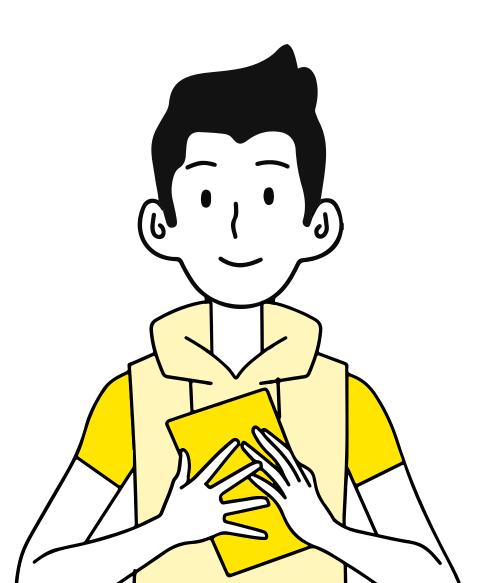
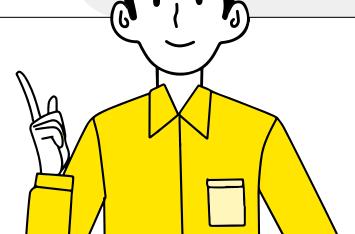
## Santé publique France

Parcours Al Engineer



#### Notre mission:

## Nettoyer et explorer les données pour la faisabilité d'une application.



### Déroulement

1 Traiter le jeu de données.

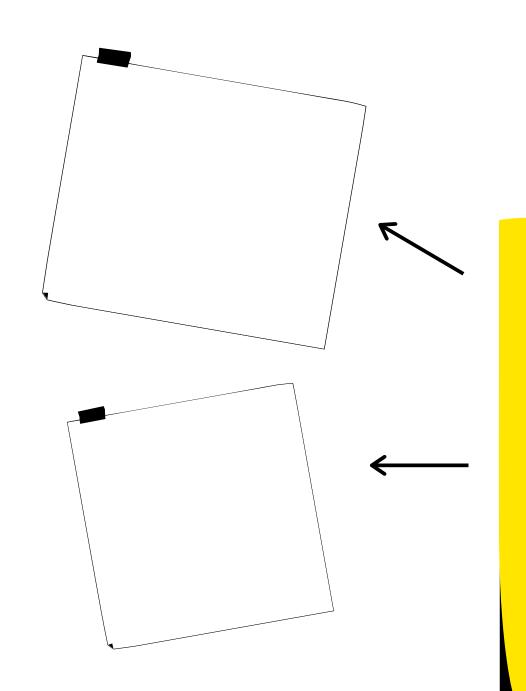
3

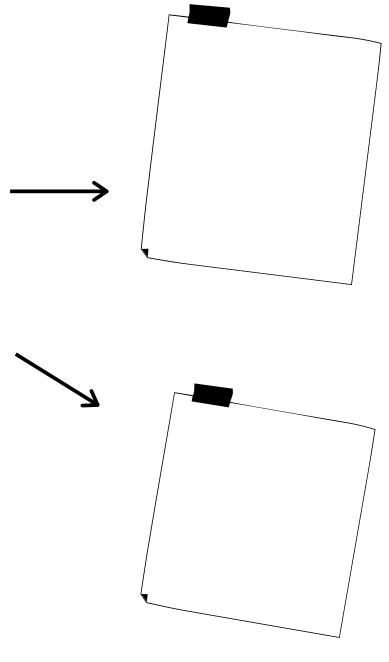
Observation - Corrélation -

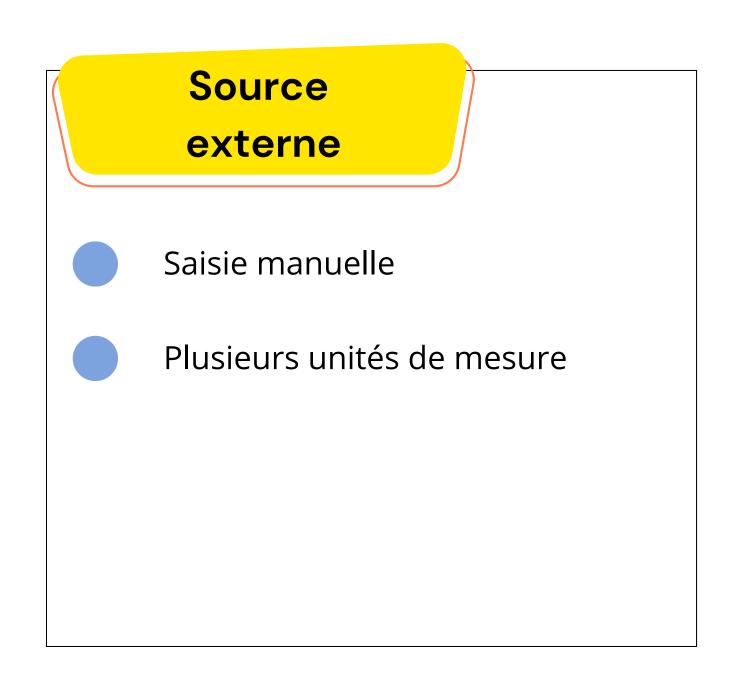
2 Observation -tendance-

4

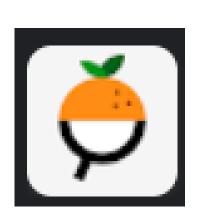
Rapport d'exploration.











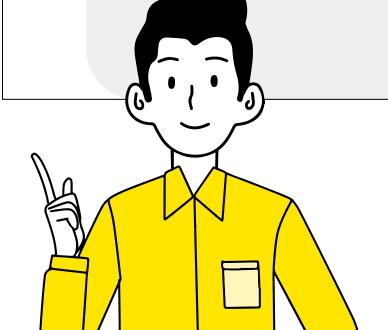


Jeu de données 320772 rows x 162 columns International Nombreuse colonnes 'doublons' Nombreuse colonnes pour 100g'

**Piste** d'analyse allergens nutrition\_grade\_fr carbon-footprint-100g additives countries

Question.

# Il y a-t-il un lien entre le Nutri Score et les allergènes ?

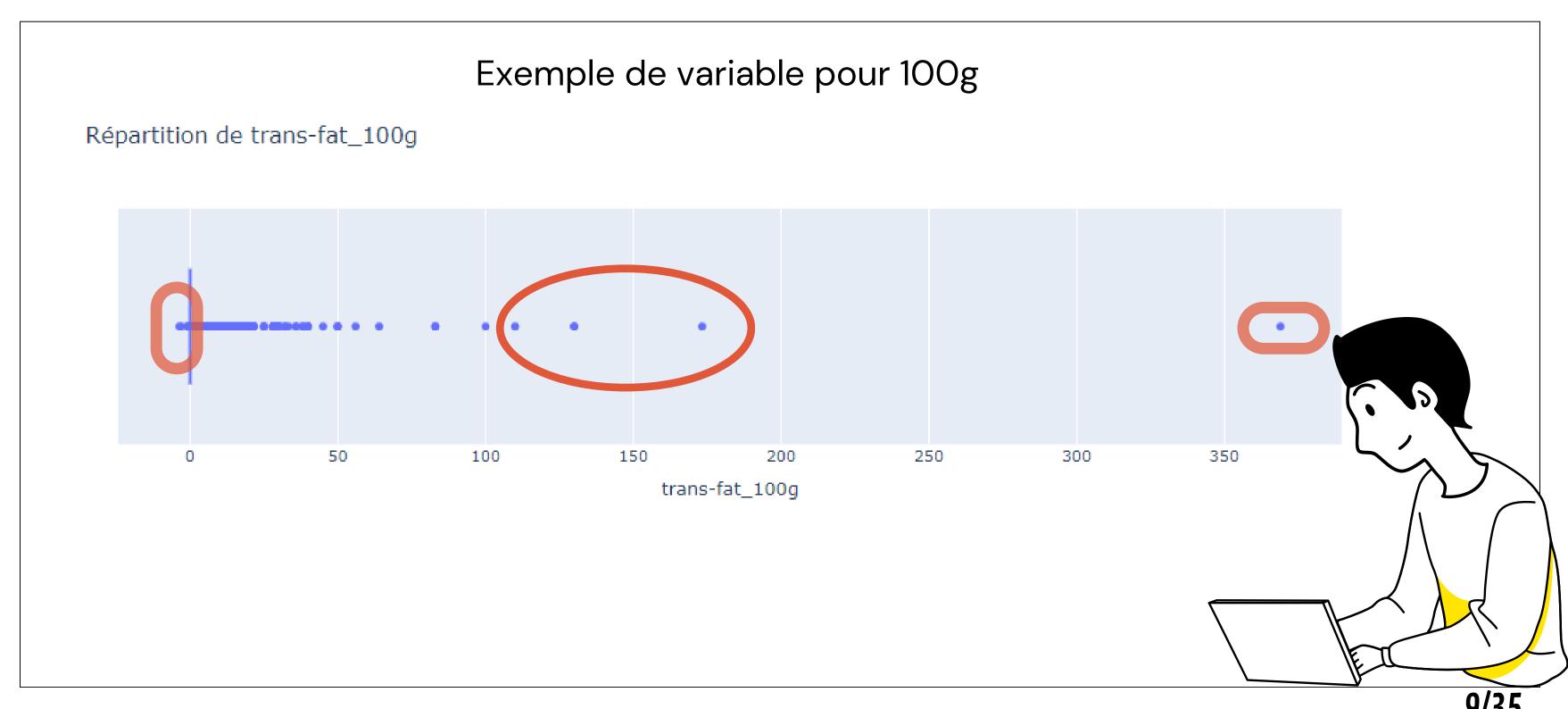


#### Principe RGPD

- 1. Licéité, loyauté et transparence
- 2. Limitation des finalités
- 3. Minimisation des données
- 4. Exactitude des données
- 5. Limitation de la conservation



Nettoyage et complétions



Nettoyage et complétions

#### Solution pour les variables "pour 100g"

```
# liste de mots-clés que vous souhaitez rechercher dans les noms de colonnes
keywords = ['_100g']

# Parcourez les colonnes du DataFrame
for col_name in data.filter(like='_100g'):
    # Appliquez la transformation
    data[col_name] = data[col_name].apply(lambda x: x / 1000 if x > 100 else 0 if x < 0 else x)</pre>
```

#### Suppression

#Suppression des variables à 100% de Val Manquante data = data.dropna(axis=1, how='all')



10/35

Nettoyage

Doublons via 'code'

```
1]: dataFrance.duplicated('code').sum()
1]: 4
```

```
Traitement des doublons :

-On va dans un premier temps fusionner les doublons pour remplacer un maximum de valeur manquante
-Ensuite on conserve : -1 Celui qui possède déjà un résultat nutriscore.

-2 Ou l'individu qui a le plus de valeur à la fin de la fusion
-3 La date de création la plus récente.
```



Liste de nettoyage rapide

- Simplification des informations sur l'huile de palme
- Suppression des lignes sans le nom du produit
- Suppression des lignes de boissons alcoolisées
- Fusion des colonnes sur le même sujet
- Suppression des colonnes après fusion et obsolètes
- Création de la variable DataFrance

Calcul de Nutri-Score avec du machine learning -Préparation-

```
for col in dataFrance.columns:
    #Si le type de la colonne est un objet (chaîne de caractères), remplacer les valeurs manquantes par 'Abs'
                                                                                       1 - Pré-remplissage
   if dataFrance[col].dtype == '0': # '0' est le code pour le type d'objet
       dataFrance[col] = dataFrance[col].fillna('Abs')
    #Si le type de la colonne est un entier ou un flottant, remplacer les valeurs manquantes par 0
   if dataFrance[col].dtype in ['int64', 'float64']: # 'int64' et 'float64' sont les types d'entier et de flottant
       dataFrance[col] = dataFrance[col].fillna(0)
def convert nutrition grade(grade):
   mapping = {'a': 1, 'b': 2, 'c': 3, 'd': 4, 'e': 5}
   return mapping.get(grade, None) # Retourne None si La note n'est pas tropée - Numérisation
# Appliquer la fonction de conversion à la colonne 'nutrition_grade_fr'
dataFrance['nutrition_grade_fr'] = dataFrance['nutrition_grade_fr'].apply(convert_nutrition_grade)
# Séparation du jeu de données en deux parties
dataFrance with y = dataFrance.dropna(subset=['nutrition grade fr'])
dataFrance_without_y = dataFrance[dataFrance['nutrition_grade_fr'].isna()]
                                                3 - Création de variable
# Définition de X train et y train
X_train = dataFrance_with_y[['salt_100g', 'proteins_100g', 'fiber_100g', 'sugars_100g', 'carbohydrates_100g', 'saturated-f
y_train = dataFrance_with_y['nutrition_grade_fr']
```

Calcul de Nutri-Score avec du machine learning -Application-

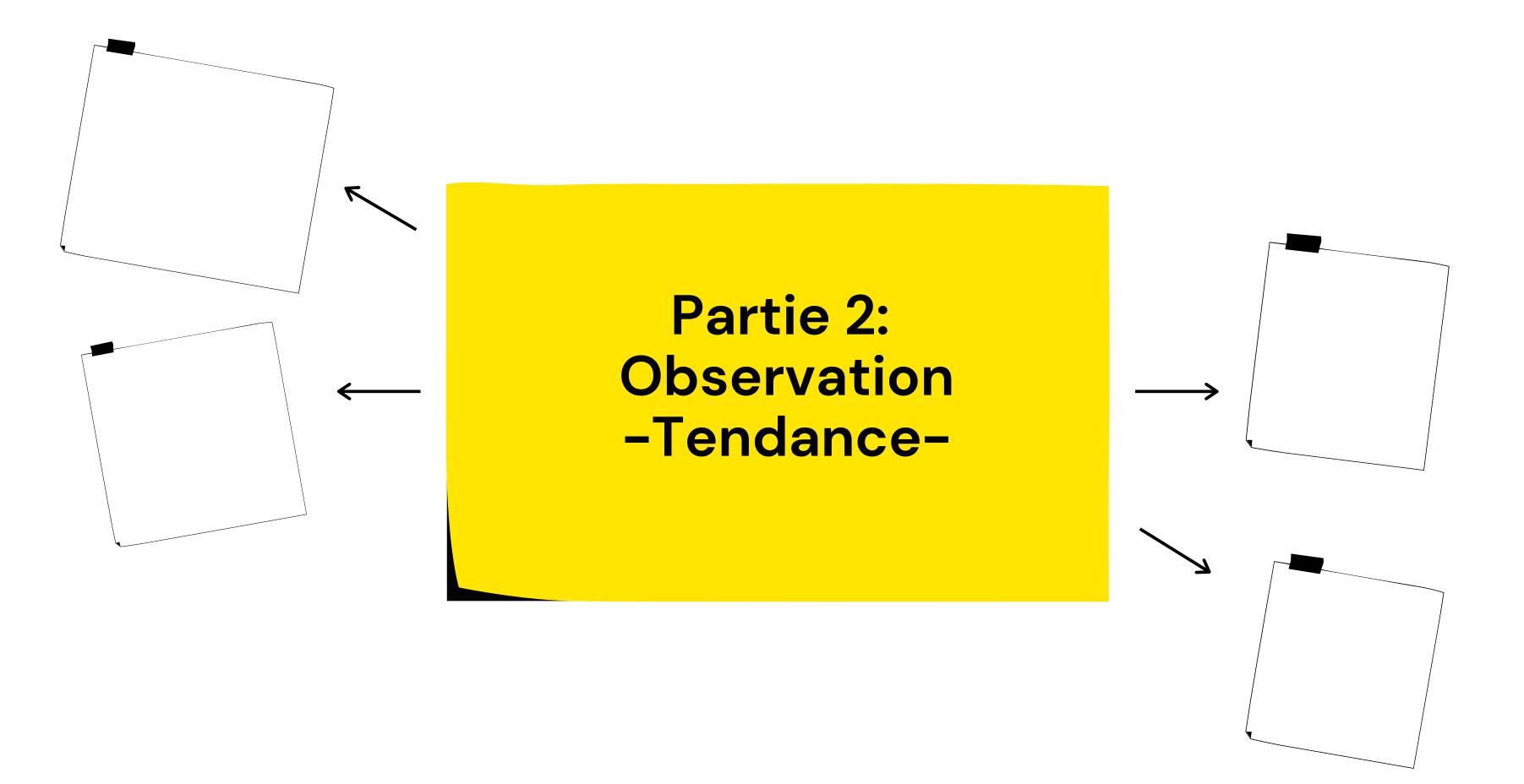
```
from sklearn.impute import SimpleImputer
                                                              1 - Utilisation d'un modèle
from sklearn.linear model import LinearRegression
# Création de l'imputer
imputer = SimpleImputer(strategy='mean') # Vous pouvez choisir 'median' ou une autre stratégie si nécessaire
# Application de l'imputation sur X train
X_train_imputed = imputer.fit_transform(X_train)
# Entraînement du modèle sur les données imputées
                                              2 - Entrainement et application
model = LinearRegression()
model.fit(X train imputed, y train)
# Utiliser le modèle pour prédire les valeurs manquantes de y
X missing = dataFrance without y[['salt 100g', 'proteins 100g', 'fiber 100g', 'sugars 100g', 'carbohydrates 100g', 'saturated-fat 100g', 'fat 100g', 'en
X_missing_imputed = imputer.transform(X_missing) # Utiliser le même imputer pour garantir la cohérence
predicted y = model.predict(X missing imputed)
# Utiliser le modèle pour prédire les valeurs manquantes de y
predicted y = model.predict(X missing imputed)
# Arrondir les prédictions au chiffre supérieur et appliquer les contraintes de plage
predicted_y_rounded = np.ceil(predicted_y)
predicted_y_constrained = np.clip(predicted_y_rounded, 1, 5)

B - Harmonise et valide
# Ajouter les prédictions ajustées à l'ensemble de données sans y
dataFrance_without_y['nutrition_grade_fr'] = predicted_y_constrained
# Combiner les deux parties du jeu de données
dataFrance = pd.concat([dataFrance_with_y, dataFrance_without_y])
```

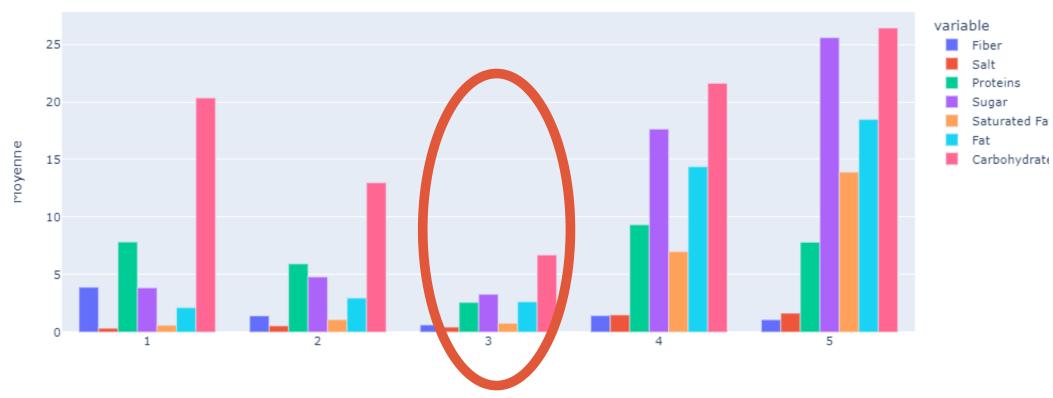
Travail sur les allergènes

Objectif: Savoir le nombre d'allergènes d'un produit.

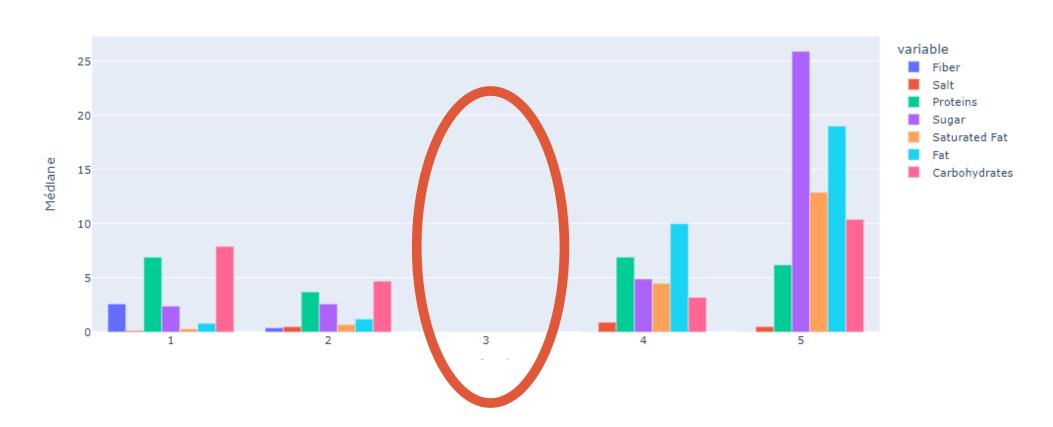
- Mettre en minuscule.
- Supprimer les espaces.
- Remplacer les termes inappropriés : 'creme' 'lait'.
- Supprimer les termes non remplaçable : 'couscous' .
- Supprimer les répétitions.
- Standardiser l'énumération ', ' (virgule + espace).
- Création de Nbr\_Allergènes et calcul.
- Supprimer les lignes avec + de 15.



Moyenne des valeurs nutritionnelles par catégorie de notation nutritionnelle



Médiane des valeurs nutritionnelles par catégorie de notation nutritionnelle



Moyennes ou Médianes ?

## Moyenne



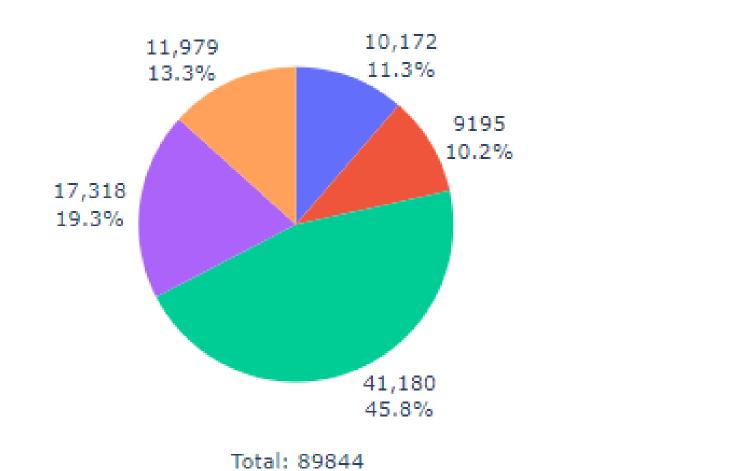
#### La situation

Nous avons tous les produits vendu en France

Nous avons calculé tous les Nutri-Scores

Nous avons connaissons les allergènes et leurs nombres pour chaque produits.

### Répartition des Nutri-Scores





Nutri-Score A

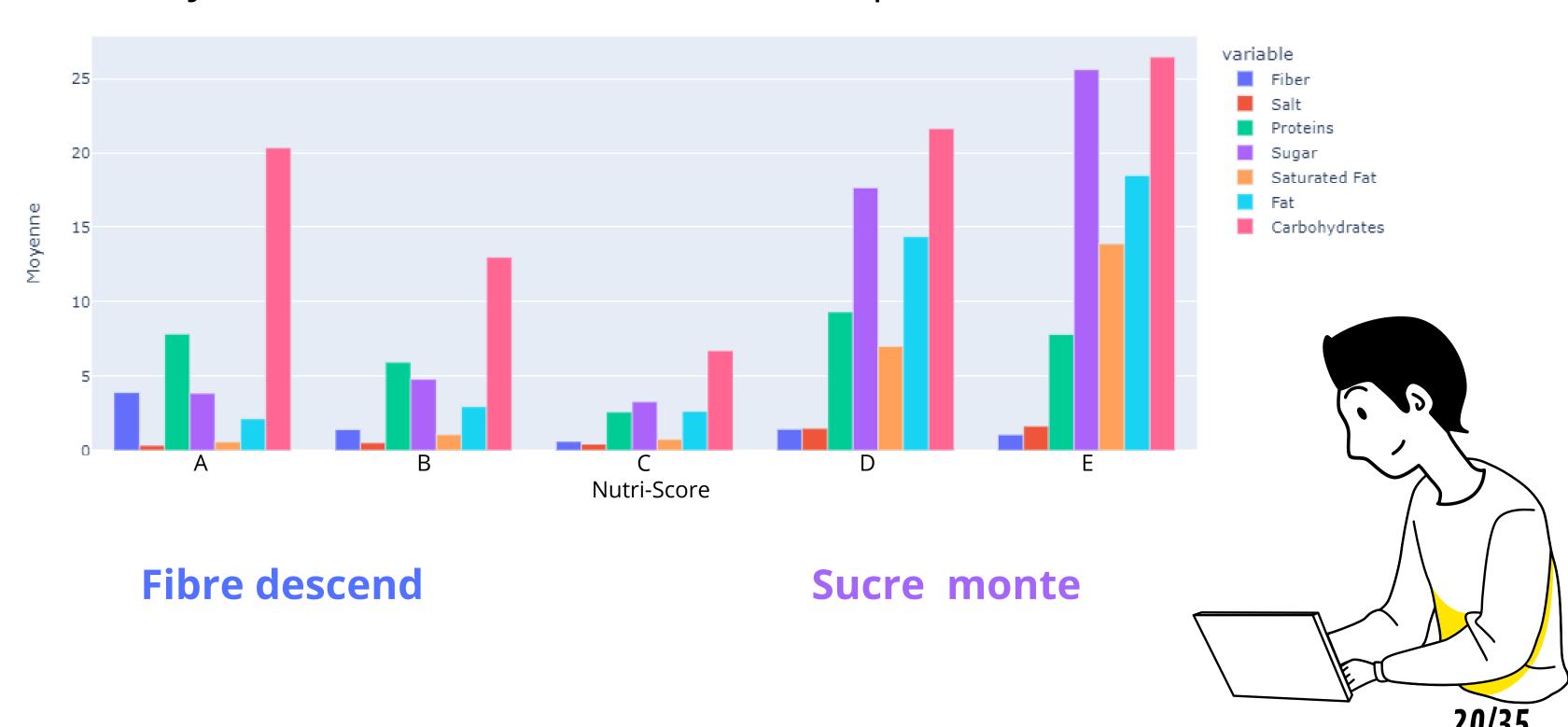
Nutri-Score B

Nutri-Score C

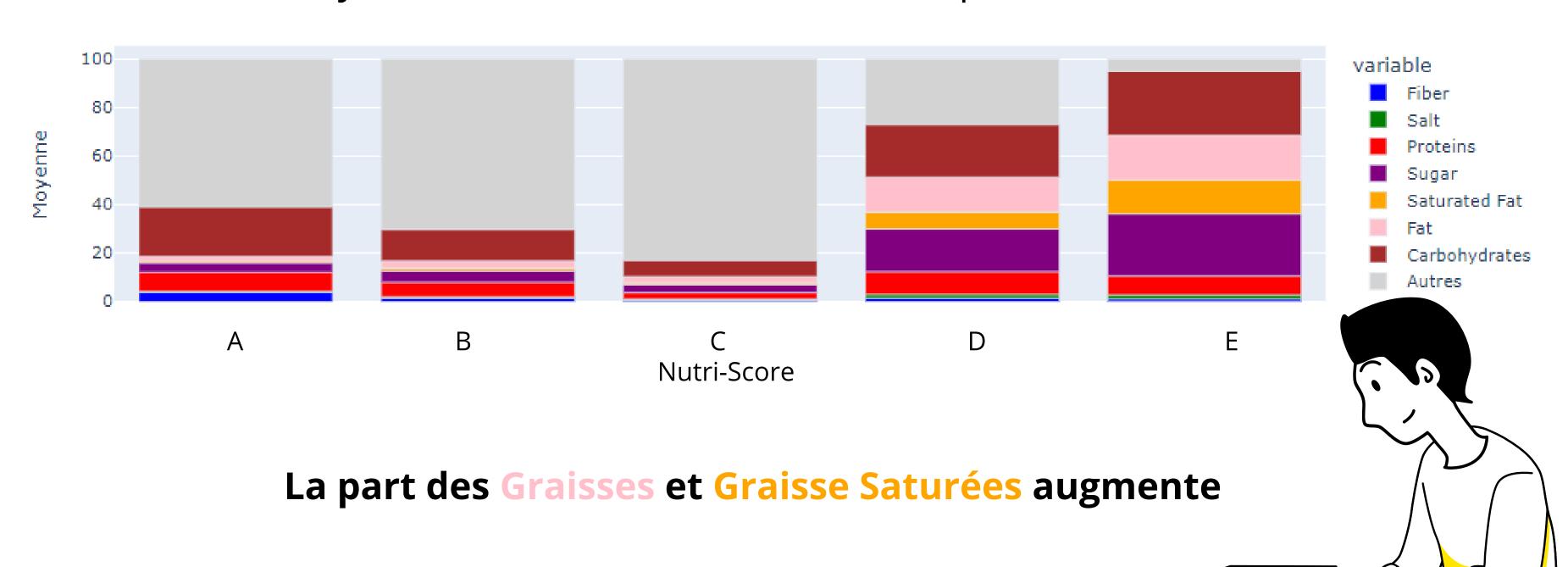
Nutri-Score D

Nutri-Score E

Moyennes des valeurs nutritionnelles par Nutri-Score

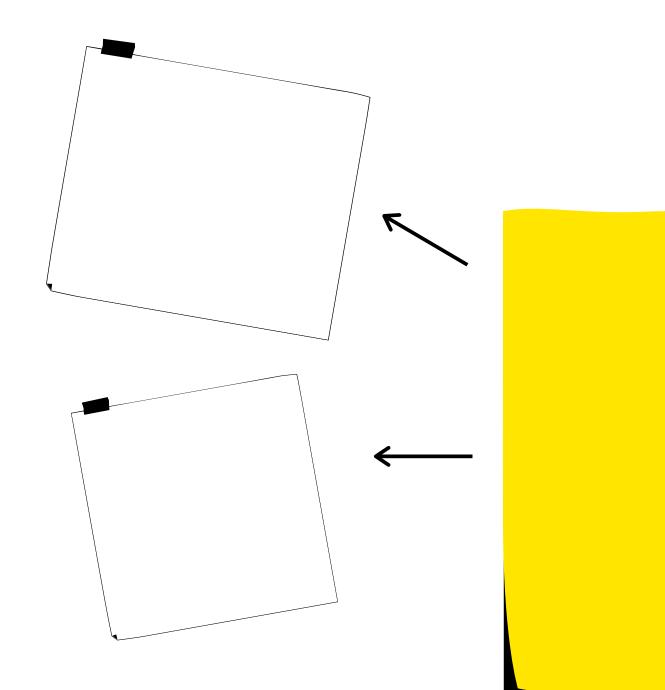


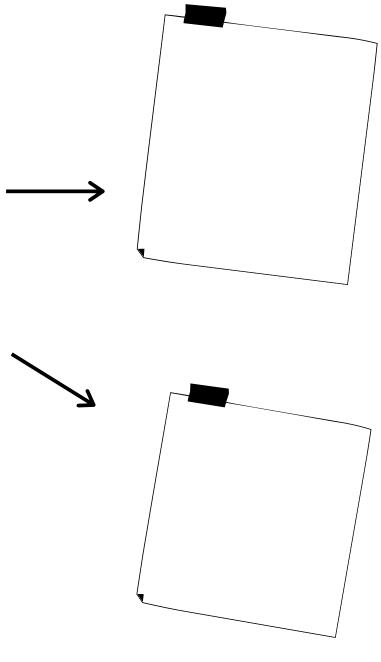
Moyennes des valeurs nutritionnelles par Nutri-Score



Proportion de présence d'allergènes dans chaque Nutri-Score







Mise en place

names names	= dataFrance.nutrition_grade_fr
2	5.0
9	4.0
10	4.0
14	5.0
16	4.0

Par rapport au

Nutri-Score

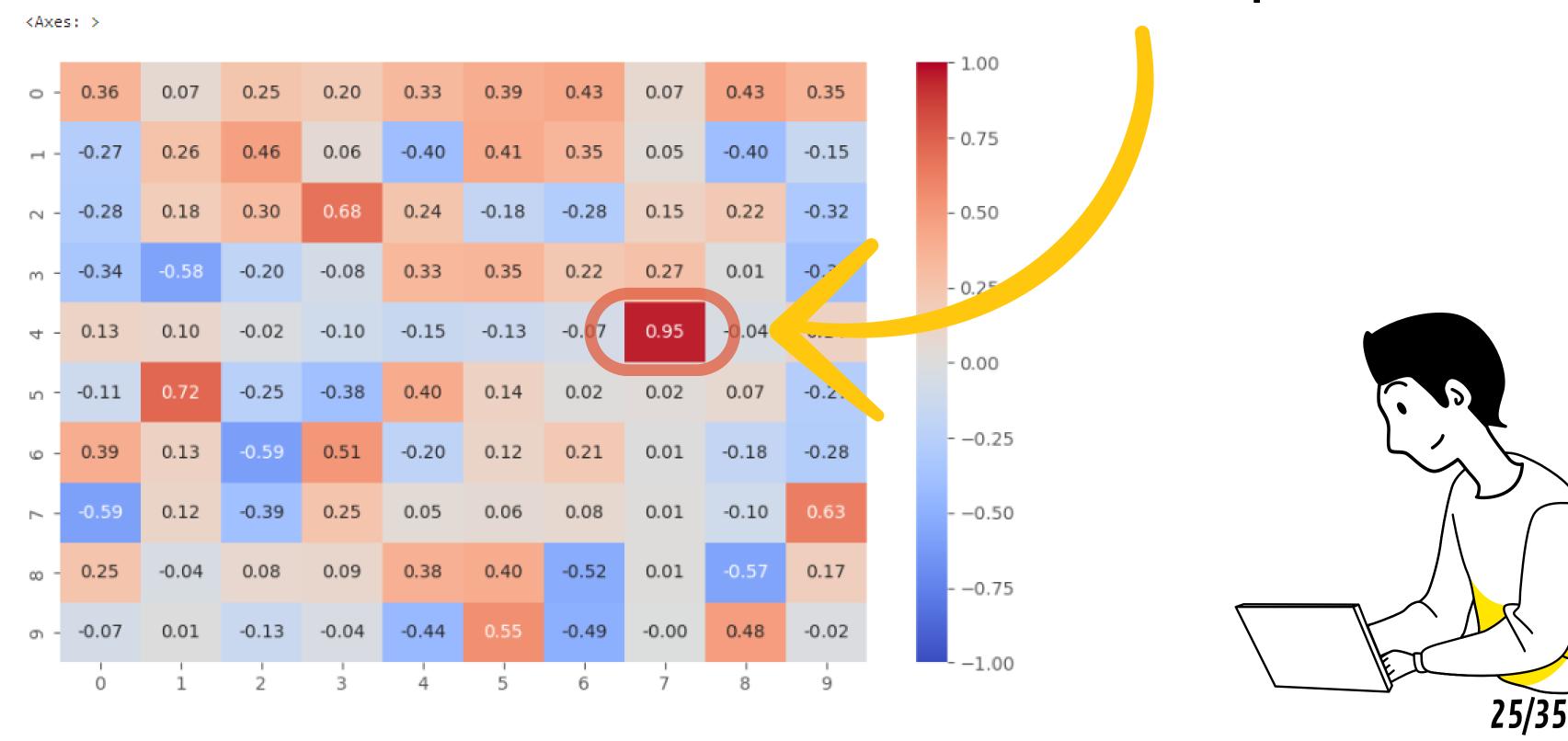
du produit

salt_100g	proteins_100g	fiber_100g	sugars_100g	saturated-fat_100g	fat_100g	energy-100g	carbohydrates_100g	Nbr_Allergenes
0.09652	2.5	2.5	57.5	12.5	20.0	1883.0	70.0	0
0.01000	0.6	0.9	87.7	0.8	0.0	1753.0	0.0	0
0.00300	9.5	3.9	50.3	2.9	0.0	2406.0	0.0	0
0.02540	0.0	0.0	10.4	0.0	0.0	177.0	10.4	0
0.80000	7.5	1.4	1.0	11.0	0.0	1079.0	0.0	0
0.00000	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0

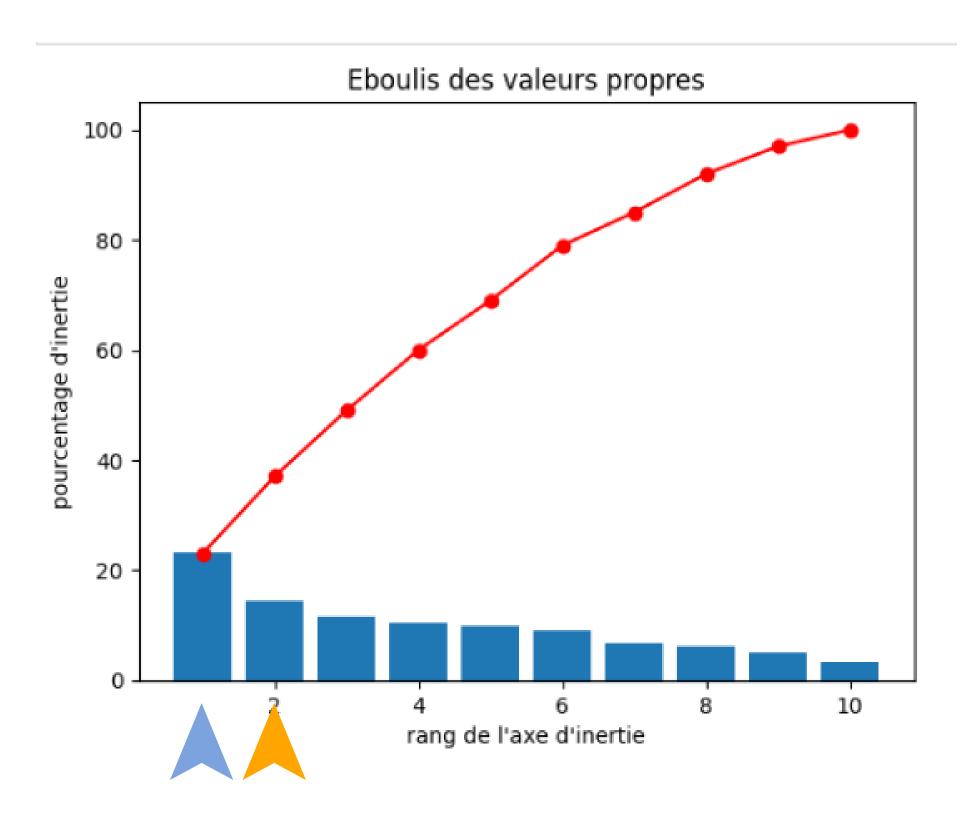
A quoi correspond : sel, protéine, fibre, graisse, graisse saturée, énergie, glucide et d'allergènes

Observation de tendance.

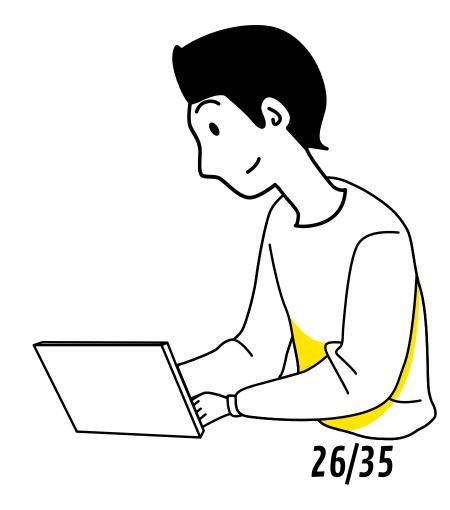
#### Corrélation forte et positive



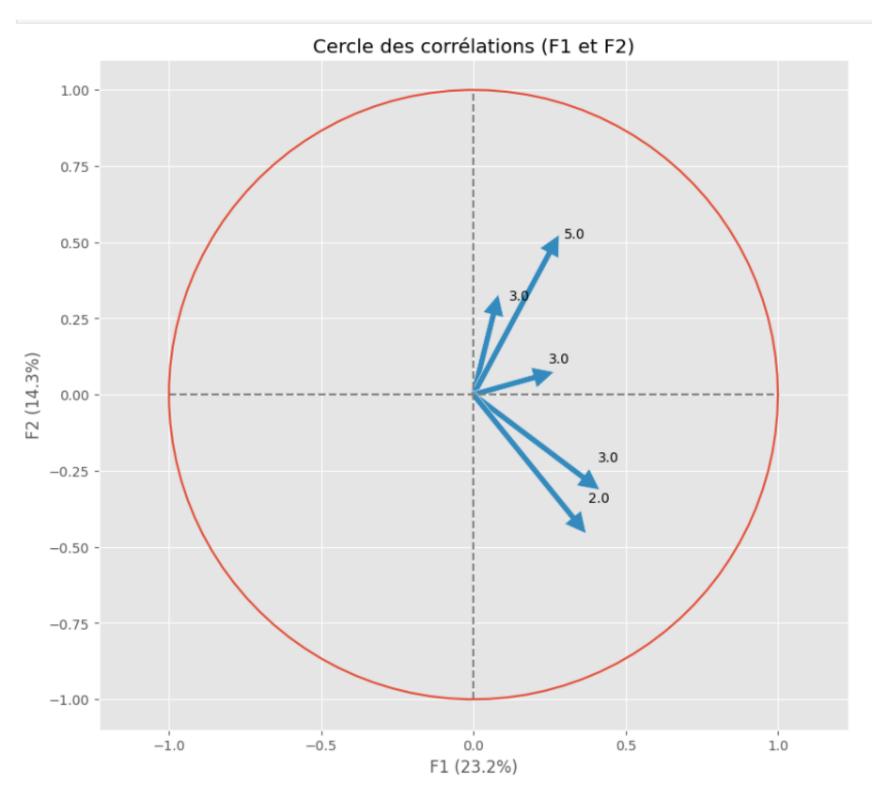
Observation de tendance.



La majorité de l'information est capturée par le premier composant principal (F1) et, dans une moindre mesure, par le deuxième (F2).



Observation de tendance.



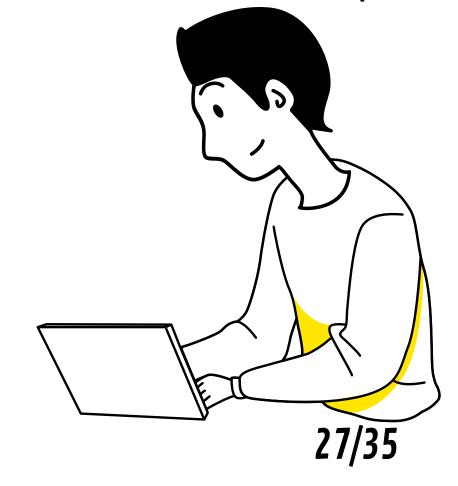


Les variables corrèlent avec F1

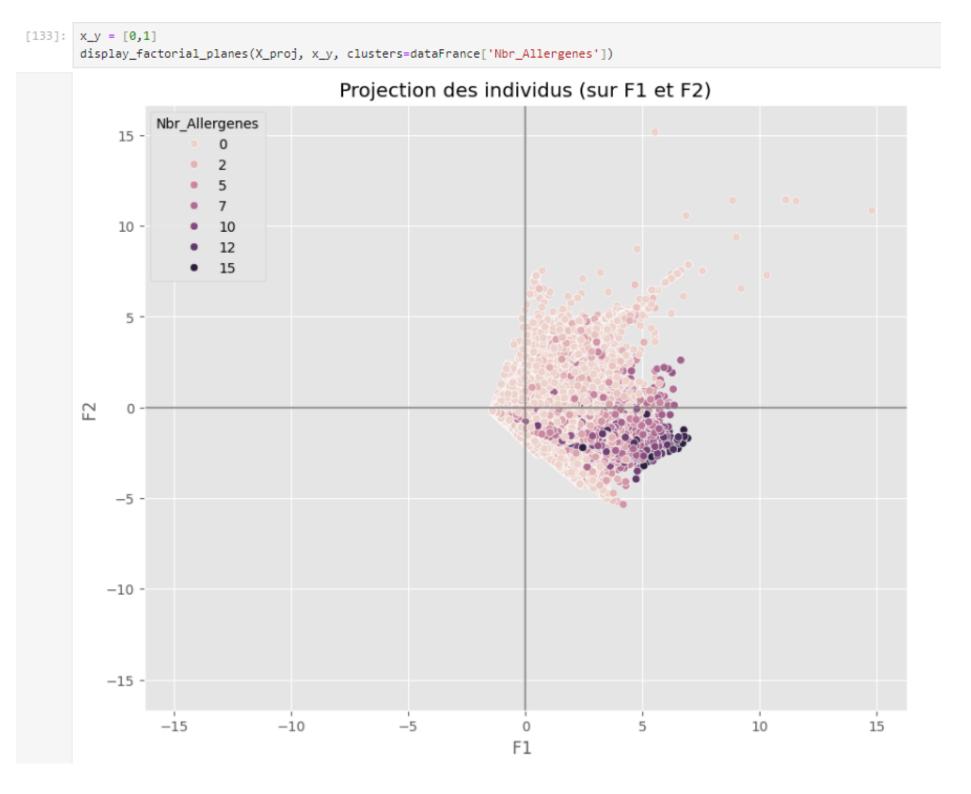
#### **MAIS**

- 0 -

Attention aux noms identiques



Observation de tendance.



Nuage, principalement autour de F1

+ il y a d'allergène, + le nuage se dirige en bas à droite. (F2)

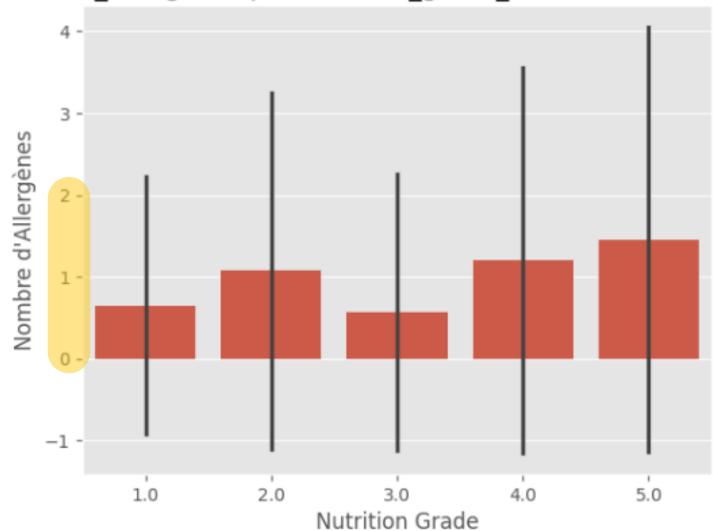


Observation de tendance.



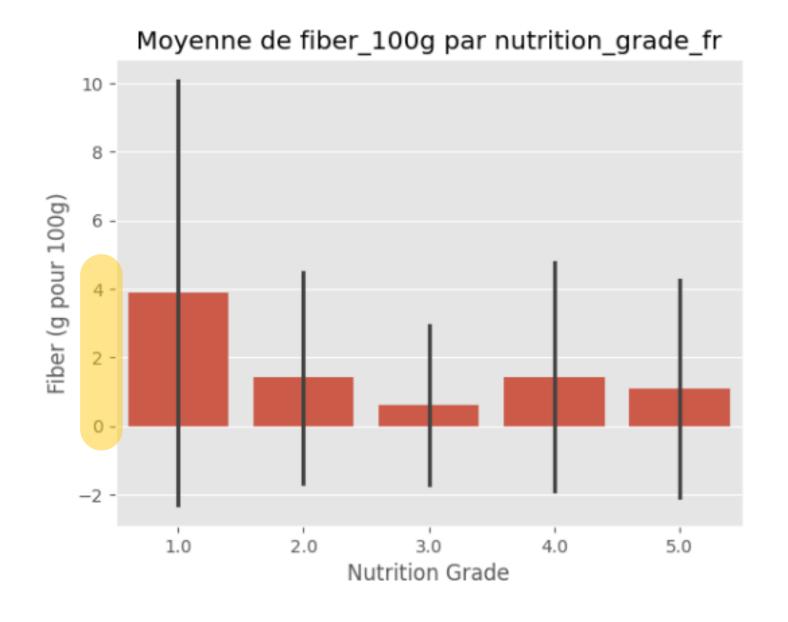
Observation de tendance.

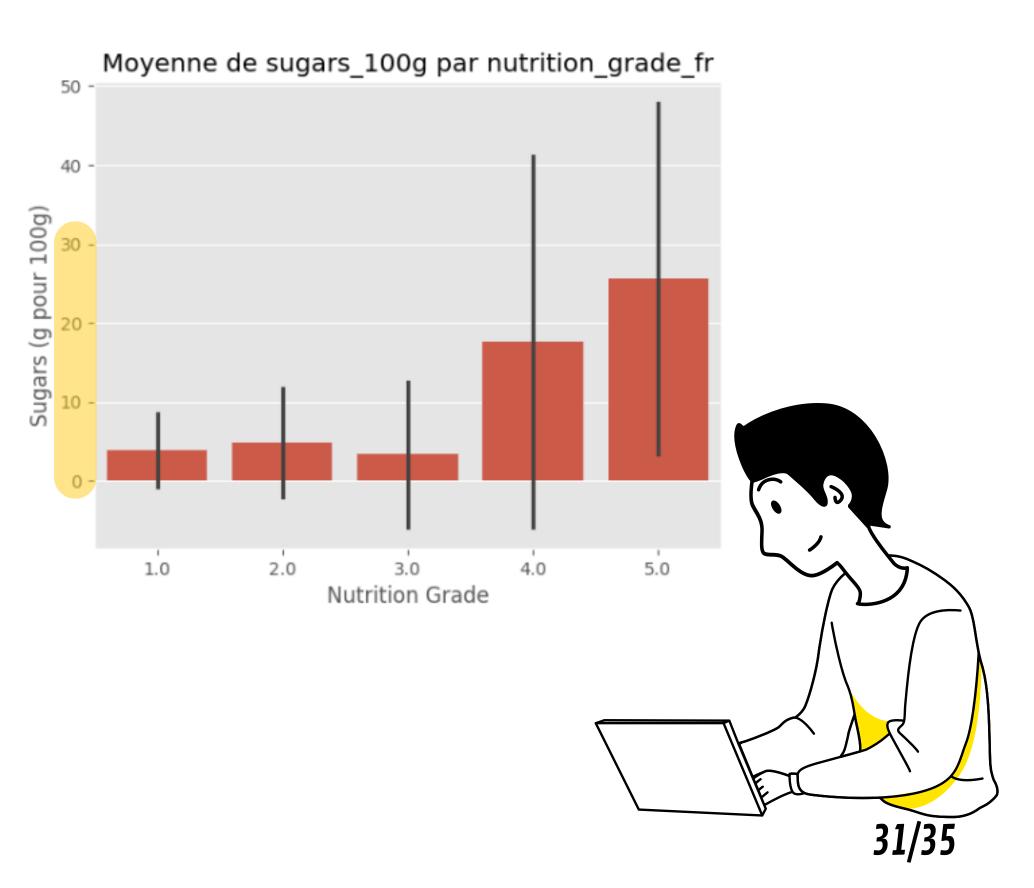
Moyenne de Nbr\_Allergenes par nutrition\_grade\_fr avec intervalles de confiance





Observation de tendance.





Observation de tendance.

						====	
Dep. Variable:	Nbr_Al	llergenes	R-squared:		0	.028	
Model:		OLS	Adj. R-squar	red:	0	.028	
Method:	Least	Squares	F-statistic:	:	65	56.5	
Date:	Thu, 18	Jan 2024	Prob (F-stat	tistic):		0.00	
Time:		15:51:10	Log-Likeliho	ood:	-1.9018	e+05	
No. Observations:		89844	AIC:		3.804	e+05	
Df Residuals:		89839	BIC:		3.804	e+05	
Df Model:		4					
Covariance Type:	r	nonrobust					
		coef	std err	t	P> t	[0.025	0.975
Intercept		0.6509	0.020	32.670	0.000	0.612	0.69
C(nutrition_grade_fr	)[T.2.0]	0.4202	0.029	14.533	0.000	0.364	0.47
C(nutrition_grade_fr	)[T.3.0]	-0.0828	0.022	-3.720	0.000	-0.126	-0.039
C(nutrition_grade_fr	)[T.4.0]	0.5485	0.025	21.852	0.000	0.499	0.598
C(nutrition_grade_fr	)[T.5.0]	0.8049	0.027	29.708	0.000	0.752	0.85
						====	
Omnibus:			Durbin-Watso		1		
Prob(Omnibus):				(JB):			
Skew:			Prob(JB):			0.00	
Kurtosis:		13.406	Cond. No.		8	8.03	

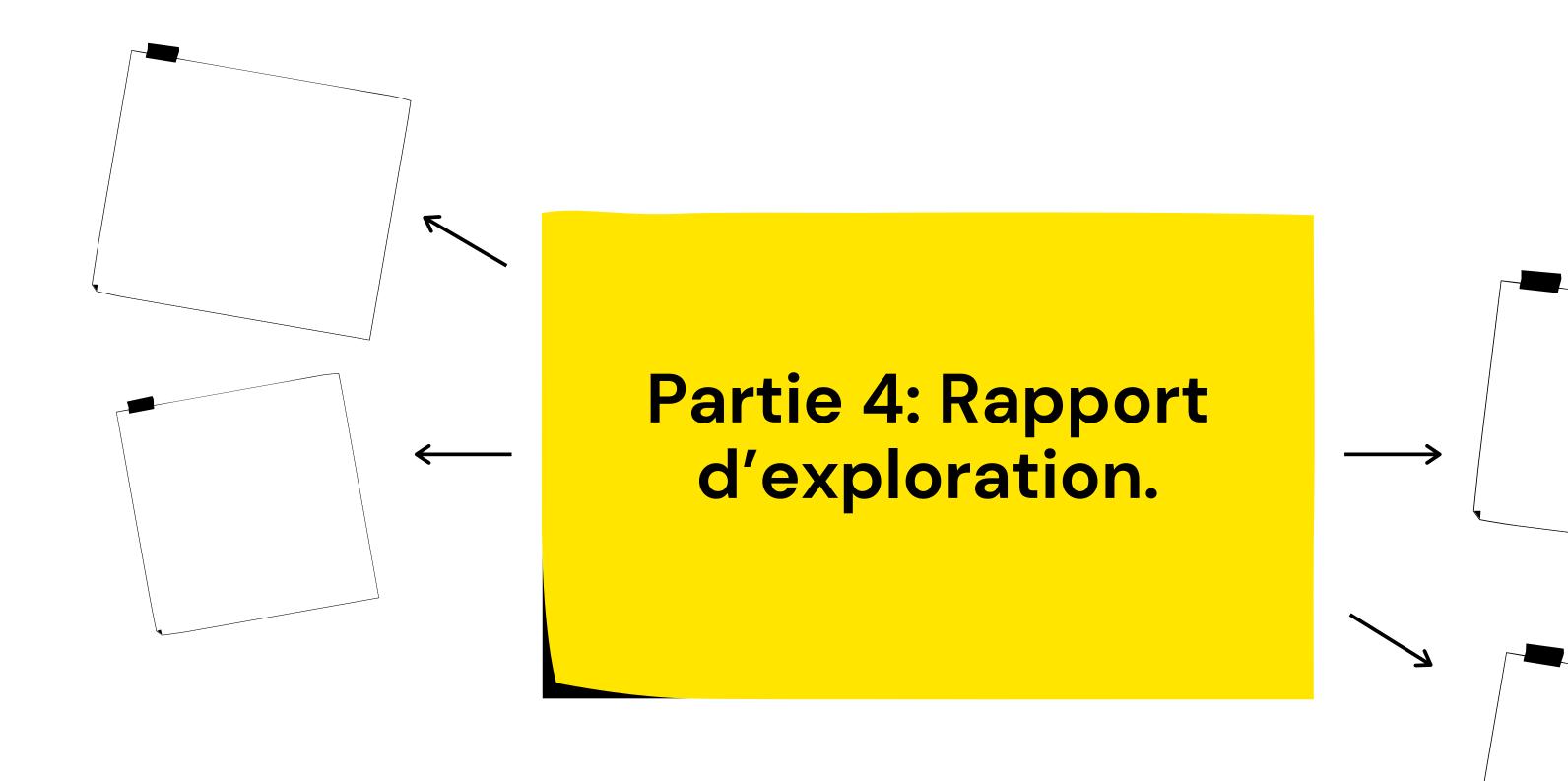
#### Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

# Afficher les résultats du test de Tukey
print(tukey.summary())

Multiple Comparison of Means - Tukey HSD, FWER=0.05

group1	group2	meandiff	p-adj	lower	upper	reject
1.0	2.0	0.4202	0.0	0.3413	0.4991	True
1.0	3.0	-0.0828	0.0019	-0.1435	-0.0221	True
1.0	4.0	0.5485	0.0	0.4801	0.617	True
1.0	5.0	0.8049	0.0	0.731	0.8788	True
2.0	3.0	-0.503	0.0	-0.5662	-0.4398	True
2.0	4.0	0.1283	0.0	0.0576	0.199	True
2.0	5.0	0.3847	0.0	0.3087	0.4607	True
3.0	4.0	0.6313	0.0	0.5817	0.681	True
3.0	5.0	0.8877	0.0	0.8308	0.9446	True
4.0	5.0	0.2564	0.0	0.1912	0.3215	True



### Partie 4: Rapport d'exploration.

L'application est-elle faisable ?

Oui

Les premières pistes ergonomiques:

- Liste déroulante
- Pré-remplissage
- Utilisation d'IA d'analyse d'image

Pistes pour de futurs analyses du futurs jeu de donnés:

- Nombres d'ingrédients
- Croisement JdD Additif
- Simplification
- Mise en avant de l'allergènes



## Merci!

