



ZAKŁAD SYSTEMÓW ZŁOŻONYCH

Wydział Elektrotechniki i Informatyki

ul. Wincentego Pola 2, 35-959 Rzeszów,

tel. 17 865 1340

zsz.prz.edu.pl



WYDZIAŁ
ELEKTROTECHNIKI
I INFORMATYKI
POLITECHNIKI RZESZOWSKIEJ

Pracownia Problemowa II

Projekt

Gromadzenie danych z Internetu dotyczących ofert
pracy w oraz ich analiza.

Przemysław Skubel, 160466

Inżynieria i analiza danych, rok V, sem. IX

Rzeszów 2024r.



Spis treści

1. Czym jest Web Scraping.....	2
2. Opis Projektu.....	3
3. Zbieranie danych	3
3.1. Just Join IT	3
3.2. No Fluff Jobs.....	10
3.3. FlexJobs	14
4. Analiza zebranych danych.....	18
Wnioski.....	21
Źródła	22

1. Czym jest Web Scraping

Web scraping to proces automatycznego zbierania danych z witryn internetowych. Jest to technika, która polega na ekstrakcji informacji z stron internetowych poprzez analizę ich struktury HTML, CSS lub innych formatów danych. Web scraping jest szeroko stosowany w różnych dziedzinach, takich jak analiza danych, zbieranie informacji marketingowej, monitorowanie cen online, badania konkurencji i wiele innych.



2. Opis Projektu

W projekcie do zbierania danych wykorzystam bibliotekę BeautifulSoup, jest to biblioteka w języku Python do parsowania dokumentów HTML i XML. Jest często używana do ekstrakcji danych z witryn internetowych, przetwarzania dokumentów HTML/XML i manipulowania strukturą drzewa dokumentu. BeautifulSoup umożliwia łatwe nawigowanie po drzewie dokumentu HTML, wyszukiwanie, filtrowanie i modyfikowanie danych.

3. Zbieranie danych

Każda ze stron, z której będę zbierać informacje o dostępnych ofertach pracy w IT jest tego samego typu. Na stronie tytułowej znajdują się oferty, które możemy filtrować za pomocą dostępnych filtrów poprzez wybranie technologii, lokalizacji, widełek zarobków i tym podobnych.

3.1. Just Join IT

Pierwszą stroną z jakiej pobiorę oferty będzie strona Job Offers IT - Just Join IT. Oto jak przedstawia się strona główna pierwszej wyszukiwarki



ZARŁAD SYSTEMÓW ZŁOŻONYCH

Wydział Elektrotechniki i Informatyki
ul. Wincentego Pola 2, 35-959 Rzeszów,
tel. 17 865 1340
eez.prz.edu.pl



WYDZIAŁ
ELEKTROTECHNIKI
I INFORMATYKI
POLITECHNIKI RZESZOWSKIEJ

justjoin.it #1 Job Board for tech industry in Europe

Job offers Top Companies Geek Post a job Sign in

Search Location JS HTML PHP Ruby Python Java Net Scala C Mobile Testing DevOps Admin UX/UI PM Game Analytics Security Data Go Support ERP Archite... Other

Offers with salary All offers 17 108 offers Remote Default

Work: All offers - 17 108 offers

- STATSCORE** Senior C#/.NET Developer (Robots) 21 600 - 27 600 PLN New
STATSCORE Sp. z o.o. Katowice .NET Core C# Cloud
- P&G** Platform Engineer Undisclosed Salary New
Procter & Gamble Warsaw Software Development CI/CD Azure
- BROWN BROTHERS HARRIMAN** Senior Test Automation Engineer Undisclosed Salary New
Brown Brothers Harriman Kraków TestNG Selenium Java
- DSV** Senior Network Specialist Undisclosed Salary New
DSV ISS Warszawa ACI Cisco Data Center
- B BRAUN** SAP BASIS Specialist Undisclosed Salary New
B. Braun Poznań SAP SAP HANA Fiori
- luxoft** C#/.Net Developer with CyberArk Undisclosed Salary New
Luxoft Poland Kraków C# .NET C# CyberArk
- Zeta Labs** Senior Backend Engineer 30 000 - 40 000 PLN New
Warszawa SQL fastapi Python
- Google** Multidisciplinary Prompt Engineer 20 000 - 25 000 PLN New

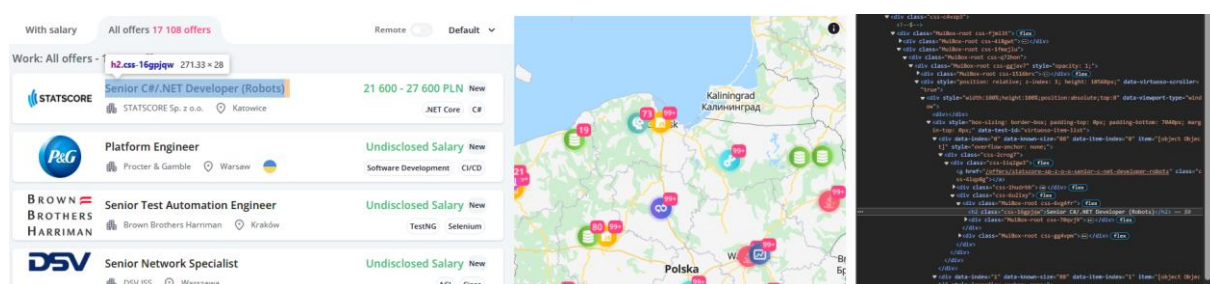
Na początku potrzebujemy tylko adres URL naszej strony.

```
url = 'https://justjoin.it/'
page = requests.get(url)
soup = BeautifulSoup(page.text, "html.parser")
print(soup.prettify())
```

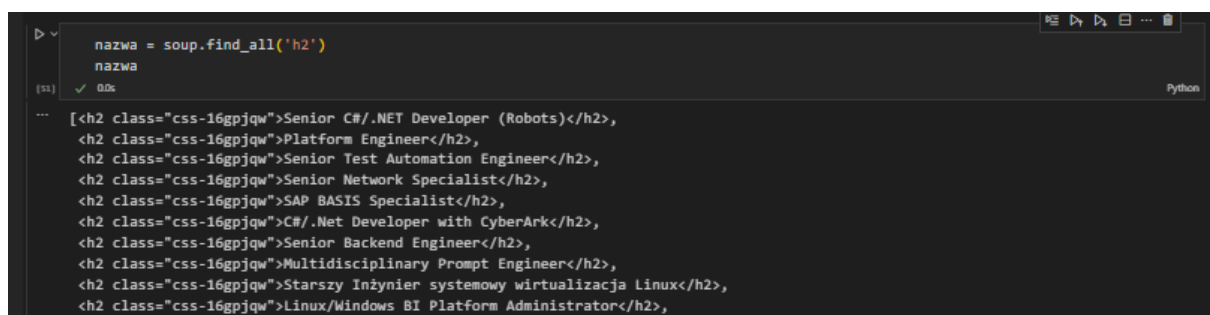
```
<!DOCTYPE html>
<html>
<head>
<meta charset="utf-8">
<link href="https://api.justjoin.it" rel="preconnect">
<link href="https://public.justjoin.it" rel="preconnect">
<link as="font" crossorigin="" href="https://public.justjoin.it/fonts/open-sans/open-sans-v34-latin-latin-ext-regular.woff2" rel="prel
<link as="font" crossorigin="" href="https://public.justjoin.it/fonts/open-sans/open-sans-v34-latin-latin-ext-600.woff2" rel="prel
<link as="font" crossorigin="" href="https://public.justjoin.it/fonts/open-sans/open-sans-v34-latin-latin-ext-700.woff2" rel="prel
<meta content="The first map of the labor market in the IT sector. We want to simplify the search process to minimum." property="de
<meta content="https://justjoin.it/gfx/justjoinit/og-fallback-image.png" property="og:image"/>
<meta content="1200" property="og:image:width"/>
<meta content="630" property="og:image:height"/>
<meta content="https://justjoin.it" property="og:url"/>
<meta content="Just Join IT" property="og:title"/>
<meta content="website" property="og:type"/>
<meta content="Just Join IT: #1 Job Board for tech industry in Europe" name="description" property="og:description"/>
<meta content="summary_large_image" name="twitter:card"/>
<meta content="https://justjoin.it/gfx/justjoinit/og-fallback-image.png" property="twitter:image"/>
<meta content="https://justjoin.it" property="twitter:site"/>
<meta content="Just Join IT" property="twitter:title"/>
<meta content="Just Join IT: #1 Job Board for tech industry in Europe" property="twitter:description"/>
<meta content="width=device-width,initial-scale=1,maximum-scale=1,user-scalable=0,shrink-to-fit=no,interactive-widget=resizes-cont
<link href="/_next/static/media/favicon_jjit_16x16.17c0797f.png" rel="icon" sizes="16x16" type="image/png"/>
<link href="/_next/static/media/favicon_jjit_32x32.5fa64152.png" rel="icon" sizes="32x32" type="image/png"/>
...
{"props":{"pageProps":{"tenantConfig":{"privacyPolicyUrl":{"rjNew":[{"url":"https://public.justjoin.it/Polityka%2BPrywatności%2B31
</script>
</body>
</html>
```



Następnie na naszej stronie włączamy opcję „Zbadaj” i wyszukujemy interesujące nas treści. W tym przypadku na początku będziemy chcieli zebrać informacje na temat stanowiska.



Możemy zauważyć, że nazwy stanowisk są wpisane w akapicie `<h2>` w kodzie HTML, zatem możemy wykorzystać to do ich wyciągnięcia z tak zwanej „Zupy”:



Następnie naszym celem jest wydobyć informacje o nazwie stanowiska z poszczególnego elementu listy za pomocą Regexu. Regex jest to wzorzec znaków, który służy nam do wyciągania potrzebnych nam znaków z łańcucha znaków. Do stworzenia Regexu pomogła mi strona regex101: build, test, and debug regex.



```
tekst = str(nazwa)
pattern = r'>\s*(.*?)\s*<'
matches = re.findall(pattern, tekst)
macthes
stanowiska = []
for indeks, i in enumerate(matches):
    if indeks % 2 == 0:
        stanowiska.append(i)
stanowiska

[('Manager ds. Zarządzania',
'SASE Engineer',
'PHP Software Developer',
'Python Developer - TALENT POOL',
'IT Project Manager - Banking 🧑‍💻',
'Java Tech Lead',
'Microsoft Power Platform Engineer',
'Entry ServiceNow Specialist',
'Administrator ds. Utrzymania Systemów Linux',
```

Za pomocą prostej pętli i regexu wyciągamy interesujące nas informacje. Identyfikujemy postępujemy z kolejnymi informacjami, które będziemy chcieli zawrzeć w naszym zbiorze.

Za każdym razem jednak zmieniamy sposób wyszukiwania tj. czy informacje znajdują się w nagłówku `` lub `<div>`, jeżeli znajdują się w jednym z nich, będziemy musieli zawęzić zbiór informacji, do najczęściej klasy, w której znajdują się potrzebne dane.

```
lokalizacja = soup.find_all('span', class_ = 'css-1o4wo1x')
lokalizacja

[<span class="css-1o4wo1x">Łódź</span>,
<span class="css-1o4wo1x">Warszawa</span>,
<span class="css-1o4wo1x">Gdańsk</span>,
<span class="css-1o4wo1x">Warszawa</span>,
<span class="css-1o4wo1x">Warsaw</span>,
```



ZAKŁAD SYSTEMÓW ZŁOŻONYCH

Wydział Elektrotechniki i Informatyki
ul. Wincentego Pola 2, 35-959 Rzeszów,
tel. 17 865 1340
zsz.prz.edu.pl



WYDZIAŁ
ELEKTROTECHNIKI
I INFORMATYKI
POLITECHNIKI RZESZOWSKIEJ

```
tekst = str(lokalizacja)
pattern = r'>\s*(.*?)\s*<'
loc = re.findall(pattern, tekst)
lokalizacja = []
for indeks, i in enumerate(loc):
    if indeks % 2 == 0:
        lokalizacja.append(i)
lokalizacja

[18] ✓ 0.1s

... ['tódź',
      'Warszawa',
      'Gdańsk',
      'Warszawa',
      'Warsaw',
      'Warszawa',
      'Warszawa',
```

```
widelki = soup.find_all('div', class_ = 'css-1b2ga3v')
widelki

()

pattern = r'^\b(\d{1,3})?: \d{3})*(?:\.\d{2})?\b$'

widly = []
for indeks, i in enumerate(widelki):
    w1 = str(widelki[indeks])
    w2 = re.findall(pattern, w1)
    widly.append(w2)
widly

()

Widelki = []

for indeks, i in enumerate(widly):
    l = []
    for j in range(3):
        if len(widly[indeks]) == 2:
            a = widly[indeks][j]
            b = a.replace(".", "")
            a = int(a)
            b = int(b)
            l.append(b)
        else:
            continue
    Widelki.append(l)

[18] ✓ 0.1s

Widelki
[18] ✓ 0.1s

... [[7670, 11250],
      [27000, 32000],
      []]
```

```
# currency (waluta)
waluty = soup.find_all('span', class_ = 'css-jmy9db')
waluty

()

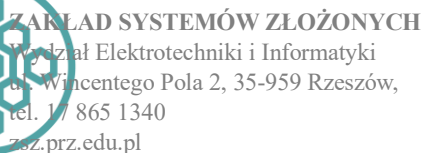
waluta = []
for indeks, i in enumerate(waluty):
    txt = str(waluty[indeks])
    pattern = r'<span class="css-jmy9db">(.*?)</span>'
    w = re.findall(pattern, txt)
    waluta.append(w)
waluta

()

waluta1 = []
for indeks, i in enumerate(waluta):
    k = waluta[indeks][0]
    waluta1.append(k)
waluta1

[61] ✓ 0.0s

... ['pln',
      'pln',
      'pln',
      'pln',
      'pln',
```



```

nazwa = soup.find_all('div', class_ = 'css-ldh1c9')
firma = []
for indeks, i in enumerate(nazwa):
    txt = str(nazwa[indeks])
    pattern = r'<span>(.*?)</span>'
    w = re.findall(pattern, txt)
    firma.append(w)
firma

[]

```

```

firma = []
for indeks, i in enumerate(firma):
    k = firma[indeks][0]
    firma.append(k)
firma

['Rossmann',
 'Codillime',
 'GetResponse S.A.',
 'Polcode',
 'ITDS',
 'ITDS',
 'ITDS']

```



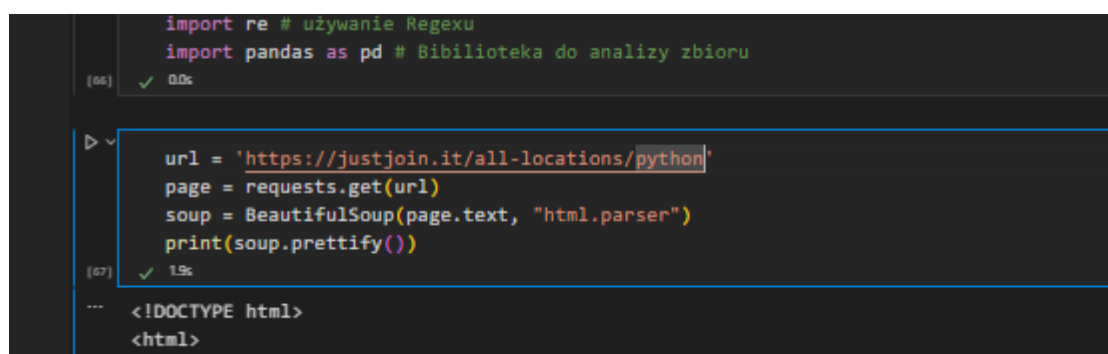
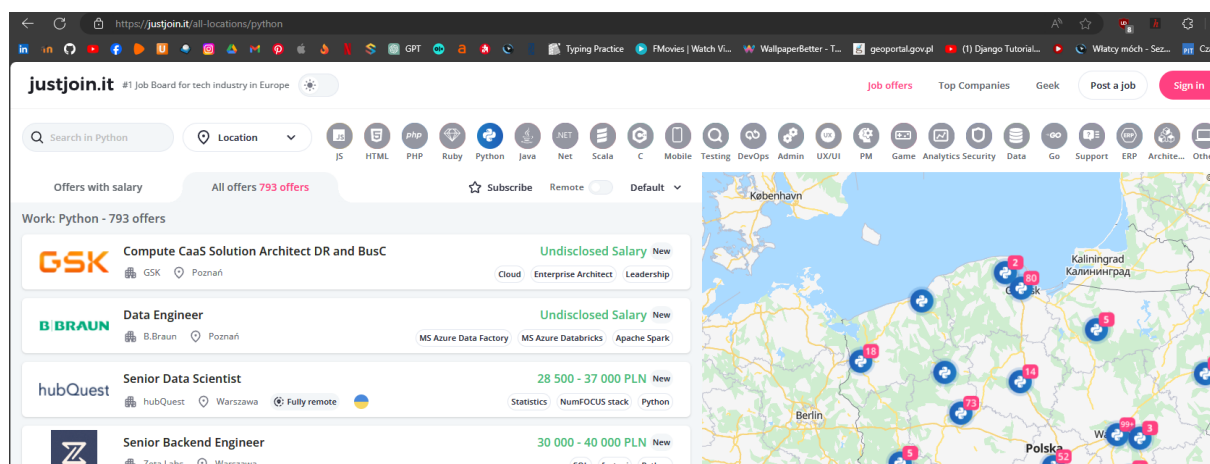

```
data = {
    'Firma' : firma,
    'Stanowisko' : stanowiska,
    'Lokalizacja' : lokalizacja,
    'Widelki' : widly,
}

df = pd.DataFrame(data) # utworzenie DataFrame
df.to_excel('JustJoinIT.xlsx', index = False) # zapis tabeli do Excel
df
```

	Firma	Stanowisko	Lokalizacja	Widelki
0	GSK	Product Manager - PAT System	Poznań	[]
1	Human4Human Recruitment	C++ Engineering Lead (Medical industry)	Kraków	[]
2	Human4Human Recruitment	Odoo and Big Data Architect	Kraków	[(32 000, 36 000)]
3	Code Poets	Developer C++	Wroclaw	[(17 000, 23 500)]
4	Human4Human Recruitment	Senior Google Analyst	Kraków	[(20 000, 24 000)]
...
95	DSV ISS	Senior Back-end Java Developer	Warszawa	[]
96	Schibsted	Engineering Manager .Net	Kraków	[]
97	Schibsted	Engineering Manager	Kraków	[]
98	DRÄXLMAIER	V&V Engineer - Cybersecurity Testing	Gliwice	[]
99	ITDS	GoLang Developer - Analytics Consulting Team -...	Warsaw	[(19 000, 22 000)]

100 rows x 4 columns

Tak przedstawia się tabela zebranych przeze mnie danych ze strony JustJoinIT. Warto wspomnieć o tym, że strona oferuje filtrowanie ofert, które następnie jest zapisywane w adresie URL strony. Zatem, jeśli byśmy chcieli zebrać informacje tylko i wyłącznie informacje odnośnie ofert w np. języku Python, musimy tylko zaznaczyć na stronie technologię Python, a następnie zamienić początkowy adres URL na otrzymany po dodaniu filtra.



3.2. No Fluff Jobs

Drugą z kolei stroną będzie strona [Portal z ofertami pracy IT | Praca IT dla Ciebie | No Fluff Jobs](#) tej znajdują się poza ofertami pracy dla programistów również stanowiska w np. Human Resources lub stanowisk stricte dla inżynierów.



Oto jak przedstawia się strona tytułowa:

The screenshot shows the homepage of the 'NO FLUFF JOBS' website. The header includes navigation links: 'OFERTY PRACY', 'PROFILE FIRM', 'MYSALARY', 'DLA PRACODAWCY', 'PUBLIKUJ', 'ZALOGUJ SIĘ', and a language selector 'PL, POL, PLN'. A main banner states 'Szanujemy Twój czas. 100% ogłoszeń z widelkami.' with a search bar showing '18110 ofert z widelkami'. Below the banner are filters for 'IT', 'HR', 'Marketing', 'Sprzedaż', 'Finanse', 'Pozostałe', 'Lokalizacja', 'Doświadczenie', 'Wynagrodzenie', 'Więcej', and 'Mapa'. A section titled 'Szykujesz się do zmiany pracy?' offers a 'Pomóż My Salary' button. The 'OFERTY DNIA' section lists four job offers:

Job Title	Salary Range	Location
Principal Autosar Developer	21 300 – 37 300 PLN	Zdalnie +1
Principal Software Engineer	23 750 – 30 083 PLN	Kraków
Senior TAX Manager (CIT/M&A)	17 000 – 21 000 PLN	Warszawa
Specjalista / Specjalistka w dziale HelpDesk	5 300 – 8 800 PLN	Kraków

Wszystkie operacje zbierania danych są analogiczne jak przy zbieraniu informacji ze strony Just Join IT z małymi różnicami w kodach regexu. No Fluff Jobs wymagało również więcej operacji na pętlach i głębszego wyciągania danych z kodu HTML, ponieważ informacje na tej stronie są inaczej rozmieszczone i ułożone w bardziej skomplikowanym i złożonym formacie.



ZARŁAD SYSTEMÓW ZŁOŻONYCH

Wydział Elektrotechniki i Informatyki
ul. Wincentego Pola 2, 35-959 Rzeszów,
tel. 17 865 1340
zoz.prz.edu.pl



WYDZIAŁ
ELEKTROTECHNIKI
I INFORMATYKI
POLITECHNIKI RZESZOWSKIEJ

```
NoFluffJobs.py:nb x
NoFluffJobs.py:nb > url = 'https://nofluffjobs.com/pl'
+ Code + Markdown | Run All | Restart | Clear All Outputs | Go To | Variables | Outline ...

from bs4 import BeautifulSoup # możliwość ściągnięcia kodu HTML
import requests # wysłanie requestu do strony o jej kod HTML
import re # używanie Regexu
import pandas as pd # Biblioteka do analizy zbioru

url = 'https://nofluffjobs.com/pl'
page = requests.get(url)
soup = BeautifulSoup(page.text, "html.parser")
print(soup.prettify())

job_name_soup = soup.find_all('h3', class_ = "posting-title__position text-truncate color-main ng-star-inserted")

txt = str(job_name_soup)
pattern = r">\s*(.*?)\s*<"
matches = re.findall(pattern, txt)
matches
job_names = []
for indeks, i in enumerate(matches):
    if indeks % 2 == 0:
        job_names.append(i)
job_names
```

```
salary = soup.find_all('span', class_ = "text-truncate badgy salary lg:tw-btn tw-text-ink lg:tw-btn-secondary-outline tw-text-xs lg:tw-py-0.5 lg:tw-px-2 ng-star-inserted")
salaries = []
for i in salary:
    s = i.get_text()
    salaries.append(s)

# salaries

requirements = soup.find_all('div', class_ = "lg:tw-w-[50%] lg:tw-pr-3 tw-flex-[0_0_50%]")
requirements

misa = []

for i in requirements:
    noodle = i.find_all('span', class_ = "lg:tw-text-gray-60 lg:tw-border-2 lg:tw-border-gray-ddd tw-text-xs tw-lowercase lg:tw-py-0.5 lg:tw-px-2 tw-text-gray-60")
    # print(vd)
    bowl = []
    for j in noodle:
        carrot = str(j)
        # print(j)
        spoon = re.findall(pattern, carrot)
        # print(spoon)
        bowl.append(spoon)
    misa.append(bowl)

# misa
```



```
NoFluffJobs.py  
+ Code + Markdown | ▶ Run All | ⏏ Restart | 🧹 Clear All Outputs | 📄 Variables | 📄 Outline | ... Python 3.10.0  
  
company = soup.find_all('h4', class_ = 'tw-text-gray-60 company-name tw-w-[50%] lg:tw-auto tw-mb-0 !tw-text-xs !lg:tw-text-sm tw-font-semibold  
lg:tw-font-normal')  
pattern = '>\s(?:.*?)\s<'  
companies = []  
for i in company:  
    i = str(i)  
    carrot = re.findall(pattern, i)  
    companies.append(carrot)  
  
# companies  
  
city = soup.find_all('span', class_ = 'tw-text-ellipsis tw-inline-block tw-overflow-hidden tw-whitespace-nowrap tw-max-w-[100px] md:tw-max-w-[200px]  
tw-text-right')  
cities = []  
pattern = '>\s(?:.*?)\s<'  
for i in city:  
    i = str(i)  
    carrot = re.findall(pattern, i)  
    cities.append(carrot)  
#cities  
  
df = {  
    'Firma': companies,  
    'Stanowisko': job_names,  
    'Widélki': salaries,  
    'Wymagania': misa,  
    'Lokalizacja': cities  
}  
  
df = pd.DataFrame(df)  
df.to_excel('NoFluffJobs.xlsx', index = False)
```

	Firma	Stanowisko	Widélki	Wymagania	Lokalizacja
0	[State Street]	Senior UI/UX Developer	17 416 – 30 083 PLN	[[Design], [ui], [ui design], [ux design]]	[Kraków]
1	[Alsendo Sp. z o. o.]	Remote Data/ Business Analyst (data provisioning)	21 840 – 23 520 PLN	[[Data], [data mapping], [csv], [business anal...]]	[Warszawa]
2	[LeasingTeam Professional]	Logisztikai SAP koordínátor-6355	10 315 – 12 608 PLN	[[Logistyka], [sap], [węgierski], [angielski]]	[Warszawa]
3	[State Street]	Principle DBA Engineer	23 750 – 30 083 PLN	[[Backend], [relational database], [operating ...]]	[Kraków]
4	[Neontr]	HR Specialist with German	10 500 – 14 280 PLN	[[HR], [niemiecki]]	[Zdálanie]
5	[7N]	Senior DevOps Developer	23 750 – 30 083 PLN	[[DevOps], [python], [powershell], [perl]]	[Zdálanie]
6	[State Street]	Specjalista / Specjalistka w dziale HelpDesk	5 300 – 8 800 PLN	[[IT Support], [windows], [macos], [windows se...]]	[Gdańsk]
7	[Ringier Axel Springer Tech]	Senior .NET Developer	16 000 – 20 000 PLN	[[Fullstack], [.net], [.net core], [c#]]	[Kraków]
8	[Inter Cars S.A.]	Remote Senior PHP Developer (Magento)	10 000 – 20 000 PLN	[[Backend], [php], [magento], [html]]	[Warszawa]
9	[Forward Thinking Systems Polska]	Senior .NET Developer	16 000 – 23 000 PLN	[[Backend], [.net], [c#], [azure devops]]	[Zdálanie]
10	[CTHINGS.CO]	Software Development Engineer	20 000 – 25 000 PLN	[[Backend], [golang], [java], [sql]]	[Warszawa]
11	[Exorigo-Upos S. A.]	Senior PHP Developer	15 000 – 18 000 PLN	[[Backend], [php], [oop], [symfony]]	[Zdálanie]
12	[Hitachi Energy]	Mid Java/Kotlin/Scala Developer	11 424 – 17 136 PLN	[[Backend], [java], [spring], [git]]	[Kraków]
13	[HL Tech]	Frontend Engineer (React, Nextjs, AWS)	15 000 – 21 000 PLN	[[Frontend], [react], [html], [css]]	[Zdálanie]
14	[LegalZoom]	Senior Frontend Developer	20 000 – 30 000 PLN	[[Frontend], [react], [javascript], [typescript]]	[Kraków]

Podobnie jak w przypadku strony Just Join IT, aby pobierać informacje tylko o stanowiskach z filtrami, które nas interesują, musimy wybrać te filtry i następnie wprowadzić zmieniony URL do programu.



Jeżeli będziemy chcieli zbierać dane masowo, niestety nie jesteśmy w stanie tego zrobić. Klikając na dole strony „załaduj więcej ofert”, kod URL zmienia się na: <https://nofluffjobs.com/pl/?page=2>, co mogłoby nam umożliwić ustawienie pętli i zapisywanie bardzo dużej ilości danych. Niestety nie jesteśmy w stanie tego zrobić. W pliku „proba1.py”, zastosowałem taki kod, który zapisuje dane z pierwszych 5 stron, ale dane są takie same w każdym arkuszu. Powodem tego jest to, że strona nie generuje więcej kodu HTML, dopóki nie klikniemy banneru „Załaduj więcej”, a przejście bezpośrednio z linku do strony np. 4, przekierowuje nas na sam początek strony.

3.3. FlexJobs

Trzecią stroną, z której zbiorę informację jest [Computer & IT Jobs - Remote Work From Home & Online | FlexJobs](#) która zawiera oferty z branży IT jak i innych branż inżynierskich. Strona jest anglojęzyczna (oferty USA) i oferuje stanowiska w większości pracy zdalnej. Wybrałem tą stronę, ponieważ przy pobieraniu danych, widnieje jej opis, co pozwala nam na lepszą analizę tych danych.

Podobnie jak przy scrapowaniu poprzednich stron, cała technika jest analogiczna.



ZARŁĄD SYSTEMÓW ZŁOŻONYCH

Wydział Elektrotechniki i Informatyki
ul. Wincentego Pola 2, 35-959 Rzeszów,
tel. 17 865 1340
zez.prz.edu.pl



WYDZIAŁ
ELEKTROTECHNIKI
I INFORMATYKI
POLITECHNIKI RZESZOWSKIEJ

```
from bs4 import BeautifulSoup # możliwość ściągnięcia kodu HTML
import requests # wysłanie requestu do strony o jej kod HTML
import re # używanie Regexu
import pandas as pd # biblioteka do analizy zbioru

url = 'https://www.flexjobs.com/remote-jobs/computer-it'
page = requests.get(url)
soup = BeautifulSoup(page.text, "html.parser")
print(soup.prettify())

job_name_soup = soup.find_all('a', class_ = 'job-title job-link')
job_name_soup

jobs = []
pattern = r"><(.*)><"
for indeks, i in enumerate(job_name_soup):
    i = str(i)
    job_name = re.findall(pattern, i)
    a = []
    for i in job_name:
        if i != "":
            a.append(i)
        else:
            continue
    jobs.append(a)
jobs

location = soup.find_all('div', class_ = 'col pe-0 job-locations text-truncate')
location
```

```
job_loc = []
pattern = r"><(.*)><"
for indeks, i in enumerate(location):
    i = str(i)
    job_name = re.findall(pattern, i)
    a = []
    for i in job_name:
        if i != "":
            a.append(i)
        else:
            continue
    job_loc.append(a)
job_loc

desc_soup = soup.find_all('div', class_ = 'job-description')
desc_soup

job_description = []
pattern = r"><(.*)><"
for indeks, i in enumerate(desc_soup):
    i = str(i)
    job_name = re.findall(pattern, i)
    a = []
    for i in job_name:
        if i != "":
            a.append(i)
        else:
            continue
    job_description.append(a)
job_description

data = {
    'Stanowisko': jobs,
    'Lokalizacja': job_loc,
    'Opis': job_description
}

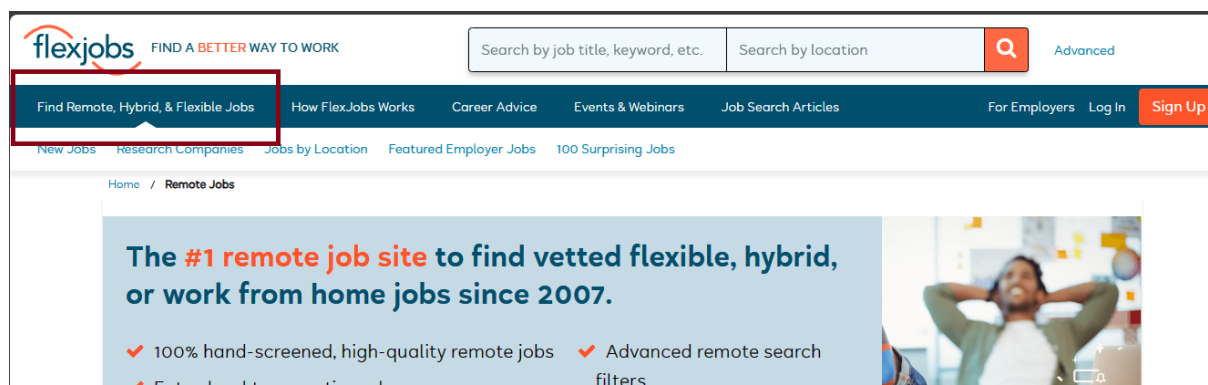
df = pd.DataFrame(data)
```



```
df = pd.DataFrame(data)
df.to_excel('FlexJobs.xlsx', index = False)
```

	Stanowisko	Lokalizacja	Opis
0	[Senior JDE E1 Developer]	[Thornton, CO]	[Develop, modify and maintain application syst...
1	[Golang Core Developer]	[Bloomfield, CT]	[Design, develop, and maintain Golang-based mi...
2	[Senior Salesforce Developer]	[Mexico]	[Architect and develop complex technical solut...
3	[Application Developer, Workday Integration]	[Mexico]	[Design, build, deliver, and support technical...
4	[Senior Software Development Engineer, Big Data]	[Mexico]	[Design, implement, and take ownership of miss...
5	[Director, IT Product Management]	[Oakland, CA, Long Beach, CA, El Dorado Hills,...	[Strategically define and design technical sol...
6	[Senior Manager, Machine Learning Engineering]	[US National]	[Lead team of engineers in building scalable L...
7	[Customer Data Platform Technologist]	[US National]	[Play a key role in configuring, integrating, ...
8	[Business Systems Analyst II - Remote]	[Raleigh, NC]	[Perform business systems analysis for Core ba...
9	[Senior Staff Technical Product Manager]	[MA, or Work from Anywhere]	[Responsible for leading, implementing, and ma...
10	[Intune Administrator]	[US National]	[An SCCM/Intune Administrator is needed for a ...
11	[Lead Configuration Management Specialist]	[Chantilly, VA]	[Responsible for developing and maintaining co...
12	[Manager, Infrastructure Core Engineering]	[Dublin, OH]	[Manage a team of Systems Engineers and Admini...
13	[Analytics Specialist - Claims Litigation - Le...]	[Madison, WI]	[Analyze business results and external market ...
14	[Data Engineer III]	[Dublin, OH]	[Design and build new data integrations. Work ...
15	[Lead Data Engineer]	[Dublin, OH]	[Lead and deliver enterprise projects for the ...

Strona Flexjobs jest o wiele bardziej rozbudowana pod względem różnorodności dostępnych ofert pracy. Ja wybrałem zawody z branży IT jak widać po kodzie URL. Możemy jednak łatwo to zmienić przechodząc na stronie do zakładki „Find Remote, Hybrid & Flexible Jobs”



Aby następnie wybrać kierunek, który nas interesuje.



All Job Categories

Account Management Jobs	▼	Environmental & Green Jobs	Nonprofit & Philanthropy Jobs	▼
Accounting & Finance Jobs	▼	Event Planning Jobs	Operations Jobs	
Administrative Jobs	▼	Fashion & Beauty Jobs	Product Jobs	
Advertising & PR Jobs	▼	Food & Beverage Jobs	Project Management Jobs	
Animals & Wildlife Jobs		Government & Politics Jobs	Research Jobs	▼
Art & Creative Jobs	▼	Graphic Design Jobs	Retail Jobs	
Bilingual Jobs	▼	HR & Recruiting Jobs	Sales Jobs	▼
Business Development Jobs		Human Services Jobs	Science Jobs	▼
Call Center Jobs	▼	Insurance Jobs	Software Development Jobs	▼
Communications Jobs		International Jobs	Sports & Fitness Jobs	▼
Computer & IT Jobs	▼	Internet & Ecommerce Jobs	Telemarketing Jobs	
Consulting Jobs	▼	Legal Jobs	Transcription Jobs	
Customer Service Jobs		Manufacturing Jobs	Translation Jobs	
Data Entry Jobs		Marketing Jobs	Travel & Hospitality Jobs	
Editing Jobs	▼	Math & Economics Jobs	Web Design Jobs	

Po kliknięciu w wybrany obszar naszego zainteresowania, kopiujemy kod URL i zamieniamy go w naszym programie.

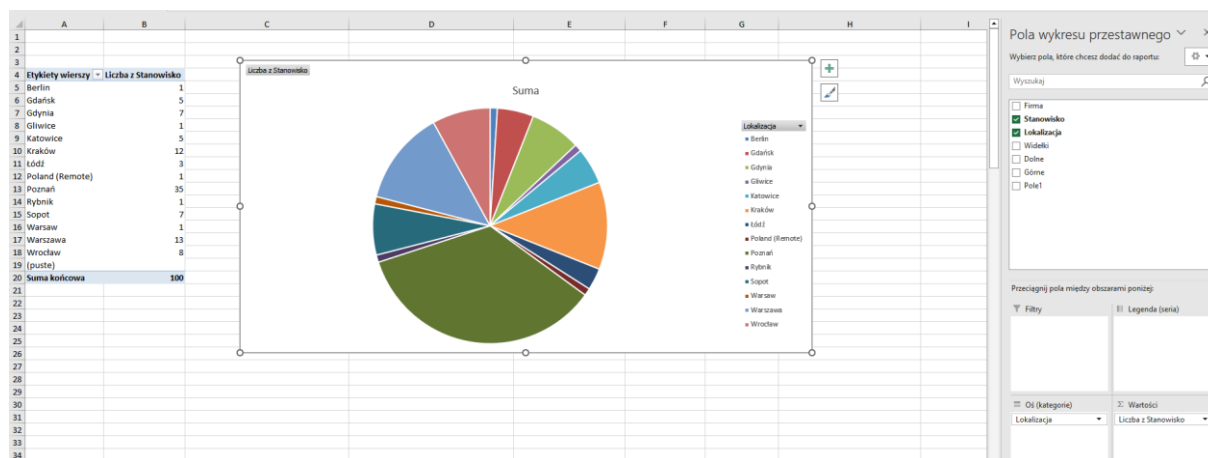


4. Analiza zebranych danych

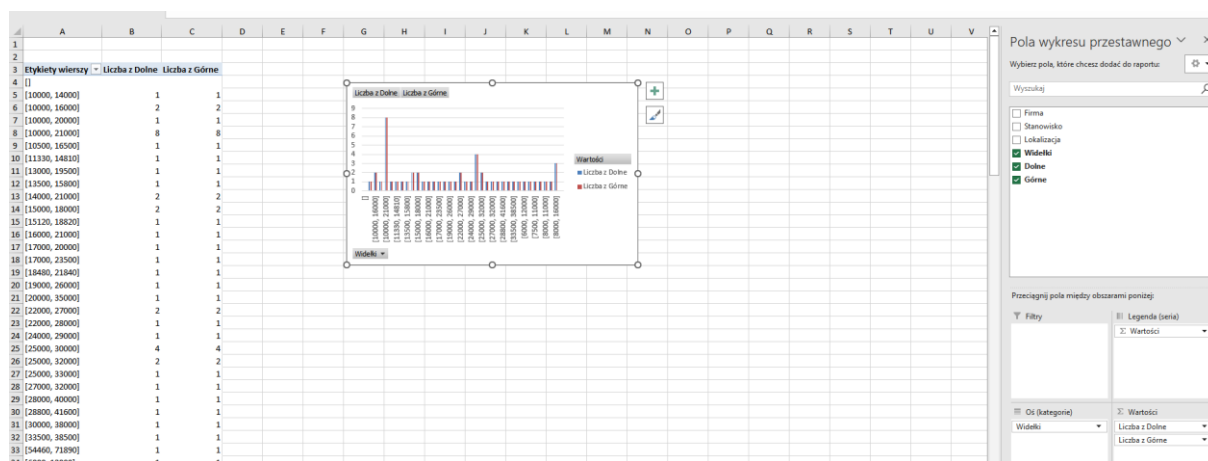
Analizę przeprowadzę na 100 próbkach uzyskanych ze strony JustJoinIT. Oto jak przedstawiają się dane w arkuszu Excel. Widełki podzieliłem na dwie górne i dolne.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
	Firma	Stanowisko	Lokalizacja	Widełki	Górne	Dolne									
1	Lufthansa Systems	Integration Systems Engineer (GCP)	Gdynia	[10000, 16000]	10000	16000									
2	Code Poets	Developer C++	Wrocław	[17000, 23500]	17000	23500									
3	Capgemini Polska	Database Administrator	Kraków	[]											
4	GSK	Veeva CRM Test Lead	Poznań	[]											
5	GSK	Senior Engineer ServiceNow Integrations	Poznań	[]											
6	GSK	Compute CaaS Solution Architect DR and BusC	Poznań	[]											
7	Qodeca	Graphic Designer	Poland (Remote)	[8000, 10000]	8000	10000									
8	GSK	Compute CaaS Solution Architect DR and BusC	Poznań	[]											
9	Antal	Principal SAP Finance Solution Architect (Capex)	Warszawa	[]											
10	GSK	Digital Engineering Architect	Poznań	[]											
11	GSK	Digital Engineering Architect	Poznań	[]											
12	GSK	Compute CaaS Solution Architect DR and BusC	Poznań	[]											
13	Insys Video Technologies	Technical Writer	Warszawa	[7500, 11000]	7500	11000									
14	ARCHE CONSULTING	DCS System & Cyber Security Specialist	Gliwice	[10000, 14000]	10000	14000									
15	Lufthansa Systems	Java and Reactive Programming	Sopot	[10000, 21000]	10000	21000									
16	GSK	Veeva CRM Test Lead	Poznań	[]											
17	GSK	Principle DevOps Engineer	Poznań	[]											
18	Lufthansa Systems	Application Security Engineer (DH)	Gdynia	[10000, 21000]	10000	21000									
19	Insys Video Technologies	Junior Backend Developer (.net)	Łódź	[6700, 10700]	6700	10700									
20	GSK	Compute CaaS Solution Architect DR and BusC	Poznań	[]											
21	Nortal	Technology Team Lead	Kraków	[25000, 30000]	25000	30000									
22	GSK	Principle DevOps Engineer	Poznań	[]											
23															

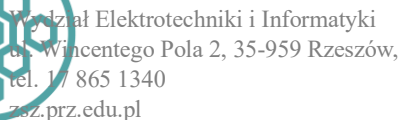
Następnie nasze dane przedstawiłem w postaci tabeli przestawnej (PowerPivot), aby analiza była dobrze widoczna. Na samym początku sprawdzimy z jakiego miasta jest najwięcej ofert.



Przy pomocy wykresu kołowego, możemy zauważyć, że najwięcej ofert z aktualnie zscrapowanych pochodzi z Poznania, następna jest Warszawa i na trzecim miejscu Kraków. Popatrzmy na widełki



Ofert z widełkami na poziomie 10000 – 21000 PLN jest najwięcej.



Najwięcej zaś można zarobić w Warszawie. Największa górna granica zarobków znalazła się właśnie w Warszawie i jest na poziomie 71890 PLN miesięcznie na stanowisku FullStack Engineer w firmie BlockLabs .Najmniejsze widełki również zaczynają się w Warszawie na poziomie 6000 PLN na stanowisku Administratora Systemów IT w Locon Sp. z o.o..

[illegible]



Podsumowanie i wnioski

Projekt zakładał zebranie danych ze stron oferujących zatrudnienie w różnych branżach. Udało mi się zscrapować i przeanalizować zebrane przeze mnie dane z trzech stron. Technologia w jakiej pracowałem o Python oraz Excel.

Po analizie ofert możemy łatwo wywnioskować, że dobrze płatną pracę w IT możemy znaleźć najłatwiej w stolicy w Warszawie oraz największych miastach Polski. Zdecydowałem się na branżę IT z powodu, że mnie najbardziej interesowała. Technologia jaką wykorzystałem niestety nie potrafi scrapować stron, w których zostało użyte bardzo dużo JavaScript, co zawęży zakres stron. Technika scrapowania jest bardzo przydatna, jest używana przez np. Urząd Statystyczny w Rzeszowie, co zostało nam pokazane na laboratoriach z pracownikiem tego urzędu.



Źródła

- [regex101: build, test, and debug regex](#)
- <https://nofluffjobs.com/>
- justjoin.it
- [Computer & IT Jobs - Remote Work From Home & Online | FlexJobs](#)
- <https://chat.openai.com/>