# An Introduction to Web Scraping with Python and DataCamp

Olga Scrivner, Research Scientist, CNS, CEWIT

WIM, February 23, 2018

Materials: DataCamp.com

◎ Review: Importing files

◎ Accessing Web

◎ Review: Processing text

◎ Practice, practice, practice!

## Credits

◎ Hugo Bowne-Anderson - *Importing Data in Python (Part 1 and Part 2)*

◎ Jeri Wieringa - *Intro to Beautiful Soup*

# Importing Files

◎ Text files are structured as a sequence of lines

◎ Each line includes a sequence of characters

◎ Each line is terminated with a special character **End of Line**

*Exercice 1.* Link special characters with their definition

| Special Characters | Definition |
| --- | --- |
| \ | A new line or line breaker |
| \t | A carriage return |
| \n | A separator using a tabulation |
| \r | A way to use quotes in the string |
| # | An escape Character |
| \r\n | Comment - python will ignore everything after |
| """ | A carriage return followed by the end-of-line (Window system) |
| '' or "" | A list, a sequence of strings |
| [] | A string, a sequence of characters |

| Special Characters | Definition |
| --- | --- |
| \ | An escape Character |
| \t | A separator using a tabulation |
| \n | A new line or line breaker |
| \r | A carriage return |
| # | Comment - python will ignore everything after |
| \r\n | A carriage return followed by the end-of-line (Window system) |
| """ | A way to use quotes in the string |
| '' or "" | A string, a sequence of characters |
| [] | A list, a sequence of strings |

# Modes

◎ **Reading Mode**

  ○ **'r'**

◎ **Writing Mode**

  ○ **'w'**

◎ **Reading Mode**

  ○ **'r'**

◎ **Writing Mode**

  ○ **'w'**

Quiz question: Why do we use quotes with **'r'** and **'w'**?

◎ **Reading Mode**

    ○ **'r'**

◎ **Writing Mode**

    ○ **'w'**

Quiz question: Why do we use quotes with **'r'** and **'w'**?
Answer: **'r'** and **'w'** are one-character strings

◎ Open File - **open(name, mode)**

> *file_object = open("filename", "mode")*

- **name = 'filename'**
- **mode = 'r' or mode = 'w'**

# Open New File

```
file = open("testfile.txt","w")

file.write("Hello World")
file.write("This is our new text file")
file.write("and this is another line.")
file.write("Why? Because we can.")

file.close()
```

# Open New File

```
file = open("testfile.txt","w")

file.write("Hello World")
file.write("This is our new text file")
file.write("and this is another line.")
file.write("Why? Because we can.")

file.close()
```

Hello World
This is our new text file
and this is another line.
Why? Because we can.

# Read File

◎ Read the Entire File - **filename.read()**

> *file = open("testfile.text", "r")*
> *print (file.read())*

◎ Read ONE Line - **filename.readline()**
> *print (file.readline())*    - Return the FIRST line
> *print (file.readline(3))*    - Return the THIRD line

◎ Read lines - **filename.readlines()**

*['Hello World', 'This is our new text file', 'and this is another line.', 'Why? Because we can.']*

## Read File

◎ Read the Entire File - **filename.read()**

*file = open("testfile.text", "r")*
*print (file.read())*

◎ Read ONE Line - **filename.readline()**
*print (file.readline())*    - Return the FIRST line
*print (file.readline(3))*    - Return the THIRD line

◎ Read lines - **filename.readlines()**

*['Hello World', 'This is our new text file', 'and this is another line.', 'Why? Because we can.']*

What type of object and what is the length of this object?

# Python Libraries

# Import Modules (Libraries)

◎ **Beautiful Soup**

◎ **urllib**

◎ More in next slides ...

For installation - `https://programminghistorian.org/lessons/intro-to-beautiful-soup`

To use external functions (**modules**), we need to import them:

1. Declare it at the top of the code
2. Use **import**
3. Call the module

```python
import random          Name of the module

                       Call a function from the module
for i in range(10):
    print(random.randint(1, 25))
```

To refer and import a specific function from the module

1. Declare it at the top pf the code

2. Use **from import**



**A Module Name**          **A Specific Functon**

```
from random import randint
```

3. Call the **randint** function from **random** module:

**random.randint()**

# How to Import Packages with Modules

1.  Install via a **terminal** or **console**
    - Type **command prompt** in window search
    - Type **terminal** in Mac search

# How to Import Packages with Modules

1. Install via a **terminal** or **console**
   - Type **command prompt** in window search
   - Type **terminal** in Mac search

2. Check your Python Version

```
●●●                🏠 olgascrivner — -bash — 80×24
Last login: Thu Feb 22 21:16:09 on ttys000
Olgas-MacBook-Pro-2:~ olgascrivner$ python --version
```

3. Click **return/enter**

```
Last login: Thu Feb 22 21:16:09 on ttys000
Olgas-MacBook-Pro-2:~ olgascrivner$ python --version
Python 3.6.0 :: Anaconda custom (x86_64)
```

**pip** or **pip3** - a tool for installing Python packages

If you installed Python from source, with an installer from python.org, or via Homebrew you should already have pip. If you're on Linux and installed using your OS package manager, you may have to install pip separately, see Installing pip/setuptools/wheel with Linux Package Managers.

To check if pip is installed:

To install the latest version of "SomeProject":

```
pip install 'SomeProject'
```

https://packaging.python.org/tutorials/installing-packages/

To install the latest version of "SomeProject":

```
pip install 'SomeProject'
```

# Web Scraping Workflow

## Web Concept

1. Import the necessary modules (functions)
2. Specify URL
3. Send a REQUEST
4. Catch RESPONSE
5. Return HTML as a STRING
6. Close the RESPONSE

# URLs

# URLs

1. URL - Uniform/Universal Resource Locator
2. A URL for web addresses consists of two parts:
   2.1 Protocol identifier - http: or https:
   2.2 Resource name - datacamp.com

## URLs

1. URL - Uniform/Universal Resource Locator
2. A URL for web addresses consists of two parts:
   2.1 Protocol identifier - http: or https:
   2.2 Resource name - datacamp.com
3. HTTP - HyperText Transfer Protocol
4. HTTPS - more secure form of HTTP
5. Going to a website = sending HTTP request (GET request)
6. HTML - HyperText Markup Language

## URLLIB package

Provide interface for getting data across the web. Instead of file names we use URLS

Step 1 Install the package urllib (**pip install urllib**)

Step 2 Import the function **urlretrieve** - to RETRIEVE urls during the REQUEST

Step 3 Create a variable **url** and provide the url link

**url = 'https:somepage'**

Step 4 Save the retrieved document locally

Step 5 Read the file
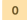
DataCamp.com - create a free account using IU email

1. Log in
2. Select Groups



3. Select **RBootcampIU** - see Jennifer if you do not see it



4. Go to Assignments and select Importing Data in Python

# Today's Practice

```
1  # Import package
2  from ____ import ____
3
4  # Import pandas
5  import pandas as pd
6              A library for data structure and analysis
7  # Assign url of file: url
8
       url = ?
9
10 # Save file locally
11
       urlretrieve(?, ?)
12
13 # Read file into a DataFrame and print its head
14 df = pd.read_csv('winequality-red.csv', sep=';')
15 print(df.head())
```
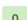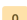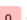
**urlretrieve** has two arguments: url (input) and file name (output)

Example: **urlretrieve(url, 'file.name')**

```
 1  # Import package
 2  from urllib.request import urlretrieve
 3
 4  # Import pandas
 5  import pandas as pd
 6
 7  # Assign url of file: url
 8  url = 'https://s3.amazonaws.com/assets.datacamp.com/production
    /course_1606/datasets/winequality-red.csv'
 9
10  # Save file locally
11  urlretrieve(url,'winequality-red.csv')
12
13  # Read file into a DataFrame and print its head
14  df = pd.read_csv('winequality-red.csv', sep=';')
15  print(df.head())
16
```

```
2  import matplotlib.pyplot as plt
3  import pandas as pd        ← A shortcut name
4
5  # Assign url of file: url
6
7
8  # Read file into a DataFrame: df
9          module.function
10         pd.read_csv(?,?)
11 # Print the head of the DataFrame
12 print(____)
```

read_csv has two arguments: url and sep (separator)

pd.head()

```
2  import matplotlib.pyplot as plt
3  import pandas as pd
4
5  # Assign url of file: url
6  url = 'https://s3.amazonaws.com/assets.datacamp.com/production
    /course_1606/datasets/winequality-red.csv'
7
8  # Read file into a DataFrame: df
9  df = pd.read_csv(url, sep=';')
10
11 # Print the head of the DataFrame
12 print(df.head())
```

read_csv has two arguments: url and sep (separator)

pd.head()

```
 1  # Import package
 2  import pandas as pd
 3
 4  # Assign url of file: url
 5
 6
 7  # Read in all sheets of Excel file: xl
 8                      xl = pd.read_excel(?, ?)
 9
10  # Print the sheetnames to the shell
11
12          Excel is a dictionary file with SHEETS as keys
13  # Print the head of the first sheet (using its name, NOT its
    index)
14
```

read_excel has two arguments: url and sheetname

To read all sheets, sheetname = None

Let's use a sheetname '1700'

```
1  # Import package
2  import pandas as pd
3
4  # Assign url of file: url
5  url = 'http://s3.amazonaws.com/assets.datacamp.com/course
   /importing_data_into_r/latitude.xls'
6
7  # Read in all sheets of Excel file: xl
8  xl = pd.read_excel(url, sheetname=None)
9
10 # Print the sheetnames to the shell
11 print(xl.keys())
12
13 # Print the head of the first sheet (using its name, NOT its
   index)
14 print(xl['1700'].head())
```

```
 1  # Import package
 2  import pandas as pd
 3
 4  # Assign url of file: url
 5
 6
 7  # Read in all sheets of Excel file: xl
 8                    xl = pd.read_excel(?, ?)
 9
10  # Print the sheetnames to the shell
11
12        Excel is a dictionary file with SHEETS as keys
13  # Print the head of the first sheet (using its name, NOT its
    index)
14
```

read_excel has two arguments: url and sheetname

To read all sheets, sheetname = None

Let's use a sheetname '1700'

Going to a website = sending HTTP request

- GET request

`urlretrieve()` performs a GET request

```
[1]: from urllib.request import urlopen, Request
[2]: url = "https://www.wikipedia.org/"
[3]: request = Request(url)
[4]: response = urlopen(request)
[5]: html = response.read()
[6]: response.close()
```

Import **request** package

```
 1  # Import packages
 2
 3
 4  # Specify the url
 5  url = "http://www.datacamp.com/teach/documentation"
 6
 7  # This packages the request: request
 8
 9
10  # Sends the request and catches the response: response
11
12
13  # Print the datatype of response
14  print(type(response))
15
16  # Be polite and close the response!
17  response.close()
```

```python
 1  # Import packages
 2  from urllib.request import urlopen, Request
 3
 4  # Specify the url
 5  url = "http://www.datacamp.com/teach/documentation"
 6
 7  # This packages the request: request
 8  request = Request(url)
 9
10  # Sends the request and catches the response: response
11  response = urlopen(request)
12
13  # Print the datatype of response
14  print(type(response))
15
16  # Be polite and close the response!
17  response.close()
```

```
1   # Import packages
2   from urllib.request import urlopen, Request
3
4   # Specify the url
5   url = "http://www.datacamp.com/teach/documentation"
6
7   # This packages the request
8   request = Request(url)
9
10  # Sends the request and catches the response: response
11
12
13  # Extract the response: html
14
15
16  # Print the html
17
```

Use **response.read()**

```
 1  # Import packages
 2  from urllib.request import urlopen, Request
 3
 4  # Specify the url
 5  url = "http://www.datacamp.com/teach/documentation"
 6
 7  # This packages the request
 8  request = Request(url)
 9
10  # Sends the request and catches the response: response
11  response = urlopen(request)
12
13  # Extract the response: html
14  html = response.read()
15
16  # Print the html
17  print(html)
```

# Return Web as a String

```
 1  # Import package
 2
 3
 4  # Specify the url: url
 5
 6
 7  # Packages the request, send the request and catch the response:
    r
 8
 9
10  # Extract the response: text
11
12
13  # Print the html
14  print(text)
15
```

Use **r.text**

```
 1  # Import package
 2  import requests
 3
 4  # Specify the url: url
 5  url = "http://www.datacamp.com/teach/documentation"
 6
 7  # Packages the request, send the request and catch the respons
    r
 8  r = requests.get(url)
 9
10  # Extract the response: text
11  text = r.text
12
13  # Print the html
14  print(text)
15
```

# Scraping Web - HTML

Mix of unstructured and structured data

Structured data:

- Has pre-defined data model, or
- Organized in a defined manner

Unstructured data: neither of these properties

Mix of unstructured and structured data

Structured data:

- Has pre-defined data model, or

- Organized in a defined manner

Unstructured data: neither of these properties

- Parse and extract structured data from HTML

```
[1]: from bs4 import BeautifulSoup

[2]: import requests

[3]: url = 'https://www.crummy.com/software/BeautifulSoup/'

[4]: r = requests.get(url)

[5]: html_doc = r.text

[6]: soup = BeautifulSoup(html_doc)
```

◎ soup.title

◎ soup.get_text()

◎ soup.find_all('a')

```
1   # Import packages
2   import requests
3   from ____ import ____
4
5   # Specify url: url
6
7   # Package the request, send the request and catch the response:
    r
8
9   # Extracts the response as html: html_doc
10
11  # Create a BeautifulSoup object from the HTML: soup
12
13  # Prettify the BeautifulSoup object: pretty_soup
14
15  # Print the response
16  print(pretty_soup)
```

# Parsing HTML with BeautifulSoup

```python
# Import packages
import requests
from bs4 import BeautifulSoup

# Specify url: url
url = 'https://www.python.org/~guido/'
# Package the request, send the request and catch the response: r
r = requests.get(url)
# Extracts the response as html: html_doc
html_doc = r.text
# Create a BeautifulSoup object from the HTML: soup
soup = BeautifulSoup(html_doc)
# Prettify the BeautifulSoup object: pretty_soup
pretty_soup = soup.prettify()
# Print the response
print(pretty_soup)
```

```
 1  # Import packages
 2  import requests
 3  from bs4 import BeautifulSoup
 4  # Specify url: url
 5  url = 'https://www.python.org/~guido/'
 6  # Package the request, send the request and catch the response: r
 7  r = requests.get(url)
 8  # Extract the response as html: html_doc
 9  html_doc = r.text
10  # Create a BeautifulSoup object from the HTML: soup
11  # Get the title of Guido's webpage: guido_title
12  # Print the title of Guido's webpage to the shell
13  # Get Guido's text: guido_text
14  # Print Guido's text to the shell
15  print(guido_text)
16
```

soup.title

soup.get_text()

```python
 2  import requests
 3  from bs4 import BeautifulSoup
 4  # Specify url: url
 5  url = 'https://www.python.org/~guido/'
 6  # Package the request, send the request and catch the response: r
 7  r = requests.get(url)
 8  # Extract the response as html: html_doc
 9  html_doc = r.text
10  # Create a BeautifulSoup object from the HTML: soup
11  soup = BeautifulSoup(html_doc)
12  # Get the title of Guido's webpage: guido_title
13  guido_title = soup.title
14  # Print the title of Guido's webpage to the shell
15  print(guido_title)
16  # Get Guido's text: guido_text
17  guido_text = soup.get_text()
18  # Print Guido's text to the shell
19  print(guido_text)
```

◎ HTML tag - <a>

◎ find_all('a')

◎ Collect all href: link.get('href')

```
 2  import requests
 3  from bs4 import BeautifulSoup
 4  # Specify url
 5  url = 'https://www.python.org/~guido/'
 6  # Package the request, send the request and catch the response: r
 7  r = requests.get(url)
 8  # Extracts the response as html: html_doc
 9  html_doc = r.text
10  # create a BeautifulSoup object from the HTML: soup
11  soup = BeautifulSoup(html_doc)
12  # Print the title of Guido's webpage
13  print(soup.title)
14  # Find all 'a' tags (which define hyperlinks): a_tags
15  
16  # Print the URLs to the shell
17  for ____ in ____:
18      ____
```

44

```python
 2  import requests
 3  from bs4 import BeautifulSoup
 4  # Specify url
 5  url = 'https://www.python.org/~guido/'
 6  # Package the request, send the request and catch the response: r
 7  r = requests.get(url)
 8  # Extracts the response as html: html_doc
 9  html_doc = r.text
10  # create a BeautifulSoup object from the HTML: soup
11  soup = BeautifulSoup(html_doc)
12  # Print the title of Guido's webpage
13  print(soup.title)
14  # Find all 'a' tags (which define hyperlinks): a_tags
15  a_tags = soup.find_all('a')
16  # Print the URLs to the shell
17  for link in a_tags:
18      print(link.get('href'))
```