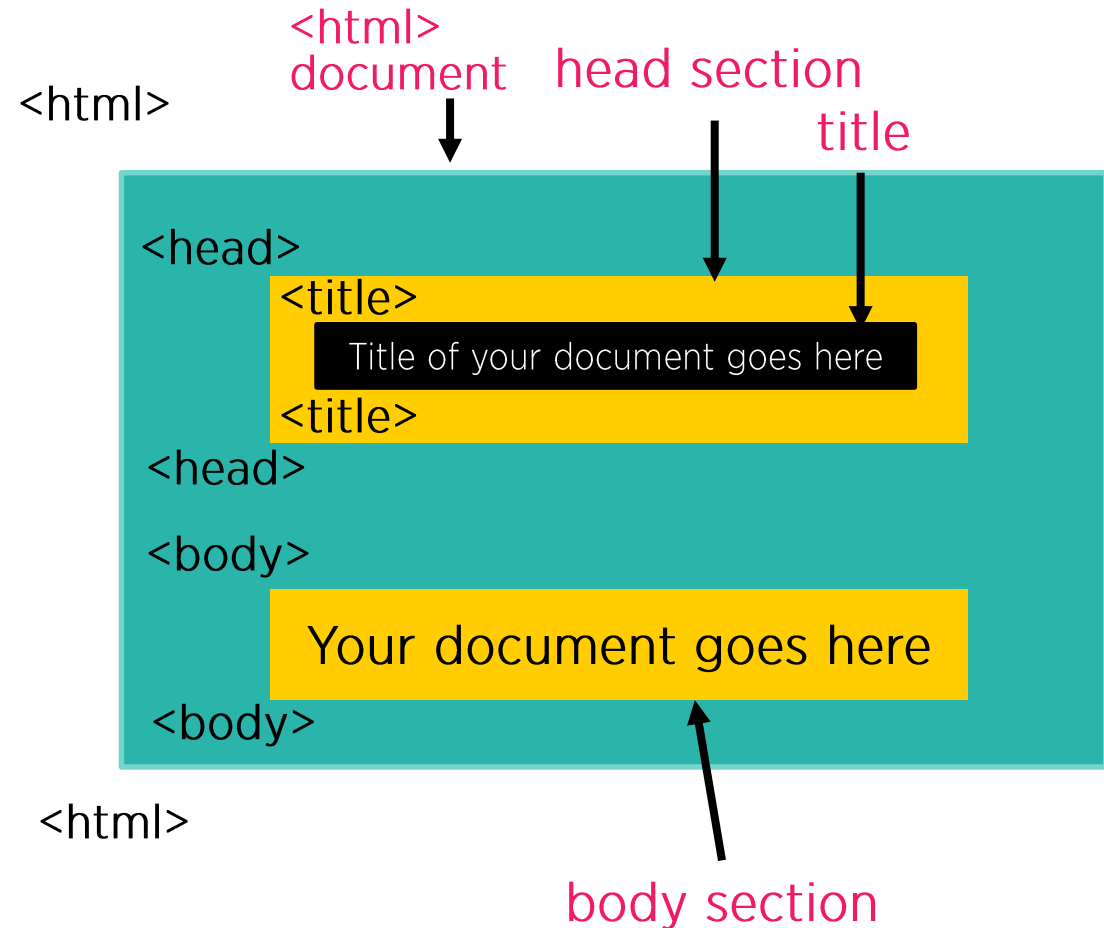# I
# Introduction to Web Scraping

# Web scraping

Web scraping is a technique to automatically access and extract large amounts of information from a website, which can save a huge amount of time and effort.

# Sample structure of a HTML

- o HTML-Stands for Hypertext Markup Language
- o Describes the structure of Web pages using markup language
- o HTML elements are the building blocks of HTML pages
- o HTML elements are represented by HTML tags

<html>

<html> document     head section
title

<head>
<title>

Title of your document goes here

<title>
<head>
<body>

Your document goes here

<body>

<html>

body section

# HTML tags explained

o   The <!DOCTYPE html> declaration defines this document to be HTML5

o   The <html> element is the root element of an HTML page

o   The <head> element contains title of document

o   The <title> element specifies a title for the document

o   The <body> element contains the visible page content

o   The <h1> element defines a large heading

o   The <p> element defines a paragraph

o   The <b> element defines the bold of text

# HTML headings

o    Headings are title of HTML pages

o    Headings are defined with the *<h1>* to *<h6>* tags.

o    h1 is the largest heading tag while h6 is the smallest heading tag.

o    *<h1>* headings should be used for main headings, followed by *<h2>* headings, then the less important *<h3>,* and so on

# **Paragraphs and breaking lines**

Paragraph tag

Most browsers render (process) this with blank lines between each paragraph
*<p>* text *</p>*

Line break tag

Used when the webmaster wants a carriage return but doesn't want a blank line to follow
*<br>*

# Lists – Unordered Lists

o   HTML supplies several list elements. Most list elements are composed of one or more *<li>* (List Item) *</li>* elements.

o    UL : Unordered List: Items in this list start with a list mark such as a bullet. Browsers will usually change the list mark in nested lists.

```
<ul>
        <li>Item One </li>
        <li>Item Two </li>
        <li>Item Three </li>
        <li>Item Four </li>
</ul>
```

Unordered list –

- Item One
- Item Two
- Item Three
- Item Four

# Lists – Ordered Lists

o Ordered List: Items in this list are numbered automatically by the browser.

```
<ol>

        <li> Item One </li>
        <li> Item Two </li>
        <li> Item Three </li>
        <li> Item Four </li>

</ol>
```

o You have the choice of setting the TYPE Attribute to one of five numbering styles.

| Type | Numbering Styles | |
|------|------------------|--------------|
| 1 | Arabic numbers | 1,2,3, …… |
| a | Lower alpha | a, b, c, …… |
| A | Upper alpha | A, B, C, …… |
| i | Lower roman | i, ii, iii, …… |
| I | Upper roman | I, II, III, ……. |

# Lists – Definition Lists

This kind of list is different from the others. Each item in a DL consists of one or more Definition Terms (DT elements), followed by one or more Definition Description (DD elements).

```
<dl>
        <dt>List Name One
                <dd>This is where information about List Name One would go</dd>
        </dt>
        <dt>List Name Two
                <dd>This is where information about List Name Two would go</dd>
        </dt>
</dl>
```

# Links

o The anchor tag *<a>* is used to link one document to another or from one part of a document to another part of the same document.

o Basic Links:

<a href="http://www.stanford.edu/">Stanford University</a>

o Inter-document Links:

<a href="#spot">Point to 'spot' in this document</a>

o Defining a point in a document:

<a name="spot">Spot</a>

o Email links:

<a href="mailto:someone@somehost.com">Email someone@somehost.com</a>

# HTML images

o     Use the HTML <img> element to define an image

o     Use the HTML src attribute to define the URL of the image

o     Use the HTML alt attribute to define an alternate text for an image, if it cannot be displayed

o     Use the HTML width and height attributes to define the size of the image

o     Note: Loading images takes time. Large images can slow down your page. Use images carefully.

# HTML class & id

o   You can create customized classes by adding the class attribute to HTML tags:

o   The syntax for creating a class is:
        &lt;h4 class="class_name" /h4&gt;
        h4 is the HTML header tag
        *class_name* is the name of the class

o   The id attribute must be unique; there can not be more than one tag with the same id value.

o   The syntax for creating an id is:
        &lt;h4 id="id_name" /h4&gt;
        h4 is the HTML header tag
        *id_name* is an id name assigned to the tag

# II
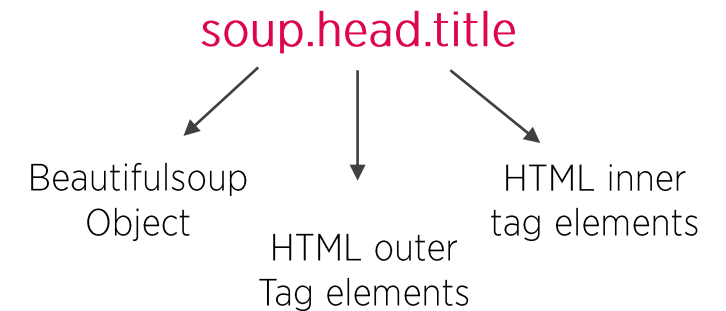# Introduction to BeautifulSoup

# Web Scraping in Python with BeautifulSoup

o  BeautifulSoup is basically an HTML and XML parser and requires additional libraries such as requests, urlib2 to open URLs and store the result

o  BeautifulSoup is relatively easy to understand, used for simple programming tasks with efficiency and can get smaller tasks done in no time

o  BeautifulSoup is slower than Scrapy if you do not use multiprocessing

o  If the project does not require much logic, BeautifulSoup is good for the job, but if you require much customization such as proxies, managing cookies, and data pipelines, Scrapy is the best option.

# BeautifulSoup Functions

```
URL = 'https://en.wikipedia.org/wiki/List_of_game_engines'
      content = requests.get(URL)
      soup = BeautifulSoup(content.text, 'html.parser')
soup.head.title           # returns <title>Head's title</title>
soup.find_all('a')         # if we want all a tags in doc.html
soup.find('a' , href = True)['href']  # returns http://example.com/element1
soup.get_text()
```

soup.head.title

Beautifulsoup
Object

HTML outer
Tag elements

HTML inner
tag elements

# BeautifulSoup Example

```html
<html>
  <head>
    <title> A simple example page
    </title>
  </head>
<body>
    <div> <p class="inner-text first-item" id="first"> First paragraph. </p> </div>
      <p class="outer-text first-item" id="first">
      <b> First outer paragraph. </b>
      </p>
      <p class="outer-text">
      <b> Second outer paragraph. </b> </p>
  </body>
</html>
```

# BeautifulSoup Example

`soup.find_all(id = "first")`

[<p class="inner-text first-item" id="first"> First paragraph. </p>, <p class="outer-text first-item" id="first"> <b> First outer paragraph. </b>]

`soup.find_all('p', class_='outer-text')`

[<p class="outer-text first-item" id="second"> <b> First outer paragraph. </b> </p>, <p class="outer-text"> <b> Second outer paragraph. </b> </p>]

## Using CSS

o BeautifulSoup objects support searching a page via CSS selectors using the select method.

`soup.select("div p")`

[<p class="inner-text first-item" id="first"> First paragraph. </p>]

# Extracting information

```python
page = requests.get("http://forecast.weather.gov/MapClick.php?lat=37.7772&lon=-122.4168")
soup = BeautifulSoup(page.content, 'html.parser')
seven_day = soup.find(id = "seven-day-forecast")
forecast_items = seven_day.find_all(class_ = "tombstone-container")
tonight = forecast_items[0]
print(tonight.prettify())
```

# BeautifulSoup Extracting information

```html
<html>
  <head>
    <title>A simple example page
    </title>
  </head>
 <body>
    <div class="tombstone-container">
      <p class="period-name">
      Tonight
      </p>
      <p class="short-desc"> Mostly Clear </p>
      <p class="temp temp-low"> Low: 49 °F </p>
    </div>
  </body>
</html>
```

# Extracting information

```
period = tonight.find(class_= 'period-name').get_text()
short_desc = tonight.find(class_= 'short-desc').get_text()
temp = tonight.find(class_= 'temp').get_text()
print(period)
print(short_desc)
print(temp)


'''result

Tonight

Mostly Clear

Low: 49 °F '''
```

# III
# Introduction to Scrapy

# What is Scrapy?

**Scrapy** is a Python framework for large scale web scraping. It gives you all the tools you need to efficiently extract data from websites, process them as you want, and store them in your preferred structure and format.

o   With Scrapy you can:

- Fetch millions of data efficiently

- Run it on server

- Run spider in multiple processes

# Web Scraping in Python with Scrapy

o Scrapy is the complete package for downloading web pages, processing them and save it in files and databases

o Scrapy is a powerhouse for web scraping and offers a lot of ways to scrape a web page. It requires more time to learn and understand how Scrapy works but once learned,  eases the process of making web crawlers and running them from just one line of  command.

o Scrapy can get big jobs done very easily. It can crawl a group of URLs in no more than a minute depending on the size of the group and does it very smoothly asynchronously (non-blocking) for concurrency.
Information: **Synchronous** means that you have to wait for a job to finish to start a new job while **Asynchronous** means you can move to another job before the previous job has finished

DATA SCIENCE ACADEMY

# Scrapy Installation

o   If you are using Anaconda :
   *conda install -c conda-forge scrapy*

o   Alternatively, use pip for Windows (Linux, Mac) from command prompt:
   *pip install scrapy*
   or
   *python –m pip install scrapy*

o    For jupyter notebook:
   *!pip install scrapy*

# Scrapy Shell

o Similar to Python shell Scrapy provides a shell of its own that you can use to experiment.

o To start the **scrapy shell** type the following in your command line
   *scrapy shell*
   (If you are using anaconda, you can write it at the anaconda prompt as well.)

o You can use the Scrapy shell to see what components the web page returns.

# **Fetch command**

o   To run Scrapy crawler you have to use **fetch** command in the Scrapy shell:

fetch("https://www.reddit.com/r/gameofthrones/")

o   A crawler or spider goes through a webpage downloading its text and  metadata.

o   The crawler returns a response which can be viewed by using the view(response) command on shell:

view(response)

o   This command will open the downloaded page in your default  browser.

o   You can view the raw HTML script by using the following command in Scrapy  shell:

print(response.text)

# Using CSS Selectors for Extraction

o   You can extract this using the element attributes or the css selector like class, id, etc.

o   **response.css(..)** is a function that helps extract content based on css selector passedto it

- Let's say, we have the following **<a>** tag with a class "product" and the text containing the product name:

    <a class = "product" href = "reddit.com/.../" title = "CHUWI Original" target = "blank">
    CHUWI Hi9 Air Tablet PC Android 0
    </a>

# Using CSS Selectors for Extraction

You can extract the text using the element attributes or the css selector like classes:

response.css(".product::text").extract_first()

**extract_first()** returns the first element that matches the selector.

"." is used with the product because it's a css . Use "::" text to tell your scraper to extract only  text content of the matching elements. If you want to extract all the product names use **extract():**

response.css(".product::text").extract()

# Example

```html
<html>
  <head>
    <base href="http://example.com/"/>
    <title>Example website</title>
  </head>
<body>
  <div id="images">
    <a href="image1.html"> Name: My image 1 </br><img src="image1_thumb.jpg"/></a>
    <a href="image2.html"> Name: My image 2 </br><img src="image2_thumb.jpg"/></a>
    <a href="image3.html"> Name: My image 3 </br><img src="image3_thumb.jpg"/></a>
    <a href="image4.html"> Name: My image 4 </br><img src="image4_thumb.jpg"/></a>
    <a href="image5.html"> Name: My image 5 </br><img src="image5_thumb.jpg"/></a>
  </div>
</body>
</html>
```

# Using CSS Selectors for Extraction

To extract the textual data you must use .**get()** or .**getall()** methods.
**get()** – always returns a single result; If there are more than 1 matches then it returns the first one;

If there are no matches it returnsNone  get(default='not-found') sets default return value None to "not-found"  **getall()** – returns a list with results:

response.css('title::text').get() – Output will be:
'Example website'
response.css('img').xpath('@src').getall()
['image1_thumb.jpg','i
mage2_thumb.jpg',
'image3_thumb.jpg',
'image4_thumb.jpg',
'image5_thumb.jpg']

# Using CSS Selectors for Extraction

extract_first() commonly used in previous Scrapy versions (i.e., alies of .get())

extract() commonly used in previous Scrapy versions (i.e., alies of .getall())
.attrib can be used instead if *@src*; it returns attributes for the first matching element:

1. response.css('img').attrib['src']   =>   'image1_thumb.jpg'   is equivalent to
2. response.css('img::attr(src)').get() => 'image1_thumb.jpg'    is equivalent to
3. response.css('img').xpath('@src').get() => 'image1_thumb.jpg'
4. response.css('base').attrib['href']  =>  'http://example.com/' is equivalent to
5. response.css('base::attr(href)').get() =>  'http://example.com/' is equivalent to
6. response.css('base').xpath('@href').get() =>  'http://example.com/'

# Using CSS Selectors for Extraction

```
<html>
  <head>
        <item>
        <title> My Website </title>
      </item>
  </head>
  <body>
    <span>Hello world!!!</span>
    <div class ="links">
      <a href ="one.html"> Link 1 <img src ="image1.jpg"/></a>
      <a href ="two.html"> Link 2 <img src = "image2.jpg"/></a>
      <a href ="three.html"> Link 3 <img src ="image3.jpg"/></a>
    </div>
    <dc:creaor><![CDATA[Qss, Bootcamp]]></dc:creator>
  </body>
</html>
```

# Using XPath for Extraction

o   XPath (XML Path Language) is a query language for selecting nodes from an XML document. You can navigate through an XML document using XPath.

o   Similar to **response.css(..)** , the function **response.xpath(..)** in scrapy to deal with Xpath  This will show you all the code under the **<html>** tag:

<p style="text-align:center">response.xpath('/html').extract()</p>

/ means direct child of the node
// means descendent nodes

o   If you want to get the **<div>** tags under the html tag you will write:

<p style="text-align:center">response.xpath('/html//div').extract()</p>

o   Behind the scenes, Scrapy uses Xpath to navigate to HTML document items. The CSS  selectors you used above are also converted to XPath, but in many cases, CSS is very  easy to use

# Using XPath for Extraction

○   Let's extract the title of the first post.

response.xpath("//item/title").extract_first()

# result
u'<title> My Website </title>'

# Using XPath for Extraction

o   To filter out the output to only get the text content of the title:

  - response.xpath("//item/title/text()").extract_first()

o   The output:

  -  u' My Website'

o   Notice that **text()** here is equivalent of **::text** from CSS selectors.

o   XPath **//item/title/text()** here you are basically saying find the element "item" and extract the "text" content of its sub element "title".

# Using XPath for Extraction

o  If you want to get all **<div>** tags instead of using /html tag you can do:

        response.xpath("//div").extract()

o  You can further filter your nodes that you start from and reach your desired nodes by using  attributes and their values:

        response.xpath("//div[@class='links']/span[@class='text']").extract()

o  Use /text() to extract all text inside nodes:

        response.xpath("//div[@class='links']/span[@class='text']/text()").extract()

Bütün hüquqlar qorunur.

# Using XPath for Extraction

o   "creator" tag has some text "dc:" because of which it can't be extracted using XPath and the author name itself is crowded with "![CDATA.." irrelevant text.

o   These are just XML namespaces and so we'll ask Scrapy to remove the namespace:

<dc:creaor><![CDATA[Qss, Bootcamp]]></dc:creator>

# Dealing with namespaces in XML Removing namespaces

o   Removing namespaces:

     response.selector.remove_namespaces()

o   Now when you try extracting the author name , it will work :

     response.xpath("//item/creator/text()").extract_first()

o   Output : u' Qss, Bootcamp'

# Site for Web Scraping

o For practical purpose, we will go through an easy example of how to automate downloading cars' names from the web site "turbo.az"

o First go to this web-site:
   https://turbo.az

# Web Scraping

o   Read through the website's Terms and Conditions to understand how you can legally use the data. Most sites prohibit you from using the data for commercial purposes.

o   Make sure you are not downloading data at too rapid a rate because this may break the website. You may potentially be blocked from the site as well.
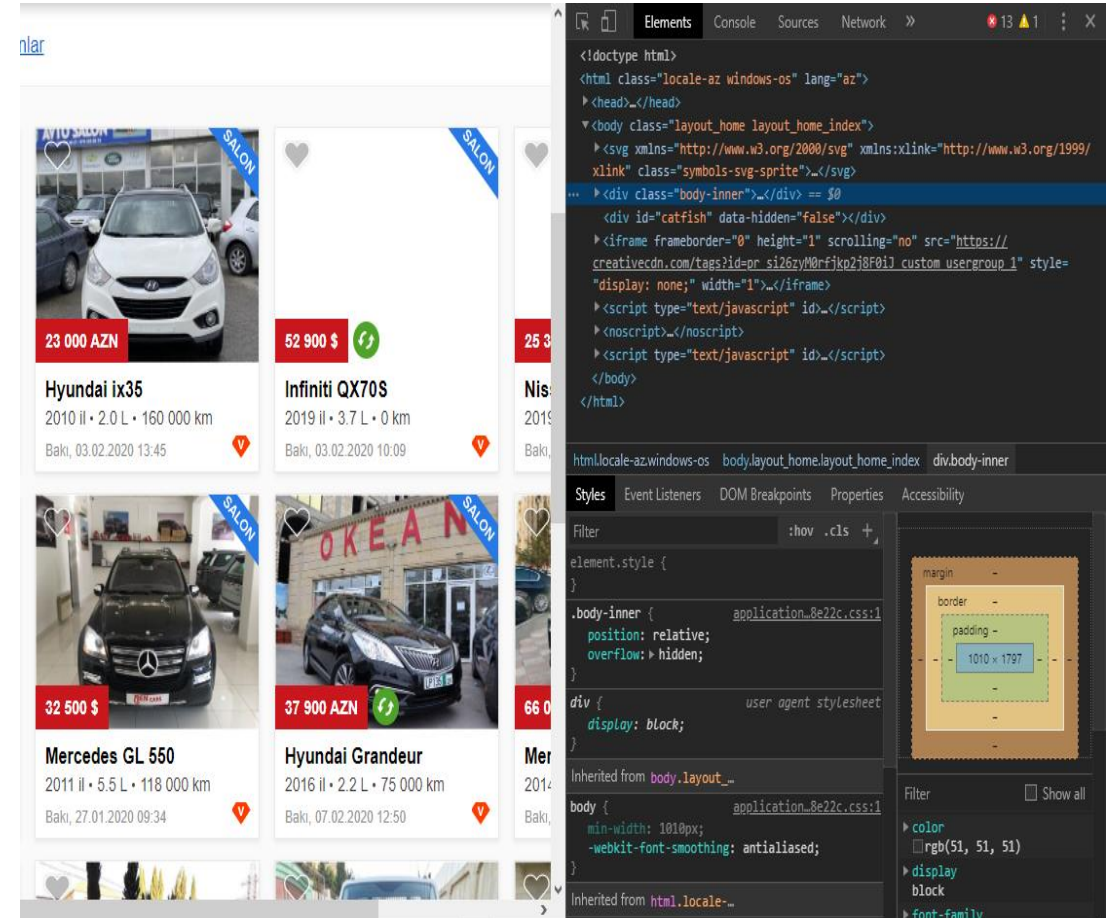
Important libraries:

```
import requests
from bs4 import BeautifulSoup
import pandas as pd
```

# Web Scraping

o   Go turbo.az website

o   Figure out where are our cars' names inside the multiple levels of HTML tags. There is a lot of code on a website page and we want to find the relevant pieces of code that contains our data.

o   Right click and click on "Inspect". This allows you to see the raw code behind the site.

o   Click on arrow and then click on an area of the site itself, the code for that particular item will be highlighted in the console.

o   We will define "https://turbo.az" as url.

url= "https://turbo.az"

Bütün hüquqlar qorunur.

# Web Scraping

o Next we parse the html with BeautifulSoup so that we can work with a nicer, nested BeautifulSoup data structure. If you are interested in learning more about this library, check out the [BeatifulSoup documentation.](#)

o Structure of Beautiful Soup

      Beautiful Soup(response.text, "html.parser")

# Code Structure

```python
# Set the URL you want to webscrape from
    url = "LINK"

# Connect to the url
    response = requests.get(url)

# Parse HTML and save to BeautifulSoup object
    soup = BeautifulSoup(response.text, "html.parser")

# To download the whole data set, let's do a loop through all tags
    fin_names = []
    names = car_soup.findAll('p',attrs={"class":"products-name"})
    for car_name in names:
        fin_names.append(car_name.text)
    print(fin_names)

# Download data as csv
    cars = pd.DataFrame(fin_names)
    cars_csv = cars.to_csv('turbo_cars_names.csv')
```

# Scrapy

o   Installation:  *pip install scrapy*

o   Documentation link: https://docs.scrapy.org/en/latest/intro/install.html

o   Scrapy is an application framework for crawling web sites and extracting structured data which can be used for a wide range of useful applications, like data mining, information processing or historical archival.