

# Extrakcia informácií z webu pre analýzu dát pomocou Pythonu

Metódy inžinierskej práce 2023/2024

Štefan Kučerák

`xkucera@fiit.stuba.sk`

Ústav informatiky, informačných systémov a softvérového inžinierstva  
Fakulta informatiky a informačných technológií  
Slovenská technická univerzita v Bratislave

26. november 2023

# Úvod

- automatizované zbieranie dát z web stránok
- časová náročnosť

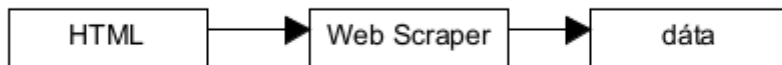
# Prehľad

1 Princíp

2 Príklady využitia

3 Príklad

# Princíp extrakcie dát



Nejaká poznámka k obrázku, možno zdroj. . .

# Knižnice

Knižnica	Popis
Requests <sup>1</sup>	HTTP požiadavky a to POST aj GET
Beautiful Soup <sup>2</sup>	táto knižnica obsahuje nástroje na html alebo xml
Scrapy <sup>3</sup>	je rýchly framework pre prehľadávanie webu na vysokej úrovni
Selenium <sup>4</sup>	je knižnica, ktorá je web driver

---

<sup>1</sup> <https://pypi.org/project/requests/>

<sup>2</sup> <https://pypi.org/project/beautifulsoup4/>

<sup>3</sup> <https://pypi.org/project/Scrapy/>

<sup>4</sup> <https://pypi.org/project/selenium/>

# Príklady využitia

Všetky z týchto projektov využili Python a niektorú z uvedených knižníc:

- Phishing Web Page Detection using web Scraping
- Food Genie, Recipe Search Algorithm using Web Scraping
- NEWSONE- AN AGGREGATION SYSTEM FOR NEWS USING WEB SCRAPING METHOD

# Jednoduchý príklad

Vypísanie názvu stránky pomocou knižnice BeautifulSoup

## ■ Program

```
import requests
from bs4 import BeautifulSoup

page = requests.get("https://www.google.com/")
soup = BeautifulSoup(page.content, "html.parser")

print(soup.title.string)
```

## ■ Výstup

Google

## Zhodnotenie a ďalšia práca

- Každá prezentácia musí byť nejako uzavretá
- Ale vždy je čo robiť ďalej...