# Scraping 2: libraries

**SCRAPING FOR JOURNALISTS**

How to grab data from hundreds of sources, put it in a form you can interrogate – and still hit deadlines

PAUL BRADSHAW

**ONLINE JOURNALISM BLOG**
A conversation.
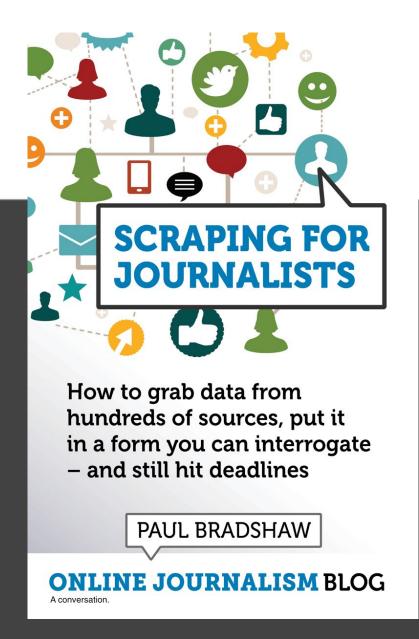
Paul Bradshaw
**Leanpub.com/scrapingforjournalists**

# What we'll cover

- What are **libraries** in Python - and why you need to know
- How to **import** libraries in a Python notebook in Google Colab

# Libraries

- A library is a **collection of recipes (functions)** and other stuff that someone has created for a particular type of problem
- Make it possible to 'stand on the shoulders of giants' & use code created by others
- E.g. the **Beautiful Soup (bs4)** library is a collection of tools for solving scraping problems
- And **requests** is a library for fetching URLs
- **Pandas** is a library for data analysis
- **Matplotlib** is a library for visualisation

```python
import requests
from bs4 import BeautifulSoup

def fetch_content(url):
    # Send an HTTP GET request to the URL
    response = requests.get(url)

    # Check if the request was successful
    if response.status_code == 200:
        # Parse the HTML content using BeautifulSoup
        soup = BeautifulSoup(response.content, 'html.parser')

        # Find the first <h1> tag and extract its text
        h1_tag = soup.find('h4')
        if h1_tag:
            data = h1_tag.text
        else:
            data = "No <h1> tag found"

        return data
    else:
        print("Failed to fetch content from the URL.")
        return None
```

# Spot the libraries

# Libraries… in Colab

- (Some) libraries need **installing** first
- (All) libraries need **importing**

# (How do you know?)

Trial and error...

```
import scraperwiki
```

```
---------------------------------------------------------------------------
ModuleNotFoundError                       Traceback (most recent call last)
<ipython-input-2-71791e80ea22> in <module>()
----> 1 import scraperwiki

ModuleNotFoundError: No module named 'scraperwiki'

---------------------------------------------------------------------------
NOTE: If your import is failing due to a missing package, you can
manually install dependencies using either !pip or !apt.

To view examples of installing some common dependencies, click the
"Open Examples" button below.
---------------------------------------------------------------------------
```

OPEN EXAMPLES    SEARCH STACK OVERFLOW

```
#install the library
!pip install scraperwiki
#import the library
import scraperwiki
```

# `import pandas as pd`

- A library can be **renamed** at the same time as it is imported (typically with shorter names for convenience)
- ...because when you use a function from a library you need to name the library

# `from bs4 import BeautifulSoup`

- Sometimes you'll find code where only part of a library is imported (just one function)
- In this case the name of the library is **bs4** but we only want to use **BeautifulSoup**
- You don't need to know any of this for the code to work!

# Using a library

- When you use a **function** from a library you name the library and the function, with a period joining them:
- **requests.**get(fullurl)
- **pandas.**DataFrame(columns=["title"])

...or if renamed when imported:
**pd.**DataFrame(columns=["title"])

# Hold on — functions?

# Functions = recipes

- A **function** is a name for a recipe. Used in Excel, e.g. SUM, AVERAGE, VLOOKUP
- A function is always followed by parentheses to 'pass' any ingredients, e.g. =SUM(A1:A10)
- requests.**get(fullurl)**
- pd.**DataFrame(columns=["title"])**

# Recap

- A library is (pre-)installed, and imported:

```
!pip install scraperwiki
import scraperwiki
import requests
```

- Functions (recipes) from that library are joined by a period and followed by parentheses:

```
html = requests.get("http://blah.com")
```

# Try it now:

- Create a notebook and import the libraries we will need:
  - **import requests**
  - **from bs4 import BeautifulSoup**
  - **import pandas as pd**

- Tip: If you get an error, ask Bard/ChatGPT what you might have done wrong