**CHULA ΣNGINEERING**
Foundation toward Innovation
**COMPUTER**

**Data♥ind**
Data Mining Group
Machine Intelligence and Knowledge Discovery Lab
Chulalongkorn University

# Web Scraping using Python Python for Data Analytics

**Peerapon Vateekul, Ph.D.**
Department of Computer Engineering,
Faculty of Engineering, Chulalongkorn University
Peerapon.v@chula.ac.th

# Outlines

- What is **Web scraping?**

- Usage

- Advantages

- Workflows

- Libraries

- BeautifulSoup

- Collecting and Parsing a Web Page

- Pulling Text From a Web Page

- Example

Reference:
1. Manoj Pandey, https://slides.com/manojp/introws#/1
2. Paul Schreiber, Web Scraping with Python Slide, NICAR 2015
3. https://www.digitalocean.com/community/tutorials/how-to-scrape-web-pages-with-beautiful-soup-and-python-3  v

# What is **Web scraping?**

- **Web scraping** is a technique for gathering data or information on web pages.

- You could revisit your favorite web site every time it updates for new information.

- Or you could write a **web scraper** to have it do it for you!

- It is a method to extract data from a website that does not have an **API** or we want to extract a LOT of data which we can not do through an API due to **rate limiting**.

- Through web scraping we can extract any data which we can see while browsing the web

# Usage

**WEB SCRAPING IN REAL LIFE**

- Extract product information.

- Extract job postings and internships.

- Extract offers and discounts from deal-of-the-day websites.

- Crawl forums and social websites.

- Extract data to make a search engine.

- Gathering weather data etc.

# + Advantages

- Web Scraping is **not rate limited.**

- **Anonymously** access the website and gather data.

- Some websites do not have an API.

- Some data is not accessible through an API

- **and many more !**

# + Workflows

- **Web Scraping follows this workflow:**
    - **Get** the website - using **HTTP** library
    - **Parse** the html document - using any parsing library
    - **Store** the results - either a database, csv, text file, etc
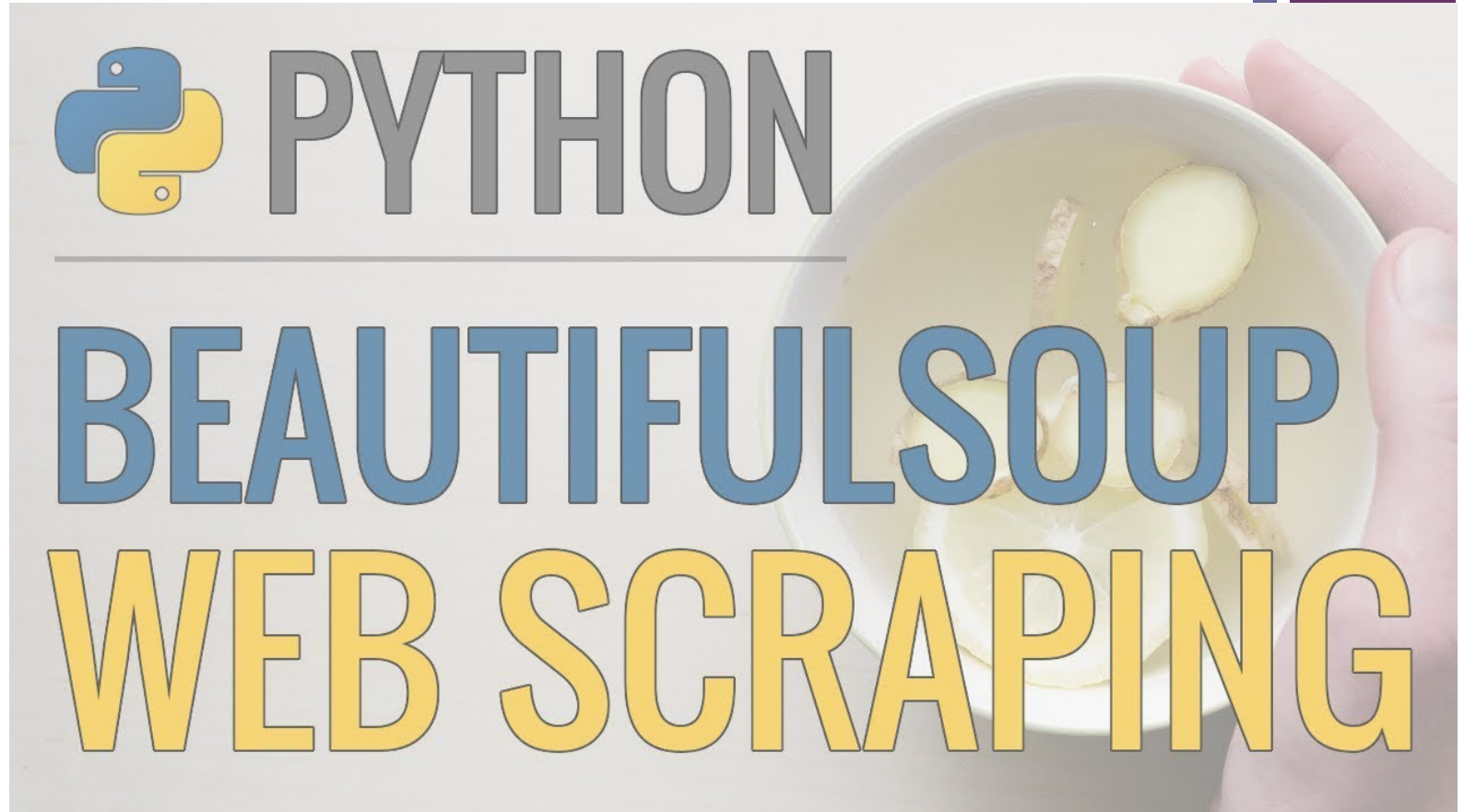
# + Libraries

- **BeautifulSoup**

- Lxml

- Scrapy

- Selenium

# Http Libraries: Get Web Content (html)

- Requests

```
r = requests.get('https://www.google.com').html
```

- urllib/urllib2

```
html = urllib2.urlopen('http://python.org/').read()
```

# Parsing Libraries: Get Title

- beautifulsoup

```
tree = BeautifulSoup(html_doc)
tree.title
```

- lxml

```
tree = lxml.html.fromstring(html_doc)
title = tree.xpath('/title/text()')
```

- re

```
title = re.findall('<title>(.*?)</title>', html_doc)
```

# BeautifulSoup

- A beautiful API

- Very easy to use

- Can handle broken markup

- Purely in Python

```python
soup = BeautifulSoup(html_doc)
last_a_tag = soup.find("a", id="link3")
all_b_tags = soup.find_all("b")
```

# + Collecting and Parsing a Web Page

■ The next step we will need to do is collect the URL of the first web page with Requests. We'll assign the URL for the first page to the variable page by using the method requests.get().

```
nga_z_artists.py

import requests
from bs4 import BeautifulSoup


# Collect first page of artists' list
page = requests.get('https://web.archive.org/web/20121007172955/https://www.nga.gov/col
```

page = requests.get('https://web.archive.org/web/20121007172955/https://www.nga.gov/collection/anZ1.htm')

# Collecting and Parsing a Web Page (cont.)

- We'll now create a BeautifulSoup object, or a parse tree.

- This object takes as its arguments the page.text document from Requests (the content of the server's response) and then parses it from Python's built-in html.parser.

```
nga_z_artists.py

import requests
from bs4 import BeautifulSoup


page = requests.get('https://web.archive.org/web/20121007172955/https://www.nga.gov/col


# Create a BeautifulSoup object
soup = BeautifulSoup(page.text, 'html.parser')
```
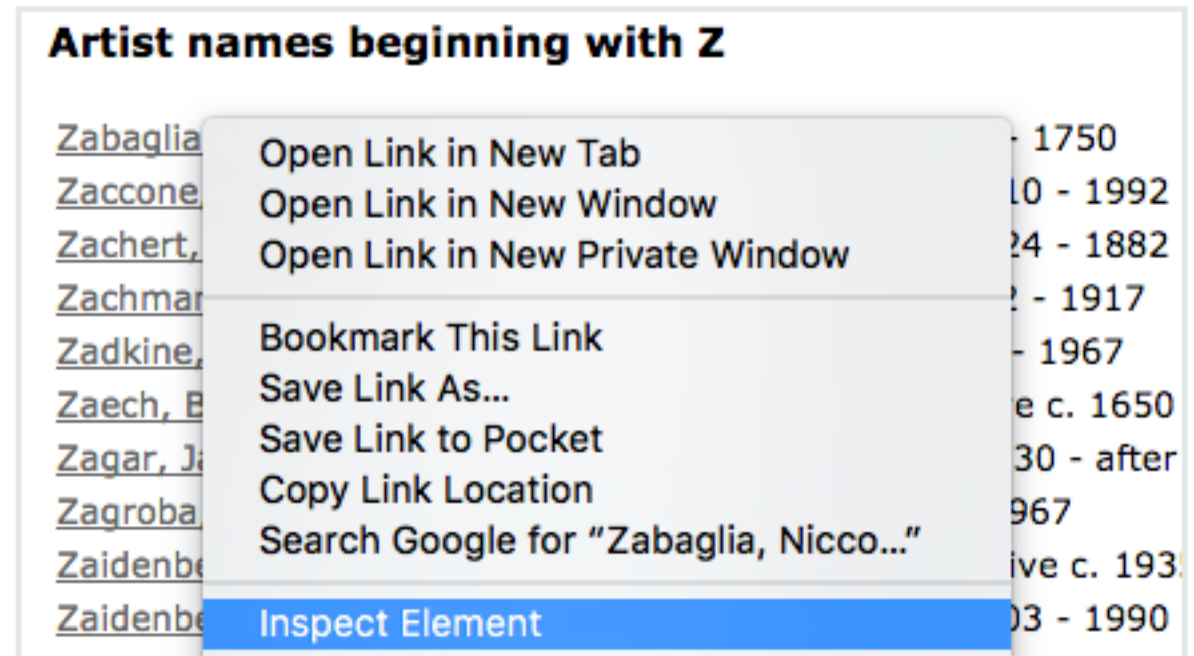
- With our page collected, parsed, and set up as a BeautifulSoup object, we can move on to collecting the data that we would like.
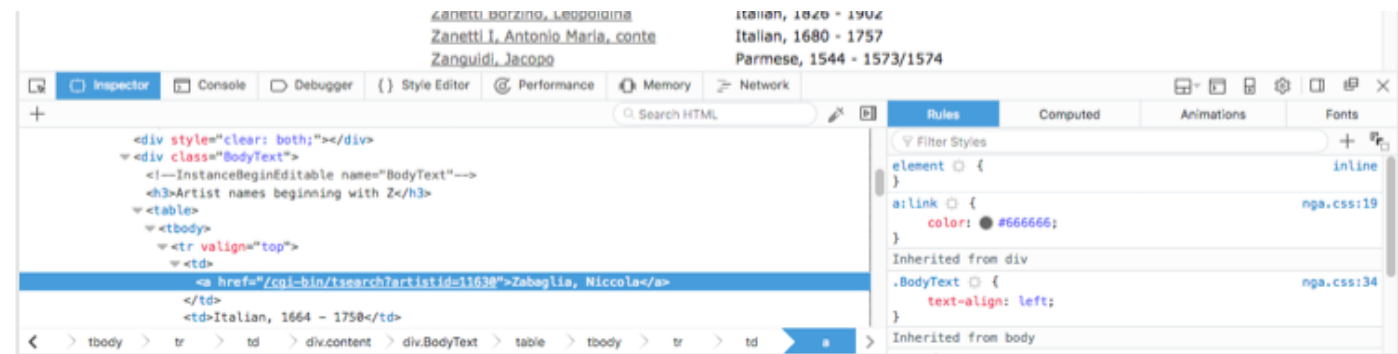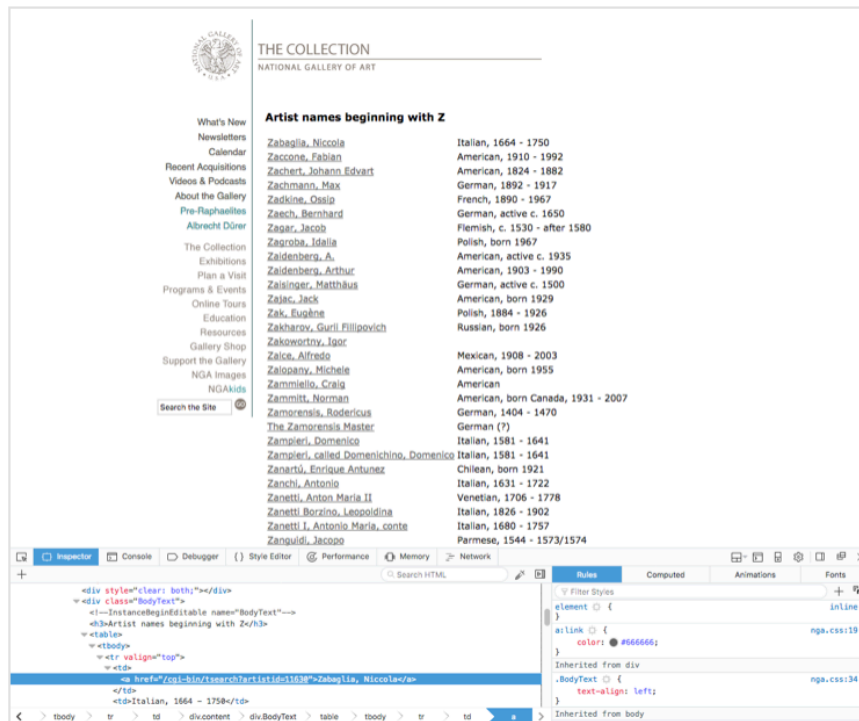
# Pulling Text From a Web Page

- For this project, we'll collect artists' names and the relevant links available on the website. You may want to collect different data, such as the artists' nationality and dates. Whatever data you would like to collect, you need to find out how it is described by the DOM of the web page.

- To do this, in your web browser, **right-click — or CTRL + click on macOS** — on the first artist's name, **Zabaglia, Niccola**. Within the context menu that pops up, you should see a menu item similar to **Inspect Element** (Firefox) or **Inspect** (Chrome).

**Artist names beginning with Z**

Zabaglia
Zaccone
Zachert,
Zachman
Zadkine,
Zaech, B
Zagar, Ja
Zagroba
Zaidenbe
Zaidenbe

Open Link in New Tab
Open Link in New Window
Open Link in New Private Window

Bookmark This Link
Save Link As...
Save Link to Pocket
Copy Link Location
Search Google for "Zabaglia, Nicco..."

Inspect Element

1750
10 - 1992
24 - 1882
2 - 1917
- 1967
e c. 1650
30 - after
967
ive c. 193
03 - 1990

# Pulling Text From a Web Page (cont.)

- Once you click on the relevant **Inspect** menu item, the tools for web developers should appear within your browser.

- We want to look for the class and tags associated with the artists' names in this list.

+

# Any Questions?