

# Analyse des données de la Boutique Bottleneck

David GAMARD  
Formation OpenClassrooms Data  
Analyst

# Missions

---

- ► Nettoyage des données
- ► Rapprochement des données de l'ERP et du CMS
- ► Calcul du chiffre d'affaire par produits
- ► Calcul du chiffre d'affaire total
- ► Analyse des prix des produits

# Ressources

---

- ► Données venant de l'ERP
- ► Données venant du CMS
- ► Table de liaison

# Outils

## Outils

- ▶ Langage de programmation : Python



- ▶ Logiciel : Jupyter



# Nettoyage des données

## Première vu des données WEB

```
RangeIndex: 0 to 1512
Data columns (total 28 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   sku              1428 non-null    object 
 1   virtual          1513 non-null    int64  
 2   downloadable     1513 non-null    int64  
 3   rating_count     1513 non-null    int64  
 4   average_rating   1430 non-null    float64
 5   total_sales      1430 non-null    float64
 6   tax_status        716 non-null    object  
 7   tax_class         0 non-null     float64
 8   post_author       1430 non-null    float64
 9   post_date         1430 non-null    object  
 10  post_date_gmt    1430 non-null    object  
 11  post_content      0 non-null     float64
 12  post_title        1430 non-null    object  
 13  post_excerpt      716 non-null    object  
 14  post_status       1430 non-null    object  
 15  comment_status    1430 non-null    object  
 16  ping_status       1430 non-null    object  
 17  post_password     0 non-null     float64
 18  post_name         1430 non-null    object  
 19  post_modified     1430 non-null    object  
 20  post_modified_gmt 1430 non-null    object  
 21  post_content_filtered 0 non-null    float64
 22  post_parent       1430 non-null    float64
 23  guid              1430 non-null    object  
 24  menu_order        1430 non-null    float64
 25  post_type          1430 non-null    object  
 26  post_mime_type    714 non-null    object  
 27  comment_count     1430 non-null    float64
dtypes: float64(19), int64(2), object(7)
memory usage: 40.1 MB
```

| sku                 | virtual | downloadable | rating_count | average_rating | total_sales | tax_status | tax_class | post_author | post_date           | post_name                                      | post_modified       |
|---------------------|---------|--------------|--------------|----------------|-------------|------------|-----------|-------------|---------------------|--|---------------------|
| bon-cadeau-25-euros | 0       | 0            | 0            | 0.0            | 10.0        | taxable    | Nan       | 1.0         | 2018-01-13 23:00    | bon-cadeau-25-euros                            | 2019-05-01 14:13:57 |
| 1                   | 15298   | 0            | 0            | 0.0            | 6.0         | taxable    | Nan       | 2.0         | 2018-02-08 12:56:32 | plene-jean-villa-saint-joseph-preco-2019       | 2019-12-30 09:30:29 |
| 2                   | 15299   | 0            | 0            | 0.0            | 0.0         | taxable    | Nan       | 2.0         | 2018-02-08 13:49:41 | plene-jean-villa-saint-joseph-2017             | 2019-12-21 09:03:17 |
| 3                   | 15300   | 0            | 0            | 0.0            | 0.0         | taxable    | Nan       | 2.0         | 2018-02-09 14:08:39 | plene-jean-villa-croze-hauterives-acconciacce- | 2020-08-26 18:15:03 |
| 4                   | 19914   | 0            | 0            | 0.0            | 3.0         | taxable    | Nan       | 2.0         | 2018-02-09 14:01:05 | plene-jean-villa-saint-joseph-garmino-2018     | 2020-01-04 10:35:01 |
| 5                   | 1508    | 16135        | 0            | 0              | 0.0         | 5.0        | Nan       | Nan         | 2.0                 | moufles-le-bihan-ame-2019                      | 2020-09-26 17:35:03 |
| 6                   | 1509    | 15091        | 0            | 0              | 0.0         | 0.0        | Nan       | Nan         | 2.0                 | cannier-lamédy-jurancio-2018                   | 2020-09-26 17:35:02 |
| 7                   | 1510    | 15087        | 0            | 0              | 0.0         | 0.0        | Nan       | Nan         | 2.0                 | jamet-cote-roche-frérottes-vézoulées-2019      | 2020-08-14 18:15:03 |
| 8                   | 1511    | 13127-1      | 0            | 0              | 0.0         | 0.0        | Nan       | Nan         | 2.0                 | clou-du-mont-saint-châtelain-du-pape-2019      | 2020-07-20 17:09:06 |
| 9                   | 1512    | 16230        | 0            | 0              | 0.0         | 0.0        | Nan       | Nan         | 2.0                 | domane-nicolas-saint-vérand-blancières         | 2020-08-13 10:45:03 |

## Constats :

- ▶ Erreurs lexical ( bon-cadeau-25-euros et 13127-1)
- ▶ Beaucoup de valeurs nulles (NaN) et de zéro
- ▶ Des dates au format datetime

# Nettoyage des données

## Suppressions des colonnes inutiles

```
le data set de web contient 1513 lignes
le data set de web contient 29 colonnes
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1513 entries, 0 to 1512
Data columns (total 29 columns):
 #   Column            Non-Null Count  Dtype  
0   sku               1428 non-null    object 
1   virtual           1513 non-null    int64  
2   downloadable      1513 non-null    int64  
3   rating_count      1513 non-null    int64  
4   average_rating    1430 non-null    float64
5   total_sales       1430 non-null    float64
6   tax_class          710 non-null    object 
7   post_author        0 non-null     float64
8   post_content       1430 non-null    float64
9   post_date          1430 non-null    datetime64[ns]
10  post_content_filtered  0 non-null     float64
11  post_content_type 0 non-null     float64
12  product_type      1429 non-null    object 
13  post_title         1430 non-null    object 
14  post_excerpt       710 non-null    float64
15  post_status         1430 non-null    object 
16  comment_status     1430 non-null    object 
17  ping_status        1430 non-null    object 
18  post_password      0 non-null     float64
19  post_name          1430 non-null    object 
20  post_modified      1430 non-null    datetime64[ns]
21  post_modified_gmt  1430 non-null    datetime64[ns]
22  post_content_filtered  0 non-null     float64
23  post_parent        1430 non-null    float64
24  guid               1430 non-null    object 
25  menu_order         1430 non-null    float64
26  post_type          1430 non-null    object 
27  post_mime_type     714 non-null    object 
28  comment_count      1430 non-null    float64
dtypes: datetime64[ns](4), float64(10), int64(3), object(12)
memory usage: 346.9+ KB
```

Nous décidons alors de garder les uniquement les valeurs qui nous serviront, à savoir les colonnes liées à la vente, et aux quantités

#Selon vous, quelles sont les colonnes à conserver ?

```
df_web = df_web.drop([
    'average_rating', 'post_parent', 'menu_order', 'comment_count',
    'tax_status', 'post_author', 'post_date', 'post_date_gmt',
    'post_status', 'comment_status', 'ping_status', 'post_name',
    'post_modified', 'post_modified_gmt', 'guid',
    'post_mime_type', 'tax_class', 'rating_count', 'downloadable', 'post_content', 'post_password', 'post_content_filtered', 'post_content_filtered'
], axis='columns', errors='ignore')
```

```
df_web.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1513 entries, 0 to 1512
Data columns (total 6 columns):
 #   Column            Non-Null Count  Dtype  
0   sku               1428 non-null    object 
1   total_sales       1430 non-null    float64
2   product_type      1429 non-null    object 
3   post_title         1430 non-null    object 
4   post_excerpt       716 non-null    object 
5   post_type          1430 non-null    object 
dtypes: float64(1), object(5)
memory usage: 71.1+ KB
```

- ▶ Beaucoup de colonnes comporte plusieurs lignes vide (83 lignes)
- ▶ 4 Colonnes comportant que des valeurs nulles (NaN)
- ▶ 3 Colonnes comportant beaucoup de valeurs nulle (NaN)

# Nettoyage des données

## Suppression des valeurs anomalies sur la clé

```
# affichage Anomalie 1 :  
df_web[df_web["sku"] == 'bon-cadeau-25-euros']  
  
# choix apporté je conserve cette appellation bon-cadeau-25-euros et mais je supprime l'un des deux doublons pour garder l'authentique "prod"  
df_web = df_web.drop(index=842)  
df_web[df_web["sku"] == 'bon-cadeau-25-euros']  
  
:   sku  total_sales  product_type      post_title      post_excerpt  post_type  
1387  bon-cadeau-25-euros       7.0        NaN  Bon cadeau de 25€  <span style="color: #a85253;"><strong>Parlons ...</strong></span>  product  
  
# affichage Anomalie 2 :  
df_web[df_web["sku"] == "13127-1"]  
  
# choix apporté : je supprime l'un des deux doublons pour garder l'authentique "product"  
df_web = df_web.drop(index=1117)  
  
df_web[df_web["sku"] == "13127"]  
  
:   sku  total_sales  product_type  post_title  post_excerpt  post_type
```

- ▶ Recherche des 'sku' ne respectant pas la règle de codification

# Nettoyage des données

## Sélection des données « product » et non « attachment »

Quelques variables, dont "post\_type", diffèrent. Cette dernière prend soit la valeur "product", soit la valeur "attachment". Nous n'allons conserver que les lignes "product". Malgré tout, par mesure de sécurité, vérifions si "post\_type" ne prend pas d'autre(s) valeur(s).

```
df_web['post_type'].unique()  
array(['attachment', 'product', nan], dtype=object)  
  
La vérification est concluante, on peut supprimer les lignes.  
  
df_web = df_web.loc[df_web['post_type'] != 'attachment']  
  
Voilà maintenant à quoi ressemble le dataframe.
```

|      | sku   | total_sales | product_type | post_title  | post_excerpt                                      | post_type |
|------|-------|-------------|--------------|---|---|-----------|
| 2    | 14692 | 5.0         | Vin          | Château Fonréaud Bordeaux Blanc Le Cygne 2016     | <div>Grâce à la complémentarité des 3 cépages ... | product   |
| 4    | 15328 | 2.0         | Vin          | Agnès Levet Côte Rôtie Maestria 2017              | <span style="float: none; background-color: tr... | product   |
| 6    | 16515 | 10.0        | Vin          | Château Turcaud Bordeaux Rouge Cuvée Majeure 2018 | <div id="wrapper"><n<div id="container-wraper..." | product   |
| 8    | NaN   | NaN         | NaN          | NaN   | NaN   | NaN       |
| 11   | 16585 | 15.0        | Vin          | Xavier Frissant Touraine Sauvignon 2019           | Un joli sauvignon frais et minéral, avec d'int... | product   |
| ...  | ...   | ...         | ...          | ...   | ...   | ...       |
| 1503 | 13074 | 4.0         | Vin          | Château de Vaudieu Châteauneuf-du-Pape L'Avenu... | "L'Avenue" est issue d'une parcelle de vieux g... | product   |
| 1505 | 16322 | 0.0         | Vin          | Moulin de Gassac IGP Pays d'Hérault Guilhem Ro... | Belle complexité aromatique alliant fruits rou... | product   |
| 1507 | 12385 | 10.0        | Vin          | Parés Baltà Penedès Electio 2013                  | Une cuvée produite avec une très vieille vigne... | product   |
| 1508 | 16326 | 5.0         | Vin          | Camin Larreda Jurançon Moelleux Au Capcéu 2018    | Sur le millésime 2017, Au Capceu du domaine Ca... | product   |
| 1509 | 15662 | 15.0        | Vin          | Chermette Domaine du Vissoux Beaujolais Griott... | C'est le Beaujolais typique : fruité, frais, g... | product   |

799 rows × 6 columns

- ▶ Tout les produits ont des doublons
- ▶ Chaque produit est attaché à une pièce joint (image,etc...)

- ▶ Reconstruction du DataFrame avec les données intéressante pour notre étude

# Nettoyage des données

## Traitement des valeurs nulles (NaN)

```
On constate encore qu'il y a des doublons pour cette colonne, voyons voir la nature de ces derniers
df_web[df_web["sku"].duplicated(keep=False)]
```

| sku  | total_sales | product_type | post_title | post_excerpt | post_type |
|------|-------------|--------------|------------|--------------|-----------|
| 8    | Nan         | Nan          | Nan        | Nan          | Nan       |
| 20   | Nan         | Nan          | Nan        | Nan          | Nan       |
| 30   | Nan         | Nan          | Nan        | Nan          | Nan       |
| 37   | Nan         | Nan          | Nan        | Nan          | Nan       |
| 41   | Nan         | Nan          | Nan        | Nan          | Nan       |
| ...  | ...         | ...          | ...        | ...          | ...       |
| 1384 | Nan         | Nan          | Nan        | Nan          | Nan       |
| 1429 | Nan         | Nan          | Nan        | Nan          | Nan       |
| 1432 | Nan         | Nan          | Nan        | Nan          | Nan       |
| 1445 | Nan         | Nan          | Nan        | Nan          | Nan       |
| 1457 | Nan         | Nan          | Nan        | Nan          | Nan       |

83 rows x 6 columns

```
Il s'agit de valeurs nulles , les lignes completes affichent des NaN nous pouvons donc les supprimer
df_web = df_web.dropna(subset=["sku"])

Vérifions maintenant une nouvelle fois.
```

```
df_web["sku"].isna().sum()
0
```

Notre Clé ne possède maintenant plus de doublons nous pouvons continuer vers l'Analyse Exploratoire du fichier liaison

- ▶ Remplacement de zéro par des valeur (NaN)
- ▶ Suppression des lignes contenant que des valeurs nulles (NaN)
- ▶ Passage de 1513 lignes à 799 lignes dans le DataFrame

# Nettoyage des données

## Première vu des données ERP

```
RangeIndex: 825 entries, 0 to 824
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   product_id  825 non-null    int64  
 1   onsale_web   825 non-null    int64  
 2   price        825 non-null    float64 
 3   stock_quantity  825 non-null  int64  
 4   stock_status  825 non-null    object  
dtypes: float64(1), int64(3), object(1)
```

|     | product_id | onsale_web | price | stock_quantity | stock_status |
|-----|------------|------------|-------|----------------|--------------|
| 0   | 3847       | 1          | 24.2  | 0              | outofstock   |
| 1   | 3849       | 1          | 34.3  | 0              | outofstock   |
| 2   | 3850       | 1          | 20.8  | 0              | outofstock   |
| 3   | 4032       | 1          | 14.1  | 0              | outofstock   |
| 4   | 4039       | 1          | 46.0  | 0              | outofstock   |
| ... | ...        | ...        | ...   | ...            | ...          |
| 820 | 7203       | 0          | 45.0  | 30             | instock      |
| 821 | 7204       | 0          | 45.0  | 9              | instock      |
| 822 | 7247       | 1          | 54.8  | 23             | instock      |
| 823 | 7329       | 0          | 26.5  | 14             | instock      |
| 824 | 7338       | 1          | 16.3  | 45             | instock      |

825 rows × 5 columns

Constats :

- ▶ Aucune valeur nulles (NaN)
- ▶ Pas d'erreur lexical

# Nettoyage des données

## Suppression des valeurs nulles sur la clé

```
print(dataerp.isnull().sum())
```

```
product_id      0  
onsale_web     0  
price          0  
stock_quantity  0  
stock_status    0  
dtype: int64
```

```
dataerp.loc[dataerp['product_id'].duplicated(keep=False),:]
```

| product_id | onsale_web | price | stock_quantity | stock_status |
|------------|------------|-------|----------------|--------------|
|------------|------------|-------|----------------|--------------|

- ▶ Aucune valeur nulles (NaN)

- ▶ Aucun doublon

# Nettoyage des données

## Première vu des données de liaison

```
#Consulter le nombre de colonnes
print(f"le data set de liaison contient {df_liaison.shape[1]} colonnes")
#La nature des données dans chacune des colonnes
print(f" la nature des données dans chacune des colonnes est : ")
print (df_liaison.dtypes,"\n")
#Le nombre de valeurs présentes dans chacune des colonnes
print("Nombre de valeurs non nulles dans chaque colonne :")
print(df_liaison.count(), "\n")

le data set de liaison contient 2 colonnes
 la nature des données dans chacune des colonnes est :
id_web      object
product_id   int64
dtype: object

Nombre de valeurs non nulles dans chaque colonne :
id_web      734
product_id   825
dtype: int64
```



- ▶ Aucune valeur nulles (NaN) dans la colonne ‘product\_id’
- ▶ Des valeurs nulles (NaN) dans la colonne ‘id\_web’

# Nettoyage des données

## Correction du nom de la colonne érroné

---

| # | Column | Non-Null Count | Dtype  |
|---|--------|----------------|--------|
| 0 | id_web | 1428 non-null  | object |



```
df_liaison = df_liaison.rename(columns = {'id_web':'sku'})  
df_liaison
```



| # | Column | Non-Null Count | Dtype  |
|---|--------|----------------|--------|
| 0 | sku    | 1428 non-null  | object |

# Nettoyage des données

## Vérification et traitement des valeurs nulles (NaN) et des doublons

```
#Les valeurs de la colonne "id_web" OU "product_id" sont-elles toutes uniques?  
df_liaison.isna().sum()
```

```
sku      91  
product_id      0  
dtype: int64
```

Il manque 91 valeurs pour "sku". Les produits concernés ne pourront pas être rapprochés de leur(s) vente(s), il est donc inutile de les conserver.

```
df_liaison = df_liaison.dropna().reset_index(drop=True)
```

Pour finir, nous nous assurons que les codes "product\_id" et "sku" ne contiennent pas de doublons.

```
# Test des doublons sur le code produit et le SKU  
test_doublons = df_liaison.duplicated(subset=['product_id', 'sku']).sum()  
  
print(f'Il y a {test_doublons} doublons dans les codes produits et les SKU.')  
  
del test_doublons
```

Il y a 0 doublons dans les codes produits et les SKU.

- ▶ Suppression des 91 valeurs nulles (NaN)

- ▶ Aucun doublon

# Fusion de données

## Premier merge

- ▶ Premier merge entre les données de liaison et les données ERP

```
#Fusion des fichiers df_erp et df_liaison
df_merge = pd.merge((df_erp), (df_liaison), on = ['product_id'], how = 'outer' ,indicator=True)
```

|     | product_id | onsale_web | price | stock_quantity | stock_status | purchase_price | stock_quantity_corrigé | sku | _merge        |
|-----|------------|------------|-------|----------------|--------------|----------------|------------------------|-----|---------------|
| 19  | 4055       | 0          | 86.1  | 0              | outofstock   | 37.88          |                        | 0   | NaN left_only |
| 49  | 4090       | 0          | 73.0  | 0              | outofstock   | 33.79          |                        | 0   | NaN left_only |
| 50  | 4092       | 0          | 47.0  | 0              | outofstock   | 25.25          |                        | 0   | NaN left_only |
| 119 | 4195       | 0          | 14.1  | 0              | outofstock   | 7.36           |                        | 0   | NaN left_only |
| 131 | 4209       | 0          | 73.5  | 0              | outofstock   | 33.01          |                        | 0   | NaN left_only |
| ... | ...        | ...        | ...   | ...            | ...          | ...            | ...                    | ... | ...           |
| 814 | 7196       | 0          | 31.0  | 55             | instock      | 31.20          |                        | 55  | NaN left_only |
| 815 | 7200       | 0          | 31.0  | 6              | instock      | 15.54          |                        | 6   | NaN left_only |
| 816 | 7201       | 0          | 31.0  | 18             | instock      | 16.02          |                        | 18  | NaN left_only |
| 817 | 7203       | 0          | 45.0  | 30             | instock      | 23.48          |                        | 30  | NaN left_only |
| 818 | 7204       | 0          | 45.0  | 9              | instock      | 24.18          |                        | 9   | NaN left_only |

# Fusion de données

## Deuxième merge

- Deuxième merge entre le premier merge et les données WEB

| #Fusionnez les datasets df_merge et df_web   |      |       |                |              |                |                        |     |             |              |  |   |           |           |
|--|------|-------|----------------|--------------|----------------|------------------------|-----|-------------|--------------|--|---|-----------|-----------|
| df_final = pd.merge((df_merge), (df_web), on = ['sku'], how = 'outer', indicator=True) |      |       |                |              |                |                        |     |             |              |  |   |           |           |
| df_final   | eb   | price | stock_quantity | stock_status | purchase_price | stock_quantity_corrige | sku | total_sales | product_type | post_title                                       | post_excerpt                                      | post_type | _merge    |
|  | 0.0  | 8.6   | 26.0           | instock      | 4.22           | 26.0                   | 38  | 10.0        | Vin          | Emile Boeckel Crémant Brut Blanc de Blancs       | Ce Crémant est vif et délicat, gourmand et cro... | product   | both      |
|  | 0.0  | 41.0  | 11.0           | instock      | 20.12          | 11.0                   | 41  | 6.0         | Vin          | Marcel Windholtz Eau de Vie de Marc de Gewurz... | Les eaux de vie naissent d'une subtile alchimi... | product   | both      |
|  | 0.0  | 39.0  | 123.0          | instock      | 24.86          | 123.0                  | 304 | 8.0         | Champagne    | Champagne Gosset Grande Réserve                  | Le nez, ouvert et expressif, évoque les fruits... | product   | both      |
|  | 0.0  | 59.9  | 13.0           | instock      | 27.18          | 13.0                   | 523 | 0.0         | Cognac       | Cognac Normandin Mercier VFC                     | Issus des meilleurs crus de Grande et de Petit... | product   | both      |
|  | 0.0  | 22.5  | 76.0           | instock      | 13.78          | 76.0                   | 531 | 8.0         | Champagne    | Champagne Petit Lebrun & Fils Blanc de Bla...    | Cuvée bien équilibrée à la fois vive et souple... | product   | both      |
| ...  | ...  | ...   | ...            | ...          | ...            | ...                    | ... | ...         | ...          | ...  | ...   | ...       | ...       |
| 0.0  | 31.0 | 55.0  | instock        | 31.20        | 55.0           | NaN                    | NaN | NaN         | NaN          | NaN  | NaN   | NaN       | left_only |
| 0.0  | 31.0 | 6.0   | instock        | 15.54        | 6.0            | NaN                    | NaN | NaN         | NaN          | NaN  | NaN   | NaN       | left_only |
| 0.0  | 31.0 | 18.0  | instock        | 16.02        | 18.0           | NaN                    | NaN | NaN         | NaN          | NaN  | NaN   | NaN       | left_only |
| 0.0  | 45.0 | 30.0  | instock        | 23.48        | 30.0           | NaN                    | NaN | NaN         | NaN          | NaN  | NaN   | NaN       | left_only |
| 0.0  | 45.0 | 9.0   | instock        | 24.18        | 9.0            | NaN                    | NaN | NaN         | NaN          | NaN  | NaN   | NaN       | left_only |

# Analyse des données

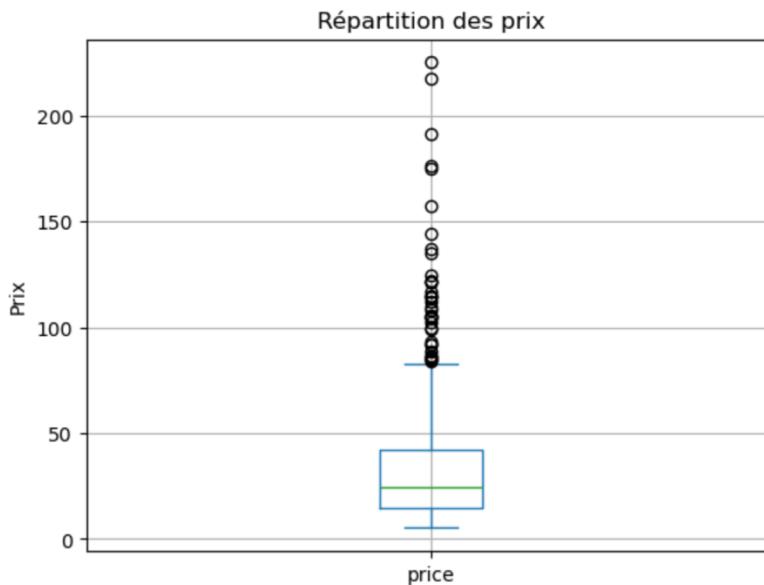
## Calcul du CA par produit et du CA total

| #Créez une colonne calculant le CA par article                                     |              |                |                        |       |             |              |                 |  |   |         |                |       |
|--|--------------|----------------|------------------------|-------|-------------|--------------|-----------------|--|---|---------|----------------|-------|
| #Calculez la somme de la colonne "ca_par_article"                                  |              |                |                        |       |             |              |                 |  |   |         |                |       |
| ca_total = df_final["ca_par_article"].sum()<br>ca_col_total = round (ca_total , 2) |              |                |                        |       |             |              |                 |  |   |         |                |       |
| print (f" le chiffre d'affaire total est de {ca_total} euros ")                    |              |                |                        |       |             |              |                 |  |   |         |                |       |
| #Ce résultat correspond au chiffre d'affaire du site web<br>df_final               |              |                |                        |       |             |              |                 |  |   |         |                |       |
| le chiffre d'affaire total est de 143680.1 euros                                   |              |                |                        |       |             |              |                 |  |   |         |                |       |
| quantity   | stock_status | purchase_price | stock_quantity_corrigé | sku   | total_sales | product_type | post_title      | post_excerpt                                     | post_type   | _merge  | ca_par_article |       |
| 26.0   | instock      | 4.22           |                        | 26.0  | 38          | 10.0         | Vin             | Emile Boeckel Crémant Brut Blanc de Blancs       | Ce Crémant est vif et délicat, gourmand et cro... | product | both           | 86.0  |
| 11.0   | instock      | 20.12          |                        | 11.0  | 41          | 6.0          | Vin             | Marcel Windholtz Eau de Vie de Marc de Gewurz... | Les eaux de vie naissent d'une subtile alchimi... | product | both           | 246.0 |
| 123.0  | instock      | 24.86          |                        | 123.0 | 304         | 8.0          | Champagne       | Champagne Gosset Grande Réserve                  | Le nez, ouvert et expressif, évoque les fruits... | product | both           | 312.0 |
| 13.0   | instock      | 27.18          |                        | 13.0  | 523         | 0.0          | Cognac          | Cognac Normandin Mercier VFC                     | Issus des meilleurs crus de Grande et de Petit... | product | both           | 0.0   |
|  |              |                |                        |       |             |              | Champagne Petit | Cuvée bien                                       |   |         |                |       |

Chiffre d'affaire total  
143680.1 euros €

# Analyse des données

## Analyse des prix des produits



- Moyenne : 32.35
- Mediane : 24.4

- Q1 : 14.6
- Q3 : 42.0
- IQR : 27.4
- Min = 5.2
- Max = 84.0

- ▶ 50% des prix sont compris entre 14€ et 42€
- ▶ 75% des prix sont compris entre 5,2€ et 42€
- ▶ Les Outliers sont supérieurs à 84€



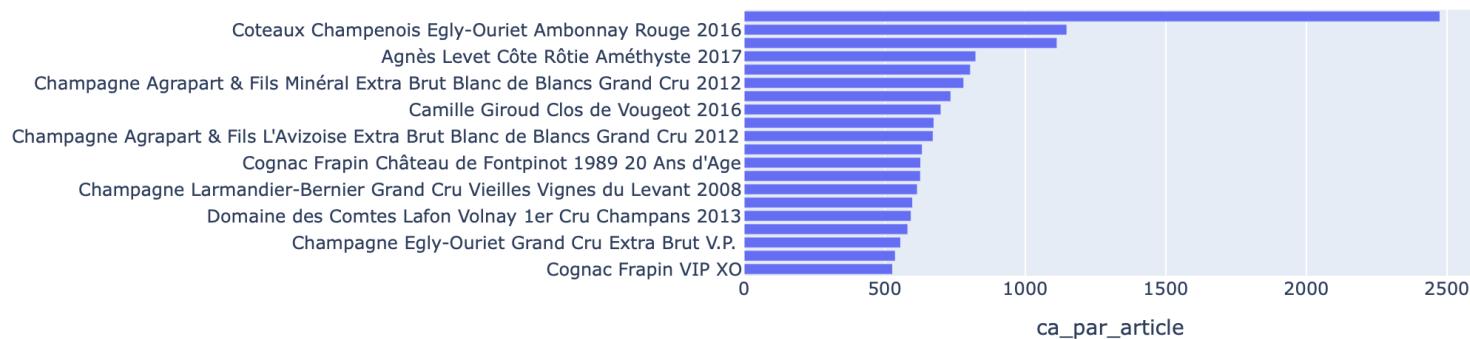
# Analyse des données

## Analyse des prix des produits



Palmarès des 20 premiers articles en CA

post\_title



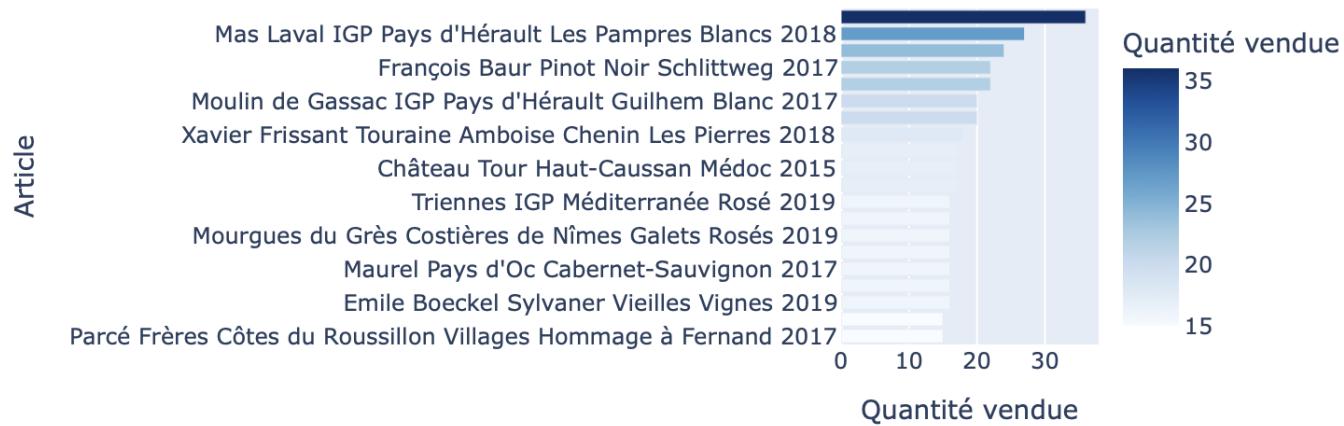
Nombre d'articles représentant 80% du CA : 434 Ces articles représentent 52.67% du catalogue total.

# Analyse des données

## Analyse des Quantités vendues des produits



Top 20 articles par quantité vendue



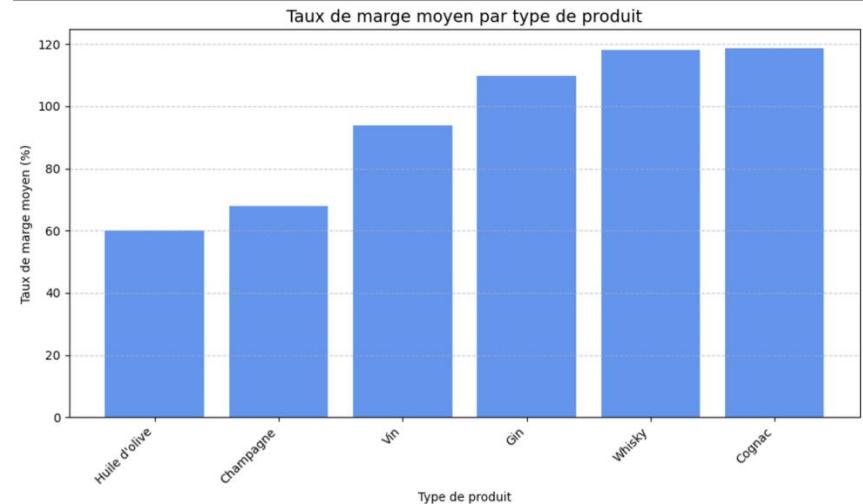
# Analyse des données

## Taux de marge moyen par type de produit

1. ✓ Les spiritueux (Cognac, Whisky, Gin) sont les produits les plus rentables. Il peut être judicieux d'en renforcer la promotion ou les stocks.

2. ⚠ Huile d'olive et Champagne ont une marge relativement faible. Cela peut indiquer :

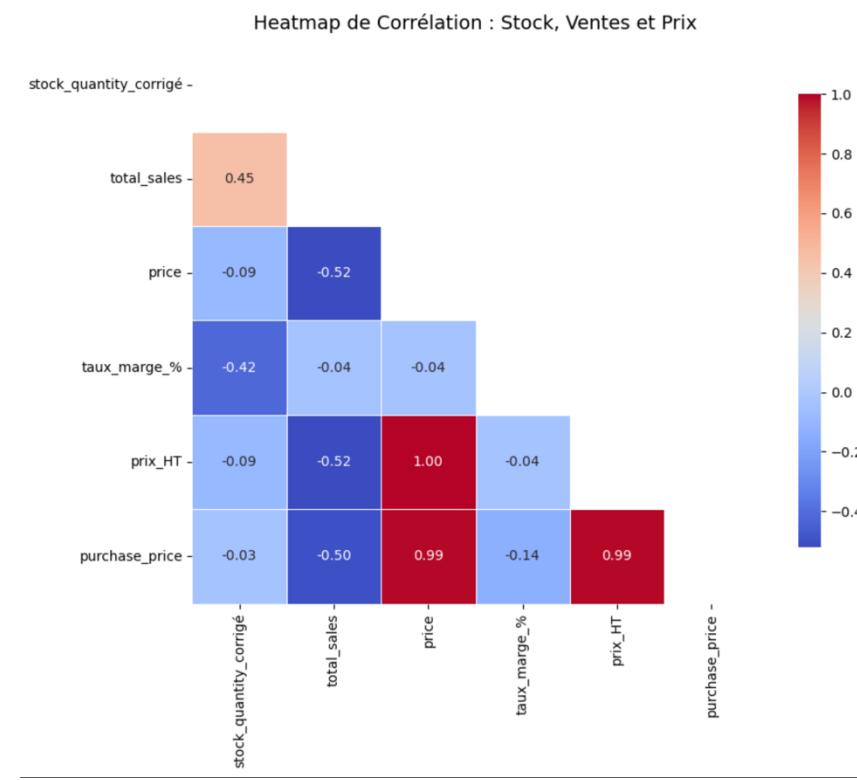
- Des coûts d'achat élevés.
- Une fixation des prix de vente trop basse.
- Une faible valeur ajoutée ou concurrence accrue.



3. ⟳ Le vin est intermédiaire : rentable mais peut être optimisé

# Analyse des données

## Heatmap de corrélation



# Synthèse sur le nettoyage

- Fichier WEB :
  - ► Beaucoup de données manquantes
  - ► Des colonnes non-utilisé
  - ► Des doublon pour chaque produits
  - ► Des erreurs lexical à revoir
- Fichier ERP :
  - ► Faire attention au cohésion entre la quantité et le statut des stocks
- Fichier Liaison :
  - ► Quelques valeurs nulles

# Synthèse sur l'analyse

- ► Chiffre d'affaire de 143680.1€
- ► 50% des prix sont compris entre 14€ et 42€
- ► Les Outlier sont des valeurs Extrême
- ► 51,46% des articles du catalogue produisent 80% du chiffre d'affaires de la boutique
- ► Les produits du type Gin , Whisky et Cognac sont les plus rentables