

Scraping Data from a Real Website + Pandas

In [1]:

```
from bs4 import BeautifulSoup
import requests
```

In [2]:

```
url = 'https://en.wikipedia.org/wiki/List_of_largest_companies_in_the_United_States_'
page = requests.get(url)
soup = BeautifulSoup(page.text, 'html')
```

In [3]:

```
print(soup)
```

```
<!DOCTYPE html>
<html class="client-nojs vector-feature-language-in-header-enabled vec
tor-feature-language-in-main-page-header-disabled vector-feature-stick
y-header-disabled vector-feature-page-tools-pinned-disabled vector-fea
ture-toc-pinned-enabled vector-feature-main-menu-pinned-disabled vecto
r-feature-limited-width-clientpref-1 vector-feature-limited-width-cont
ent-enabled vector-feature-zebra-design-disabled vector-feature-custom
-font-size-clientpref-disabled" dir="ltr" lang="en">
<head>
<meta charset="utf-8"/>
<title>List of largest companies in the United States by revenue - Wik
ipedia</title>
<script>(function(){var className="client-js vector-feature-language-i
n-header-enabled vector-feature-language-in-main-page-header-disabled
vector-feature-sticky-header-disabled vector-feature-page-tools-pinned
-disabled vector-feature-toc-pinned-enabled vector-feature-main-menu-p
inned-disabled vector-feature-limited-width-clientpref-1 vector-featur
e-limited-width-content-enabled vector-feature-zebra-design-disabled v
ector-feature-custom-font-size-clientpref-disabled";var cookie=document
```

In [5]:

```
soup.find_all('table')[1]
```

Out[5]:

```
<table class="wikitable sortable">
<caption>
</caption>
<tbody><tr>
<th>Rank
</th>
<th>Name
</th>
<th>Industry
</th>
<th>Revenue <br/>(USD millions)
</th>
<th>Revenue growth
</th>
<th>Employees
</th>
<th>Headquarters
</th></tr>
```

In [6]:

```
soup.find('table', class_ = 'wikitable sortable')
```

Out[6]:

```
<table class="wikitable sortable">
<caption>
</caption>
<tbody><tr>
<th>Rank
</th>
<th>Name
</th>
<th>Industry
</th>
<th>Revenue <br/>(USD millions)
</th>
<th>Revenue growth
</th>
<th>Employees
</th>
<th>Headquarters
</th></tr>
```

In [53]:

```
table = soup.find_all('table')[1]
```

In [54]:

```
print(table)
```

```
<table class="wikitable sortable">
<caption>
</caption>
<tbody><tr>
<th>Rank
</th>
<th>Name
</th>
<th>Industry
</th>
<th>Revenue <br/>(USD millions)
</th>
<th>Revenue growth
</th>
<th>Employees
</th>
<th>Headquarters
</th></tr>
<tr>
...

```

In [24]:

```
world_titles = table.find_all('th')
```

In [25]:

```
world_titles
```

Out[25]:

```
[<th>Rank
</th>,
<th>Name
</th>,
<th>Industry
</th>,
<th>Revenue <br/>(USD millions)
</th>,
<th>Revenue growth
</th>,
<th>Employees
</th>,
<th>Headquarters
</th>]
```

In [26]:

```
world_table_titles = [title.text.strip() for title in world_titles]
```

```
print(world_table_titles)
```

```
['Rank', 'Name', 'Industry', 'Revenue (USD millions)', 'Revenue growth', 'Employees', 'Headquarters']
```

In [27]:

```
import pandas as pd
```

In [28]:

```
df = pd.DataFrame(columns = world_table_titles)
```

```
df
```

Out[28]:

Rank	Name	Industry	Revenue (USD millions)	Revenue growth	Employees	Headquarters
------	------	----------	------------------------	----------------	-----------	--------------

In [33]:

```
column_data = table.find_all('tr')
```

In [48]:

```
for row in column_data[1:]:  
    row_data = row.find_all('td')  
    individual_row_data = [data.text.strip() for data in row_data]  
  
    length = len(df)  
    df.loc[length] = individual_row_data
```

In [55]:

df

Out[55]:

	Rank	Name	Industry	Revenue (USD millions)	Revenue growth	Employees	Headquarters
0	1	Walmart	Retail	611,289	6.7%	2,100,000	Bentonville, Arkansas
1	2	Amazon	Retail and Cloud Computing	513,983	9.4%	1,540,000	Seattle, Washington
2	3	Exxon Mobil	Petroleum industry	413,680	44.8%	62,000	Spring, Texas
3	4	Apple	Electronics industry	394,328	7.8%	164,000	Cupertino, California
4	5	UnitedHealth Group	Healthcare	324,162	12.7%	400,000	Minnetonka, Minnesota
...
195	96	Best Buy	Retail	46,298	10.6%	71,100	Richfield, Minnesota
196	97	Bristol-Myers Squibb	Pharmaceutical industry	46,159	0.5%	34,300	New York City, New York
197	98	United Airlines	Airline	44,955	82.5%	92,795	Chicago, Illinois
198	99	Thermo Fisher Scientific	Laboratory instruments	44,915	14.5%	130,000	Waltham, Massachusetts
199	100	Qualcomm	Technology	44,200	31.7%	51,000	San Diego, California

200 rows × 7 columns

In [52]:

df.to_csv(r'/Users/abigailmoore/Documents/CV Projects/Companies.csv', index = False)

In []: