

# Projet Final Data

# Membres du groupe

- BRIMESSE Wail GRP-1
- BENAHMED Rayan GRP-2
- FEDDILA Yanis GRP-2
- HASSAOUI Zaid GRP-2
- SKALLI Yanis GRP-3

# Sommaire

- I- Introduction
- II- Travail réalisé
- III- Base de données
- IV- Code et résultat
- V- Conclusion

# Introduction

Dans le cadre de ce projet, nous avons réalisé une solution complète de traitement et d'analyse de données allant du scrapping d'un site web jusqu'à la visualisation des résultats après avoir nettoyé et stocké les informations dans une base de données.

Le site pour ce projet est **Books to Scrape**, une plateforme permettant de tester des techniques de scrapping sur des données de livres.

Nous avons listé plusieurs informations clés à récupérer pour chaque livre :

- Le titre
- Le prix (TTC et HT)
- La disponibilité
- Le stock
- La taxe
- La notation en étoiles
- La catégorie

Nous avons donc réalisé les étapes dans cet ordre :

1. Scrapping des données
2. Nettoyage et structuration des données
3. Stockage des données dans une base de données
4. Analyse et virtualisation des données

Adopter cette structure étape par étape nous a permis de gagner en efficacité et en qualité de travail, tout en assurant une organisation claire et une meilleure gestion des données.

# Travail réalisé

## 1 - Scraping des données

L'extraction des données a été réalisée à l'aide de **Requests** et **BeautifulSoup**, permettant d'envoyer des requêtes HTTP et de parser le HTML des pages. Un script a été mis en place pour parcourir jusqu'à 30 pages du site et récupérer les données de chaque livre.

## 2 - Nettoyage et Structuration des données

Une fois les données brutes extraites, elles ont été nettoyées et organisées à l'aide de **Pandas**. Plusieurs étapes ont été mises en place :

- Suppression des espaces inutiles et caractères spéciaux dans les prix
- Conversion des prix en valeurs numériques pour faciliter les analyses
- Gestion des valeurs manquantes

## 3 - Stockage des données dans une base de données

Pour garantir une meilleure gestion des données et leur réutilisation, elles ont été stockées dans une base de données **MySQL**. Une table dédiée a été créée avec les colonnes correspondantes, et un script Python utilisant MySQL Connector a permis d'insérer les données proprement dans la base.

## 4 - Analyse et Visualisation des données

Nous avons récupéré les données de la base **MySQL livre**, ajouté une colonne **vente**, et nettoyé les valeurs (prix, stock, étoiles).

**Visualisations principales :**

- **Répartition des livres par notation (étoiles), catégorie et intervalle de prix.**
- **Analyse des stocks et des ventes par notation, catégorie et prix.**

Ces graphiques permettent de **visualiser les tendances des ventes et du stock**, aidant à optimiser l'approvisionnement et la stratégie de vente.

# Base de données

La base de données **sprint\_data** contient une table principale livres, qui stocke les informations extraites et nettoyées des livres du site **Books to Scrape**.

## Structure de la Table Livres

```
CREATE TABLE livres (  
id int(11) NOT NULL,  
titre varchar(255) NOT NULL,  
prix decimal(10,2) NOT NULL,  
disponibilite varchar(50) NOT NULL,  
etoiles int(11) DEFAULT NULL CHECK (etoiles between 0 and 5),  
prix_ht decimal(10,2) NOT NULL,  
prix_ttc decimal(10,2) NOT NULL,  
taxe decimal(10,2) NOT NULL,  
categorie varchar(100) NOT NULL,  
stock int(11) NOT NULL,  
vente int(11) DEFAULT NULL );
```

## Remplissage de la Base

Nous avons intégré dans la base de données les données récupérées sur le site grâce aux commandes **SQL** suivantes :

### **LOAD DATA INFILE**

```
'C:/Users/Rayan/PycharmProjects/PythonProject/.venv/books_toscraper_cleaned.csv'
```

### **INTO TABLE livres**

```
FIELDS TERMINATED BY ';' ;
```

```
LINES TERMINATED BY '\n' IGNORE 1 ROWS
```

```
(titre, prix, disponibilite, etoiles, prix_ht, prix_ttc, taxe, categorie, stock);
```

## Objectif

- Stocker les données de manière organisée après extraction et nettoyage.
- Faciliter le traitement des données.
- Préparer les données pour des visualisations ou exportations.

# Code et Résultat

## 1 – Scraping des données

Code :

```
scrap.py x
1 import requests
2 from bs4 import BeautifulSoup
3 import pandas as pd
4
5 # URL de base pour accéder aux différentes pages
6 BASE_URL = "https://books.toscrape.com/catalogue/page-{}.html"
7
8 # Headers pour éviter d'être bloqué par le site
9 HEADERS = {
10     "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/119.0.0.0 Safari/537.36"
11 }
12
13 # Conversion des étoiles en notation numérique
14 STAR_RATINGS = {
15     "One": 1,
16     "Two": 2,
17     "Three": 3,
18     "Four": 4,
19     "Five": 5
20 }
21
22 # Dictionnaire pour stocker les données des livres
23 books = {
24     "Titre": [],
25     "Prix": [],
26     "Disponibilité": [],
27     "Étoiles": [],
28     "Price (excl. tax)": [],
29     "Price (incl. tax)": [],
30     "Tax": [],
31     "Catégorie": [],
32     "Stock": []
33 }
34
35 # Limite du nombre de pages à scraper
36 max_pages = 30
37
38 # Parcourir les pages une par une
39 for page in range(1, max_pages + 1):
40     url = BASE_URL.format(page)
41     response = requests.get(url, headers=HEADERS)
42
43     if response.status_code != 200:
44         print(f"Problème d'accès à la page {page}, arrêt du scraping.")
45         break # Si une page ne charge pas, on ne continue pas
46
47     soup = BeautifulSoup(response.text, 'html.parser')
48
49     # Récupérer tous les livres affichés sur la page
50     items = soup.find_all("article", class_="product_pod")
51
52     if not items:
53         print(f"Arrêt du scraping : seulement {page - 1} pages trouvées.")
54         break # Si aucune donnée trouvée, inutile de continuer
55
56     for item in items:
57         try:
58             title = item.find("h3").find("a")["title"]
59             price = item.find("p", class_="price_color").text.strip()
```

```

59     price = item.find("p", class_="price_color").text.strip()
60     availability = item.find("p", class_="instock availability").text.strip()
61
62     # Récupération des étoiles (notation de 1 à 5)
63     star_class = item.find("p", class_="star-rating")["class"]
64     star_rating = STAR_RATINGS.get(star_class[1], 0) # Si pas trouvé, on met 0
65
66     # Accéder à la page du produit pour récupérer les prix détaillés
67     product_link = "https://books.toscrape.com/catalogue/" + item.find("h3").find("a")["href"]
68     product_response = requests.get(product_link, headers=HEADERS)
69     product_soup = BeautifulSoup(product_response.text, "html.parser")
70
71     price_excl_tax = product_soup.find("th", string="Price (excl. tax)").find_next("td").text.strip()
72     price_incl_tax = product_soup.find("th", string="Price (incl. tax)").find_next("td").text.strip()
73     tax = product_soup.find("th", string="Tax").find_next("td").text.strip()
74
75     # Récupérer la catégorie
76     category = product_soup.find("ul", class_="breadcrumb").find_all("a")[2].text.strip()
77
78     # Récupérer le nombre disponible dans le stock
79     stock = product_soup.find("p", class_="instock availability").text.strip().replace("In stock (", "").replace(" available)", "")
80
81     # Ajouter les données extraites dans le dictionnaire
82     books["Titre"].append(title)
83     books["Prix"].append(price)
84     books["Disponibilité"].append(availability)
85     books["Étoiles"].append(star_rating)
86     books["Price (excl. tax)"].append(price_excl_tax)
87     books["Price (incl. tax)"].append(price_incl_tax)
88     books["Tax"].append(tax)
89     books["Catégorie"].append(category)
90     books["Stock"].append(stock)
91
92     except Exception as e:
93         print(f"Erreur lors de l'extraction d'un livre : {e}")
94         continue
95
96 # Enregistrement des données dans un fichier CSV
97 df = pd.DataFrame(books)
98 df.to_csv("books_toscrape.csv", index=False, encoding='utf-8', sep=";")
99
100 print("Scraping terminé ! Données enregistrées dans books_toscrape.csv")

```



## Résultat : Extrait de books\_toscrape.csv

	A	B	C	D	E	F	G	H	I
1	Titre	Prix	Disponibilit	%toiles	Price (excl. t	Price (incl. t	Tax	Catégorie	Stock
2	A Light in the	£51.77	In stock	3	£51.77	£51.77	£0.00	Poetry	22
3	Tipping the V	£53.74	In stock	1	£53.74	£53.74	£0.00	Historical Fic	20
4	Soumission	£50.10	In stock	1	£50.10	£50.10	£0.00	Fiction	20
5	Sharp Object	£47.82	In stock	4	£47.82	£47.82	£0.00	Mystery	20
6	Sapiens: A Bi	£54.23	In stock	5	£54.23	£54.23	£0.00	History	20
7	The Requiem	£22.65	In stock	1	£22.65	£22.65	£0.00	Young Adult	19
8	The Dirty Litt	£33.34	In stock	4	£33.34	£33.34	£0.00	Business	19
9	The Coming	£17.93	In stock	3	£17.93	£17.93	£0.00	Default	19
10	The Boys in t	£22.60	In stock	4	£22.60	£22.60	£0.00	Default	19
11	The Black Me	£52.15	In stock	1	£52.15	£52.15	£0.00	Poetry	19
12	Starving Hea	£13.99	In stock	2	£13.99	£13.99	£0.00	Default	19
13	Shakespeare	£20.66	In stock	4	£20.66	£20.66	£0.00	Poetry	19
14	Set Me Free	£17.46	In stock	5	£17.46	£17.46	£0.00	Young Adult	19
15	Scott Pilgrim	£52.29	In stock	5	£52.29	£52.29	£0.00	Sequential A	19
16	Rip it Up and	£35.02	In stock	5	£35.02	£35.02	£0.00	Music	19
17	Our Band Co	£57.25	In stock	3	£57.25	£57.25	£0.00	Music	19
18	Olio	£23.88	In stock	1	£23.88	£23.88	£0.00	Poetry	19
19	Mesaerion: T	£37.59	In stock	1	£37.59	£37.59	£0.00	Science Ficti	19
20	Libertarianis	£51.33	In stock	2	£51.33	£51.33	£0.00	Politics	19
21	It's Only the	£45.17	In stock	2	£45.17	£45.17	£0.00	Travel	19
22	In Her Wake	£12.84	In stock	1	£12.84	£12.84	£0.00	Thriller	19
23	How Music V	£37.32	In stock	2	£37.32	£37.32	£0.00	Music	19
24	Foolproof Pr	£30.52	In stock	3	£30.52	£30.52	£0.00	Food and Dri	19
25	Chase Me (P	£25.27	In stock	5	£25.27	£25.27	£0.00	Romance	19
26	Black Dust	£34.53	In stock	5	£34.53	£34.53	£0.00	Romance	19
27	Birdsong: A	£54.64	In stock	3	£54.64	£54.64	£0.00	Childrens	19
28	America's Cr	£22.50	In stock	3	£22.50	£22.50	£0.00	Default	19
29	Aladdin and	£53.13	In stock	3	£53.13	£53.13	£0.00	Default	19
30	Worlds Elsev	£40.30	In stock	5	£40.30	£40.30	£0.00	Nonfiction	18
31	Wall and Pie	£44.18	In stock	4	£44.18	£44.18	£0.00	Art	18
32	The Four Agr	£17.66	In stock	5	£17.66	£17.66	£0.00	Spirituality	18
33	The Five Love	£31.05	In stock	3	£31.05	£31.05	£0.00	Nonfiction	18
34	The Elephant	£23.82	In stock	5	£23.82	£23.82	£0.00	Thriller	18
35	The Bear and	£36.89	In stock	1	£36.89	£36.89	£0.00	Childrens	18
36	Sophie's Wor	£15.94	In stock	5	£15.94	£15.94	£0.00	Philosophy	18
37	Pennv Mavbr	£33.29	In stock	3	£33.29	£33.29	£0.00	Default	18
	< >	books_toscrape		+					

## 2 – Nettoyage des données

Code :

```
nettoyage.py ×
1 import pandas as pd
2
3 # Charger les données avec le bon séparateur et encodage
4 def charger_donnees(fichier):
5     try:
6         df = pd.read_csv(fichier, sep=";", encoding="utf-8")
7         # Vérifier les valeurs manquantes
8         print("Valeurs manquantes par colonne :")
9         print(df.isnull().sum())
10        return df
11    except Exception as e:
12        print(f"Erreur lors du chargement des données : {e}")
13        return None
14
15
16 # Nettoyer les données
17 def nettoyer_donnees(df):
18     if df is None or df.empty:
19         return None
20
21     try:
22         # Colonnes numériques à nettoyer
23         colonnes_numeriques = ["Prix", "Price (excl. tax)", "Price (incl. tax)", "Tax", "Stock"]
24
25         for col in colonnes_numeriques:
26             if col in df.columns:
27                 # Supprimer les symboles non numérique et convertir en nombres
28                 df[col] = df[col].astype(str).str.replace("Â£", "", regex=False)
29                 df[col] = df[col].str.replace("[^0-9.]", "", regex=True) # Garder uniquement les chiffres et les points
30                 df[col] = pd.to_numeric(df[col], errors="coerce") # Convertir en nombres, avec gestion des erreurs
31
32                 # Remplacer les valeurs manquantes par la moyenne de la colonne
33                 moyenne = df[col].mean()
34                 df[col] = df[col].fillna(moyenne)
35
36         # Nettoyage de la colonne "Stock" en extrayant uniquement le nombre disponible
37         if "Stock" in df.columns:
38             df["Stock"] = df["Stock"].astype(str).str.extract("(\\d+)").astype(float)
39
40         # Convertir la colonne "Stock" en entiers pour supprimer le ".0"
41         if "Stock" in df.columns:
42             df["Stock"] = df["Stock"].astype(int)
43
44         return df
45     except Exception as e:
46         print(f"Erreur lors du nettoyage des données : {e}")
47         return None
48
49
50 # Sauvegarder les données nettoyées
51 def sauvegarder_donnees(df, fichier_sortie):
52     if df is None or df.empty:
53         print("Erreur lors de la sauvegarde : aucune donnée à enregistrer.")
54         return
55
56     try:
57         # Vérifier les valeurs manquantes avant la sauvegarde
58         print("Valeurs manquantes après nettoyage :")
59         print(df.isnull().sum())
```

```

59     print(df.isnull().sum())
60
61     df.to_csv(fichier_sortie, sep=";", encoding="utf-8", index=False)
62     print("Donnée sauvegardée avec succès !")
63 except Exception as e:
64     print(f"Erreur lors de la sauvegarde : {e}")
65
66
67 #Exécution du programme
68 if __name__ == "__main__":
69     fichier_entree = "books_toscrape.csv"
70     fichier_sortie = "books_toscrape_cleaned.csv"
71
72     df = charger_donnees(fichier_entree)
73     df = nettoyer_donnees(df)
74     sauvegarder_donnees(df, fichier_sortie)

```

## Résultat : le csv nettoyé :

	A	B	C	D	E	F	G	H	I
1	Titre	Prix	Disponibilité	Annotations	Price (excl. t	Price (incl. t	Tax	Catégorie	Stock
2	A Light in the	51.77	In stock	3	51.77	51.77	0.0	Poetry	22
3	Tipping the V	53.74	In stock	1	53.74	53.74	0.0	Historical Fic	20
4	Soumission 50.1		In stock	1	50.1	50.1	0.0	Fiction	20
5	Sharp Object	47.82	In stock	4	47.82	47.82	0.0	Mystery	20
6	Sapiens: A Bi	54.23	In stock	5	54.23	54.23	0.0	History	20
7	The Requiem	22.65	In stock	1	22.65	22.65	0.0	Young Adult	19
8	The Dirty Littl	33.34	In stock	4	33.34	33.34	0.0	Business	19
9	The Coming	17.93	In stock	3	17.93	17.93	0.0	Default	19
10	The Boys in t	22.6	In stock	4	22.6	22.6	0.0	Default	19
11	The Black M	52.15	In stock	1	52.15	52.15	0.0	Poetry	19
12	Starving Hea	13.99	In stock	2	13.99	13.99	0.0	Default	19
13	Shakespeare	20.66	In stock	4	20.66	20.66	0.0	Poetry	19
14	Set Me Free	17.46	In stock	5	17.46	17.46	0.0	Young Adult	19
15	Scott Pilgrim	52.29	In stock	5	52.29	52.29	0.0	Sequential A	19
16	Rip it Up and	35.02	In stock	5	35.02	35.02	0.0	Music	19
17	Our Band Co	57.25	In stock	3	57.25	57.25	0.0	Music	19
18	Olio	23.88	In stock	1	23.88	23.88	0.0	Poetry	19
19	Mesaerion: T	37.59	In stock	1	37.59	37.59	0.0	Science Ficti	19
20	Libertarianis	51.33	In stock	2	51.33	51.33	0.0	Politics	19
21	It's Only the	45.17	In stock	2	45.17	45.17	0.0	Travel	19
22	In Her Wake	12.84	In stock	1	12.84	12.84	0.0	Thriller	19
23	How Music V	37.32	In stock	2	37.32	37.32	0.0	Music	19
24	Foolproof Pr	30.52	In stock	3	30.52	30.52	0.0	Food and Dri	19
25	Chase Me (P	25.27	In stock	5	25.27	25.27	0.0	Romance	19
26	Black Dust	34.53	In stock	5	34.53	34.53	0.0	Romance	19
27	Birdsong: A	54.64	In stock	3	54.64	54.64	0.0	Childrens	19
28	America's Cr	22.5	In stock	3	22.5	22.5	0.0	Default	19
29	Aladdin and	53.13	In stock	3	53.13	53.13	0.0	Default	19
30	Worlds Elsev	40.3	In stock	5	40.3	40.3	0.0	Nonfiction	18
31	Wall and Pie	44.18	In stock	4	44.18	44.18	0.0	Art	18
32	The Four Agr	17.66	In stock	5	17.66	17.66	0.0	Spirituality	18
33	The Five Love	31.05	In stock	3	31.05	31.05	0.0	Nonfiction	18
34	The Elephant	23.82	In stock	5	23.82	23.82	0.0	Thriller	18
35	The Bear anc	36.89	In stock	1	36.89	36.89	0.0	Childrens	18
36	Sophie's Wor	15.94	In stock	5	15.94	15.94	0.0	Philosophy	18
37	Pennv Mavhr	33.29	In stock	3	33.29	33.29	0.0	Default	18

## 3 – Base de données

### Code : Insertion des livres dans la base de données

```
28
29
30 * CREATE TABLE 'livres' (
31   'id' int(11) NOT NULL,
32   'titre' varchar(255) NOT NULL,
33   'prix' decimal(10,2) NOT NULL,
34   'disponibilite' varchar(50) NOT NULL,
35   'etoiles' int(11) DEFAULT NULL CHECK ('etoiles' between 0 and 5),
36   'prix_ht' decimal(10,2) NOT NULL,
37   'prix_ttc' decimal(10,2) NOT NULL,
38   'taxe' decimal(10,2) NOT NULL,
39   'categorie' varchar(100) NOT NULL,
40   'stock' int(11) NOT NULL,
41   'vente' int(11) DEFAULT NULL,
42 ) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_general_ci;
43
44 --
45 -- Déchargement des données de la table 'livres'
46 --
47
48 * INSERT INTO 'livres' ('id', 'titre', 'prix', 'disponibilite', 'etoiles', 'prix_ht', 'prix_ttc', 'taxe', 'categorie', 'stock', 'vente') VALUES
49 (1024, 'A Light in the Attic', 51.77, 'In stock', 3, 51.77, 51.77, 0.00, 'Poetry', 22, 18),
50 (1025, 'Tipping the Velvet', 53.74, 'In stock', 1, 53.74, 53.74, 0.00, 'Historical Fiction', 20, 20),
51 (1026, 'Soumission', 50.10, 'In stock', 1, 50.10, 50.10, 0.00, 'Fiction', 20, 20),
52 (1027, 'Sharp Objects', 47.82, 'In stock', 4, 47.82, 47.82, 0.00, 'Mystery', 20, 20),
53 (1028, 'Sapiens: A Brief History of Humankind', 54.23, 'In stock', 5, 54.23, 54.23, 0.00, 'History', 20, 20),
54 (1029, 'The Requiem Red', 22.65, 'In stock', 1, 22.65, 22.65, 0.00, 'Young Adult', 19, 21),
55 (1030, 'The Dirty Little Secrets of Getting Your Dream Job', 33.34, 'In stock', 4, 33.34, 33.34, 0.00, 'Business', 19, 21),
56 (1031, 'The Coming Woman: A Novel Based on the Life of the Infamous Feminist, Victoria Woodhull', 17.93, 'In stock', 3, 17.93, 17.93, 0.00, 'Default', 19, 21),
57 (1032, 'The Boys in the Boat: Nine Americans and Their Epic Quest for Gold at the 1936 Berlin Olympics', 22.60, 'In stock', 4, 22.60, 22.60, 0.00, 'Default', 19, 21),
58 (1033, 'The Black Maria', 52.15, 'In stock', 1, 52.15, 52.15, 0.00, 'Poetry', 19, 21),
59 (1034, 'Starving Hearts (Triangular Trade Trilogy, #1)', 13.99, 'In stock', 2, 13.99, 13.99, 0.00, 'Default', 19, 21),
60 (1035, 'Shakespeare's Sonnets', 20.66, 'In stock', 4, 20.66, 20.66, 0.00, 'Poetry', 19, 21),
61 (1036, 'Set Me Free', 17.46, 'In stock', 5, 17.46, 17.46, 0.00, 'Young Adult', 19, 21),
62 (1037, 'Scott Pilgrim's Precious Little Life (Scott Pilgrim #1)', 52.29, 'In stock', 5, 52.29, 52.29, 0.00, 'Sequential Art', 19, 21),
63 (1038, 'Rip It Up and Start Again', 35.02, 'In stock', 5, 35.02, 35.02, 0.00, 'Music', 19, 21),
64 (1039, 'Our Band Could Be Your Life: Scenes from the American Indie Underground, 1981-1991', 57.25, 'In stock', 3, 57.25, 57.25, 0.00, 'Music', 19, 21),
65 (1040, 'Olio', 23.88, 'In stock', 1, 23.88, 23.88, 0.00, 'Poetry', 19, 21),
66 (1041, 'Mesaeron: The Best Science Fiction Stories 1880-1840', 37.59, 'In stock', 1, 37.59, 37.59, 0.00, 'Science Fiction', 19, 21),
67 (1042, 'Libertarianism for Beginners', 51.33, 'In stock', 2, 51.33, 51.33, 0.00, 'Politics', 19, 21),
68 (1043, 'It's Only the Himalayas', 45.17, 'In stock', 2, 45.17, 45.17, 0.00, 'Travel', 19, 21),
69 (1044, 'In Her Wake', 12.84, 'In stock', 1, 12.84, 12.84, 0.00, 'Thriller', 19, 21),
70 (1045, 'How Much Water?', 37.52, 'In stock', 2, 37.52, 37.52, 0.00, 'Music', 19, 21),
71 (1046, 'Foodproof Preserving: A Guide to Small Batch Jams, Jellies, Pickles, Condiments, and More: A Foolproof Guide to Making Small Batch Jams, Jellies, Pickles, Condiments, and More: A Foolproof Guide to Making Small Batch Jams, Jellies, Pickles, Condiments, and More', 30.52, 'In stock', 3, 30.52, 30.52, 0.00, 'Food and Drink', 19, 21),
72 (1047, 'Chase Me (Paris Nights #2)', 25.27, 'In stock', 5, 25.27, 25.27, 0.00, 'Romance', 19, 21),
73 (1048, 'Black Out', 34.53, 'In stock', 5, 34.53, 34.53, 0.00, 'Romance', 19, 21),
74 (1049, 'Hikings: A Story in Pictures', 54.04, 'In stock', 3, 54.04, 54.04, 0.00, 'Children', 19, 21),
75 (1050, 'America's Cradle of Quakerhicks: Western Pennsylvania's Football Factory from Johnny Unitas to Joe Montana', 22.58, 'In stock', 3, 22.58, 22.58, 0.00, 'Default', 19, 21),
76 (1051, 'Aladdin and His Wonderful Lamp', 53.13, 'In stock', 3, 53.13, 53.13, 0.00, 'Default', 19, 21),
```

### Résultat :

id	titre	prix	disponibilite	etoiles	prix_ht	prix_ttc	taxe	categorie	stock
1024	A Light in the Attic	51.77	In stock	3	51.77	51.77	0.00	Poetry	22
1025	Tipping the Velvet	53.74	In stock	1	53.74	53.74	0.00	Historical Fiction	20
1026	Soumission	50.10	In stock	1	50.10	50.10	0.00	Fiction	20
1027	Sharp Objects	47.82	In stock	4	47.82	47.82	0.00	Mystery	20
1028	Sapiens: A Brief History of Humankind	54.23	In stock	5	54.23	54.23	0.00	History	20
1029	The Requiem Red	22.65	In stock	1	22.65	22.65	0.00	Young Adult	19
1030	The Dirty Little Secrets of Getting Your Dream Job	33.34	In stock	4	33.34	33.34	0.00	Business	19
1031	The Coming Woman: A Novel Based on the Life of the...	17.93	In stock	3	17.93	17.93	0.00	Default	19
1032	The Boys in the Boat: Nine Americans and Their Epi...	22.60	In stock	4	22.60	22.60	0.00	Default	19
1033	The Black Maria	52.15	In stock	1	52.15	52.15	0.00	Poetry	19
1034	Starving Hearts (Triangular Trade Trilogy, #1)	13.99	In stock	2	13.99	13.99	0.00	Default	19
1035	Shakespeare's Sonnets	20.66	In stock	4	20.66	20.66	0.00	Poetry	19
1036	Set Me Free	17.46	In stock	5	17.46	17.46	0.00	Young Adult	19
1037	Scott Pilgrim's Precious Little Life (Scott Pilgr...	52.29	In stock	5	52.29	52.29	0.00	Sequential Art	19
1038	Rip It Up and Start Again	35.02	In stock	5	35.02	35.02	0.00	Music	19
1039	Our Band Could Be Your Life: Scenes from the Amer...	57.25	In stock	3	57.25	57.25	0.00	Music	19
1040	Olio	23.88	In stock	1	23.88	23.88	0.00	Poetry	19
1041	Mesaeron: The Best Science Fiction Stories 1800-1...	37.59	In stock	1	37.59	37.59	0.00	Science Fiction	19
1042	Libertarianism for Beginners	51.33	In stock	2	51.33	51.33	0.00	Politics	19
1043	It's Only the Himalayas	45.17	In stock	2	45.17	45.17	0.00	Travel	19
1044	In Her Wake	12.84	In stock	1	12.84	12.84	0.00	Thriller	19

## 4 – Analyse et Visualisation

Code :

```
analyse.py
1 import mysql.connector
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import numpy as np
6
7 # Connexion a la bdd mysql
8 try:
9     conn = mysql.connector.connect(
10         host="127.0.0.1",      # @serv
11         user="root",
12         password="",
13         database="livre"      # Nom bdd
14     )
15     print("Connexion à la base de données réussie !")
16 except mysql.connector.Error as err:
17     print(f"Erreur de connexion à la base de données : {err}")
18     exit()
19
20 # Récup data table livre
21 cursor = conn.cursor()
22 query = "SELECT * FROM livres"
23 cursor.execute(query)
24
25 # Convertir les résultats en DataFrame pandas
26 columns = [col[0] for col in cursor.description] # recuperation des noms des col
27 data = cursor.fetchall() # recuperation des lignes
28 df = pd.DataFrame(data, columns=columns)
29
30 # Ajout d la colonne vente au DataFrame
31 df['vente'] = 40 - df['stock'] # Calcul de la colonne vente
32
33 # MAJ de la bdd MySQL avec la nouvelle col vente
34 try:
35     # Ajout la colonne vente à la table livres si elle n'existe pas déjà
36     cursor.execute("ALTER TABLE livres ADD COLUMN IF NOT EXISTS vente INT")
37
38     # MAJ chaque ligne avec la valeur de vente
39     for index, row in df.iterrows():
40         update_query = f"UPDATE livres SET vente = {row['vente']} WHERE id = {row['id']}"
41         cursor.execute(update_query)
42
43     # Valider les modif dans la bdd
44     conn.commit()
45     print("Colonne 'vente' ajoutée et mise à jour dans la base de données !")
46 except mysql.connector.Error as err:
47     print(f"Erreur lors de la mise à jour de la base de données : {err}")
48 finally:
49     # Fermer le curseur et la connexion
50     cursor.close()
51     conn.close()
52
53 # Vérif et convertir les types de données
54 df['etoiles'] = df['etoiles'].astype(int)
55 df['stock'] = df['stock'].astype(int)
56 df['prix'] = df['prix'].astype(float)
57 df['vente'] = df['vente'].astype(int)
58
59 # Palette de couleurs commune
```

```

116 plt.show()
117 plt.close()
118
119 # Figure 4 : Somme des livres en stock par catégorie (trié par ordre alphabétique)
120 stock_par_catégorie = df.groupby('catégorie')['stock'].sum().reset_index()
121
122 stock_par_catégorie = stock_par_catégorie.sort_values(by='catégorie')
123
124 plt.figure(figsize=(12, 6))
125 sns.barplot(x='catégorie', y='stock', data=stock_par_catégorie, hue='catégorie', palette=palette, legend=False)
126 plt.title('Somme des livres en stock par catégorie')
127 plt.xlabel('Catégorie')
128 plt.ylabel('Somme du stock')
129 plt.xticks(rotation=45, ha='right')
130 plt.tight_layout()
131 plt.show()
132 plt.close()
133
134 # Figure 4bis : Somme des ventes par catégorie (trié par ordre alphabétique)
135 ventes_par_catégorie = df.groupby('catégorie')['vente'].sum().reset_index()
136
137 ventes_par_catégorie = ventes_par_catégorie.sort_values(by='catégorie')
138
139 plt.figure(figsize=(12, 6))
140 sns.barplot(x='catégorie', y='vente', data=ventes_par_catégorie, hue='catégorie', palette=palette, legend=False)
141 plt.title('Somme des ventes par catégorie')
142 plt.xlabel('Catégorie')
143 plt.ylabel('Somme des ventes')
144 plt.xticks(rotation=45, ha='right')
145 plt.tight_layout()
146 plt.show()
147 plt.close()
148
149 # Figure 5 : Nombre de livres par intervalle de prix (graphique en barres)
150 # Creer des intervalle de prix en excluant les deux premier
151 bins = np.arange(10, df['prix'].max() + 5, 5) # Commence à 10 au lieu de 0
152 df['Prix Intervalle'] = pd.cut(df['prix'], bins=bins)
153
154 # Compte le nombre de livres par intervalle de prix
155 prix_counts = df['Prix Intervalle'].value_counts().sort_index().reset_index()
156 prix_counts.columns = ['Prix Intervalle', 'count']
157
158 # Filtre pour exclure les intervalles vides (au cas où des fois que)
159 prix_counts = prix_counts.dropna()
160
161 plt.figure(figsize=(12, 6))
162 sns.barplot(x='Prix Intervalle', y='count', data=prix_counts, hue='Prix Intervalle', palette=palette, legend=False)
163 plt.title('Nombre de livres par intervalle de prix')
164 plt.xlabel('Intervalle de prix')
165 plt.ylabel('Nombre de livres')
166 plt.xticks(rotation=45, ha='right')
167 plt.tight_layout()
168 plt.show()
169 plt.close()

```

```

59 # Palette de couleurs commune
60 palette = 'viridis'
61
62 # Figure 1 : Nombre de livres par nombre d'étoiles (graphique en barres)
63 etoiles_counts = df['etoiles'].value_counts().sort_index().reset_index()
64 etoiles_counts.columns = ['etoiles', 'count']
65
66 plt.figure(figsize=(10, 6))
67 sns.barplot(x='etoiles', y='count', data=etoiles_counts, hue='etoiles', palette=palette, legend=False)
68 plt.title('Nombre de livres par nombre d\'étoiles')
69 plt.xlabel('Nombre d\'étoiles')
70 plt.ylabel('Nombre de livres')
71 plt.xticks(rotation=45, ha='right')
72 plt.tight_layout()
73 plt.show()
74 plt.close()
75
76 # Figure 2 : Somme des livres en stock par nombre d'étoiles
77 stock_par_etoiles = df.groupby('etoiles')['stock'].sum().reset_index()
78
79 plt.figure(figsize=(10, 6))
80 sns.barplot(x='etoiles', y='stock', data=stock_par_etoiles, hue='etoiles', palette=palette, legend=False)
81 plt.title('Somme des livres en stock par nombre d\'étoiles')
82 plt.xlabel('Nombre d\'étoiles')
83 plt.ylabel('Somme du stock')
84 plt.xticks(rotation=45, ha='right')
85 plt.tight_layout()
86 plt.show()
87 plt.close()
88
89 # Figure 2bis : Somme des ventes par nombre d'étoiles
90 ventes_par_etoiles = df.groupby('etoiles')['vente'].sum().reset_index()
91
92 plt.figure(figsize=(10, 6))
93 sns.barplot(x='etoiles', y='vente', data=ventes_par_etoiles, hue='etoiles', palette=palette, legend=False)
94 plt.title('Somme des ventes par nombre d\'étoiles')
95 plt.xlabel('Nombre d\'étoiles')
96 plt.ylabel('Somme des ventes')
97 plt.xticks(rotation=45, ha='right')
98 plt.tight_layout()
99 plt.show()
100 plt.close()
101
102 # Figure 3 : Nombre de livres par catégorie (graphique en barres, trié par ordre alphabétique)
103 categorie_counts = df['categorie'].value_counts().reset_index()
104 categorie_counts.columns = ['categorie', 'count']
105
106 # Trier les catégories par ordre alphabétique
107 categorie_counts = categorie_counts.sort_values(by='categorie')
108
109 plt.figure(figsize=(12, 6))
110 sns.barplot(x='categorie', y='count', data=categorie_counts, hue='categorie', palette=palette, legend=False)
111 plt.title('Nombre de livres par catégorie ')
112 plt.xlabel('Catégorie')
113 plt.ylabel('Nombre de livres')
114 plt.xticks(rotation=45, ha='right')
115 plt.tight_layout()
116 plt.show()

```

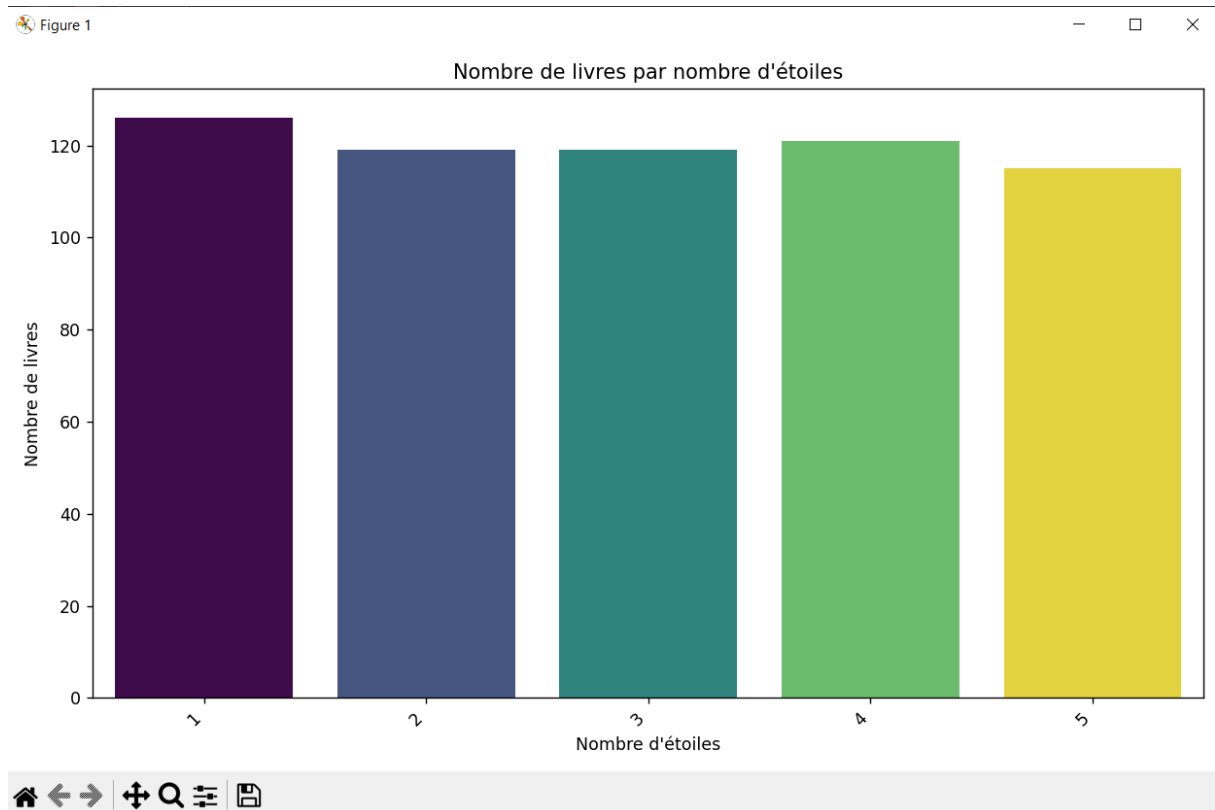


```

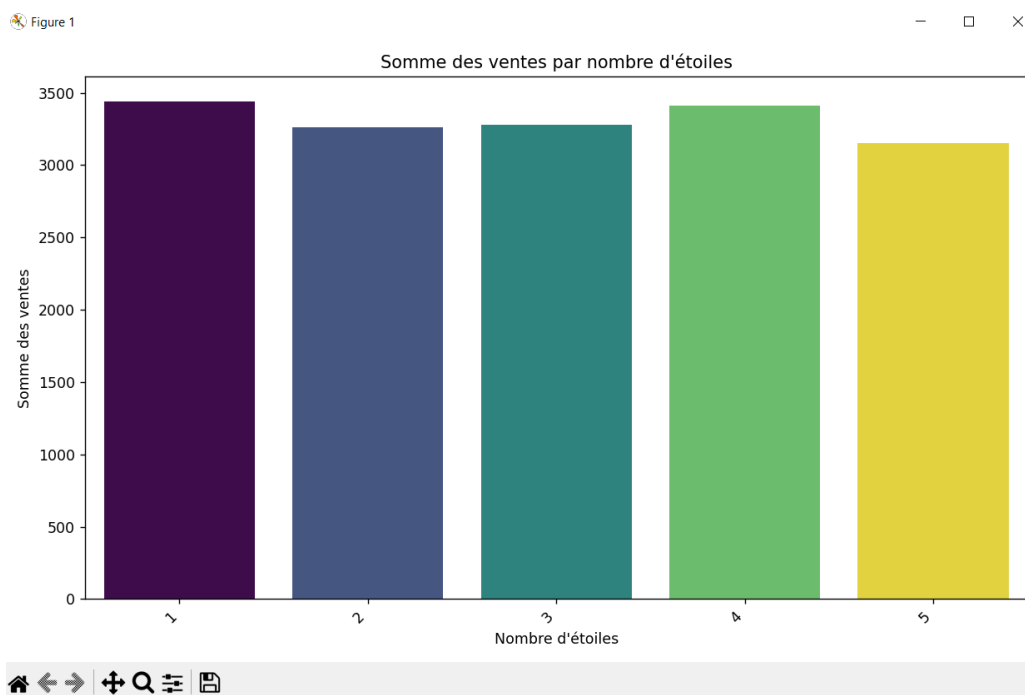
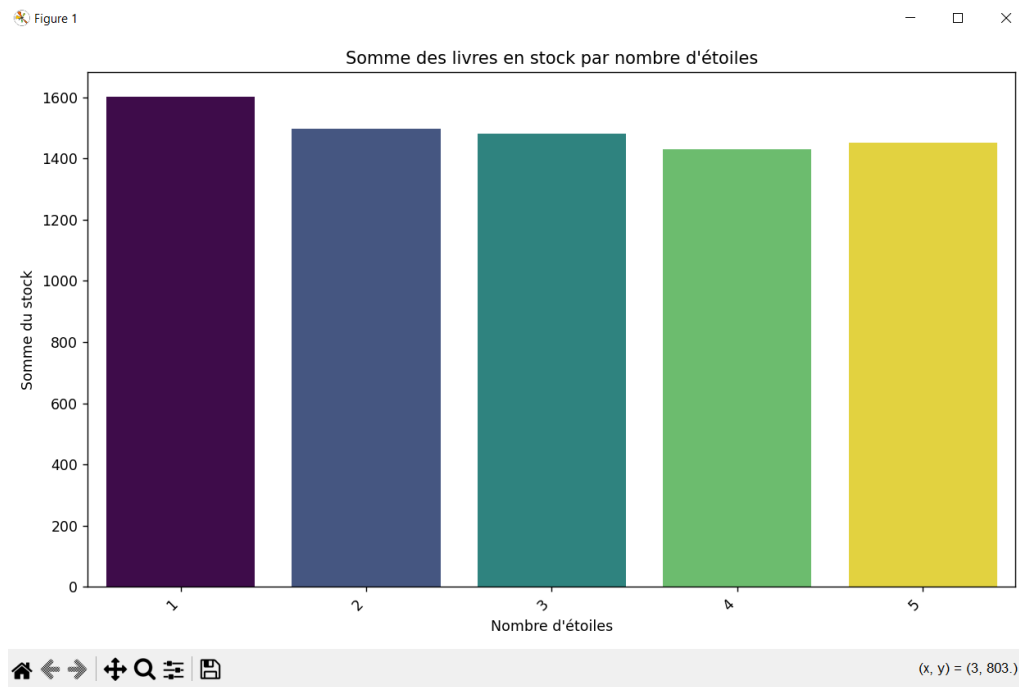
170
171 # Figure 6 : Somme des livres en stock par intervalle de prix
172 # Utiliser les mêmes intervalles que pour la Figure 5
173 stock_par_prix = df.groupby('Prix Intervalle', observed=True)['stock'].sum().reset_index()
174
175 # Filtrer pour exclure les intervalles vides (au cas où)
176 stock_par_prix = stock_par_prix.dropna()
177
178 plt.figure(figsize=(12, 6))
179 sns.barplot(x='Prix Intervalle', y='stock', data=stock_par_prix, hue='Prix Intervalle', palette=palette, legend=False, dodge=False)
180 plt.title('Somme des livres en stock par intervalle de prix')
181 plt.xlabel('Intervalle de prix')
182 plt.ylabel('Somme du stock')
183 plt.xticks(rotation=45, ha='right')
184 plt.tight_layout()
185 plt.show()
186 plt.close()
187
188 # Figure 6bis : Somme des ventes par intervalle de prix
189 ventes_par_prix = df.groupby('Prix Intervalle', observed=True)['vente'].sum().reset_index()
190
191 plt.figure(figsize=(12, 6))
192 sns.barplot(x='Prix Intervalle', y='vente', data=ventes_par_prix, hue='Prix Intervalle', palette=palette, legend=False, dodge=False)
193 plt.title('Somme des ventes par intervalle de prix')
194 plt.xlabel('Intervalle de prix')
195 plt.ylabel('Somme des ventes')
196 plt.xticks(rotation=45, ha='right')
197 plt.tight_layout()
198 plt.show()
199 plt.close()

```

**Résultats :** A savoir, on part du principe que chaque livre a été acheté par le site au fournisseur en 40 exemplaires !

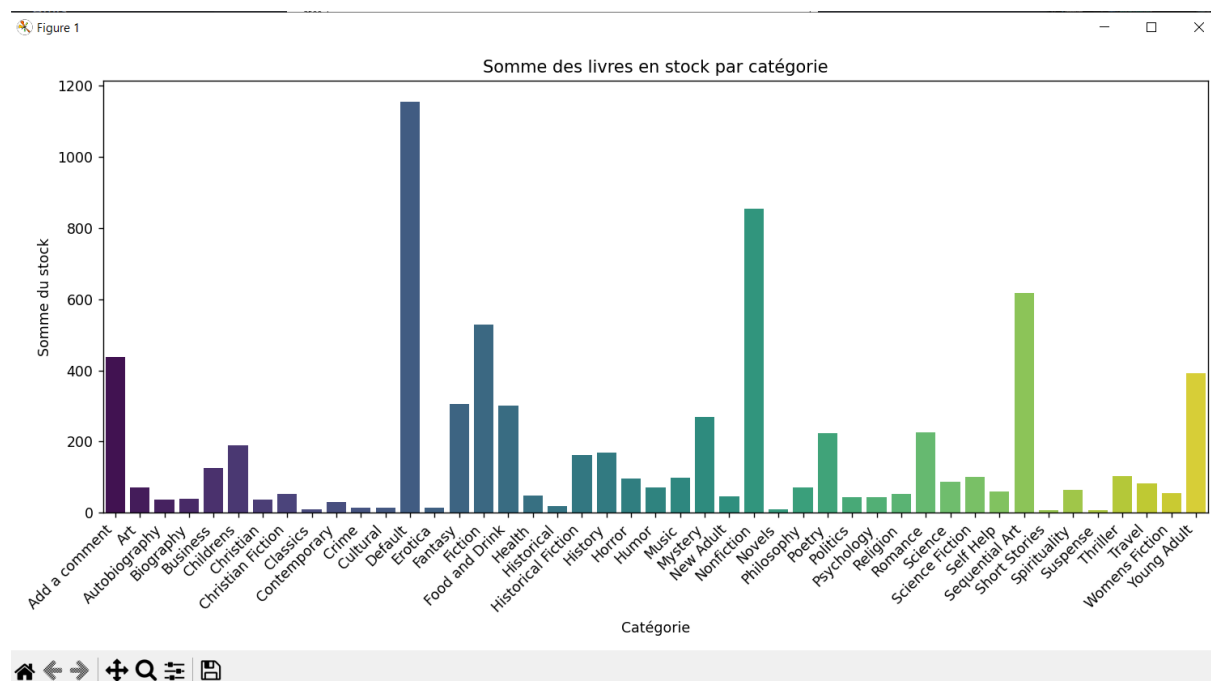


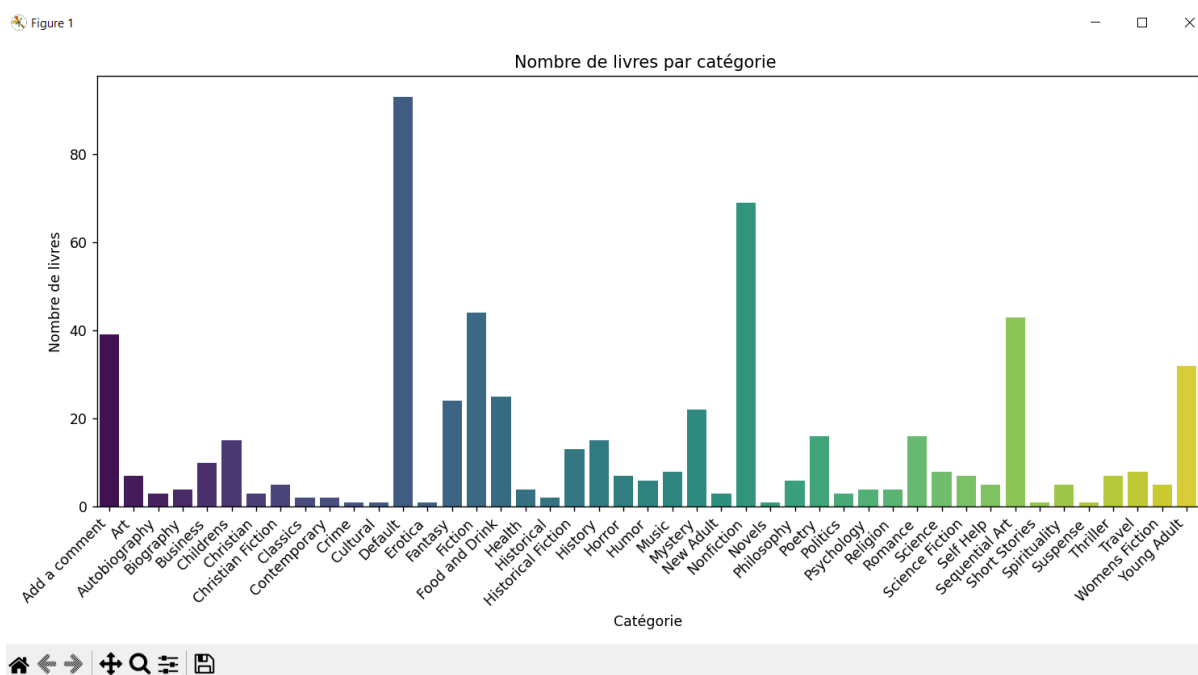
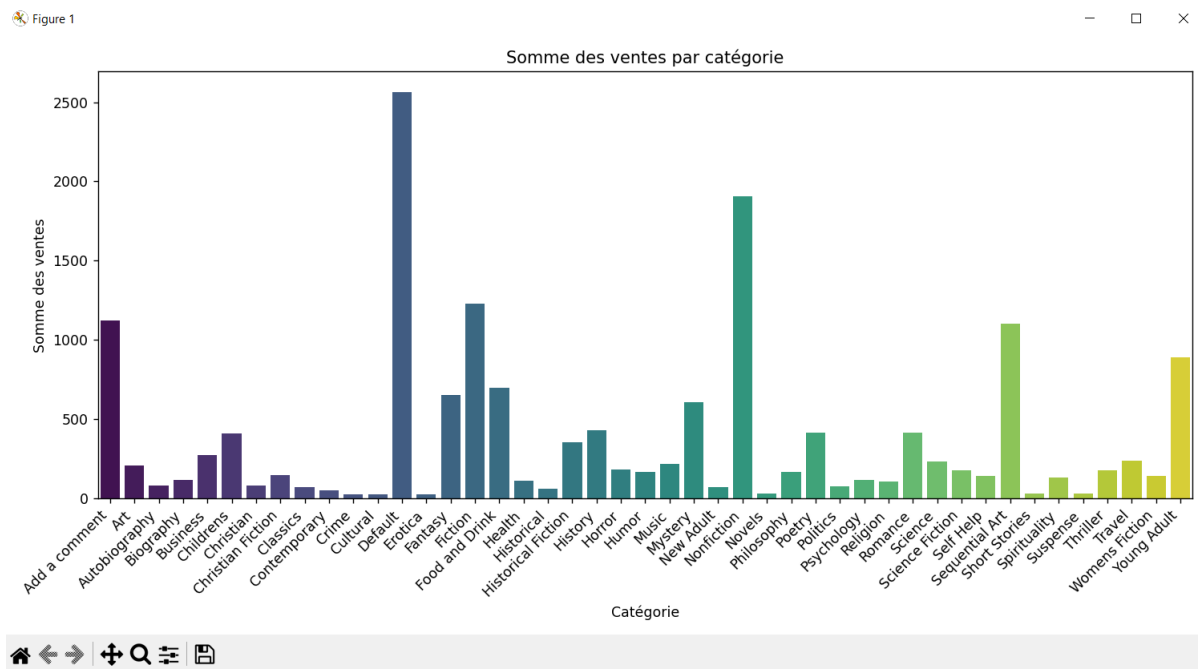




Ici, le premier graphique nous montre qu'il y a environ 125 livres (unique) qui ont 1 étoile et environ 115 livres (unique) qui ont 5 étoiles. On peut donc partir du principe qu'au total la somme de la commande passé de livre était de  $125 \times 40 (=5000)$  (pour les livres 1 étoile et de  $115 \times 40 (=4500)$  pour les livres 5 étoiles, (il y a donc là un ratio de 9/10ieme) et si on va sur le graphique numéro 3 nous pouvons simplement voir que les ordres de grandeur sont les mêmes que sur le graphique 1. On peut quand

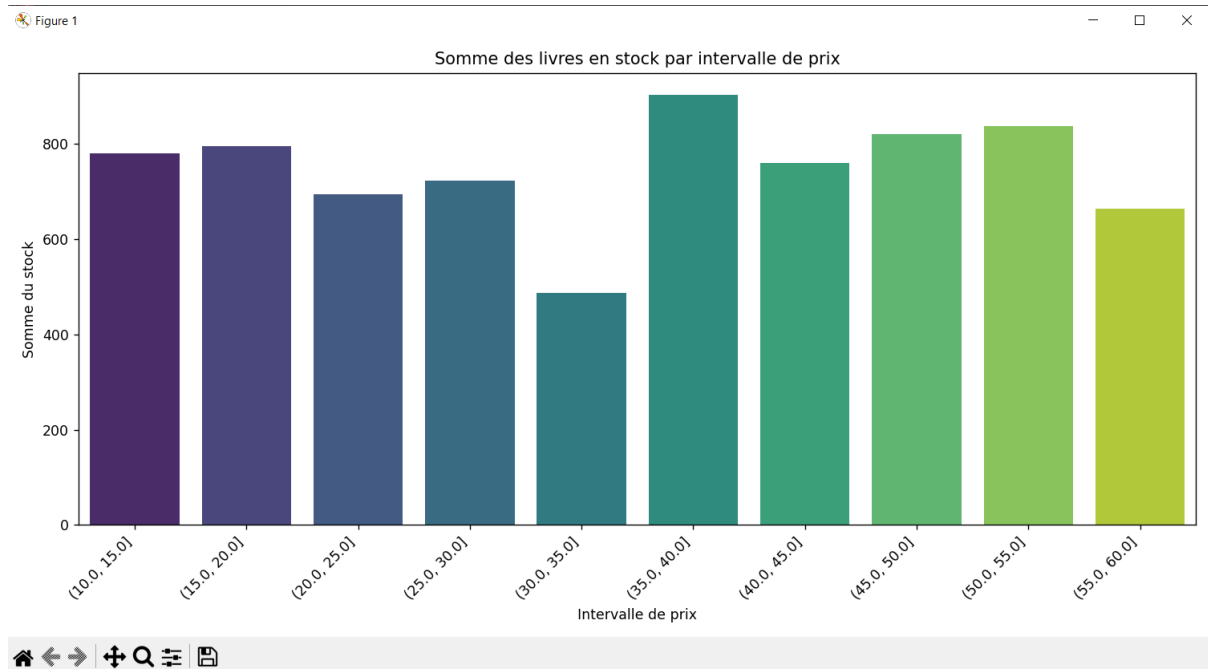
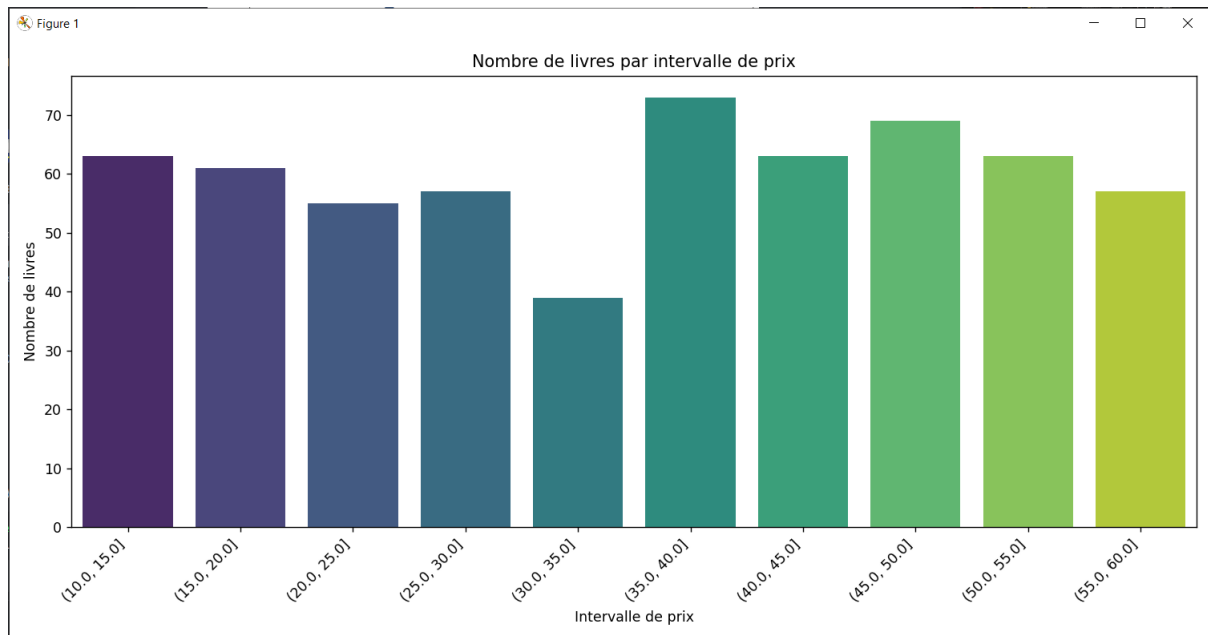
même dire que les livres 4 étoiles se sont sûrement un peu mieux vendu que les autres d'après le graphique numéro 2 car on sait que les livres ont été acheté dans les mêmes quantité peu importe le nombre d'étoile, et l'ordre de grandeur du graphique numéro ne se réplique pas sur le deuxième graphique donc on sait que les 4 étoiles se sont un peu mieux vendu, et la raison pour laquelle cette interprétation ne saute pas aux yeux sur le graphique 3, c'est parce qu'on peut croire que l'ordre de grandeur est respecté mais si on regarde attentivement le premier graphique le nombre de livre unique est de 125 pour les 1 étoile et de 120 pour les 4 étoiles ce qui nous fait un ratio de 0,96 ( $120/125$ ) alors qu'au niveau du graphique 2, on voit environ 1600 livres en stock pour les 1 étoile et 1400 livres en stock pour les 4 étoiles donc un ratio de 0,87 ! Il y a donc eu une légère surperformance au niveau de vente de ces livres 4 étoiles, mais rien de vraiment significatif. En ayant l'œil, on peut directement le voir sur le graphique 3 car les 4 étoiles ont légèrement surperformé les attentes mais ce n'est pas évident à avoir, il faut avoir le coup d'œil !

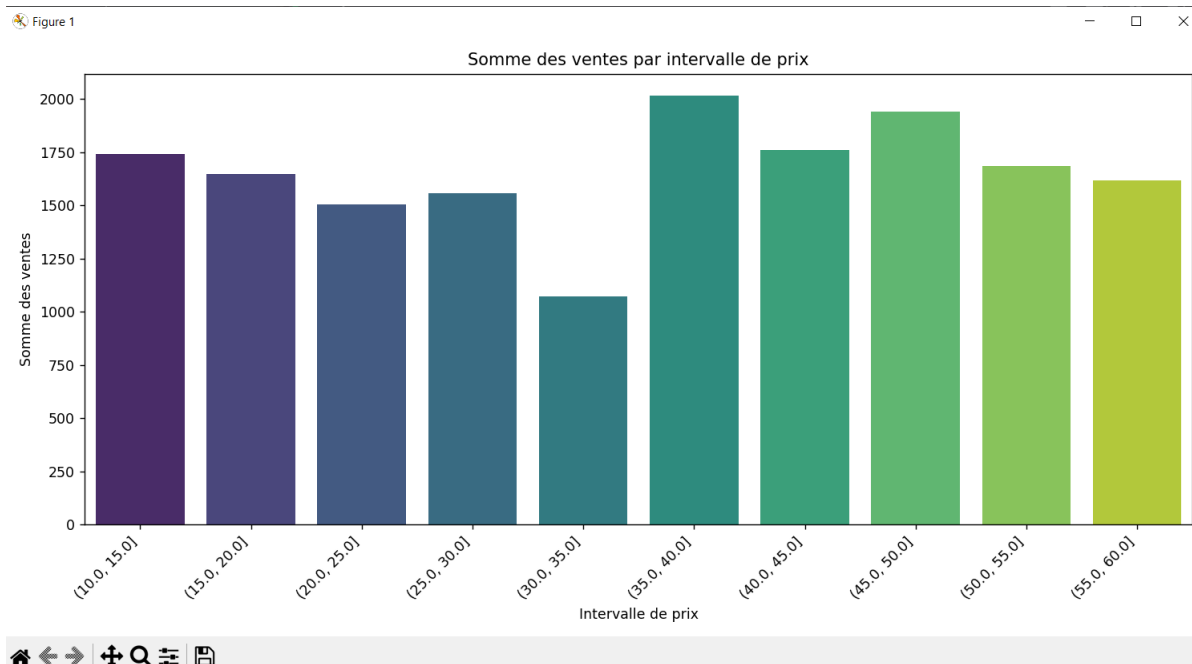




Au niveau du graphique 1 on peut voir que certaines catégories, comme "Default", "Young Adult" et "Philosophy", sont beaucoup plus représentées que d'autres, ce qui peut indiquer des tendances éditoriales. En ce qui concerne le graphique 2, les stocks sont plus élevés pour certaines catégories populaires, ce qui suggère que l'offre suit la demande du marché de la part de camarade de chez bookstoscape ! Et d'après le graphique 3, on observe que les catégories ayant le plus de ventes

correspondent souvent aux plus représentées il y a donc corrélation entre les volumes d'achats, les ventes et donc très logiquement les stocks.





On peut commencer par analyser ce graphique 3 qui démontre clairement que les livres avec un prix plus élevé se sont vendus dans de plus grande quantité en général (prix supérieur à 35) mais que cela correspond aussi au volume de livre unique et donc d'exemplaire au global de ces tranches de prix respective !

Sur cette séquence de graphique on peut voir 3 tranches intéressante de 40 £ à 55 £. Car sachant le nombre d'exemplaire unique de 40 £ à 45 £ et de 50 £ à 55 £ à 65 livres on sait qu'il y a eu 2600 livres en tout des ces catégories respectives achetés donc moins que dans la catégorie des livres qui ont un prix compris entre 50 £ et 55 £ qui eux sont au nombre de 70 exemplaires uniques donc  $70 \times 40 = 2800$  exemplaires sans considération du livre. En l'occurrence on peut voir sur le graphique 2 que les stocks des livres de 50 £ à 55 £ sont plus élevé que celle des 40 £ à 45 £ et cela correspond parfaitement avec le graphique 3 qui montre bien que les ventes sont moins élevées pour une quantité initiale de livre pourtant supérieur.

### IMPORTANT !

Les analyses faites ici ne disent clairement rien de concret car les données entrées sur le site sont purement fictives, elles ne sont représentatives d'aucun comportement humain par conséquent on peut voir que peu importe la manière dont on souhaite analyser ces données, elles ont été entrées de manière très lisse dans ce site fictif par conséquent il est difficile de dégager de quelconque tendance !

# Conclusion

Ce projet nous a permis de mettre en place un processus complet de scraping, nettoyage, stockage et analyse de données. En partant de l'extraction automatisée des livres sur **Books to Scrape**, nous avons structuré et stocké ces informations dans une base de données **MySQL**, garantissant ainsi une gestion efficace et une exploitation optimale des données.

Grâce au nettoyage des valeurs extraites et à la conversion des types, nous avons pu obtenir une base de données fiable, évitant ainsi les erreurs lors de l'analyse. Enfin, les visualisations graphiques générées ont permis d'identifier des tendances intéressantes concernant les prix ou la notation moyenne des livres.

Ce projet illustre bien l'importance d'une approche méthodique et structurée dans la gestion des données.

En plus des aspects techniques, ce travail nous a apporté une meilleure compréhension des défis liés à la collecte et au traitement des données en conditions réelles, ainsi qu'une méthodologie efficace pour exploiter ces informations de manière pertinente.

Enfin, ce projet nous aura permis de constituer une base solide pour nos futures missions en entreprise, où la gestion et l'analyse des données joueront un rôle clé.