**Topic: Web Scraping 27: Setup & Test Applications (WebScraper & Wikipedia)**

*Speaker: | Notebook: Django: Automating Common Tasks*

---



Web scraping is an automatic method to obtain large amounts of data from websites (see full documentation here)

1. To start, we need the libraries BEAUTIFULSOUP (which we previously download) and REQUEST Django library. Install these packages if these are NOT present in your REQUIREMENTS.TXT

$ pip install requests

login.html    views.py emails    views.py image_compression    h

≡ requirements.txt

```
1   amqp==5.2.0
2   asgiref==3.8.1
3   async-timeout==4.0.3
4   beautifulsoup4==4.12.3
5   billiard==4.2.0
6   celery==5.4.0
7   certifi==2024.7.4
8   charset-normalizer==3.3.2
9   click==8.1.7
10  click-didyoumean==0.3.1
11  click-plugins==1.1.1
12  click-repl==0.3.0
13  colorama==0.4.6
14  crispy-bootstrap5==2024.2
15  Django==4.2.14
16  django-anymail==11.1
17  django-ckeditor==6.7.1
18  django-crispy-forms==2.3
19  django-js-asset==2.2.0
20  idna==3.7
21  kombu==5.4.0
22  pillow==10.4.0
23  prompt_toolkit==3.0.47
24  python-dateutil==2.9.0.post0
25  python-decouple==3.8
26  redis==5.0.8
27  requests==2.32.3
28  six==1.16.0
29  soupsieve==2.6
30  sqlparse==0.5.1
31  typing_extensions==4.12.2
32  tzdata==2024.1
33  urllib3==2.2.2
34  vine==5.1.0
35  wcwidth==0.2.13
36
```

AUTOMATINGCOMMONTASKS
- emails
  - forms.py
  - models.py
  - tasks.py
  - tests.py
  - urls.py
  - views.py
  - env
- image_compression
  - __pycache__
  - migrations
  - __init__.py
  - admin.py
  - apps.py
  - forms.py
  - models.py
  - tests.py
  - urls.py
  - views.py
  - media
  - Resources
  - static
- templates
  - dataentry
  - emails
    - send-email.html
    - track_dashboard.html
    - track_stats.html
  - image_compression
    - compress.html
  - alerts.html
  - base.html
  - home.html
  - login.html
  - register.html
  - uploads
- .env
- .gitignore
- db.sqlite3
- manage.py
- requirements.txt

OUTLINE

No symbols found in document 'requirements.txt'

PROBLEMS   OUTPUT   DEBUG CONSOLE   **TERMINAL**   PORTS   AZURE

2. Use the website, WEBSCRAPERS.IO for the TEST SITES.



Use the TABLE PLAYGROUND for testing.

## Test Sites

Semantically correct tables
Tables without the thead tag
Tables with multiple header rows

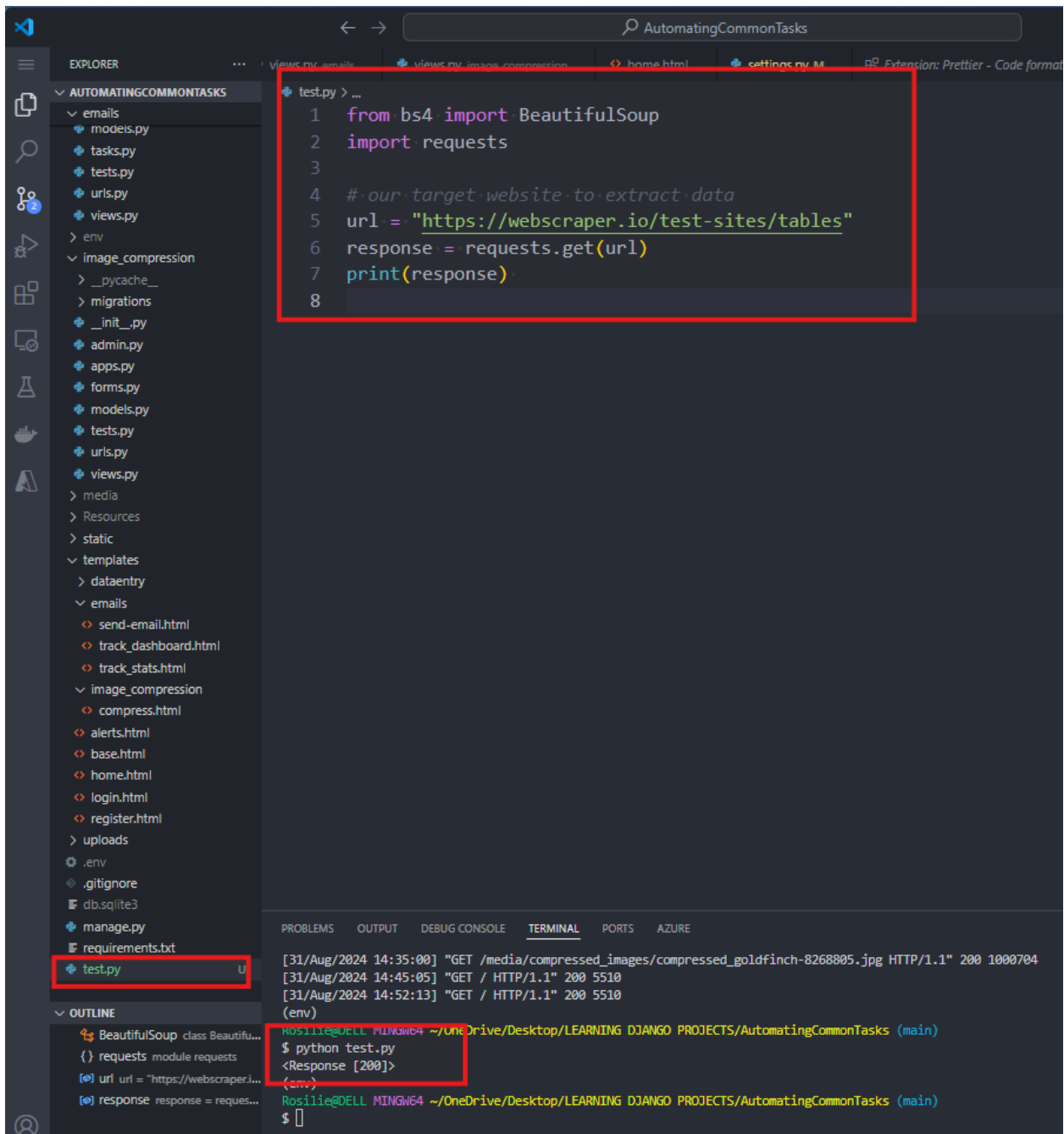## Table playground

You can train using Table selector here.

### Semantically correct table with thead and tbody

Table selector automatically detects header and data rows.

| # | First Name | Last Name | Username |
|---|-----------|-----------|----------|
| 1 | Mark | Otto | @mdo |
| 2 | Jacob | Thornton | @fat |
| 3 | Larry | the Bird | @twitter |

| # | First Name | Last Name | Username |
|---|-----------|-----------|----------|
| 4 | Harry | Potter | @hp |
| 5 | John | Snow | @dunno |
| 6 | Tim | Bean | @timbean |

3. In the root directory, create a new file, TEST.PY and update:

4. Run this file using and it returns 200 - MEANING OK.

$ python test.py



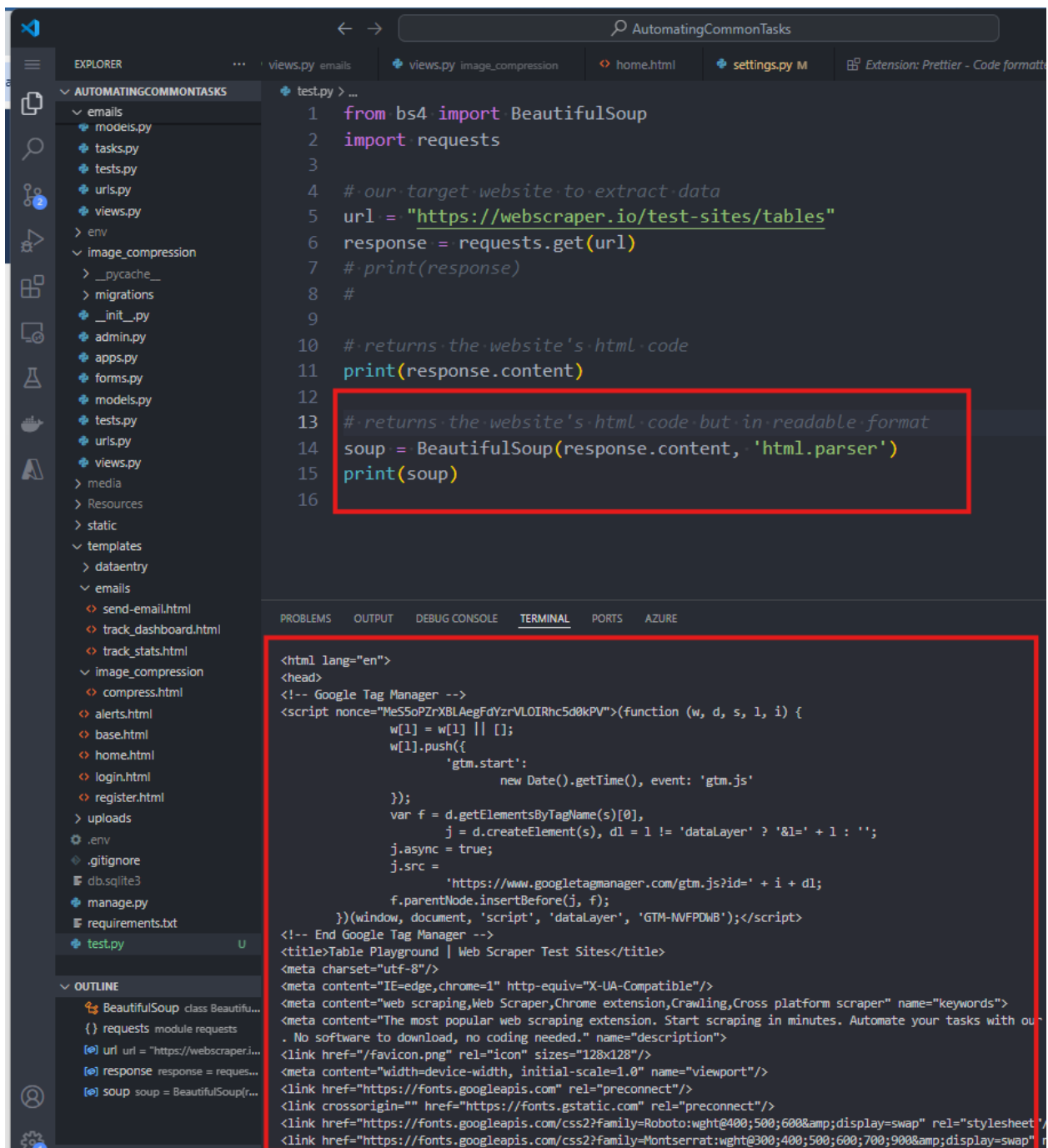5. Use this as your [guide](#) for HTTP RESPONSES:

# HTTP response status codes

HTTP response status codes indicate whether a specific HTTP request has been successfully completed. Responses are grouped in five classes:

1. Informational responses ( 100 – 199 )

2. Successful responses ( 200 – 299 )

3. Redirection messages ( 300 – 399 )

4. Client error responses ( 400 – 499 )
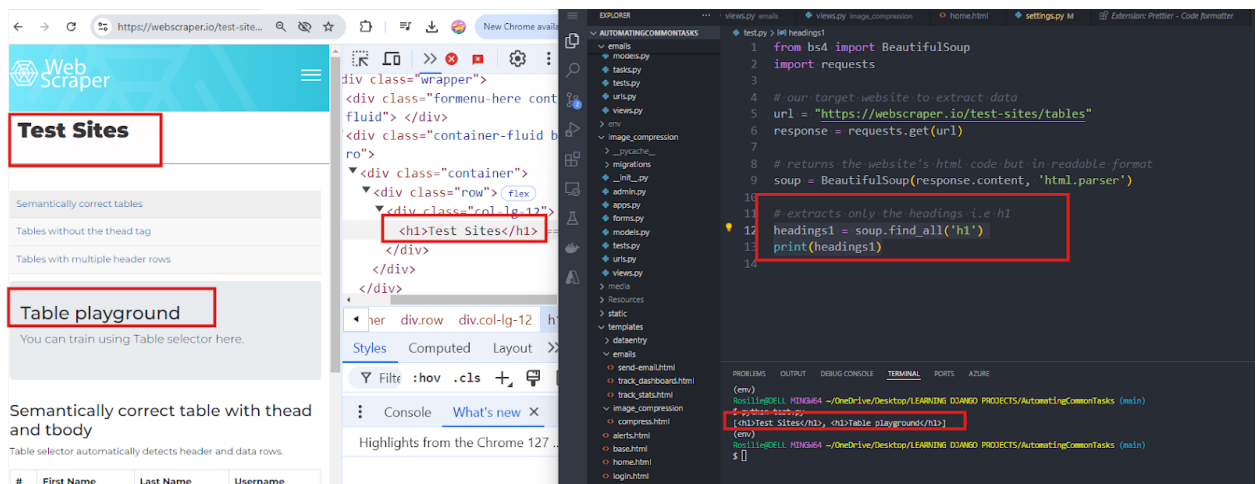
5. Server error responses ( 500 – 599 )

The status codes listed below are defined by RFC 9110 ⬀.

6. To get the website's HTML code, we type:



7. If we want to read the website's code in a more readable format, we can use BeautifulSoup.

8. To extract only certain portions of the website like H1 headings only



9. To extract the images:

10. To extract the information from the website's tables:



11. Extract only the LASTNAME OF THE SECOND TABLE:

```python
 6    response = requests.get(url)
 7
 8    # returns the website's html code but in readable format
 9    soup = BeautifulSoup(response.content, 'html.parser')
10
11    # extracts only the headings i.e h1
12    headings1 = soup.find_all('h1')
13    headings2 = soup.find_all('h2')
14    images = soup.find_all('img'[0])  # gets the first image
15
16    # extracts the tables
17    table = soup.find_all('table')[1]  # gets the 2nd table
18    # returns everything starting at index 1 (excludes 'th')
19    rows = table.find_all('tr')[1:]
20
21    last_names = []
22    for row in rows:
23        print(row.find_all('td')[2].get_text())  # gets the value only w/o tags
24
```

Terminal:

```
<td>@timbean</td>
</tr>]
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
<td>Potter</td>
<td>Snow</td>
<td>Bean</td>
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
<td>Potter</td>
Traceback (most recent call last):
  File "C:\Users\Rosilie\OneDrive\Desktop\LEARNING DJANGO PROJECTS\AutomatingCommonTasks\test.py", line 23, in <module>
    print(row.find_all('td')[2]).get_text()  # gets the value only w/o tags
AttributeError: 'NoneType' object has no attribute 'get_text'
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
Potter
Snow
Bean
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ 
```

```python
  8    # returns the website's html code but in readable format
  9    soup = BeautifulSoup(response.content, 'html.parser')
 10
 11    # extracts only the headings i.e h1
 12    headings1 = soup.find_all('h1')
 13    headings2 = soup.find_all('h2')
 14    images = soup.find_all('img'[0])  # gets the first image
 15
 16    # extracts the tables
 17    table = soup.find_all('table')[1]  # gets the 2nd table
 18    # returns everything starting at index 1 (excludes 'th')
 19    rows = table.find_all('tr')[1:]
 20
 21    last_names = []
 22    for row in rows:
 23        # adds the value only w/o tags
 24        last_names.append(row.find_all('td')[2].get_text())
 25    print(last_names)
 26
```

Terminal:

```
</tr>]
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
<td>Potter</td>
<td>Snow</td>
<td>Bean</td>
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
<td>Potter</td>
Traceback (most recent call last):
  File "C:\Users\Rosilie\OneDrive\Desktop\LEARNING DJANGO PROJECTS\AutomatingCommonTasks\test.py", line 23, in <module>
    print(row.find_all('td')[2]).get_text()  # gets the value only w/o tags
AttributeError: 'NoneType' object has no attribute 'get_text'
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
Potter
Snow
Bean
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
['Potter', 'Snow', 'Bean']
```

WEBSITE # 2: WIKIPEDIA  ON PYTHON

The goal is to categorize the data types according to MUTABLE OR IMMUTABLE

language, but may be used by external tools such as mypy to catch errors.[113][114] Mypy also supports a Python compiler called mypyc, which leverages type annotations for optimization.[115]

**Summary of Python 3's built-in types**

| Type | Mutability | Description | Syntax examples |
|---|---|---|---|
| `bool` | immutable | Boolean value | `True`<br>`False` |
| `bytearray` | mutable | Sequence of bytes | `bytearray(b'Some ASCII')`<br>`bytearray(b"Some ASCII")`<br>`bytearray([119, 105, 107, 105])` |
| `bytes` | immutable | Sequence of bytes | `b'Some ASCII'`<br>`b"Some ASCII"`<br>`bytes([119, 105, 107, 105])` |
| `complex` | immutable | Complex number with real and imaginary parts | `3+2.7j`<br>`3 + 2.7j` |
| `dict` | mutable | Associative array (or dictionary) of key and value pairs; can contain mixed types (keys and values), keys must be a hashable type | `{'key1': 1.0, 3: False}`<br>`{}` |
| `types.EllipsisType` | immutable | An ellipsis placeholder to be used as an index in NumPy arrays | `...`<br>`Ellipsis` |
| `float` | immutable | Double-precision floating-point number. The precision is machine-dependent but in practice is generally implemented as a 64-bit IEEE 754 number with 53 bits of precision.[116] | `1.33333` |
| `frozenset` | immutable | Unordered set, contains no duplicates; can contain mixed types, if hashable | `frozenset([4.0, 'string', True])` |

12. We can also use the CLASS name to find an element on the website. So, the table uses a classname 'WIKIPEDIA'

We can update our TEST.PY AS:



```python
from bs4 import BeautifulSoup
import requests

# our target website to extract data
# url = "https://webscraper.io/test-sites/tables"
url = "https://en.wikipedia.org/wiki/Python_(programming_language)"
response = requests.get(url)

# returns the website's html code but in readable format
soup = BeautifulSoup(response.content, 'html.parser')
# extracts the table using the class
datatype_table = soup.find(class_='wikitable')
print(datatype_table)

# THIS BLOCK IS FOR WEBSCRAPER WEBSITE:  https://webscraper.io/test-sit
# # extracts only the headings i.e h1
# headings1 = soup.find_all('h1')
# headings2 = soup.find_all('h2')
# images = soup.find_all('img'[0])  # gets the first image

# # extracts the tables
# table = soup.find_all('table')[1]  # gets the 2nd table
```

PROBLEMS   OUTPUT   DEBUG CONSOLE   TERMINAL   PORTS   AZURE

```
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
<table class="wikitable">
<caption>Summary of Python 3's built-in types
</caption>
<tbody><tr>
<th>Type
</th>
<th><a href="/wiki/Immutable_object" title="Immutable object">Mutability</a>
</th>
<th>Description
</th>
<th>Syntax examples
</th></tr>
<tr>
<td><code>bool</code>
</td>
<td>immutable
</td>
```

13. Each data type has a tag 'TD' inside the 'TBODY.' We need to use this info then in our TEST.PY. Each table column can be accessed using an INDEX POSITION where INDEX 0 means our first column, INDEX 1 means our second column.

14. Update TEST.PY as:



15. To remove the NEWLINE (\n), we can use the STRIP FUNCTION:

```python
11   # extracts the table using the class
12   datatype_table = soup.find(class_='wikitable')
13   body = datatype_table.find('tbody')
14   rows = body.find_all('tr')[1:]  # extracts all tr's starting at position 1
15
16   mutable_types = []
17   immutable_types = []
18
19
20   for row in rows:
21       data = row.find_all('td')
22       if data[1].get_text() == 'mutable\n':  # gets the table column # 2
23           mutable_types.append(data[0].get_text().strip())  # adds the data type
24       else:
25           immutable_types.append(data[0].get_text().strip())
26
27   print('Mutatable data types:', mutable_types)
28   print('=================')
29   print('Immutatable data types:', immutable_types)
30
31
32   # THIS BLOCK IS FOR WEBSCRAPER WEBSITE:  https://webscraper.io/test-sites/tables
33   # # extracts only the headings i.e. h1
```

PROBLEMS   OUTPUT   DEBUG CONSOLE   **TERMINAL**   PORTS   AZURE

```
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
Mutatable data types: ['bytearray\n', 'dict\n', 'list\n', 'set\n']
=================
Immutatable data types: ['bool\n', 'bytes\n', 'complex\n', 'types.EllipsisType\n', 'float\n', 'frozenset\n', 'int\n', 'types.
NoneType\n', 'types.NotImplementedType\n', 'range\n', 'str\n', 'tuple\n']
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$ python test.py
Mutatable data types: ['bytearray', 'dict', 'list', 'set']
=================
Immutatable data types: ['bool', 'bytes', 'complex', 'types.EllipsisType', 'float', 'frozenset', 'int', 'types.NoneType', 'ty
pes.NotImplementedType', 'range', 'str', 'tuple']
(env)
Rosilie@DELL MINGW64 ~/OneDrive/Desktop/LEARNING DJANGO PROJECTS/AutomatingCommonTasks (main)
$
```

16. These information are needed for STOCK MARKET ANALYSIS.