

Web Scraping com Python

```
import time
import requests
import pandas as pd
from bs4 import BeautifulSoup
from selenium import webdriver
from selenium.webdriver.firefox.options import Options
import json
```

biblioteca para
tratamento de dados

lib

para ser oter pelo firefox

* Pegar conteúdo HTML a partir da URL

```
url = " " "
```

```
option = Options()
option.headless = True
driver = webdriver.Firefox(options=option)
```

```
driver.get(url)
```

para lembrar isso

```
driver.quit()
```

```
$ python webscraping.py
```

* Aqui, ele entra e sai o navegador ...

* Agora vamos manipular o HTML ...
> 10x

```
time.sleep(10)
```

* fazer o quit p/ o fim...

* Vamos simular um clique agora dentro do elemento

```
driver.find_element_by_xpath("")
```

```
// div [@ class = 'nba-stat-table'] // table // thead // tr //  
th [@ data-field = 'PTS']"). click
```

\$ python web scraping.py

// demora por causa do sleep

```
element = driver.find_element_by_xpath ("//div [@  
class = 'nba-stat-table'] // table")
```

```
html-content = element.get_attribute ("outerHTML")
```

```
print (html-content)
```

\$ python ...

OK, trouxe todo o conteúdo HTML

* 2. Processar o conteúdo HTML (Beautiful Soup)

↓
lib de análise de HTML
e transformá-lo em
dados estruturados

```
soup = BeautifulSoup(html_content,  
                        'html.parser')
```

```
table = soup.find(name='table')
```

* 3. Estruturar conteúdo em um Data Frame - Pandas → data frame

```
df_full = pd.read_html(str(table))[0].head(10)
```

limpeza
em colunas

limitar em
10

```
df = df_full[['Unnamed: 0', 'PLAYER', 'TEAM', 'PTS']]
```

```
df.columns = ['pos', 'player', 'team', 'total']
```

```
print(df)
```

```
driver.quit()
```

```
$ python ...
```


* 4. Transformar os dados em dicionário de dados próprios

top10ranking = {}

data frame

top10ranking['points'] = df.to_dict('records')

driver.quit()

* 5. Converter o valor em um arquivo JSON

js = json.dumps(top10ranking)

~~fp~~ fp = open('ranking.json', 'w')

fp.write(js)

fp.close()

* python

criar um arquivo json → temos um ranking no formato de dicionário.