# pyladies Bratislava #10

Web scraping with requests and BeautifulSoup

# Components of a web page

- HTML — contain the main content of the page.

- CSS — add styling to make the page look nicer.

- JS — Javascript files add interactivity to web pages.

- Images — image formats, such as JPG and PNG allow web pages to show pictures.

# HTML

tags, properties, …

```
<a href="https://www.python.org" class="library-link">Python</a>
```

- **child** — a child is a tag inside another tag

- **parent** — a parent is the tag another tag is inside

- **sibiling** — a sibiling is a tag that is nested inside the same parent as another tag

# Requests library

- pip install requests
- https://realpython.com/python-requests/

```
>>> response = requests.get('https://weather.com/weather/today/l/48.15,17.11')

        <Response [200]>
```

# Requests library

- Status codes

```python
if response.status_code == 200:
    print('Success!')
elif response.status_code == 404:
    print('Not Found.')
```

- Content

```python
>>> response.content
>>> response.text
>>> response.json()
```
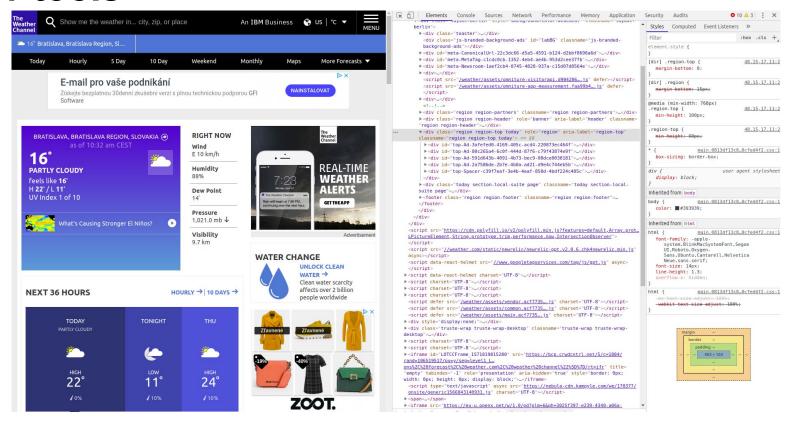
# Requests library

- Request headers and query parameters

```
response = requests.get(
    'https://api.github.com/search/repositories' ,
    params={'q': 'requests+language:python' },
    headers={'Accept': 'application/vnd.github.v3.text-match+json' },
)
```

- Other methods

```
>>> requests.post('https://httpbin.org/post' , data={'key':'value'})
>>> requests.put('https://httpbin.org/put' , data={'key':'value'})
>>> requests.delete('https://httpbin.org/delete' )
>>> requests.head('https://httpbin.org/get' )
>>> requests.patch('https://httpbin.org/patch' , data={'key':'value'})
>>> requests.options('https://httpbin.org/get' )
```

# Dev tools

# BeautifulSoup library

- pip install beautifulsoup4

- https://www.crummy.com/software/BeautifulSoup/bs4/doc/

# BeautifulSoup library

```python
from bs4 import BeautifulSoup

soup = BeautifulSoup(response,
'html.parser')

print(soup.prettify())
```

# BeautifulSoup library

```python
from bs4 import BeautifulSoup

soup = BeautifulSoup(response,
'html.parser')

print(soup.prettify())
```

```html
<html>
 <head>
  <title>
   The Dormouse's story
  </title>
 </head>
 <body>
  <p class="title">
   <b>
    The Dormouse's story
   </b>
  </p>
  <p class="story">
   Once upon a time there were three little sisters; and their names were
   <a class="sister" href="http://example.com/elsie" id="link1">
    Elsie
   </a>
   ,
   <a class="sister" href="http://example.com/lacie" id="link2">
    Lacie
   </a>
   and
   <a class="sister" href="http://example.com/tillie" id="link2">
    Tillie
   </a>
   ; and they lived at the bottom of a well.
  </p>
  <p class="story">
   ...
  </p>
 </body>
</html>
```

# BeautifulSoup library

```
soup.title
# <title>The Dormouse's story</title>
soup.title.name
# 'title'
soup.title.string
# 'The Dormouse's story'
soup.title.parent.name
# 'head'
soup.p
# <p class="title"><b>The Dormouse's story</b></p>
soup.p['class']
# 'title'
soup.a
# <a class="sister" href="http://example.com/elsie" id="link1">Elsie</a>
soup.find_all('a')
# [<a class="sister" href="http://example.com/elsie" id="link1">Elsie</a>,
#  <a class="sister" href="http://example.com/lacie" id="link2">Lacie</a>,
#  <a class="sister" href="http://example.com/tillie" id="link3">Tillie</a>]
soup.find(id="link3")
 # <a class="sister" href="http://example.com/tillie" id="link3">Tillie</a>
list(soup.p.children)
# [<b>The Dormouse's story</b>]
```

# BeautifulSoup library

soup.find_all(name, attrs)

soup.find_all("a", class_="sister")
soup.find_all("a", attrs={"class": "sister"})

soup.find(id="link2")

# BeautifulSoup library

[CSS selectors](#)

- **p a** — finds all a tags inside of a p tag

- **body p > a** — finds all a tags directly inside of a p tag inside of a body tag

- **p.outer-text** — finds all p tags with a class of outer-text

- **p#first** — finds all p tags with an id of first

- **body p.outer-text** — finds any p tags with a class of outer-text inside of a body tag

```
soup.select("div p")
```

# BeautifulSoup library

extract text from the element

```
soup.select("div p")[0].get_text()
```

# Send an email

https://www.afternerd.com/blog/how-to-send-an-email-using-python-and-smtplib/