

```
In [ ]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import missingno as msno
import plotly.express as px
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.neighbors import KNeighborsClassifier
from sklearn.decomposition import PCA
pd.options.display.max_rows = 100
pd.options.mode.copy_on_write = True
```

```
In [ ]: df = pd.read_csv('fr.openfoodfacts.org.products.csv', sep = '\t', low_memory=False)
df.describe()
```

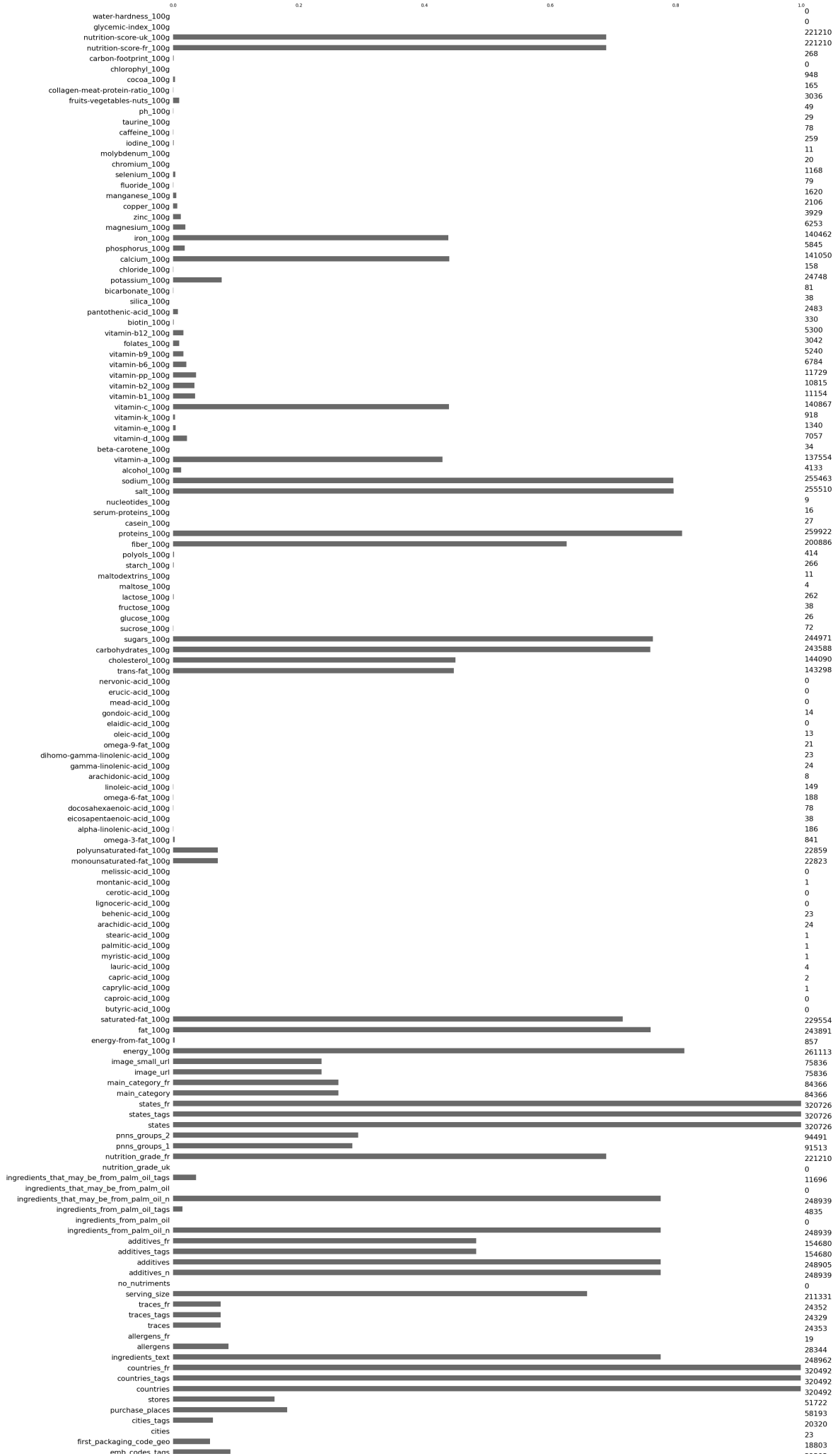
Out[]:

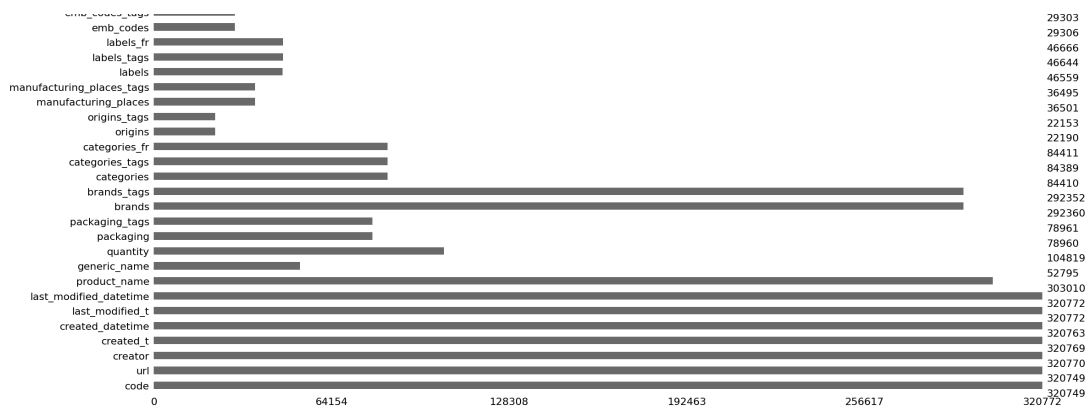
	no_nutriments	additives_n	ingredients_from_palm_oil_n	ingredients_from_palr
count	0.0	248939.000000		248939.000000
mean	NaN	1.936024		0.019659
std	NaN	2.502019		0.140524
min	NaN	0.000000		0.000000
25%	NaN	0.000000		0.000000
50%	NaN	1.000000		0.000000
75%	NaN	3.000000		0.000000
max	NaN	31.000000		2.000000

8 rows × 106 columns

```
In [ ]: %matplotlib inline
msno.bar(df)
```

Out[]: <Axes: >





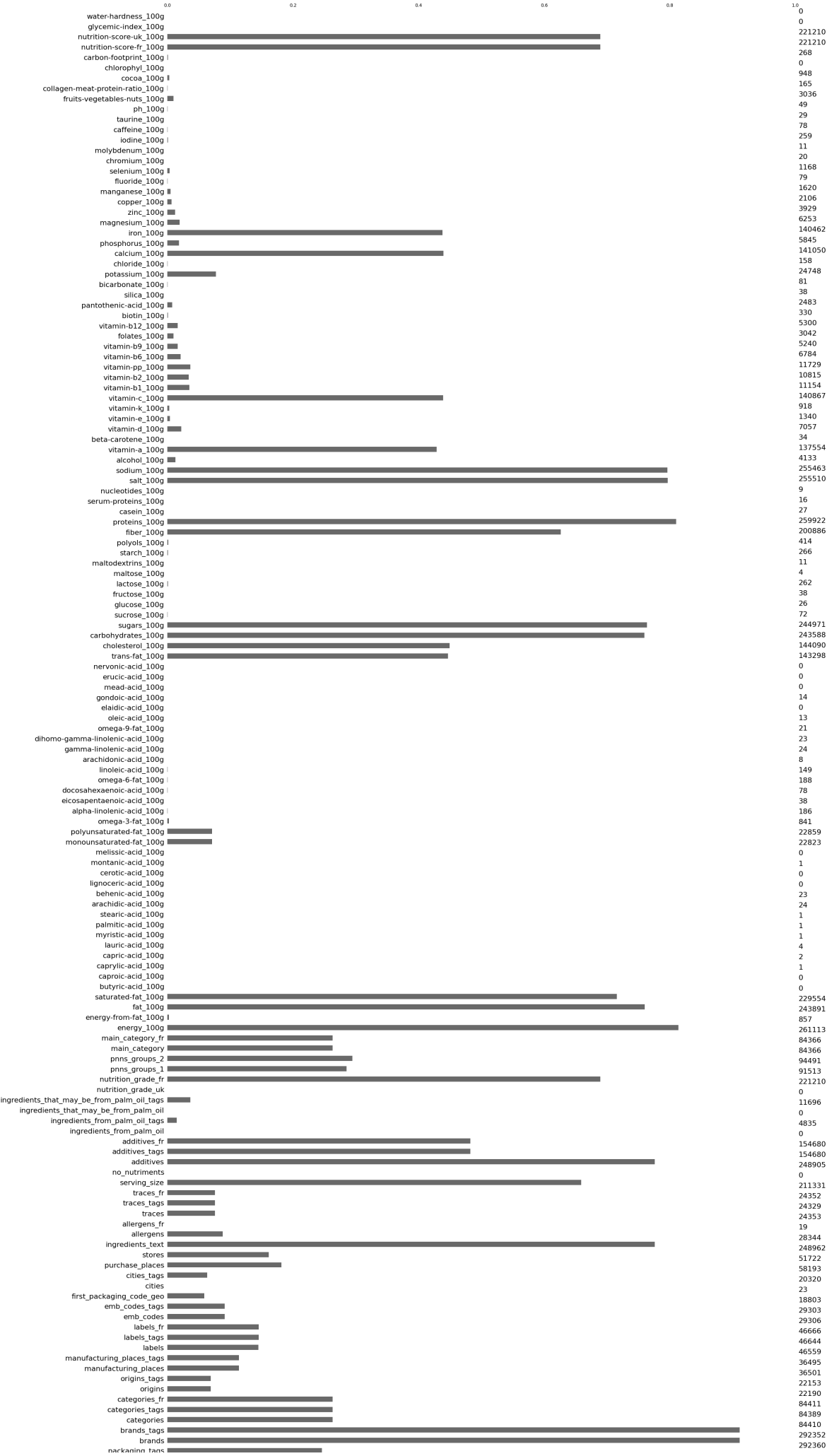
```
In [ ]: df.columns
```

```
Out[ ]: Index(['code', 'url', 'creator', 'created_t', 'created_datetime',
              'last_modified_t', 'last_modified_datetime', 'product_name',
              'generic_name', 'quantity',
              ...,
              'ph_100g', 'fruits-vegetables-nuts_100g',
              'collagen-meat-protein-ratio_100g', 'cocoa_100g', 'chlorophyl_100g',
              'carbon-footprint_100g', 'nutrition-score-fr_100g',
              'nutrition-score-uk_100g', 'glycemic-index_100g',
              'water-hardness_100g'],
              dtype='object', length=162)
```

```
In [ ]: df.drop(columns=['url', 'creator', 'created_t', 'created_datetime', 'last_modified_t',
                        'last_modified_datetime', 'product_name', 'generic_name', 'quantity',
                        'ph_100g', 'fruits-vegetables-nuts_100g', 'collagen-meat-protein-ratio_100g',
                        'cocoa_100g', 'chlorophyl_100g', 'carbon-footprint_100g', 'nutrition-score-fr_100g',
                        'nutrition-score-uk_100g', 'glycemic-index_100g', 'water-hardness_100g'],
                  inplace=True)
df.drop(columns=['code'], inplace=True)
```

```
In [ ]: msno.bar(df)
```

```
Out[ ]: <Axes: >
```





```
In [ ]: # Calculer Le nombre total de lignes
total_rows = len(df)

# Calculer Le nombre de valeurs non nulles dans chaque colonne
non_null_counts = df.notna().sum()

# Calculer Le taux de remplissage de chaque colonne
filling_rates = non_null_counts / total_rows * 100

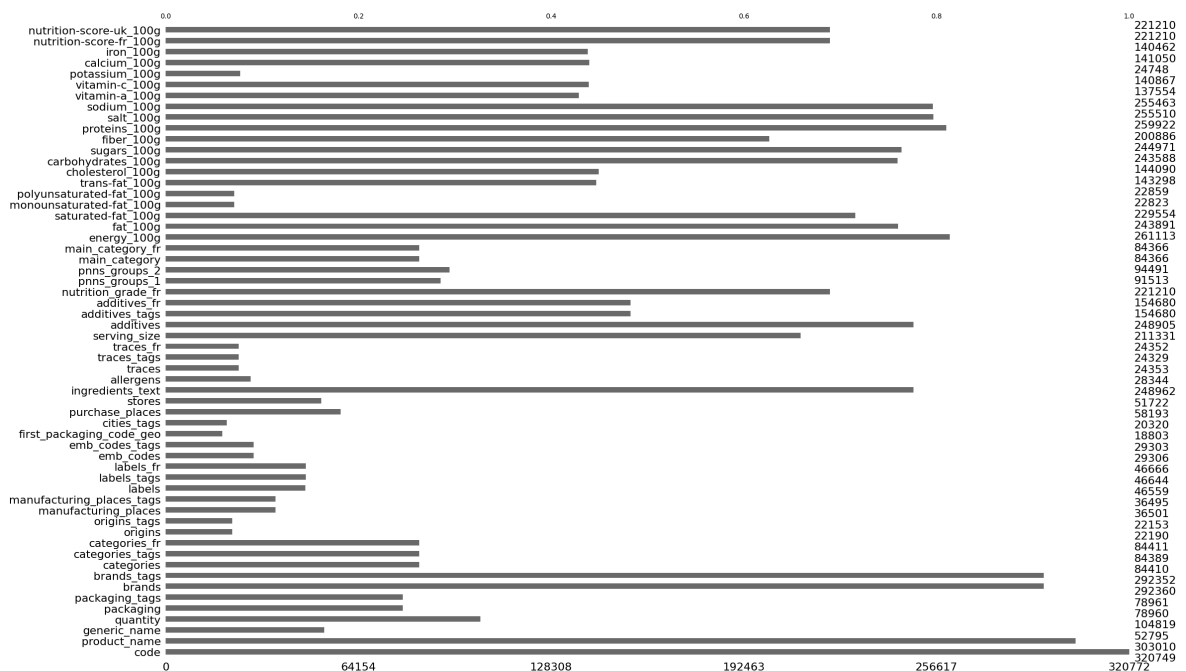
# Filtrer Les colonnes avec un taux de remplissage entre 20% et 50%
selected_columns = filling_rates[(filling_rates >= 20) && (filling_rates <= 50)].index

# Conserver uniquement Les colonnes sélectionnées
df_cleaned = df[selected_columns]
len(df_cleaned)
```

Out[]: 320772

```
In [ ]: msno.bar(df_cleaned)
```

Out[]: <Axes: >



```
In [ ]: df_cleaned.drop_duplicates(subset=['code'])
len(df_cleaned)
```

Out[]: 320772

```
In [ ]: df_cible=df_cleaned.dropna(subset=['pnns_groups_1'])
df_cible_cleaned = df_cible[df_cible['pnns_groups_1']!= 'unknown']
len(df_cible)
```

Out[]: 91513

```
In [ ]: df_pnns1 = df_cible_cleaned['pnns_groups_1']
df_pnns1
```

```
Out[ ]: 174      Fruits and vegetables
        175          Sugary snacks
        177      Cereals and potatoes
        180          Sugary snacks
        182      Cereals and potatoes
        ...
        320759  Fruits and vegetables
        320763          Beverages
        320765      Fish Meat Eggs
        320766          Salty snacks
        320769          Salty snacks
        Name: pnns_groups_1, Length: 68889, dtype: object
```

```
In [ ]: counts = df_pnns1.value_counts().reset_index()
        counts.columns = ['pnns_groups_1', 'Counts']
        print(df_pnns1.value_counts())
        # Créer Le graphique avec Plotly
        fig = px.bar(counts, x='pnns_groups_1', y='Counts', title='Value Counts de la va

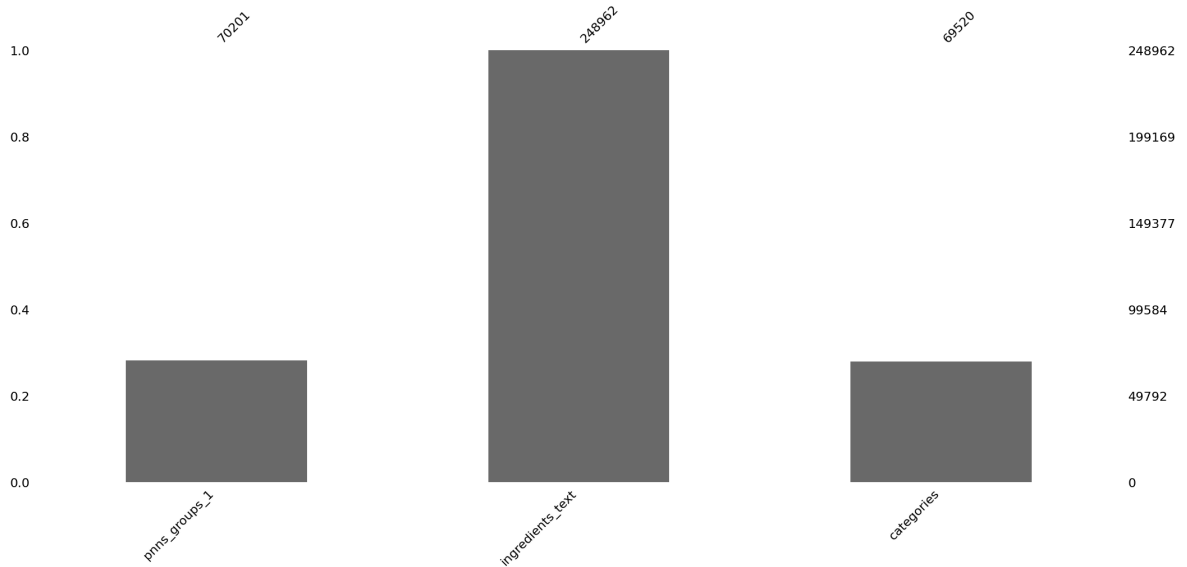
        # Afficher Le graphique
        fig.show()
```

```
pnns_groups_1
Sugary snacks      12368
Beverages          9033
Milk and dairy products  8825
Cereals and potatoes  8442
Fish Meat Eggs     8041
Composite foods    6747
Fruits and vegetables  5908
Fat and sauces     5216
Salty snacks       2809
fruits-and-vegetables  987
sugary-snacks       496
cereals-and-potatoes  16
salty-snacks        1
Name: count, dtype: int64
```

```
In [ ]: colonnes = ['pnns_groups_1', 'ingredients_text', 'categories']
        df_select = df_cleaned.dropna(subset=['ingredients_text'])
        df_travail = df_select[colonnes]

        msno.bar(df_travail)
```

```
Out[ ]: <Axes: >
```



remplacer les valeurs incohérentes

```
In [ ]: df_travail['pnns_groups_1']=df_travail.pnns_groups_1.replace({"-":" "},regex=True)
df_travail['pnns_groups_1']=df_travail['pnns_groups_1'].str.capitalize()
df_travail['pnns_groups_1'].value_counts()
```

```
Out[ ]: pnns_groups_1
Unknown          12366
Sugary snacks    10865
Milk and dairy products    7336
Cereals and potatoes    7152
Fish meat eggs    6980
Beverages        6825
Fruits and vegetables    6032
Composite foods    6025
Fat and sauces    4226
Salty snacks     2394
Name: count, dtype: int64
```

enlever les unknown

```
In [ ]: df_travail.drop(df_travail[df_travail['pnns_groups_1'] == 'Unknown'].index, inplace=True)
df_travail['pnns_groups_1'].value_counts()
```

```
Out[ ]: pnns_groups_1
Sugary snacks    10865
Milk and dairy products    7336
Cereals and potatoes    7152
Fish meat eggs    6980
Beverages        6825
Fruits and vegetables    6032
Composite foods    6025
Fat and sauces    4226
Salty snacks     2394
Name: count, dtype: int64
```

```
In [ ]: #df d'entrainement avec pnns group
df_train = df_travail.dropna(subset=['pnns_groups_1'])
#DF sans pnns group
df_missing = df_travail[df['pnns_groups_1'].isna()]
```

C:\Users\devil\AppData\Local\Temp\ipykernel_5268\1862913704.py:4: UserWarning:
Boolean Series key will be reindexed to match DataFrame index.

```
In [ ]: vectorizer = CountVectorizer(tokenizer=lambda x: x.split(','), max_features=1000

X_train = vectorizer.fit_transform(df_train['ingredients_text'])
X_missing = vectorizer.transform(df_missing['ingredients_text'])

# Encodage des Labels
y_train = df_train['pnns_groups_1']
# Affichage de la matrice
```

c:\Users\devil\anaconda3\envs\PROJET3\Lib\site-packages\sklearn\feature_extraction\text.py:523: UserWarning:

The parameter 'token_pattern' will not be used since 'tokenizer' is not None'

```
In [ ]: # Initialisation du classifieur KNN
knn = KNeighborsClassifier(n_neighbors=3)

# Entraînement du modèle
knn.fit(X_train, y_train)
```

```
Out[ ]: ▼ KNeighborsClassifier ⓘ ?
KNeighborsClassifier(n_neighbors=3)
```

```
In [ ]: # Prédiction sur les données manquantes
y_pred_missing = knn.predict(X_missing)

# Remplir les valeurs manquantes dans le DataFrame original
df_travail.loc[df_travail['pnns_groups_1'].isna(), 'pnns_groups_1'] = y_pred_missing

print("\nDataFrame avec valeurs manquantes prédites:")
print(df)
```


DataFrame avec valeurs manquantes prédites:

	code	product_name \
0	0000000003087	Farine de blé noir
1	0000000004530	Banana Chips Sweetened (Whole)
2	0000000004559	Peanuts
3	0000000016087	Organic Salted Nut Mix
4	0000000016094	Organic Polenta
...
320767	9948282780603	Tomato & ricotta
320768	99567453	Mint Melange Tea A Blend Of Peppermint, Lemon ...
320769	9970229501521	乐吧泡菜味薯片
320770	9980282863788	Tomates aux Vermicelles
320771	999990026839	Sugar Free Drink Mix, Peach Tea

	generic_name	quantity	packaging	packaging_tags \
0	NaN	1kg	NaN	NaN
1	NaN	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN
3	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN
...
320767	NaN	1	Plastique	plastique
320768	NaN	NaN	Plastique	plastique
320769	Leba pickle flavor potato chips	50 g	Plastique	plastique
320770	NaN	67g	NaN	NaN
320771	NaN	NaN	Plastique	plastique

	brands	brands_tags	categories \
0	Ferme t'y R'nao	ferme-t-y-r-nao	NaN
1	NaN	NaN	NaN
2	Torn & Glasser	torn-glasser	NaN
3	Grizzlies	grizzlies	NaN
4	Bob's Red Mill	bob-s-red-mill	NaN
...
320767	Panzani	panzani	NaN
320768	Trader Joe's	trader-joe-s	NaN
320769	乐吧	乐吧	Potato chips
320770	Knorr	knorr	NaN
320771	Market Pantry	market-pantry	NaN

	categories_tags	... ph_100g \
0	NaN	NaN
1	NaN	NaN
2	NaN	NaN
3	NaN	NaN
4	NaN	NaN
...
320767	NaN	NaN
320768	NaN	NaN
320769	en:salty-snacks,en:appetizers,en:chips-and-fri...	NaN
320770	NaN	NaN
320771	NaN	NaN

	fruits-vegetables-nuts_100g	collagen-meat-protein-ratio_100g \
0	NaN	NaN
1	NaN	NaN
2	NaN	NaN
3	NaN	NaN
4	NaN	NaN
...

320767	NaN	NaN
320768	NaN	NaN
320769	NaN	NaN
320770	NaN	NaN
320771	NaN	NaN

	cocoa_100g	chlorophyl_100g	carbon-footprint_100g	\
0	NaN	NaN	NaN	
1	NaN	NaN	NaN	
2	NaN	NaN	NaN	
3	NaN	NaN	NaN	
4	NaN	NaN	NaN	
...	
320767	NaN	NaN	NaN	
320768	NaN	NaN	NaN	
320769	NaN	NaN	NaN	
320770	NaN	NaN	NaN	
320771	NaN	NaN	NaN	

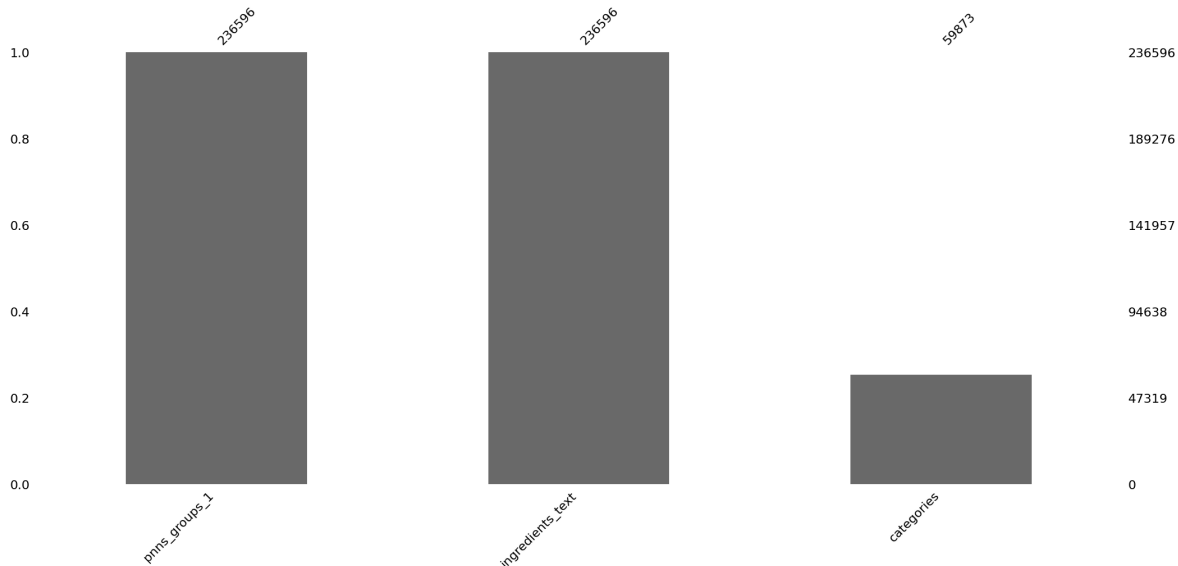
	nutrition-score-fr_100g	nutrition-score-uk_100g	glycemic-index_100g	\
0	NaN	NaN	NaN	
1	14.0	14.0	NaN	
2	0.0	0.0	NaN	
3	12.0	12.0	NaN	
4	NaN	NaN	NaN	
...	
320767	NaN	NaN	NaN	
320768	0.0	0.0	NaN	
320769	NaN	NaN	NaN	
320770	NaN	NaN	NaN	
320771	NaN	NaN	NaN	

	water-hardness_100g
0	NaN
1	NaN
2	NaN
3	NaN
4	NaN
...	...
320767	NaN
320768	NaN
320769	NaN
320770	NaN
320771	NaN

[320772 rows x 145 columns]

```
In [ ]: msno.bar(df_travail)
```

Out[]: <Axes: >



```
In [ ]: df_travail['pnns_groups_1'].value_counts()
```

```
Out[ ]: pnns_groups_1
Cereals and potatoes      75723
Beverages                 37575
Fat and sauces            36028
Sugary snacks             32117
Fruits and vegetables     13614
Milk and dairy products   12176
Fish meat eggs            11740
Salty snacks              9718
Composite foods           7905
Name: count, dtype: int64
```