

## ✓ Part 1 Web Scraping

Web Scraping is an art where one has to study the website and work according to the dynamics of that particular website.

Most common tools used for web scraping in python are demonstrated below.

1. requests <https://requests.readthedocs.io/en/latest/>
2. beautiful soup <https://beautiful-soup-4.readthedocs.io/en/latest/>
3. Selenium <https://selenium-python.readthedocs.io/>
4. Scrapy <https://docs.scrapy.org/en/latest/>

We will be working on the first three and the fourth one can be explored in the homeworks.

We will be scraping 4 websites today:

1. ScrapeThisSite
2. GeeksforGeeks
3. CNBC
4. Hoopshype

There are different techniques to be used when scraping a dynamic website vs a static website which will be discussed in the coming sections

Some websites have their APIs open and those can be used to directly fetch the data without the need of scraping the HTML or XML pages.

```
1 # installing the libraries
2 !pip install requests
3 !pip install bs4
4 !pip install selenium
5 !apt-get update
6 !apt install chromium-chromedriver
7 !cp /usr/lib/chromium-browser/chromedriver /usr/bin
```

 [Show hidden output](#)

```
1 # importing the libraries
2 import requests
3 from bs4 import BeautifulSoup
4 from selenium import webdriver
5 from selenium.webdriver.common.keys import Keys
6 import pandas as pd
7 import json
8 from google.colab import drive
9 import sys
```

```
1 drive.mount("/content/drive/")
```

 [Show hidden output](#)

```
1 url = 'https://www.scrapethissite.com/pages/forms/'
```

```
1 page = requests.get(url)
```

```
1 soup = BeautifulSoup(page.text, 'html')
```

```
1 print(soup)
```

 [Show hidden output](#)

```
1 soup.find('div')
```

 [Show hidden output](#)

```
1 soup.find_all('p')
```

 [Show hidden output](#)

```
1 soup.find('p', class_ = 'lead').text
```

↗ Show hidden output

```
1 soup.find('p', class_ = 'lead').text.strip()
```

↗ Show hidden output

```
1 soup.find_all('th')
```

↗ Show hidden output

```
1 soup.find('th').text.strip()
```

↗ Show hidden output

## ✓ My work for Task 1 and Task 5:

I created a pandas dataframe with the table available on the website. This is for multiple tasks i.e For Task 1 and Task 5. I have discussed this with the professor and she said that is fine.

```
1 import pandas as pd
2
3 rows = []
4 for i in range(1,7):
5     url = 'https://www.scrapethissite.com/pages/forms/?page_num={}&per_page=100'.format(i)
6     print(url)
7     page = requests.get(url)
8     soup = BeautifulSoup(page.text, 'html')
9     table = soup.find('table')
10    columns = [th.text.replace("\n", "").strip() for th in table.find_all('th')]
11    for tr in table.find_all('tr')[1:]:
12        temp = []
13        for td in tr.find_all('td'):
14            x = td.text.replace("\n", "").strip()
15            temp.append(x)
16        rows = rows + [temp]
17    # print(len(rows))
18    # print(len(set(tuple(row) for row in rows)))
19
20 df_from_html = pd.DataFrame(rows, columns=columns)
21 df_from_html
```

↗ [https://www.scrapethissite.com/pages/forms/?page\\_num=1&per\\_page=100](https://www.scrapethissite.com/pages/forms/?page_num=1&per_page=100)  
[https://www.scrapethissite.com/pages/forms/?page\\_num=2&per\\_page=100](https://www.scrapethissite.com/pages/forms/?page_num=2&per_page=100)  
[https://www.scrapethissite.com/pages/forms/?page\\_num=3&per\\_page=100](https://www.scrapethissite.com/pages/forms/?page_num=3&per_page=100)  
[https://www.scrapethissite.com/pages/forms/?page\\_num=4&per\\_page=100](https://www.scrapethissite.com/pages/forms/?page_num=4&per_page=100)  
[https://www.scrapethissite.com/pages/forms/?page\\_num=5&per\\_page=100](https://www.scrapethissite.com/pages/forms/?page_num=5&per_page=100)  
[https://www.scrapethissite.com/pages/forms/?page\\_num=6&per\\_page=100](https://www.scrapethissite.com/pages/forms/?page_num=6&per_page=100)

	Team Name	Year	Wins	Losses	OT Losses	Win %	Goals For (GF)	Goals Against (GA)	+ / -
0	Boston Bruins	1990	44	24		0.55	299	264	35
1	Buffalo Sabres	1990	31	30		0.388	292	278	14
2	Calgary Flames	1990	46	26		0.575	344	263	81
3	Chicago Blackhawks	1990	49	23		0.613	284	211	73
4	Detroit Red Wings	1990	34	38		0.425	273	298	-25
...	...	...	...	...	...	...	...	...	...
577	Tampa Bay Lightning	2011	38	36	8	0.463	235	281	-46
578	Toronto Maple Leafs	2011	35	37	10	0.427	231	264	-33
579	Vancouver Canucks	2011	51	22	9	0.622	249	198	51
580	Washington Capitals	2011	42	32	8	0.512	222	230	-8
581	Winnipeg Jets	2011	37	35	10	0.451	225	246	-21

582 rows × 9 columns

```
1 # getting the first URL
2 # open the URL in parallel in other tab to check the information we are extracting
3 url = "https://www.geeksforgeeks.org/python-programming-language/"
```

```
1 # hitting the url and getting the response
2 res = requests.get(url)
3 print(res.status_code)
```

 [Show hidden output](#)

```
1 # creating a soup object from the returned html page
2 sp = BeautifulSoup(res.text, "lxml")
3 sp
```

 [Show hidden output](#)

```
1 # printing it in readable format
2 print(sp.prettify())
```

 [Show hidden output](#)

```
1 # parsing title element of the page
2 print(sp.title)
3 print(sp.title.name)
4 print(sp.title.string)
5 print(sp.title.parent.name)
```

 [Show hidden output](#)

```
1 # extracting the title of the article
2 print(sp.find("div", {"class" : "article-title"}).text)
```

 [Show hidden output](#)

```
1 print(sp.find("div"))
```

 [Show hidden output](#)

```
1 # extracting the date of the article
2 date_label = sp.find('span', class_='strong', text='Last Updated : ')
3 date_span = date_label.find_next('span')
4 last_updated_date = date_span.text
5 print(f"Last Updated Date: {last_updated_date}")
6
```

 [Show hidden output](#)

```
1 # extracting the content of the article
2 # it extracts everything together, in the next sections we can see how to iteratively extract information paragraph by paragraph
3 print(sp.find("div", {"class" : "text"}).text.strip().replace("\n", " "))
```

 [Show hidden output](#)

```
1 # extracting the links found in the bottom of the article for further reading
2 for tag in sp.find("div", {"class": "gfg-similar-reads-list"}).findAll("a", href = True): print(tag.text, "\n", tag["href"], "\n")
```

 [Show hidden output](#)

```
1 # extracting the links found in the bottom of the article for further reading
2 for link in sp.find('div', class_='gfg-similar-reads-list').find_all('a'):
3     print(link.get('href'))
```

 [Show hidden output](#)

```
1 # extracting information paragraph by paragraph
2 for tag in sp.find("div", {"class": "text"}).findAll("p"): print(tag.text)
```

 [Show hidden output](#)

Almost all the geeksforgeeks articles have the same format and hence all of them can be scraped using the same code, this repeatability is useful while doing web scraping as a block of code can help one get a lot of information in a structured way

Other information can be extracted using similar methods as used in geeksforgeeks, this can be done in homework

Again same as geeksforgeeks, all articles of scrapethissite are similar in structure and hence the same code can be used to scrap through all the articles of this website

Now working with a dynamic website that has its API open.

Open the CNBC website and search for any topic. If we search for SPORTS the URL looks like this: <https://www.cnbc.com/search/?query=SPORTS&qsearchterm=SPORTS> and if we search for POLITICS the URL looks like this: <https://www.cnbc.com/search/?query=POLITICS&qsearchterm=POLITICS>

Here we can observe a pattern and we can predict what the url would look like if we search something else, this information can be used for reusability and repeatability of code.

Now we can see that there are more than 50,000 results for the topic politics but only 10 are loaded in the beginning. Once we scroll down, next 10 are loaded and so on. To load the next results when the user scrolls down, an API is hit and link to that API is found from the network section of the inspect element: [https://api.queryly.com/cnbc/json.aspx?queryly\\_key=31a35d40a9a64ab3&query=POLITICS&endindex=10&batchsize=10&callback=&showfaceted=false&timezoneoffset=420&facetedfields=formats&facetedkey=formats%7C&facetedvalue=!Press%20Release%7C&additionalindexes=4cd6f71bf22424d,937d600b0d0d4e23,3bfb e40caee7443e,626fdcd96444f28](https://api.queryly.com/cnbc/json.aspx?queryly_key=31a35d40a9a64ab3&query=POLITICS&endindex=10&batchsize=10&callback=&showfaceted=false&timezoneoffset=420&facetedfields=formats&facetedkey=formats%7C&facetedvalue=!Press%20Release%7C&additionalindexes=4cd6f71bf22424d,937d600b0d0d4e23,3bfb e40caee7443e,626fdcd96444f28)

Here playing around with the endindex and batchsize parameter we can get various results.

Now the response of this API call would be a JSON and hence the need for parsing a web page is gone when there is an open API

```
1 # getting the third url
2 url = "https://api.queryly.com/cnbc/json.aspx?queryly_key=31a35d40a9a64ab3&query=POLITICS&endindex=10&batchsize=10&callback=&showfaceted=
```

```
1 # hitting the url and getting the response
2 res = requests.get(url)
3 print(res.status_code)
```

 Show hidden output

```
1 res.text
```

 Show hidden output

```
1 # getting the description of each news article
2 for description in json.loads(res.text)["results"]: print(description["description"], "\n")
```

 Show hidden output

Finally working with Selenium

```
1 !echo | sudo add-apt-repository ppa:saiarcot895/chromium-beta
2 !sudo apt remove chromium-browser
3 !sudo snap remove chromium
4 !sudo apt install chromium-browser -qq
5 # Chromium (an open-source version of Chrome) and Chromium WebDriver (which allows Selenium to control Chromium).
6
```

 Show hidden output

```
1 !pip3 install selenium --quiet
2 !apt-get update
3 !apt install chromium-chromedriver -qq
4 !cp /usr/lib/chromium-browser/chromedriver /usr/bin
5 #Selenium requires a browser driver (in this case, chromedriver) to communicate with the browser. You're installing it using the chromium
```

 Show hidden output

 Show hidden output Show hidden output

## Salaries of players

[illegible]

```
'$48,787,676',
'$48,728,845',
'$43,827,586',
'$43,219,440',
'$43,031,940',
'$43,031,940',
'$43,031,940',
'$42,846,615',
'$42,176,400',
'$42,176,400',
'$42,176,400',
'$42,176,400',
'$41,000,000',
'$40,500,000',
'$40,338,144',
'$36,725,670',
'$36,725,670',
'$36,725,670',
'$36,637,932',
'$36,016,200',
'$36,016,200',
'$35,859,950',
'$35,859,950',
'$35,410,310',
'$35,147,000',
'$35,147,000',
'$34,848,340',
'$34,848,340',
'$34,848,340',
'$34,005,250',
'$34,005,126',
'$33,653,846',
'$33,333,333',
'$32,500,000',
'$31,666,667',
'$30,000,000',
'$30,000,000',
'$29,793,104',
'$29,651,786',
'$29,347,826',
'$29,283,801',
'$29,268,293',
'$29,000,000',
'$28,939,680',
'$27,556,817',
'$27 172 912'
```

### Task 3

#### Website 1

```
1 # Apart from that choose any 2 websites of your choice and extract meaningful and structured information from there.
2 import sys
3 from selenium.webdriver.chrome.service import Service as ChromeService
4 sys.path.insert(0, '/usr/lib/chromium-browser/chromedriver')
5 # download the selenium chromedriver executable file and paste the link in the following code
6 # this code should open a new chrome window in your machine
7 sys.path.insert(0, '/usr/lib/chromium-browser/chromedriver')
8 chrome_options = webdriver.ChromeOptions()
9 chrome_options.add_argument('--headless')
10 chrome_options.add_argument('--no-sandbox')
11 chrome_options.add_argument('--disable-dev-shm-usage')
12 chrome_service = ChromeService(
13     executable_path='/usr/lib/chromium-browser/chromedriver',
14     log_path='/dev/null' # You can change the log path as needed
15 )
16 driver = webdriver.Chrome(service=chrome_service, options=chrome_options)
17 #The ChromeService class sets up the path to the chromedriver
```

```
1 # this code should open hoops hype website in your newly opened chrome window
2 driver.get('https://www.learnpytorch.io/')
3 # After setting up the Selenium WebDriver, the browser (in headless mode) navigates to HoopsHype Salaries page
```

```
1 # getting players name list
2 main_content = driver.find_elements("xpath", '/html/body/div[3]/main/div/div[3]')
3 #uses an XPath selector to find all the HTML elements on the page that have a class attribute "name", which corresponds to the player name
```

```
1 main_content[0].text
```

➞ 'Learn PyTorch for Deep Learning: Zero to Mastery book\nWelcome to the second best place on the internet to learn PyTorch (the first being the PyTorch documentation).\nThis is the online book version of the Learn PyTorch for Deep Learning: Zero to Mastery course.\nThis course will teach you the foundations of machine learning and deep learning with PyTorch (a machine learning framework written in Python).\nThe course is video based. However, the videos are based on the contents of this online book.\nFor full code and resources see the course GitHub.\nOtherwise, you can find more about the course below.\nDoes this course cover PyTorch 2.0?\nYes. PyTorch 2.0 is an additive release to previous versions of PyTorch.\nThis means it adds new features on top of the existing baseline features of PyTorch.\nThis

```
1 # getting all list items from the unordered list
2 list_of_topics = driver.find_elements("xpath", '/html/body/div[3]/main/div/div[1]/div/div/nav/ul/li/a/span')
3
4 # # Fetching text using 'innerText'
5 for topic in list_of_topics:
6     text = topic.get_attribute('innerText').strip()
7     if text:
8         print(text)
9
```

➞ Home  
00. PyTorch Fundamentals  
01. PyTorch Workflow Fundamentals  
02. PyTorch Neural Network Classification  
03. PyTorch Computer Vision  
04. PyTorch Custom Datasets  
05. PyTorch Going Modular  
06. PyTorch Transfer Learning  
07. PyTorch Experiment Tracking  
08. PyTorch Paper Replicating  
09. PyTorch Model Deployment  
A Quick PyTorch 2.0 Tutorial  
PyTorch Extra Resources  
PyTorch Cheatsheet  
The Three Most Common Errors in PyTorch

Similarly other information such as players' salaries can also be easily extracted and can be done as homework

## ▼ Website 2

```
1 # Apart from that choose any 2 websites of your choice and extract meaningful and structured information from there.
2 import sys
3 from selenium.webdriver.chrome.service import Service as ChromeService
4 sys.path.insert(0, '/usr/lib/chromium-browser/chromedriver')
5 # download the selenium chromedriver executable file and paste the link in the following code
6 # this code should open a new chrome window in your machine
7 sys.path.insert(0, '/usr/lib/chromium-browser/chromedriver')
8 chrome_options = webdriver.ChromeOptions()
9 chrome_options.add_argument('--headless')
10 chrome_options.add_argument('--no-sandbox')
11 chrome_options.add_argument('--disable-dev-shm-usage')
12 chrome_service = ChromeService(
13     executable_path='/usr/lib/chromium-browser/chromedriver',
14     log_path='/dev/null' # You can change the log path as needed
15 )
16 driver = webdriver.Chrome(service=chrome_service, options=chrome_options)
17 #The ChromeService class sets up the path to the chromedriver
```

```
1 # this code should open hoopshype website in your newly opened chrome window
2 driver.get('https://en.wikipedia.org/wiki/Main_Page')
3 # After setting up the Selenium WebDriver, the browser (in headless mode) navigates to HoopsHype Salaries page
```

```
1 # getting players name list
2 main_content = driver.find_elements("xpath", '//*[@id="mp-tfa"]/p')
3 #uses an XPath selector to find all the HTML elements on the page that have a class attribute "name", which corresponds to the player names
4 main_content[0].text
```

➞ 'An attempted coup took place on September 13, 1964, in South Vietnam against the ruling military junta, led by Nguyễn Khánh (pictured). In the preceding month, Khánh had tried to improve his leadership by declaring a state of emergency, provoking protests and riots. He made concessions to the protesters and removed military officials linked to former President Ngô Đình Diệm, including Lâm Văn Phát and Dương Văn Đức. They responded with a coup, broadcasting their promise to revive Diệm's policies. Khánh evaded capture and rallied allies while the U.S. continued their support for his rule. Khánh forced Phát and Đức to capitulate the next morning and various coup leaders appeared at a media conference where they denied that a coup had taken place. To maintain power, Khánh tried to court support from

## ✓ Task 4

```
1 !pip install scrapy
2 !pip install pyppeteer
3 !pip install scrapy asyncio
```

 [Show hidden output](#)

```
1 import scrapy
2 from scrapy.crawler import CrawlerProcess
3
4 class BasicSpider(scrapy.Spider):
5     name = "basic"
6     start_urls = ['https://www.learnpytorch.io/']
7
8     def parse(self, response):
9         for title in response.css('h1::text'):
10             data = {'title': title.get()}
11             print(data)
12             yield data
13
14         next_page = response.css('a.next::attr(href)').get()
15         if next_page is not None:
16             yield response.follow(next_page, self.parse)
17
18 process = CrawlerProcess(settings={
19     "FEEDS": {
20         "output.json": {"format": "json"}, # Save data to a JSON file
21     },
22 })
23
24 process.crawl(BasicSpider)
25 process.start()
```





```

'log_count/DEBUG': 3,
'log_count/INFO': 11,
'memusage/max': 211329024,
'memusage/startup': 211329024,
'response_received_count': 1,
'scheduler/dequeued': 1,
'scheduler/dequeued/memory': 1,
'scheduler/enqueued': 1,
'scheduler/enqueued/memory': 1,
'start_time': datetime.datetime(2024, 9, 13, 1, 29, 35, 765763, tzinfo=datetime.timezone.utc)}
INFO:scrapy.core.engine:Spider closed (finished)
2024-09-13 01:29:36 [scrapy.core.engine] INFO: Spider closed (finished)
{'title': 'Learn PyTorch for Deep Learning: Zero to Mastery book'}

```

## ✓ Part 2: PPT Scraping

```

1 from google.colab import drive
2 drive.mount('/content/drive')

```

 Show hidden output

```

1 !pip install python-pptx

```

 Show hidden output

Extracting text from PPT using the python-pptx library, it is typically used for generating ppts from databases but we can exploit some of its features here to extract text from ppts, this is a very basic example and it can be explored further as per the need. documentation to the library: <https://python-pptx.readthedocs.io/en/latest/>

```

1 # importing library
2 from pptx import Presentation

1 # extracting texts slide wise and section wise
2 # open the PPT in parallel to check the outcome
3 prs = Presentation("/content/drive/MyDrive/Natural Language Processing/HW2/Your big idea.pptx")
4 counter_slide = 1
5 for slide in prs.slides:
6     print("slide:", counter_slide, "\n")
7     counter_content = 1
8     for shape in slide.shapes:
9         try:
10             print("content:", counter_content, shape.text, "\n")
11             counter_content += 1
12         except: continue
13     print("\n\n")
14     counter_slide += 1

```

 Show hidden output

Similarly other components of the PPT can be extracted after following the documentation as per need

## ✓ Task 6

Extracting table from PPT using pptx

```

1 # extracting texts slide wise and section wise
2 # open the PPT in parallel to check the outcome
3 prs = Presentation("/content/drive/MyDrive/Natural Language Processing/HW2/Scrape This.pptx")
4 counter_slide = 1
5 for slide in prs.slides:
6     print("slide:", counter_slide, "\n")
7     counter_content = 1
8     for shape in slide.shapes:
9         if(shape.has_text_frame):
10             print("content:", counter_content, shape.text, "\n")
11             counter_content += 1
12
13     if shape.has_table:

```

```

14     table = shape.table
15     for row_number, row in enumerate(table.rows, start=1):
16         row_data = []
17         for cell in row.cells:
18             row_data.append(cell.text)
19         print(f"Row {row_number}: {row_data}")
20
21     print()
22     counter_slide += 1

```

➡ slide: 1

content: 1 Scrape This

content: 2

slide: 2

```

Row 1: ['Team Name', 'Year', 'Wins', 'Losses', 'OT Losses', 'Win %', 'Goals For (GF)', 'Goals Against (GA)', '+ / -']
Row 2: ['Boston Bruins', '1990', '44', '24', '', '0.55', '299', '264', '35']
Row 3: ['Buffalo Sabres', '1990', '31', '30', '', '0.388', '292', '278', '14']
Row 4: ['Calgary Flames', '1990', '46', '26', '', '0.575', '344', '263', '81']
Row 5: ['Chicago Blackhawks', '1990', '49', '23', '', '0.613', '284', '211', '73']
Row 6: ['Detroit Red Wings', '1990', '34', '38', '', '0.425', '273', '298', '-25']
Row 7: ['Edmonton Oilers', '1990', '37', '37', '', '0.463', '272', '272', '0']

```

## ➤ Part 3: PDF Scrapping

[ ] ↳ 8 cells hidden

## ▼ Homework

1. As discussed in the demo above, using the example of geeksforgeeks extract the information from ScrapethisSite articles apart from those that is already demonstrated.
2. As discussed in the demo above, extract the salaries of each of the players from the hoopshype website using the example of how to extract the names.
3. Apart from that choose any 2 websites of your choice and extract meaningful and structured information from there.
4. Also explore the scrapy library to perform webscraping apart from the three discussed above in the demo
5. Pick a website that has tabular data and try to scrap it using the tools studied during the demo.

(The datasets you will be collecting for the projects would be by text extraction so make sure to extract usable structured information)

6. Extract table from a PPT using the same library.