

Microsoft

Préparez des données pour un organisme de santé publique

OpenClassrooms

Benjamin Demaille

Les données

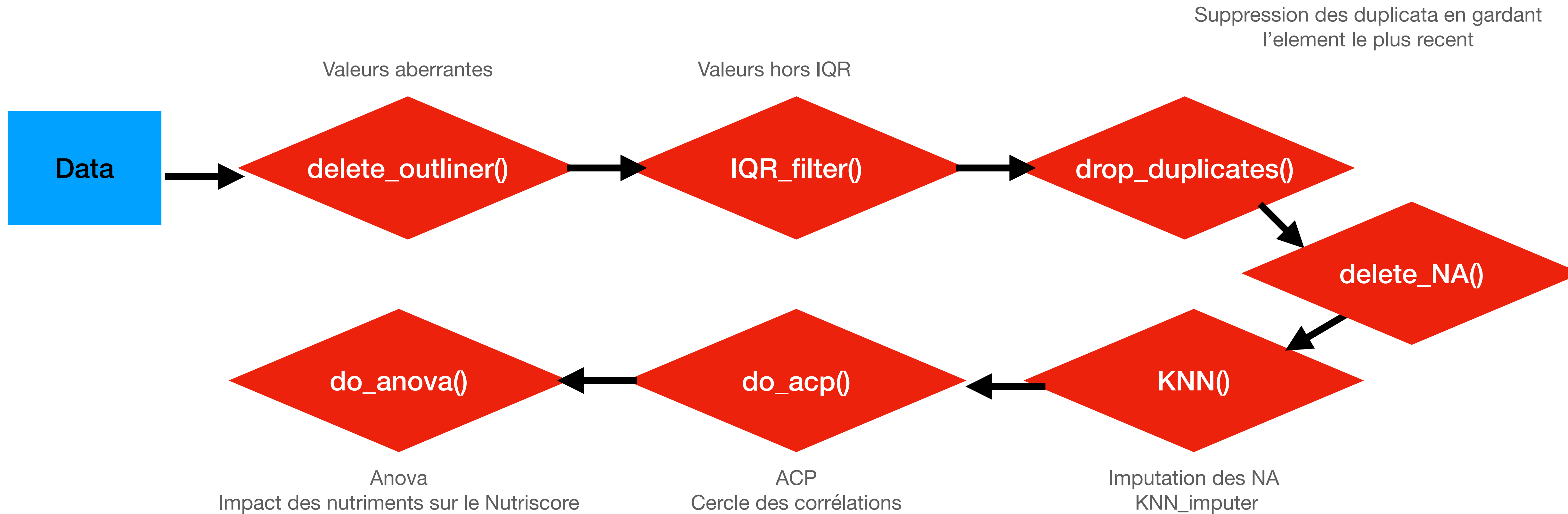
Vue générale

162 variables

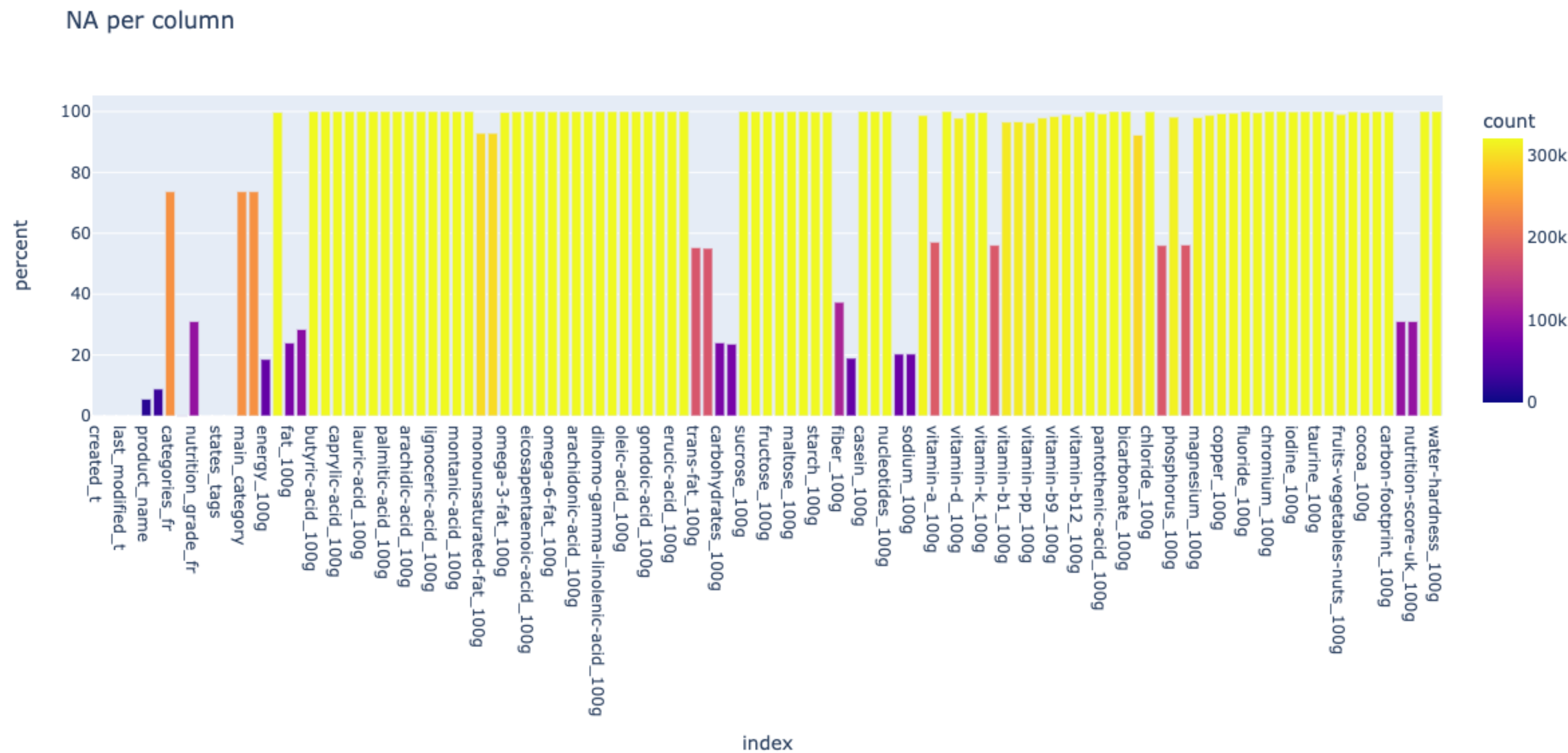
320 749
entrées

[illegible]

Pipeline



NA avant traitement



delete_outlier()

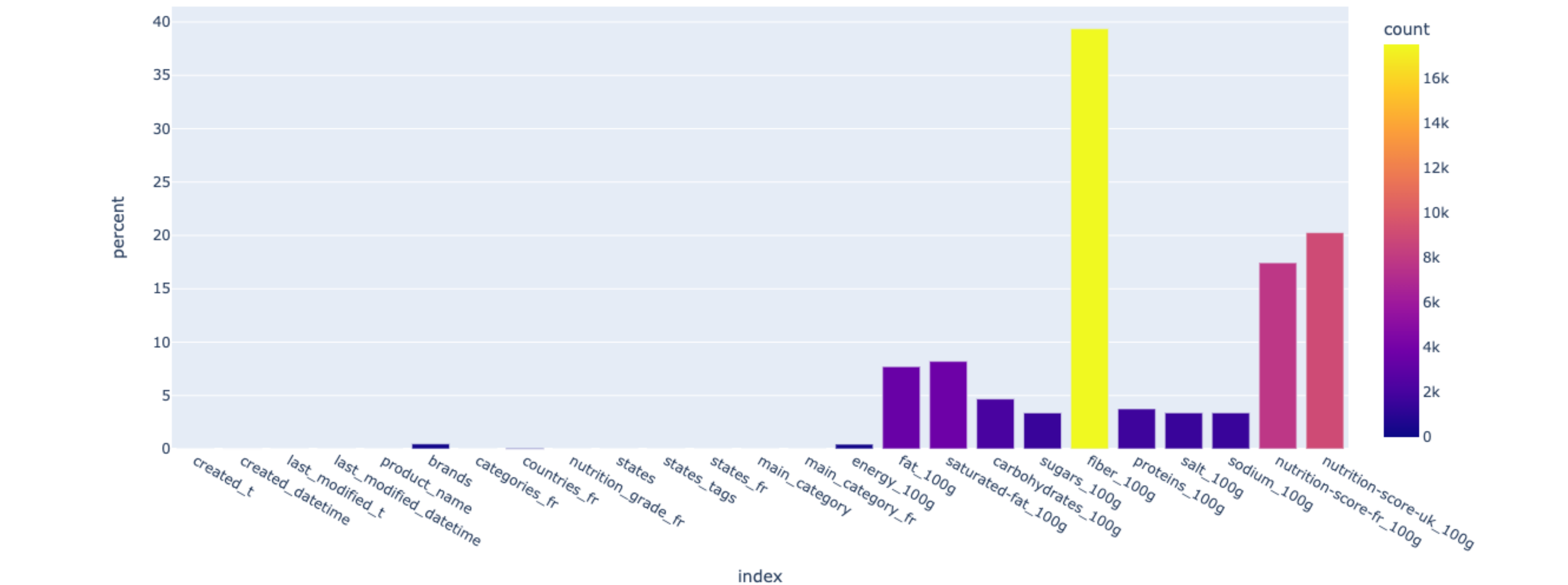
```
def delete_outliers(data: pd.DataFrame, Nutritive_columns: list):  
    for col in Nutritive_columns:  
        data.loc[:,col] = data[col].where(data[col] >= 0)  
        data.loc[:,col] = data[col].where(data[col] <= 100)  
    return(data)
```

IQR_filter()

```
def IQR_filter(s: pd.DataFrame, col_num: list, replace=np.nan):  
    for col in col_num:  
        Q1 = s[col].quantile(0.25)  
        Q3 = s[col].quantile(0.75)  
        IQR = Q3-Q1  
        s.loc[:,col] = s[col].where((s[col] > (Q1 - 1.5 * IQR)) & (s[col] < (Q3 + 1.5 * IQR)))  
    return(s)
```

NA après filtre

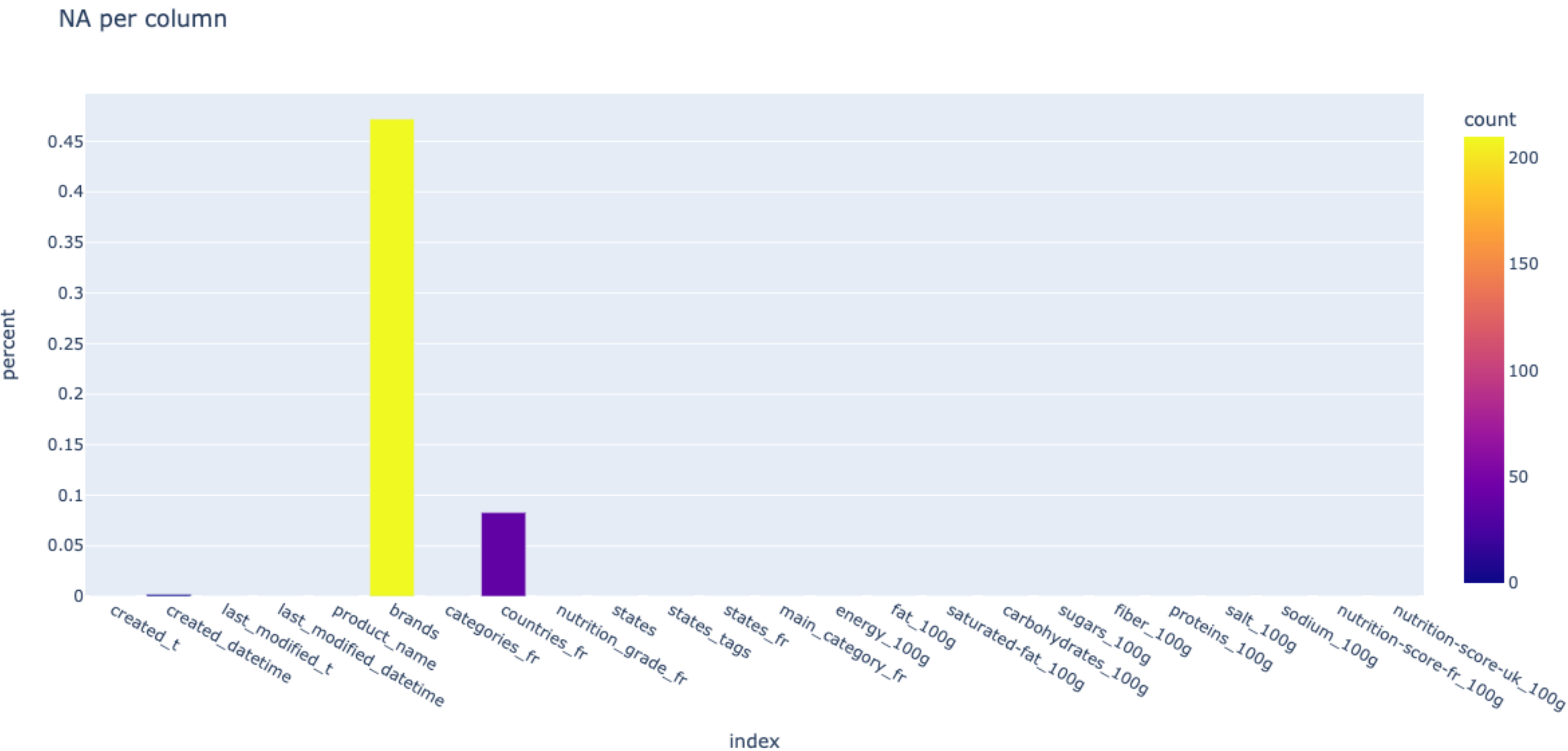
NA per column



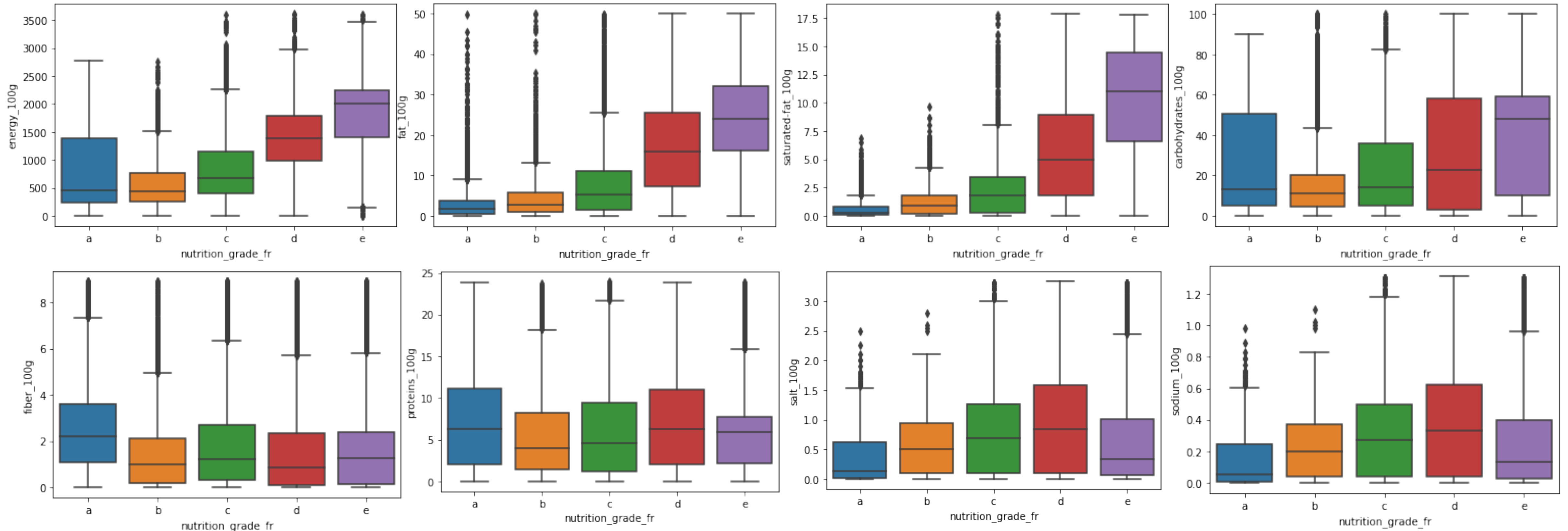
KNN()

```
def KNN(data):  
    cat = list(['created_t', 'created_datetime',  
               'last_modified_t', 'last_modified_datetime', 'product_name', 'brands',  
               'categories_fr',  
               'countries_fr', 'nutrition_grade_fr',  
               'states', 'states_tags', 'states_fr',  
               'main_category_fr'])  
    dataKNN = pd.DataFrame(KNNImputer(n_neighbors=5).fit_transform(data.filter(regex='_100g')),  
                           columns = data.filter(regex='_100g').columns)  
    data = pd.concat([data[cat].reset_index(drop=True),  
                     dataKNN.reset_index(drop = True)], axis=1)  
    return(data)
```

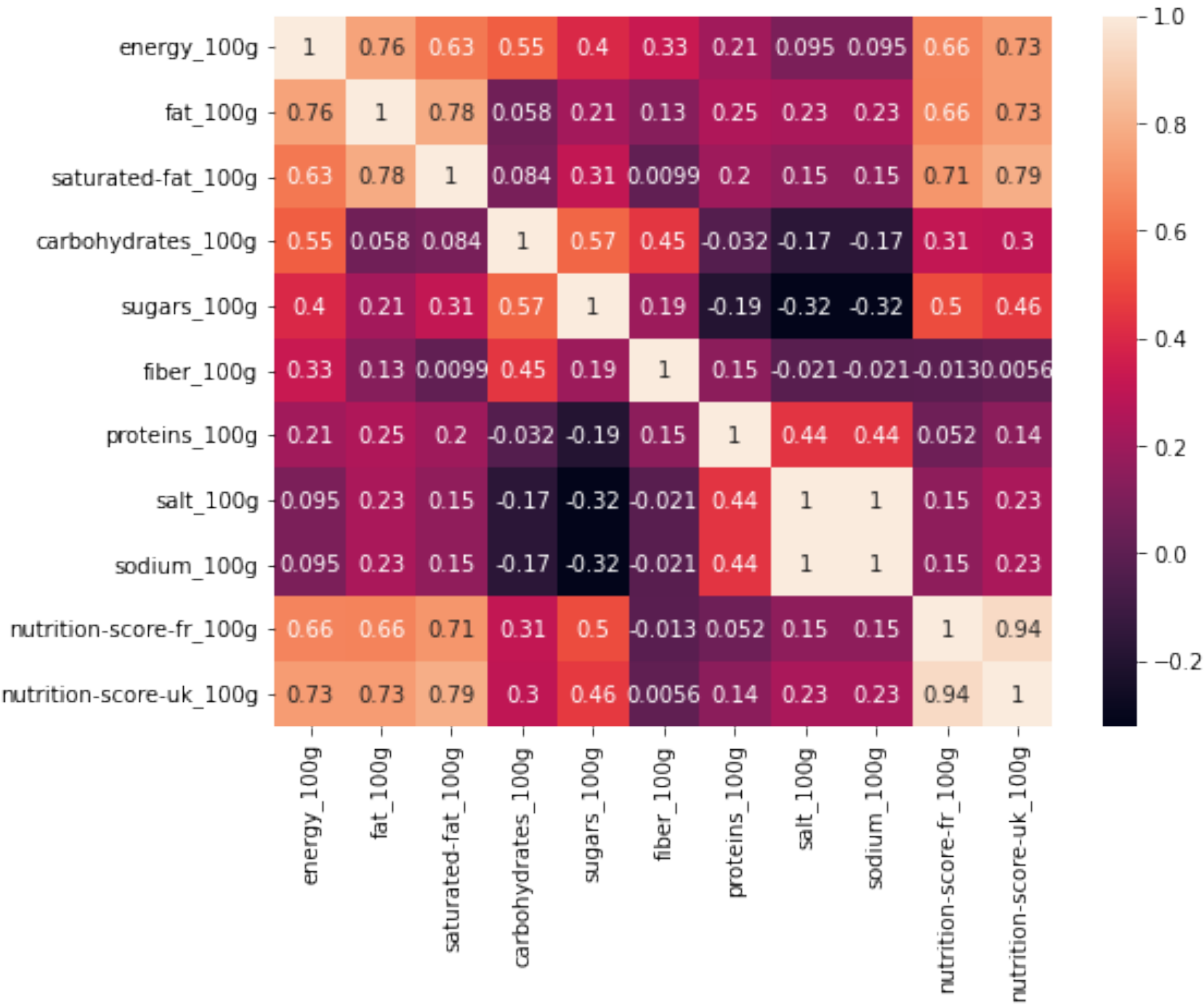

NA après KNN



Analyse univariée



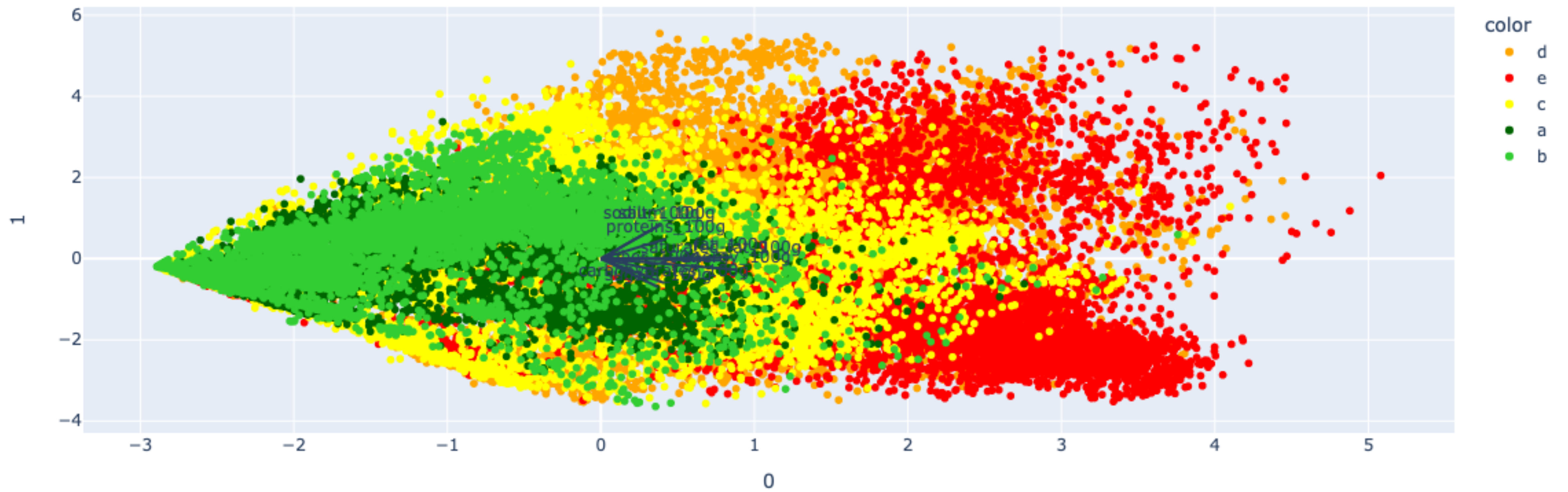
Corrélation



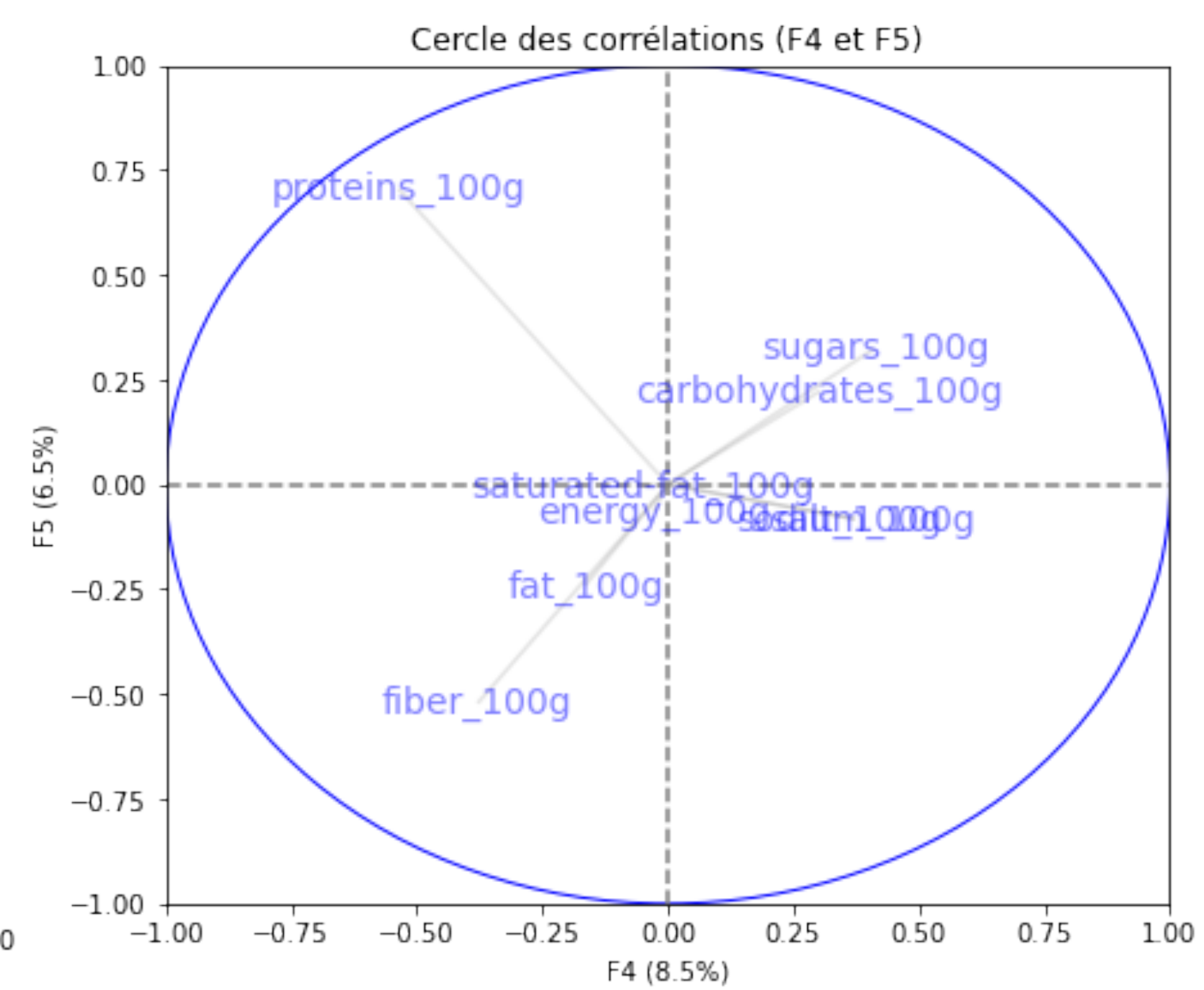
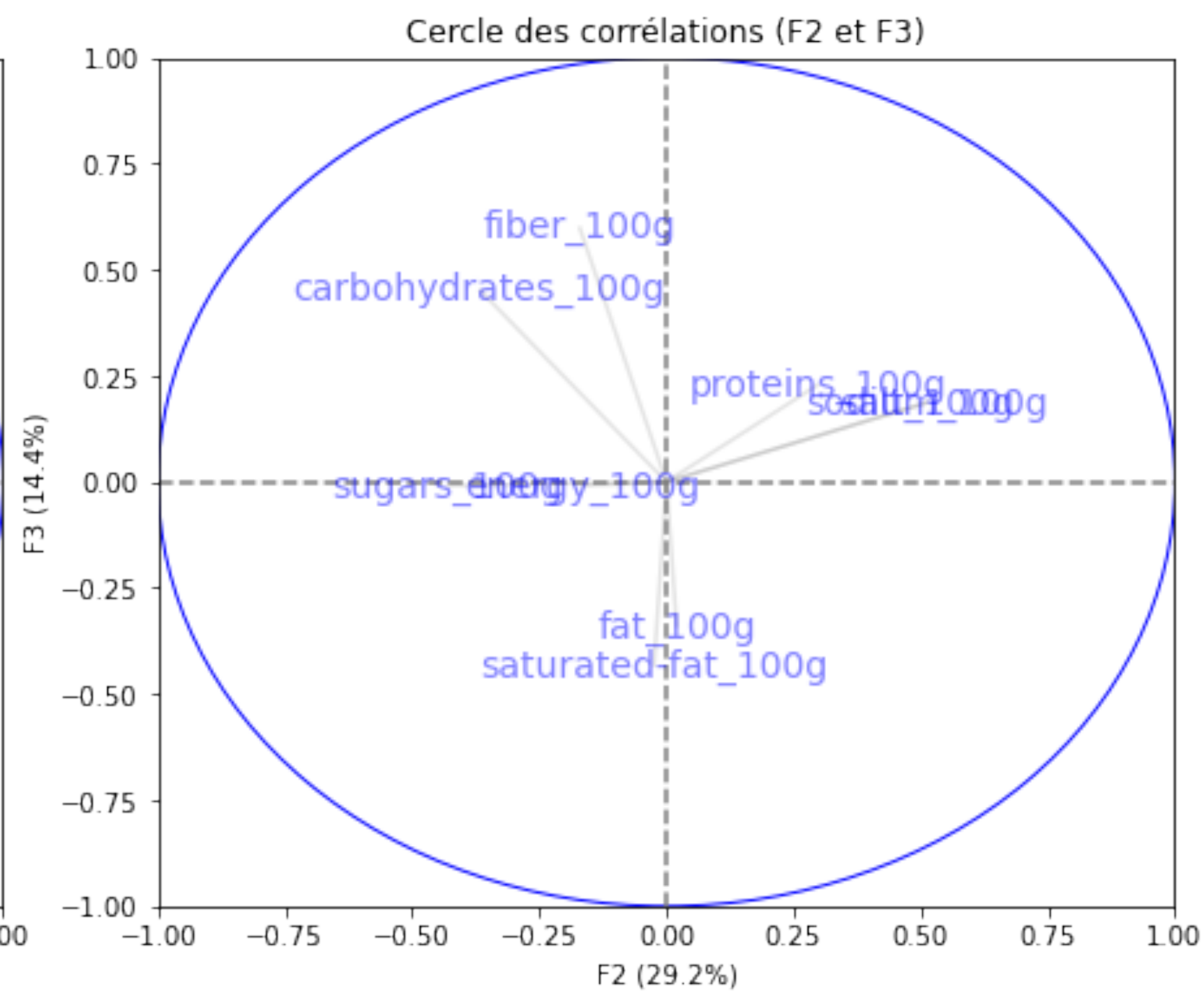
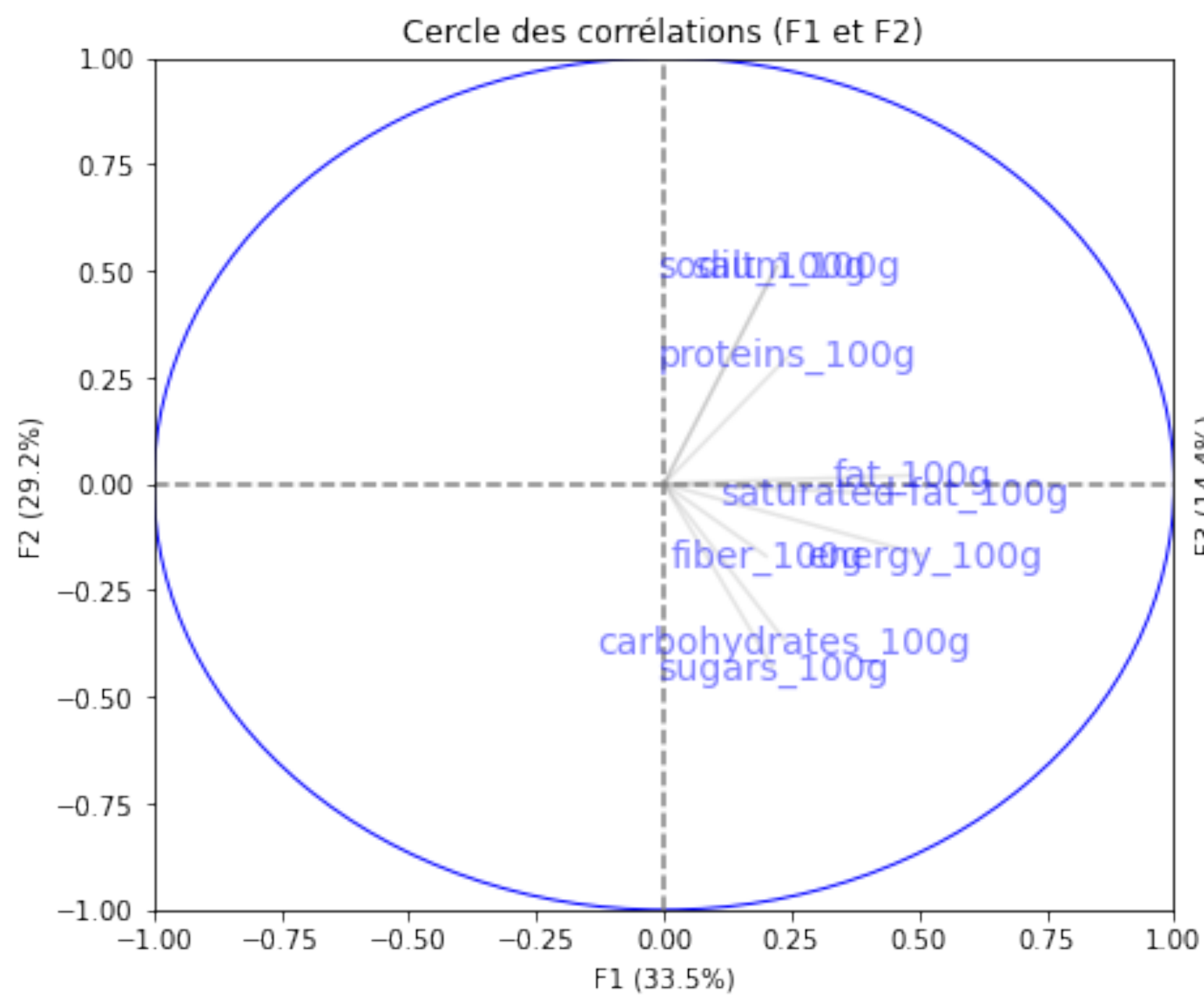
do_acp()

```
def do_pca(data):
    data2 = data.filter(regex='_100g').drop(['nutrition-score-fr_100g', 'nutrition-score-uk_100g'], axis=1)
    data2 = StandardScaler().fit_transform(data2)
    acp = PCA(n_components=data2.shape[1]).fit_transform(data2)
    #3D PCA
    n_comp = data2.shape[1]
    fig = px.scatter_3d(
        acp, x=0, y=1, z=2, color=data.nutrition_grade_fr, color_discrete_map= col,
        title= "3D PCA",
        labels={'0': 'PC 1', '1': 'PC 2', '2': 'PC 3'})
    fig.show()
    # Explained variance
    pca = PCA()
    pca.fit(data2)
    pcs = pca.components_
    exp_var_cumul = np.cumsum(pca.explained_variance_ratio_)
    px.area(
        x=range(1, exp_var_cumul.shape[0] + 1),
        y=exp_var_cumul,
        labels={"x": "# Components", "y": "Explained Variance"})
    loadings = pca.components_.T * np.sqrt(pca.explained_variance_)
    features = data.filter(regex='_100g').drop(['nutrition-score-fr_100g', 'nutrition-score-uk_100g'], axis=1).columns
    fig = px.scatter(acp, x=0, y=1, color=data['nutrition_grade_fr'], color_discrete_map= col)
    for i, feature in enumerate(features):
        fig.add_shape(
            type='line',
            x0=0, y0=0,
            x1=loadings[i, 0],
            y1=loadings[i, 1]
        )
        fig.add_annotation(
            x=loadings[i, 0],
            y=loadings[i, 1],
            ax=0, ay=0,
            xanchor="center",
            yanchor="bottom",
            text=feature,
        )
    fig.show()
    #Cercle des correlations
    display_circles(pcs, n_comp, pca, [(0,1), (1,2), (3,4)], labels=features)
```


ACP



ACP



do_anova()

```
def do_anova(data):  
    X = "nutrition_grade_fr" #qualitative  
    Y = data.filter(regex='_100g').columns #quantitative  
    anova = []  
    nutriment = []  
    for y in Y:  
        anova.append(eta_squared(data[X],data[y]))  
        nutriment.append(y)  
    return(pd.DataFrame(np.column_stack([nutriment, anova]), columns=['nutriment', 'anova']))
```

ANOVA

nutriment	anova
energy_100g	0.3163253651456699
fat_100g	0.36039822889719414
saturated-fat_100g	0.4405205482601115
carbohydrates_100g	0.05626378061387536
sugars_100g	0.22567170218930177
fiber_100g	0.04266394894237031
proteins_100g	0.011313202357831198
salt_100g	0.08407891418818202
sodium_100g	0.08407815628274008

Conclusion

- Erreur de séparateur
- Énormément de données manquante calculable si on a la formule
- Analyse par catégorie ou marque ou pays d'origine pourrai nous en apprendre plus