissage des variables sélectionnées:



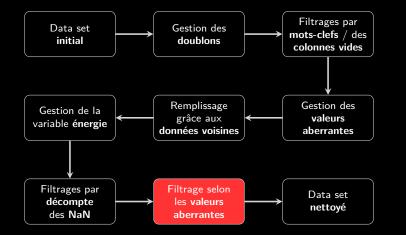








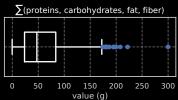
Suppression des entrées vides





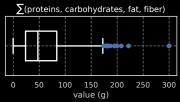
	count	mean	std	min	25%	50%	75%	max
Σ poids	244025	51.27	31.43	0	23.8	47.5	83.9	300

	count	mean	std	min	25%	50%	75%	max
$\sum$ poids	244025	51.27	31.43	0	23.8	47.5	83.9	300



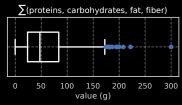
### Données statistiques sur les données brutes:

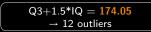
	count	mean	std	min	25%	50%	75%	max
$\sum$ poids	244025	51.27	31.43	0	23.8	47.5	83.9	300

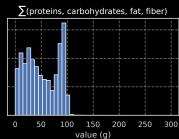


Q3+1.5\*IQ =  $\frac{174.05}{}$   $\rightarrow$  12 outliers

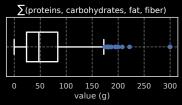
	count	mean	std	min	25%	50%	75%	max
$\sum$ poids	244025	51.27	31.43	0	23.8	47.5	83.9	300

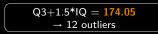


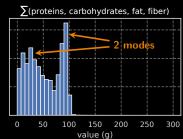




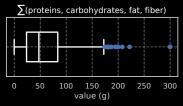
	count	mean	std	min	25%	50%	75%	max
$\sum$ poids	244025	51.27	31.43	0	23.8	47.5	83.9	300

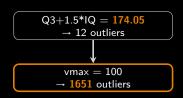


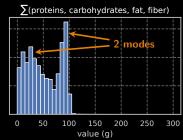




	count	mean	std	min	25%	50%	75%	max
∑ poids	244025	51.27	31.43	0	23.8	47.5	83.9	300







**EXPLORATION DES DONNÉES** 

## **EXPLORATION DES DONNÉES**

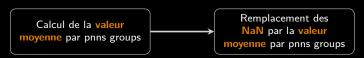
Remplissage des valeurs NaN éparses

Méthode 1: Valeur moyenne par pnns groups et par variable

Méthode 1: Valeur moyenne par pnns groups et par variable

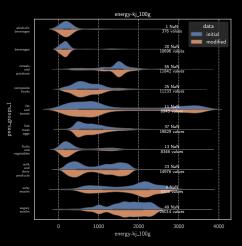
Calcul de la valeur moyenne par pnns groups

Méthode 1: Valeur moyenne par pnns groups et par variable

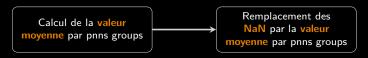


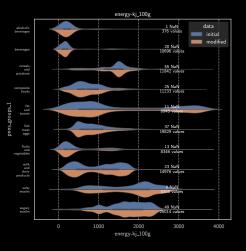
Méthode 1: Valeur moyenne par pnns groups et par variable





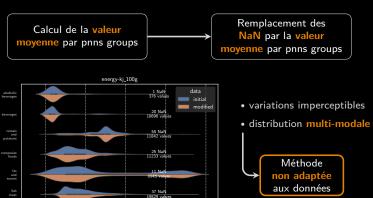
Méthode 1: Valeur moyenne par pnns groups et par variable





- variations imperceptibles
- distribution multi-modale

Méthode 1: Valeur moyenne par pnns groups et par variable



13 NaN 8346 values

23 NaN 14976 values

49 NaN

3000

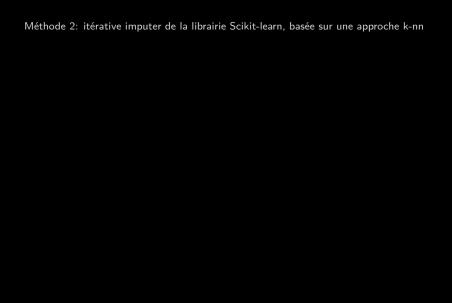
4000

2000

energy-kj\_100g

1000

milk and dairy products



Méthode 2: itérative imputer de la librairie Scikit-learn, basée sur une approche k-nn

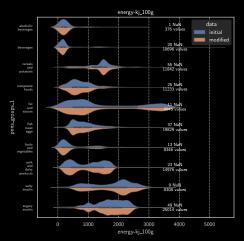
Apprentissage sur des données complètes

Méthode 2: itérative imputer de la librairie Scikit-learn, basée sur une approche k-nn



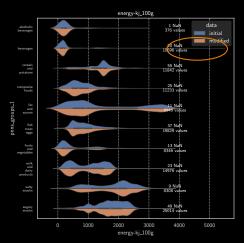
Méthode 2: itérative imputer de la librairie Scikit-learn, basée sur une approche k-nn



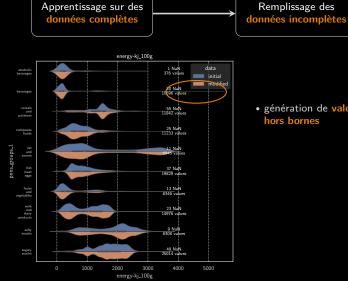


Méthode 2: itérative imputer de la librairie Scikit-learn, basée sur une approche k-nn





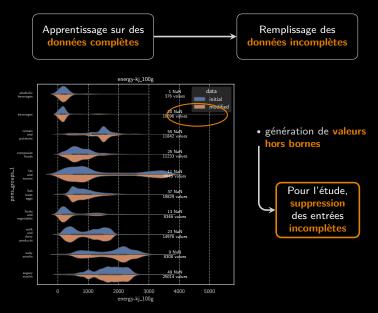
Méthode 2: itérative imputer de la librairie Scikit-learn, basée sur une approche k-nn



• génération de valeurs hors bornes

Remplissage des

Méthode 2: itérative imputer de la librairie Scikit-learn, basée sur une approche k-nn



# EXPLORATION DES DONNÉES

Prédiction du nutriscore (remplissage des NaN)

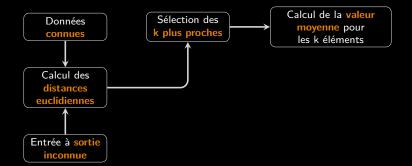
Données connues

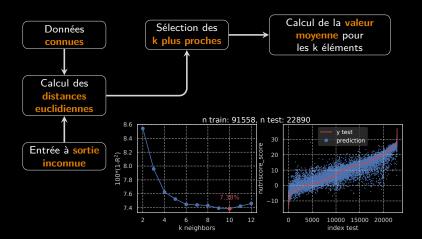
Données connues

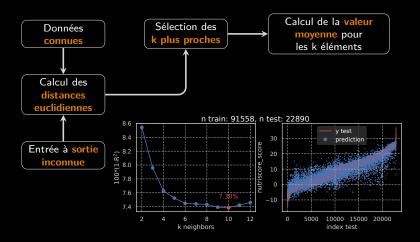
Entrée à sortie inconnue





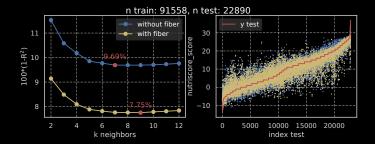




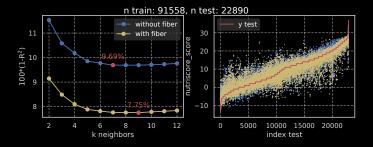


Quid de l'hypothèse NaN fibers = 0 ?

## Vérification de l'hypothèse NaN fiber = 0:



## Vérification de l'hypothèse NaN fiber = 0:



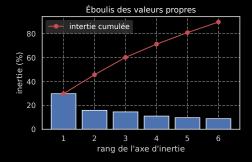
 $\rightarrow$  Hypothèse à priori validée

# **EXPLORATION DES DONNÉES**

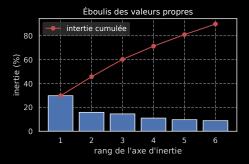
Étude de l'inertie des valeurs - Analyse en Composantes Principales

Principe: Calcul des axes en vue d'aligner les coordonnées et les principaux axes d'inertie

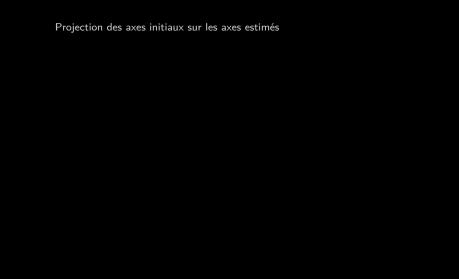
Principe: Calcul des axes en vue d'aligner les coordonnées et les principaux axes d'inertie



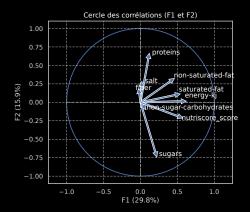
Principe: Calcul des axes en vue d'aligner les coordonnées et les principaux axes d'inertie



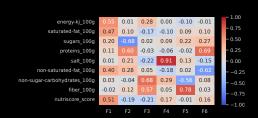
- → 90% de l'inertie sur 6 axes
- → Répartition relativement équilibrée



### Projection des axes initiaux sur les axes estimés

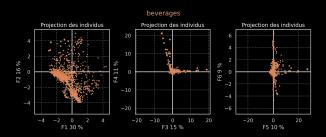


## Projection des axes initiaux sur les axes estimés



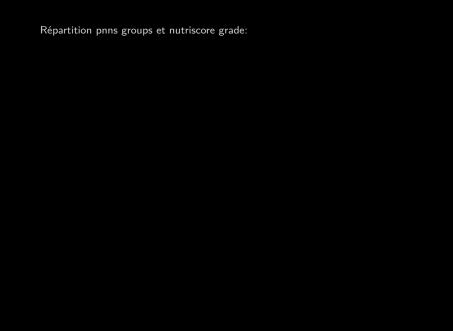
#### Projection des axes initiaux sur les axes estimés



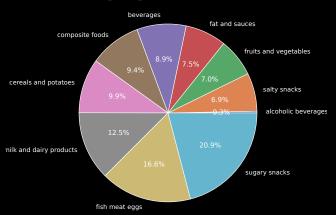


# **EXPLORATION DES DONNÉES**

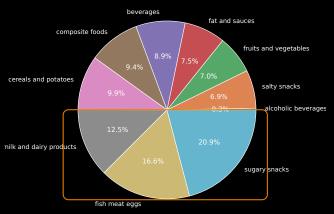
Analyse des données en lien avec l'application





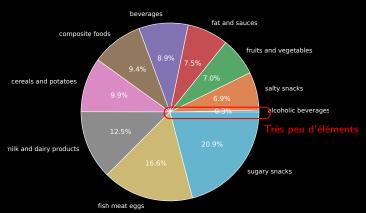


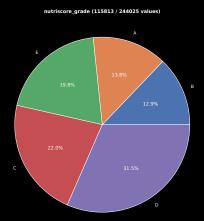


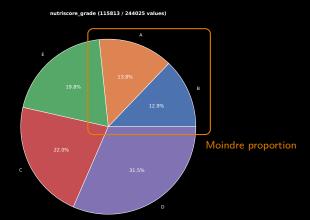


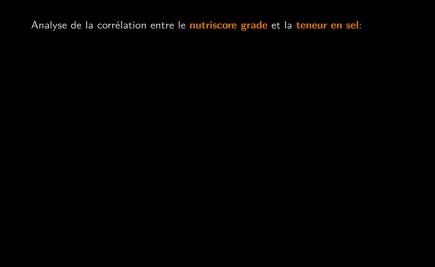
Principaux éléments



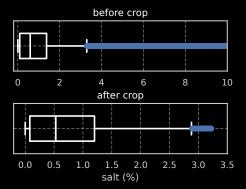




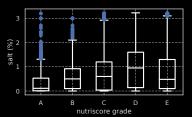


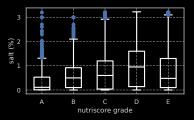


Analyse de la corrélation entre le nutriscore grade et la teneur en sel:

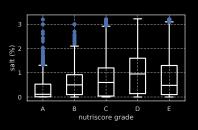


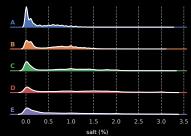
## Analyse de la corrélation entre le nutriscore grade et la teneur en sel:



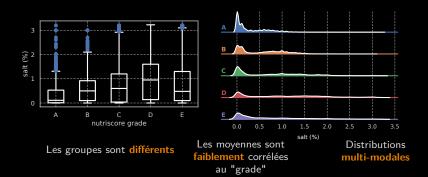


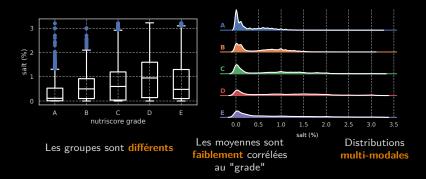
Les groupes sont différents



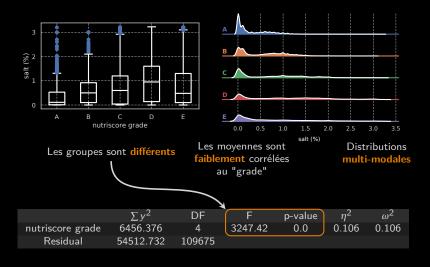


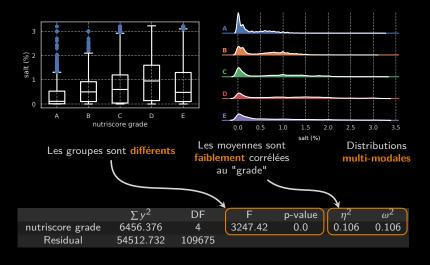
Les groupes sont différents

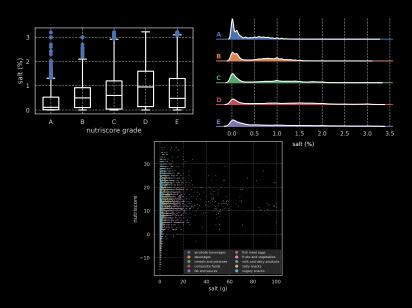


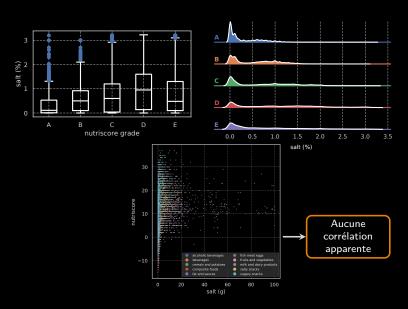


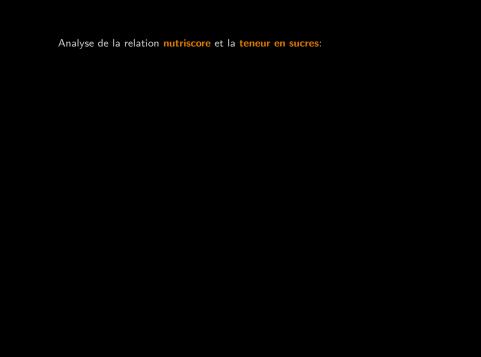
	$\sum y^2$	DF	F	p-value	$\eta^2$	$\omega^2$
nutriscore grade	6456.376	4	3247.42	0.0	0.106	0.106
Residual	54512.732	109675				



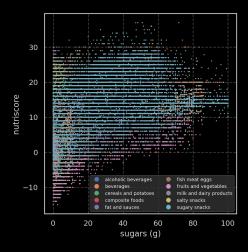




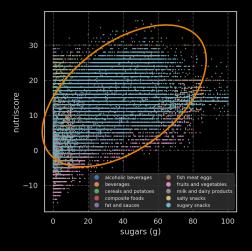




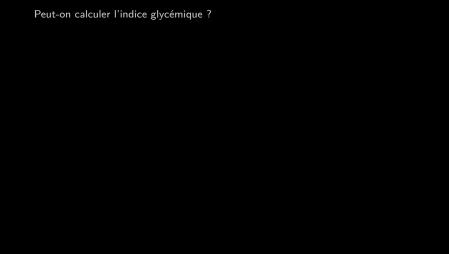
# Analyse de la relation nutriscore et la teneur en sucres:



Analyse de la relation nutriscore et la teneur en sucres:



Une corrélation est visible, mais d'autres variables influencent aussi le nutriscore



#### Peut-on calculer l'indice glycémique ?

#### Exemple de listes d'ingrédients:

- lait entier (99%); poudre de lait (1%), ferments lactiques, présure.
- jus d'orange 40% jus de pomme 40% jus d'ananas 9% purée de banane jus de raisin blanc jus de pamplemousse purée d'abricot purée de pêche.
- 1 boite de garniture: tomates fraiches (51%), eau, oignons frais, huile d'olive vierge extra, huile de colza, sel, persil, concentré de tomate, jus concentré de citron, menthe, épaississants: farine de graines de caroube et gomme guar, arômes. 1 coupelle de semoule de blé dur précuite à la vapeur.
- farine de blé, sucre, beurre 12% (lait), sirop de sucre inverti, poudres à lever: carbonates de sodium, diphosphates, lactosérum en poudre lait), lat entier en poudre, sel, émulsfiant: lécithine de soja; acidifiant: acide citrique; arôme (ocufs entiers en poudre).
- $\bullet$  palette de porc avec os 90 %, eau, sel, dextrose, conservateur : nitrite de sodium.

#### Peut-on calculer l'indice glycémique ?

#### Exemple de listes d'ingrédients:

- lait entier (99%); poudre de lait (1%), ferments lactiques, présure.
- jus d'orange 40% jus de pomme 40% jus d'ananas 9% purée de banane jus de raisin blanc jus de pamplemousse purée d'abricot purée de pêche.
- 1 boite de garniture: tomates fraiches (51%), eau, oignons frais, huile d'olive vierge extra, huile de colza, sel, persil, concentré de tomate, jus concentré de citron, menthe, épaississants: farine de graines de caroube et gomme guar, arômes. 1 coupelle de semoule de blé dur préculte à la vapeur.
- farine de blé, sucre, beurre 12% (lait), sirop de sucre inverti, poudres à lever: carbonates de sodium, diphosphates, lactosérum en poudre lait), lat entier en poudre, sel, émulsifiant: lécithine de soja; acidifiant: acide citrique; arôme (ocufs entiers en poudre).
- $\bullet$  palette de porc avec os 90 %, eau, sel, dextrose, conservateur : nitrite de sodium.

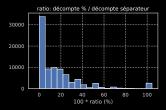
Non, car les % sont assez peu renseignés

#### Peut-on calculer l'indice glycémique ?

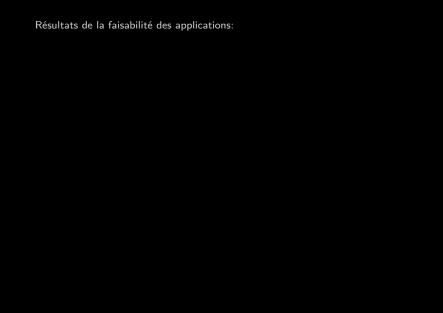
#### Exemple de listes d'ingrédients:

- lait entier (99%); poudre de lait (1%), ferments lactiques, présure.
- jus d'orange 40% jus de pomme 40% jus d'ananas 9% - purée de banane - jus de raisin blanc - jus de pamplemousse - purée d'abricot - purée de pêche.
- 1 boite de gamiture: tomates fraiches (51%), eau, olen, osen, orale, huile d'olive vierge extra, huile de colza, sel, persil, concentré de tomate, jus concentré de citron, menthe, épaississants: farine de graines de caroube et gomme guar, arômes. 1 coupelle de semoule de blé dur précuite à la vapeur.
- farine de blé, sucre, beurre 12% (lait), sirop de sucre inverti, poudres à lever: carbonates de sodium, diphosphates, lactosérum en poudre lait), lat entier en poudre, sel, émulsifiant: lécithine de soja; acidifiant: acide citrique; arôme (oeufs entiers en poudre).
- $\bullet$  palette de porc avec os 90 %, eau, sel, dextrose, conservateur : nitrite de sodium.

Non, car les % sont assez peu renseignés



# **CONCLUSIONS**





Données récupérables:

Données récupérables:

Teneur en sel

Charge calorique

Présence d'ingrédients dont allergènes

Pourcentage de macro-nutriments notamment graisses et sucres

Données récupérables:

Applications:

Teneur en sel

Charge calorique

Présence d'ingrédients dont allergènes

Pourcentage de macro-nutriments notamment graisses et sucres

Données récupérables:

Applications:

Teneur en sel

Charge calorique

Présence d'ingrédients dont allergènes

Pourcentage de macro-nutriments notamment graisses et sucres Régime pauvre en sel

Gestion de poids/IMC

Régime sans/avec ...

Données récupérables:

Applications:

Teneur en sel

Charge calorique

Présence d'ingrédients dont allergènes

Pourcentage de macro-nutriments notamment graisses et sucres Régime pauvre en sel

Gestion de poids/IMC

Régime sans/avec ...

Gestion de l'indice glycémique

Gestion équilibre acido-basique

Données récupérables: Applications:

Teneur en sel

Charge calorique

Présence d'ingrédients dont allergènes

Pourcentage de macro-nutriments notamment graisses et sucres

Régime pauvre en sel

Gestion de poids/IMC

Régime sans/avec ...

Gestion de l'indice glycémique

Gestion équilibre acido-basique

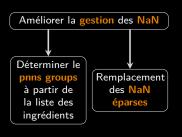
Se rapprocher de spécialistes du domaine est nécessaire pour apporter une meilleure connaissance métier (identifier/valider les contraintes alimentaires liées aux états de santé)



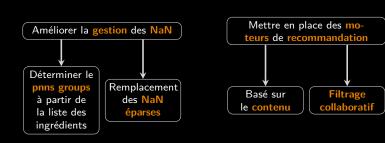
Améliorer la gestion des NaN







Mettre en place des moteurs de recommandation



# Merci pour votre attention.

knn

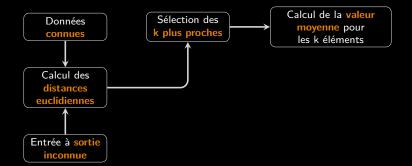
Données connues

Données connues

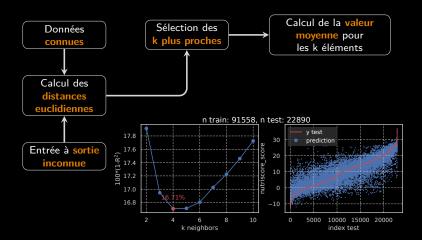
Entrée à sortie inconnue



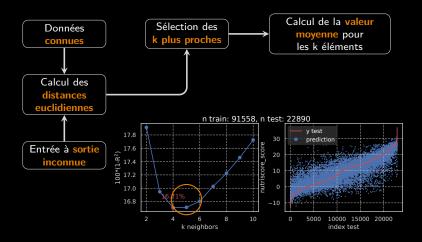




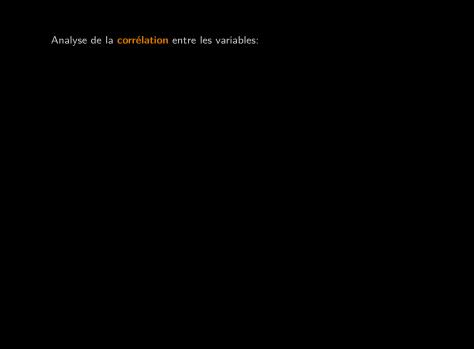
#### Principe:

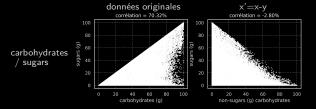


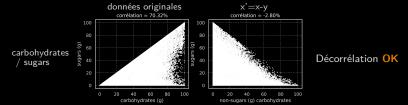
#### Principe:

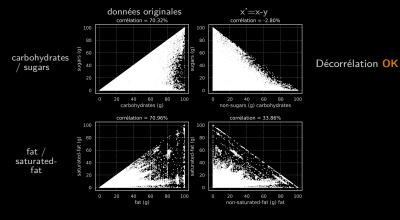


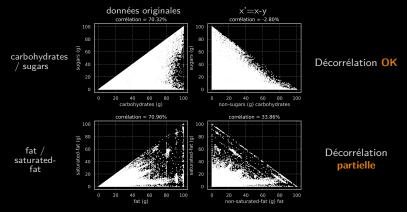
Avec les données brutes, 17% d'erreur → nécessité d'optimiser

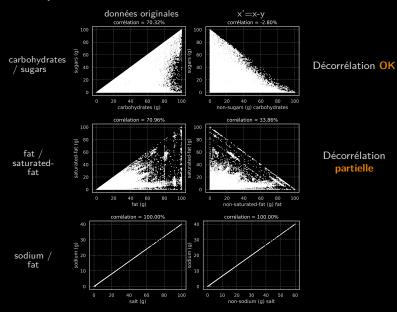


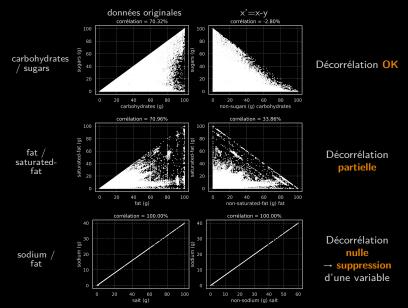






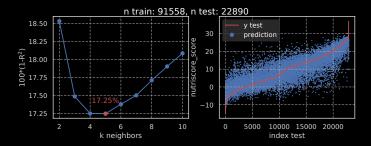




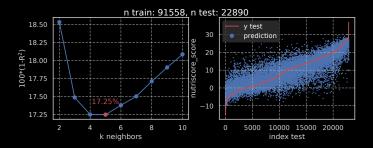


Prédiction en utilisant la base améliorée:

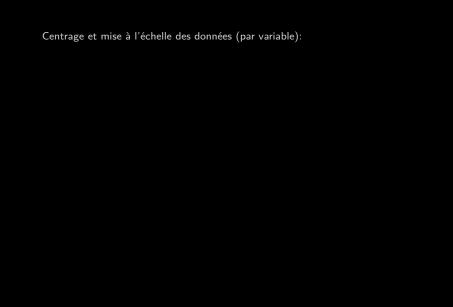
#### Prédiction en utilisant la base améliorée:



#### Prédiction en utilisant la base améliorée:



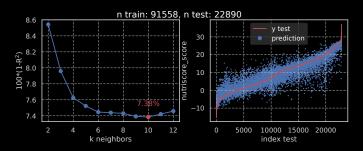
- → Amélioration non significative
- → Centrage et mise à l'échelle des données



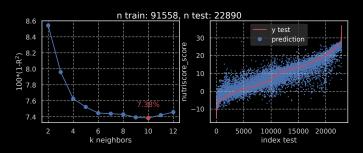






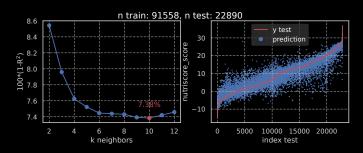






→ Nette amélioration

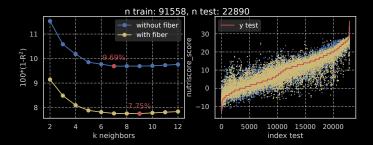




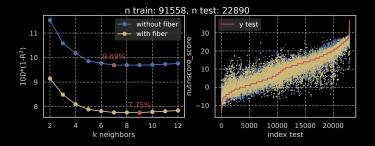
→ Nette amélioration

Quid de l'hypothèse NaN fibers = 0?

## Vérification de l'hypothèse NaN fiber = 0:

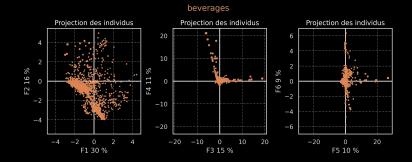


## Vérification de l'hypothèse NaN fiber = 0:

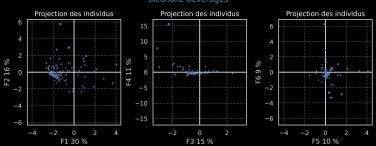


→ Hypothèse à priori validée

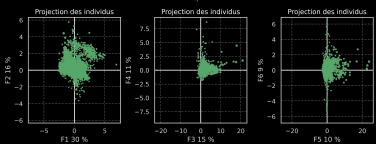
# PCA - projection

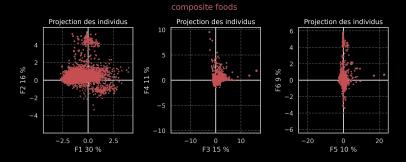


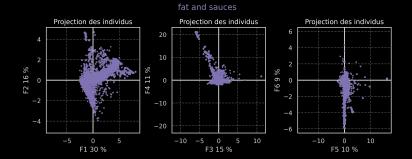
#### alcoholic heverages

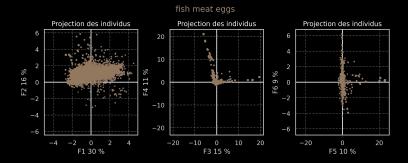


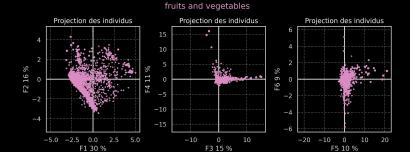


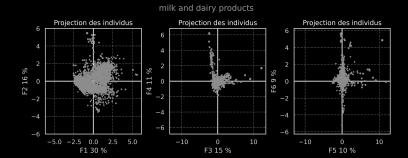


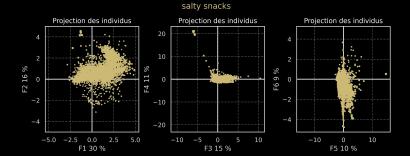


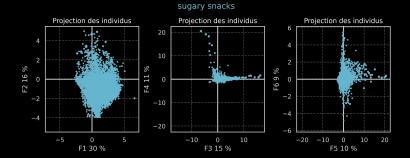












unprocessed data that should have been added to the final page this extra page been added to receive it.

ATEX now knows how many pages to expect for this document.