

Section 13

Web Scraping

Mohammed Asir Shahid

2021-08-05

Contents

1	The webbrowser Module	1
2	Downloading from the Web with the Requests Module	2
2.1	Write-binary mode: open(filename, "wb")	3
3	Parsing HTML with the Beautiful Soup Module	3

1 The webbrowser Module

```
import webbrowser

webbrowser.open("https://asir.dev")
```

Let's create a program that can open a given address on maps.

```
import webbrowser, sys, pyperclip

sys.argv # ["mapit.py", "870", "Valencia", "St."]

# Check if command line arguments were passed

if len(sys.argv) > 1:
    # ["mapit.py", "870", "Valencia", "St."] -> 870 Valencia St.
```

```
        address=" ".join(sys.argv[1:])
    else:
        address=pyperclip.paste()

webbrowser.open("https://www.google.com/maps/place/%s" % (address))
```

2 Downloading from the Web with the Requests Module

The requests module lets you easily download files from the web without having to worry about complicated network issues. The requests module is a third party module which we'll need to install on our own.

```
pip install requests
```

We can pass a URL to the `requests.get()` function in order to get the file. We can check the status code to see if it downloaded properly, if so then we'll get the status code 200. We can print out the file using `.text`. We can also see if there is an issue by calling the `raise_for_status()` method which will raise an error if we ran into any problems.

```
import requests

res=requests.get("http://automatetheboringstuff.com/files/rj.txt")

print(res.status_code)

print(len(res.text))
print(res.text[:500])

print(res.raise_for_status())

badRes=requests.get("http://automatetheboringstuff.com/files/rjuliet.txt")

print(badRes.raise_for_status())
```

2.1 Write-binary mode: open(filename, “wb”)

We can save a web page to a file using the open function. However, we must do somethings differently.

```
import requests

res=requests.get("http://automatetheboringstuff.com/files/rj.txt")
playFile=open("RomeoAndJuliet.txt","wb")

for chunk in res.iter_content(100000):
    playFile.write(chunk)

playFile.close()
```

Request module functions can be useful, but they are somewhat limited. You can only use it when you have the exact URL that you need to download. Selenium lets your Python scripts control the web browser directly.

3 Parsing HTML with the Beautiful Soup Module

Here we will learn how to write programs that pull information off of web pages. This is known as web scraping. We have a third party module called beautifulsoup which makes parsing through websites HTML much easier.

```
pip install beautifulsoup4
```

```
import bs4
```

Let's try to parse through an Amazon page and scrape the price information from that page.

```
import bs4,requests
```

```
url='https://www.amazon.in/Automate-Boring-Stuff-Python-2nd/dp/1593279922/ref=dp_ob_ti
```

```
headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 6.3; Win64; x64) AppleWebKit/537.36'}  
response = requests.get(url, headers=headers)  
print(response.raise_for_status())
```

None