

Introduction au Web Scraping

Définition et concepts

Le **web scraping** est une technique automatisée d'extraction de données depuis des sites web. BeautifulSoup est une bibliothèque Python qui facilite cette tâche en permettant de parser le HTML et d'en extraire les informations structurées.

Important

Éthique du scraping : Toujours respecter robots.txt, limiter la fréquence des requêtes, et vérifier les conditions d'utilisation des sites.

Installation des dépendances

```
1 # Installation des packages nécessaires
2 pip install requests beautifulsoup4 lxml pandas numpy
```

Rappel Complet de BeautifulSoup

Structure de base d'un script

```
1 import requests
2 from bs4 import BeautifulSoup
3 import pandas as pd
4
5 # 1. Requête HTTP
6 url = "https://example.com"
7 response = requests.get(url)
8
9 # 2. Vérification du statut
10 if response.status_code == 200:
11     # 3. Parsing du HTML
12     soup = BeautifulSoup(response.content, 'html.parser')
13
14     # 4. Extraction des données
15     # Code d'extraction ici...
16 else:
17     print(f"Erreur {response.status_code}")
```

Méthodes de sélection principales

Par nom de balise

```

1 # Premier lment d'un type
2 soup.find('div')

3

4 # Tous les lments d'un type
5 soup.find_all('div')

6

7 # Avec limite
8 soup.find_all('div', limit=5)

```

Par attributs CSS

```

1 # Par classe
2 soup.find_all('div', class_='product')
soup.select('div.product')

4

5 # Par ID
6 soup.find('div', id='header')
7 soup.select('#header')

8

9 # Combinaison
10 soup.select('div.product.special')

```

Par attributs personnalisés

```

1 # Recherche par attribut data
2 soup.find_all('div', attrs={'data-category': 'books'})

3

4 # Recherche avec expression reguli re
5 import re
6 soup.find_all('div', id=re.compile('^product_'))

```

Navigation dans l'arbre DOM

Relations parent-enfant

```

1 # Parent direct
2 element.parent

3

4 # Parents r cursifs
5 for parent in element.parents:
6     print(parent.name)

7

8 # Enfants directs
9 for child in element.children:
10    print(child)
11

```

```

12 # Tous les descendants
13 for descendant in element.descendants:
14     print(descendant)

```

Relations fraternelles

```

1 # Frere suivant
2 next_sibling = element.next_sibling
3
4 # Frere precedent
5 prev_sibling = element.previous_sibling
6
7 # Tous les freres suivants
8 for sibling in element.next_siblings:
9     print(sibling)

```

Extraction des données

Textes et contenus

```

1 # Texte direct
2 text = element.get_text()
3
4 # Texte avec separateur
5 text = element.get_text(separator=' ', strip=True)
6
7 # Contenu HTML interne
8 html_content = str(element)
9
10 # Attributs spécifiques
11 link = element['href']
12 src = element['src']
13 class_list = element.get('class', [])

```

Gestion des données manquantes

```

1 # Méthode sécurisée
2 def safe_extract(element, selector, default='N/A'):
3     found = element.select_one(selector)
4     return found.get_text(strip=True) if found else default
5
6 # Application
7 title = safe_extract(product, 'h3.title')
8 price = safe_extract(product, '.price')

```

Sélecteurs CSS avancés

```

1 # Combinaisons complexes
2 soup.select('div.products > ul > li:first-child')
3 soup.select('a[href^="https://"]') # Commence par
4 soup.select('img[src$=".jpg"]') # Termine par
5 soup.select('div:not(.excluded)') # Exclusion

```

Gestion des erreurs et bonnes pratiques

Gestion robuste des requêtes

```

1 import requests
2 from requests.adapters import HTTPAdapter
3 from requests.packages.urllib3.util.retry import Retry
4
5 def create_session_with_retries():
6     session = requests.Session()
7
8     # Stratégie de retry
9     retry_strategy = Retry(
10         total=3,
11         backoff_factor=1,
12         status_forcelist=[429, 500, 502, 503, 504],
13     )
14
15     adapter = HTTPAdapter(max_retries=retry_strategy)
16     session.mount("http://", adapter)
17     session.mount("https://", adapter)
18
19     return session

```

Respect des serveurs

```

1 import time
2 import random
3
4 def respectful_delay(min_delay=1, max_delay=3):
5     """Pause alatoire entre les requêtes"""
6     time.sleep(random.uniform(min_delay, max_delay))
7
8 def check_robots_txt(base_url):
9     """Vérifier robots.txt"""
10    try:
11        robots_url = f"{base_url}/robots.txt"
12        response = requests.get(robots_url)
13        print("Robots.txt:", response.text[:500])

```

```

14     except:
15         print("Robots.txt non accessible")

```

Exercices Challenges

Exercice 1 : Books to Scrape - Extraction complète

Site : <https://books.toscrape.com/>

Objectif : Créer un scraper qui extrait pour chaque livre :

- Titre complet et URL détaillée
- Prix (converti en float)
- Note sur 5 (convertie en numérique)
- Catégorie principale et secondaire
- Description complète (page détail)
- Stock disponible (nombre d'exemplaires)
- URL de l'image haute définition

Challenge :

- Gérer automatiquement la pagination sur toutes les pages
- Structurer les données en JSON hiérarchique
- Gérer les erreurs 404 sur les pages détail
- Sauvegarder avec timestamp dans le nom de fichier

Exercice 2 : Quotes to Scrape - Graphe de données

Site : <http://quotes.toscrape.com/>

Objectif :

- Extraire toutes les citations avec leurs métadonnées
- Pour chaque auteur, visiter sa page et extraire :
 - Biographie complète
 - Date et lieu de naissance
 - Date de décès (si applicable)
- Créer des relations : Citation → Auteur → Tags

Challenge :

- Implémenter un système de cache pour éviter les requêtes répétées
- Construire un graphe de relations entre entités
- Exporter en format GraphML ou GEXF
- Déetecter les auteurs les plus cités

Exercice 3 : Fake Jobs - Filtres avancés

Site : <https://realpython.github.io/fake-jobs/>

Objectif :

- Extraire uniquement les offres contenant "Python"
- Nettoyer et standardiser les dates
- Classer par type de contrat et localisation
- Extraire et valider les URLs d'application

Challenge :

- Créer un système de filtres dynamiques en ligne de commande
- Implémenter la détection de doublons
- Générer des statistiques par ville et type de contrat
- Sauvegarder en CSV avec encodage UTF-8

Exercice 4 : Analyse de marché livresque

Site : <https://books.toscrape.com/>

Objectif :

- Calculer le prix moyen par note et par catégorie
- Identifier les tendances de prix
- Déetecter les livres en rupture de stock
- Analyser la distribution des ratings

Challenge :

- Générer un rapport PDF avec matplotlib
- Créer des visualisations interactives
- Implémenter des alertes prix
- Calculer la corrélation note/prix

Exercice 5 : Navigation catégorielle avancée

Site : <https://books.toscrape.com/>

Objectif :

- Cartographier toute l'arborescence des catégories
- Pour chaque catégorie, calculer :
 - Nombre total de livres
 - Prix moyen, minimum et maximum moyenne pondérée
- Générer un classement des catégories

Challenge :

- Créer un arbre hiérarchique des catégories
- Déetecter les catégories sous-représentées
- Implémenter une recherche full-text
- Exporter en format arborescent (JSON nested)

Exercice 6 : Scraper résilient

Site : <http://books.toscrape.com/>

Objectif :

- Implémenter un système de retry intelligent
- Gérer les timeouts adaptatifs
- Logger détaillé des succès/échecs
- Sauvegarde automatique de la progression

Challenge :

- Reprise sur interruption
- Limitation automatique du débit
- Détection de blocage IP
- Rapport de performance temps réel

Exercice 7 : Pipeline de data cleaning

Site : <https://books.toscrape.com/>

Objectif :

- Nettoyer et standardiser tous les champs
- Valider la cohérence des données
- Corriger les encodages problématiques
- Générer un rapport de qualité

Challenge :

- Détection automatique d'anomalies
- Imputation des valeurs manquantes
- Validation croisée des données
- Métriques de qualité automatisées

Exercice 8 : Scraping multi-sources

Sites :

- <https://books.toscrape.com/>
- <http://quotes.toscrape.com/>
- <https://realpython.github.io/fake-jobs/>

Objectif :

- Architecture modulaire pour multiples sources
- Unification du format de sortie
- Gestion des schémas différents
- Comparaison des performances

Challenge :

- Système de plugins dynamique
- Configuration externe (YAML/JSON)
- Parallel processing des sources
- Agrégation intelligente des données

Exercice 9 : Authentification et sessions

Site : <http://quotes.toscrape.com/login>

Objectif :

- Gérer l'authentification complète
- Maintenir la session active
- Accéder au contenu protégé
- Gérer les cookies et tokens

Challenge :

- Automatisation du login/logout
- Refresh automatique de session
- Gestion des CAPTCHAs simples
- Sécurisation des credentials

Conclusion

Ce cours couvre l'ensemble des techniques avancées de web scraping avec BeautifulSoup. Les exercices proposés sont progressifs et couvrent des scénarios réels de complexité croissante.

Recommandations

- **Toujours tester localement** avant de scaler
- **Respecter les limites** techniques et légales
- **Documenter son code** pour le maintenir
- **Monitorer les performances** et les erreurs
- **Sauvegarder régulièrement** les données intermédiaires