

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 Príklady reálneho využitia

2 Knižnice

3 Princíp

4 Príklad

Príklady reálneho využitia

- "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"

Knižnice

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

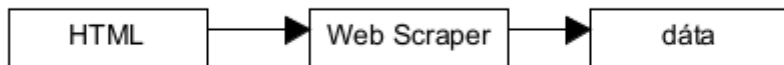
¹ <https://pypi.org/project/requests/>

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

³ <https://pypi.org/project/Scrapy/>

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

Princíp extrakcie dát



Nejaká poznámka k obrázku, možno zdroj. . .

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...