

# TP2 : Web Scraping Avancé et Analyse de Données avec Pandas

Scraping + Pandas + Visualisation

Applications Concrètes au Génie Civil

**Niveau :** École d'Ingénieurs - Première année

**Durée :** 2h30

24 octobre 2025

## 1 Introduction et rappels

### 1.1 Objectifs du TP2

Ce deuxième TP se concentre sur des techniques avancées de web scraping et l'analyse des données collectées avec Pandas. Vous allez :

- Scraper un catalogue en ligne de matériaux BTP (MarketBTP)
- Gérer la pagination pour collecter de grandes quantités de données
- **Utiliser Pandas pour analyser les données collectées**
- **Réaliser des statistiques descriptives sur des données de génie civil**
- **Visualiser les données avec Matplotlib**
- Comparer les prix et identifier les meilleures opportunités

### 1.2 Présentation du site MarketBTP

Pour ce TP, vous allez scraper un site web en ligne : **MarketBTP**, un catalogue professionnel de matériaux de construction.

**URL du site :** <http://www.malomatique.free.fr/MarketBTP/index.html>

**Caractéristiques du site :**

- 60 produits BTP répartis sur 3 pages
- 5 catégories : Béton, Acier, Maçonnerie, Isolation, Charpente
- Pour chaque produit : nom, fournisseur, prix, caractéristiques techniques, délai, région, note, disponibilité
- Données réalistes pour analyses pertinentes
- Pagination fonctionnelle : page-2.html, page-3.html

### 1.3 Vérification de l'environnement

Assurez-vous d'avoir installé les bibliothèques nécessaires :

```
1 pip install requests beautifulsoup4 lxml pandas matplotlib
```

Test rapide :

```
1 import requests
2 from bs4 import BeautifulSoup
3 import pandas as pd
4 import matplotlib.pyplot as plt
5
6 print("Environnement pret !")
7 print(f"Version Pandas : {pd.__version__}")
```

## 2 Introduction à Pandas (20 min)

### 2.1 Qu'est-ce que Pandas ?

Pandas est une bibliothèque Python puissante pour la manipulation et l'analyse de données. Elle est particulièrement utile pour :

- Organiser des données tabulaires (comme des tableaux Excel)

- Calculer des statistiques (moyennes, totaux, etc.)
- Filtrer et trier des données
- Exporter vers CSV, Excel, etc.
- Visualiser des données

## 2.2 Structure principale : le DataFrame

Un DataFrame est un tableau à deux dimensions avec des lignes et des colonnes.

```
1 import pandas as pd
2
3 # A partir d'un dictionnaire
4 data = {
5     'Materiau': ['Beton', 'Acier', 'Bois'],
6     'Prix': [95.50, 850.00, 120.00],
7     'Unite': ['m3', 'tonne', 'm3']
8 }
9
10 df = pd.DataFrame(data)
11 print(df)
12
13 # Resultat :
14 #   Materiau   Prix  Unite
15 # 0    Beton   95.50    m3
16 # 1    Acier  850.00  tonne
17 # 2    Bois   120.00    m3
```

Listing 1 – Créer un DataFrame

## 2.3 Opérations de base sur les DataFrames

### 2.3.1 Affichage et information

```
1 # Afficher les premieres lignes
2 print(df.head())
3
4 # Afficher les dernieres lignes
5 print(df.tail())
6
7 # Informations sur le DataFrame
8 print(df.info())
9
10 # Statistiques descriptives
11 print(df.describe())
12
13 # Dimensions (lignes, colonnes)
14 print(df.shape)
```

Listing 2 – Explorer un DataFrame

### 2.3.2 Sélection de données

```
1 # Selectionner une colonne
2 prix = df['Prix']
3
4 # Selectionner plusieurs colonnes
5 selection = df[['Materiau', 'Prix']]
6
7 # Selectionner des lignes par index
```

```
8 premiere_ligne = df.iloc[0]
9
10 # Filtrer des lignes selon une condition
11 materiaux_chers = df[df['Prix'] > 100]
```

Listing 3 – Selectionner des donnees

### 2.3.3 Calculs statistiques

```
1 # Moyenne
2 prix_moyen = df['Prix'].mean()
3
4 # Mediane
5 prix_median = df['Prix'].median()
6
7 # Minimum et maximum
8 prix_min = df['Prix'].min()
9 prix_max = df['Prix'].max()
10
11 # Somme
12 prix_total = df['Prix'].sum()
13
14 # Ecart-type
15 ecart_type = df['Prix'].std()
16
17 print(f"Prix moyen : {prix_moyen:.2f}")
18 print(f"Prix total : {prix_total:.2f}")
```

Listing 4 – Statistiques sur les donnees

### 2.3.4 Tri et classement

```
1 # Trier par prix croissant
2 df_trie = df.sort_values('Prix')
3
4 # Trier par prix decroissant
5 df_trie_desc = df.sort_values('Prix', ascending=False)
6
7 # Trier par plusieurs colonnes
8 df_multi = df.sort_values(['Unite', 'Prix'])
```

Listing 5 – Trier les donnees

### 2.3.5 Export de données

```
1 # Exporter en CSV
2 df.to_csv('materiaux.csv', index=False, encoding='utf-8')
3
4 # Exporter en Excel (necessite openpyxl)
5 df.to_excel('materiaux.xlsx', index=False)
6
7 # Lire un CSV
8 df_lecture = pd.read_csv('materiaux.csv')
```

Listing 6 – Exporter un DataFrame

## 2.4 Visualisation avec Matplotlib

```
1 import matplotlib.pyplot as plt
2
3 # Configuration pour affichage correct des accents
4 plt.rcParams['font.sans-serif'] = ['DejaVu Sans']
5
6 # Graphique en barres
7 df.plot(x='Materiau', y='Prix', kind='bar',
8         title='Prix des materiaux',
9         xlabel='Materiau',
10        ylabel='Prix (euros)')
11 plt.xticks(rotation=45)
12 plt.tight_layout()
13 plt.savefig('prix_materiaux.png')
14 plt.show()
15
16 # Diagramme circulaire
17 df.set_index('Materiau')['Prix'].plot(kind='pie',
18                                     autopct='%1.1f%%')
19 plt.title('Repartition des couts')
20 plt.show()
```

Listing 7 – Graphiques simples

## 3 Exercices pratiques (1h40)

### 3.1 Exercice 1 : Scraper et analyser le catalogue MarketBTP (60 min)

Site à scraper : <http://www.malomatique.free.fr/MarketBTP/index.html>

#### 3.1.1 Objectifs

Créez un script `exercice1.py` qui :

1. Scrape les 3 pages du catalogue MarketBTP
2. Extrait pour chaque produit :
  - Le type de produit (Béton, Acier, Maçonnerie, etc.)
  - Le nom du produit
  - Le fournisseur
  - Le prix (valeur numérique)
  - L'unité (m<sup>3</sup>, tonne, m<sup>2</sup>, ml, unité)
  - Les caractéristiques techniques
  - Le délai de livraison
  - La région
  - La note (nombre d'étoiles)
  - La disponibilité
3. Stocke les données dans un DataFrame Pandas
4. Effectue une analyse complète :
  - Affiche les statistiques descriptives
  - Identifie les produits les plus chers et les moins chers
  - Calcule le prix moyen par catégorie
  - Calcule le prix moyen par fournisseur

- Analyse les délais de livraison
  - Filtre les produits disponibles uniquement
5. Crée des visualisations :
    - Graphique en barres des 10 produits les plus chers
    - Histogramme de distribution des prix
    - Graphique de répartition par catégorie
    - Comparaison des prix moyens par fournisseur
  6. Exporte les données en CSV

### 3.1.2 Code de départ

```

1 import requests
2 from bs4 import BeautifulSoup
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import re
6 import time
7
8 # Configuration matplotlib pour accents
9 plt.rcParams['font.sans-serif'] = ['DejaVu Sans']
10
11 def extraire_prix(texte_prix):
12     """
13     Extrait le prix numerique d'une chaine
14     Ex: '95.50 euros/m3' -> 95.50
15     """
16     match = re.search(r'(\d+\.\d*)', texte_prix)
17     if match:
18         return float(match.group(1))
19     return 0.0
20
21 def compter_etoiles(texte_note):
22     """
23     Compte le nombre d'etoiles pleines
24     Ex: 5 etoiles pleines, 0 vides -> 5
25     """
26     # Compter les etoiles pleines (caractere Unicode U+2605)
27     return texte_note.count('\u2605')
28
29 def scraper_page(url):
30     """Scrape une page et retourne la liste des produits"""
31     response = requests.get(url)
32
33     if response.status_code != 200:
34         print(f"Erreur : {response.status_code}")
35         return []
36
37     soup = BeautifulSoup(response.text, 'html.parser')
38     produits = []
39
40     # Trouver tous les produits
41     cards = soup.find_all('div', class_='product-card')
42
43     for card in cards:
44         produit = {
45             'Type': '',
46             'Nom': '',
47             'Fournisseur': '',
48             'Prix': 0.0,

```

```

49         'Unite': '',
50         'Note': 0,
51         'Disponibilite': '',
52         'Delai': '',
53         'Region': ''
54     }
55
56     # A completer : extraire toutes les donnees
57     # Extraire le type
58     type_elem = card.find('span', class_='product-type')
59     if type_elem:
60         produit['Type'] = type_elem.text.strip()
61
62     # A completer : extraire les autres champs...
63
64     produits.append(produit)
65
66     return produits
67
68 def analyser_donnees(df):
69     """Analyse statistique des donnees"""
70     print("\n" + "="*60)
71     print("ANALYSE DES DONNEES - CATALOGUE MARKETBTP")
72     print("="*60)
73
74     # Informations generales
75     print(f"\nNombre total de produits : {len(df)}")
76     print(f"Nombre de categories : {df['Type'].nunique()}")
77     print(f"Nombre de fournisseurs : {df['Fournisseur'].nunique()}")
78
79     # Statistiques sur les prix
80     print("\n--- STATISTIQUES DES PRIX ---")
81     print(df['Prix'].describe())
82
83     # A completer : autres analyses...
84
85 def visualiser_donnees(df):
86     """Cree des graphiques"""
87
88     # 1. Top 10 des produits les plus chers
89     plt.figure(figsize=(12, 6))
90     top10 = df.nlargest(10, 'Prix')
91     plt.barh(range(len(top10)), top10['Prix'], color='steelblue')
92     plt.yticks(range(len(top10)), top10['Nom'], fontsize=9)
93     plt.xlabel('Prix (euros)', fontsize=12)
94     plt.title('Top 10 des produits les plus chers',
95             fontsize=14, fontweight='bold')
96     plt.tight_layout()
97     plt.savefig('top10_produits.png', dpi=300)
98     plt.show()
99
100     # A completer : autres graphiques...
101
102 # Programme principal
103 def main():
104     print("="*60)
105     print("COLLECTEUR ET ANALYSEUR MARKETBTP")
106     print("="*60)
107
108     base_url = 'http://www.malomatique.free.fr/MarketBTP/'
109     pages = ['index.html', 'page-2.html', 'page-3.html']
110     tous_les_produits = []
111

```

```
112 # Scraping des 3 pages
113 for i, page in enumerate(pages, 1):
114     url = base_url + page
115     print(f"\nScraping page {i}...")
116     produits = scraper_page(url)
117     tous_les_produits.extend(produits)
118     time.sleep(1) # Pause pour ne pas surcharger le serveur
119
120 # Conversion en DataFrame
121 df = pd.DataFrame(tous_les_produits)
122
123 print(f"\nTotal de produits collectes : {len(df)}")
124
125 # Nettoyage
126 df = df[df['Prix'] > 0]
127
128 # Analyse
129 analyser_donnees(df)
130
131 # Visualisation
132 visualiser_donnees(df)
133
134 # Export CSV
135 df.to_csv('marketbtp_analyse.csv', index=False, encoding='utf-8')
136 print("\nDonnees exportees dans 'marketbtp_analyse.csv'")
137
138 if __name__ == "__main__":
139     main()
```

Listing 8 – exercicel.py - Structure de base