

Importation des libraries

```
1 import requests
2 from bs4 import BeautifulSoup
3 from urllib.parse import urljoin
4 import pandas as pd
```

On récupère le code HTML de l'url

```
1 url = "https://books.toscrape.com/index.html"
2 response = requests.get(url)
3 soup = BeautifulSoup(response.text, 'html.parser')
```

On cherche toutes les catégories qu'on stock dans une liste.

```
1 categories = soup.select('ul.nav-list > li > ul > li > a')
2 categories_list = [category.text.strip() for category in categories]
```

Cette méthode spécifie la hiérarchie d'éléments à rechercher. En décomposant cela :

```
▼ <div class="side_categories">
  ::before
  ▼ <ul class="nav nav-list">
    ::before
    ▼ <li>
      ► <a href="catalogue/category/books_1/index.html">... </a>
      ▼ <ul>
        ▼ <li>
          ::marker
          ▼ <a href="catalogue/category/books/travel_2/index.html">
            " Travel " == $0
          </a>
        </li>
        ► <li>... </li>
        ► <li>... </li>
```

- `ul.nav-list` : Recherche un élément `` ayant la classe `nav-list`.
- `> li` : Sélectionne un élément `` qui est après l'élément ``.
- `> ul` : Sélectionne un élément `` qui est après l'élément ``.
- `> li` : Sélectionne un élément `` qui est après l'élément ``.
- `> a` : Sélectionne un élément `<a>` qui est après l'élément ``.

On crée une boucle qui permet de demander à l'utilisateur de saisir le nom de catégorie qu'il souhaite scraper. Il a cinq catégories à choisir, donc la boucle s'exécute cinq fois. Si une catégorie n'existe pas, soit n'est pas présente dans la liste `categories_list`, alors un message d'erreur s'affiche pour entrer une nouvelle catégorie présente dans la liste.

```
1 category_select = []
2
3 for i in range(5):
4     while True:
5         user_input = input(f"Entrez le nom de la catégorie {i + 1} à extraire : ")
6         if user_input in categories_list and user_input not in category_select:
7             category_select.append(user_input)
8             break
```

```

9         else:
10             print("La catégorie n'est pas valide. Veuillez réessayer.")

```

Dans la variable 'category_select', nous devons obtenir une liste avec cinq catégories.

Désormais, on crée une liste 'book_data' qui va stocker les éléments qu'on va scraper par livre. Pour chaque nom de catégorie, dans la variable 'category_select', on va parcourir le code HTML du lien associé à la catégorie. Cependant, les différentes catégories peuvent avoir plusieurs pages, alors on crée une condition pour exécuter le code tant qu'une nouvelle page existe.

Le code de la condition qui parcourt les pages, ainsi que la boucle qui parcourt les liens de chaque catégorie, est de la forme suivante :

```

1 book_data = []
2
3 for category_name in category_select:
4     index = categories_list.index(category_name)
5     category_link = urljoin(url, categories[index]['href'])
6
7     page_number = 1
8     while True:
9         if page_number == 1:
10             category_website = requests.get(category_link)
11         else:
12             category_website = requests.get(urljoin(category_link,
13                                                     f"page-{page_number}.html"))
14
15
16         soup = BeautifulSoup(category_website.text, 'html.parser')
17         product_containers = soup.find_all('article', class_='product_pod')
18
19         # Code qui récupérer des informations dans le code HTML de chaque page
20
21         page_number += 1
22     else:
23         break # la boucle s'arrête quand on atteint le nombre de page possible

```

Il faut parcourir l'ensemble du code HTML de chaque catégorie par page, pour obtenir le nom de chaque livre ainsi que son lien. La liste 'product_containers' stock le code HTML où se trouve chaque livre. Pour un seul livre, nous avons les informations suivantes (des informations sur l'image, nombre d'étoiles, le titre avec le lien du livre, le prix et le stock :

```

▼ <article class="product_pod"> == $0
  ▶ <div class="image_container">...</div>
  ▶ <p class="star-rating Two">...</p>
  ▼ <h3>
    <a href=".../its-only-the-himalayas_981/index.html" title="It's Only the Himalayas">
      It's Only the Himalayas</a>
    </h3>
  ▼ <div class="product_price">
    <p class="price_color">£45.17</p>
    ▼ <p class="instock availability">
      ▶ <i class="icon-ok">...</i>
        " In stock "
      </p>
      ▶ <form>...</form>
    </div>
  </article>

```

A la suite de cela, on effectue une boucle sur cette liste pour parcourir tous les livres qui se trouve sur le lien de la catégorie. Sur chaque livre, on récupère le titre, la présence d'un stock disponible, le nombre d'étoiles, lien de l'image et le lien de la page du livre de la façon suivante.

```

1  # Récupérer des informations dans le code HTML de chaque page
2  for container in product_containers:
3
4      # Titre du livre
5      title = container.find('h3').find('a')['title']
6
7      # Stocks
8      stock_available = container.find('p', class_='instock availability').text.replace('\n', '')
9
10     # Lien du livre
11     book_url = container.find('h3').find('a')['href']
12     book_url = f'https://books.toscrape.com/catalogue/{book_url}'
13     book_url = book_url.replace('..../', '')
14
15     # Lien de l'image
16     img_url = container.find('img')['src']
17     img_url = f'https://books.toscrape.com{img_url}'
18     img_url = img_url.replace('..../', '')
19
20     # Nombre d'étoiles
21     star_rating_class = container.find('p', class_='star-rating')['class'][1]
22     star_rating = {'One': 1, 'Two': 2, 'Three': 3, 'Four': 4,
23                   'Five': 5}.get(star_rating_class, 0)

```

Le nombre d'étoiles est contenu dans la 'class'. Exemple : 'star-rating two'. Pour cela, on récupère le dernier mot de la class puis on le convertit en chiffre.

De plus, dans le code HTML, on constate que dans le lien de l'image et du livre, il y a des caractères type '../' en trop pour écrit qu'il soit écrit en entier. Il faut donc les supprimer afin d'obtenir un lien qui fonctionne, en le combinant avec le début du lien de la page d'accueil.

Il est possible d'obtenir plus d'informations en accédant à la page du livre. Pour cela, on récupère le code HTML du lien du livre préalablement enregistré dans la variable 'book_url'. Puis on parcourt le code HTML, pour obtenir des informations sur l'UPC, la description, le type de produit, le prix

avec la taxe et sans, la valeur de la taxe, le nombre de livres disponibles et le nombre de revues.

La structure du code HTML se présente comme suit :

```
▼<p class="instock availability"> == $0
  ▶<i class="icon-ok">...</i>
    " In stock (19 available) "
  </p>
  ▶<p class="star-rating Two">...</p>
  <hr>
  ▶<div class="alert alert-warning" role="alert">...</div>
</div>
<!-- /col-sm-6 -->
::after
</div>
<!-- /row -->
▼<div id="product_description" class="sub-header">
  <h2>Product Description</h2>
</div>
▼<p>
  "Wherever you go, whatever you do, just . . . don't do anything stupid." –My MotherDuring her
  yearlong adventure backpacking from South Africa to Singapore, S. Bedford definitely did a few
  things her mother might classify as "stupid." She swam with great white sharks in South Africa,
  ran from lions in Zimbabwe, climbed a Himalayan mountain without training in Nepal, and wa
  "Wherever you go, whatever you do, just . . . don't do anything stupid." –My MotherDuring her
```

Nous pouvons obtenir ici le nombre de livres en stocks et la description. Pour le nombre de livres en stock, on récupère le texte de la class "instock availability" et on remplace ' In stocks (' et ') ' par rien.

```
▶<div class="sub-header">...</div> == $0
▼<table class="table table-striped">
  ▼<tbody>
    ▼<tr>
      <th>UPC</th>
      <td>a22124811bfa8350</td>
    </tr>
    ▼<tr>
      <th>Product Type</th>
      <td>Books</td>
    </tr>
    ▼<tr>
      <th>Price (excl. tax)</th>
      <td>£45.17</td>
    </tr>
    ▼<tr>
      <th>Price (incl. tax)</th>
      <td>£45.17</td>
    </tr>
    ▼<tr>
      <th>Tax</th>
      <td>£0.00</td>
    </tr>
    ▼<tr>
      <th>Availability</th>
      <td>In stock (19 available)</td>
    </tr>
    ▼<tr>
      <th>Number of reviews</th>
      <td>0</td>
    </tr>
  </tbody>
</table>
```

Pour récupérer l'UPC, le type de produit, les différents prix et le nombre de revues, nous devons parcourir les différents 'th' qui contient l'information que l'on souhaite extraire. Une fois le bon 'th' trouvé, on récupère le texte qui se trouve dans la partie 'td' juste après le 'th'.

Pour obtenir ses informations, le code généré sur python est :

```
1 # Accéder à la page du détail du livre
2 book_detail_page = requests.get(book_url)
3 detail_soup = BeautifulSoup(book_detail_page.text, 'html.parser')
4
5 # Catégorie du livre
6 category = category_name
7
8 # Description
9 description = detail_soup.find('meta', {'name': 'description'})['content']
10 description = description.replace('\n', ' ').replace('\n', ' ')
11
12 # UPC
13 upc = detail_soup.find('th', string='UPC')
14 upc = upc.find_next('td').get_text(strip=True)
15
16 # Type de produit
17 product_type = detail_soup.find('th', string='Product Type')
18 product_type = product_type.find_next('td').get_text(strip=True)
19
20 # Price sans tax
21 price_excl_tax = detail_soup.find('th', string='Price (excl. tax)')
22 price_excl_tax = price_excl_tax.find_next('td').get_text(strip=True).replace('Â£', ' ')
23
24 # Prix avec la taxe
25 price_incl_tax = detail_soup.find('th', string='Price (incl. tax)')
26 price_incl_tax = price_incl_tax.find_next('td').get_text(strip=True).replace('Â£', ' ')
27
28 # Prix de la taxe
29 tax = detail_soup.find('th', string='Tax')
30 tax = tax.find_next('td').get_text(strip=True).replace('Â£', ' ')
31
32 # Nombre de livres disponibles
33 number_available = detail_soup.find_all('p', class_='instock availability')
34 number_available = number_available[0].get_text(strip=True).replace("In stock (", "").repl
35
36 # Nombre de revues
37 number_of_reviews = detail_soup.find('th', string='Number of reviews')
38 number_of_reviews = number_of_reviews.find_next('td').get_text(strip=True)
```

Chaque variable créée récupère les informations stockées dans celle-ci. Auparavant, nous avons créé une liste 'book_data' qui doit récupérer toutes ces informations. Nous pouvons désormais attribuer toutes ces variables à notre liste 'book_data' de la façon suivante.

```
1 # Ajouter les données du livre à la liste
2 book_data.append({
3     'Title': title,
4     'Book Link': book_url,
5     'Image Link': img_url,
6     'Description': description,
7     'UPC': upc,
8     'Category': category,
```

```

9      'Product Type': product_type,
10     'Price £ (excl. tax)': price_excl_tax,
11     'Price £ (incl. tax)': price_incl_tax,
12     'Tax': tax,
13     'Stock available' : stock_available,
14     'Number of available': number_available,
15     'Number of reviews': number_of_reviews
16 })

```

Chaque fois que la boucle s'exécute, cela crée un nouvel élément à la liste. Un élément, soit un livre, est représenté de la façon suivante dans la liste :

```

{'Title': 'Private Paris (Private #10)',
 'Book Link': 'https://books.toscrape.com/catalogue/private-paris-private-10_958/index.html',
 'Image Link': 'https://books.toscrape.com/media/cache/9d/05/9d0533bae1578846d728a82913b95c26.jpg',
 'Description': "Paris is burning--and only Private's Jack Morgan can put out the fire.When Jack Morgan stops by Private's Paris office, he envisions a quick hello during an otherwise relaxing trip filled with fine food and sightseeing. But Jack is quickly pressed into duty after a call from his client Sherman Wilkerson, asking Jack to track down his young granddaughter who is on the run f Paris is burning--and only Private's Jack Morgan can put out the fire.When Jack Morgan stops by Private's Paris office, he envisions a quick hello during an otherwise relaxing trip filled with fine food and sightseeing. But Jack is quickly pressed into duty after a call from his client Sherman Wilkerson, asking Jack to track down his young granddaughter who is on the run from a brutal drug dealer.Before Jack can locate her, several members of France's cultural elite are found dead--murdered in stunning, symbolic fashion. The only link between the crimes is a mysterious graffiti tag. As religious and ethnic tensions simmer in the City of Lights, only Jack and his Private team can connect the dots before the smoldering powder keg explodes. ...more",
 'UPC': 'b12b89017878a60d',
 'Category': 'Fiction',
 'Product Type': 'Books',
 'Price £ (excl. tax)': '47.61',
 'Price £ (incl. tax)': '47.61',
 'Tax': '0.00',
 'Stock available': 'In stock',
 'Number of available': '17',
 'Number of reviews': '0'}

```

Nous pouvons désormais convertir la liste 'book_data' en dataframe de type 'csv' de la façon suivante :

```

1  # Création du dataframe avec les 5 catégories
2  data_category = pd.DataFrame(book_data)
3  data_category.to_csv('books_data_category.csv', index = False)

```