



Міністерство освіти і науки України
Національний технічний університет України
«Київський політехнічний інститут імені Ігоря Сікорського»
Фізико-технічний інститут

Лабораторна робота №4
з дисципліни
«Web - аналітика»

Виконав:
студент групи ФБ-31мп
Щур Павло
Перевірив:
Ткач В. М.

Київ-2024

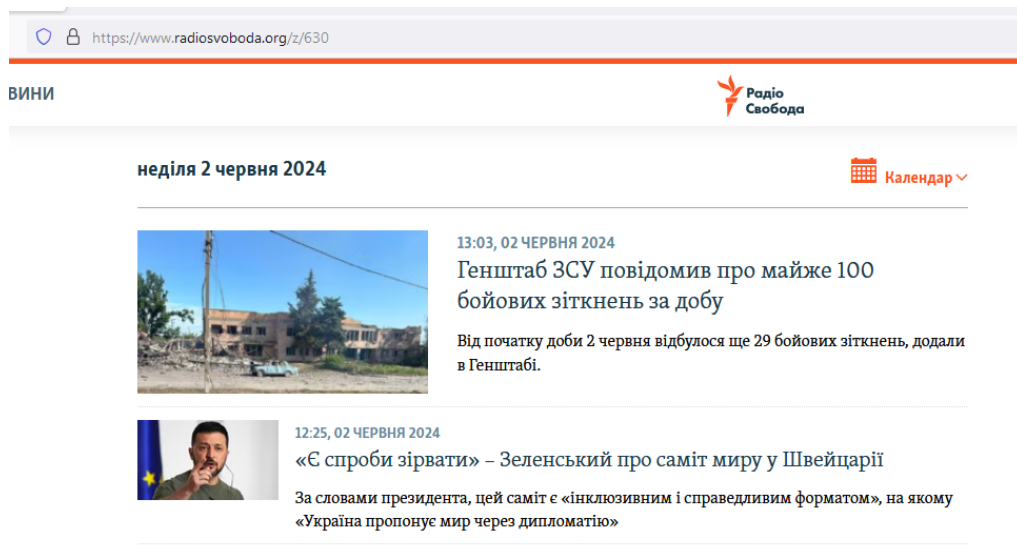
Посилання на GitHub: <https://github.com/ShchurPavlo/web-analytics-2024/tree/main/lab4>

Завдання

1. Збір великих даних із новинних ресурсів. Формування датасетів. Можна один ресурс не менше 1000 записів.

Виконання:

1) В якості піддослідного новинного ресурсу будемо використовувати сайт Радіо Свобода - <https://www.radiosvoboda.org/z/630>



Для збору даних використаємо бібліотеку для Python BeautifulSoup, а для відкриття веб-сторінки (обходу перевірки на бота) використаємо Selenium та його Chrome Driver.

2) Відкривши сторінку в режимі розробника знайдемо клас `col-xs-12 col-sm-12 col-md-12 col-lg-12 fui-grid__inner` в якому міститься інформація про новини:

```
</ul>
<ul id="ordinaryItems">
<li class="col-xs-12 col-sm-12 col-md-12 col-lg-12 fui-grid__inner">
<div class="media-block">
<a href="/a/news-zelenskyj-samit-myru-sproby-zirvaty/32975840.html" class="img-wrap img-wrap--t-spac img-wrap--size-3 img-wrap--float img-wrap--xs">
<div class="thumb thumb16_9">
<noscript class="nojs-img">

</noscript>

</div>
</a>
<div class="media-block_content media-block_content--h media-block_content--h-xs">
<span class="date date--mb date--size-3">12:25, 02 червня 2024</span>
<a href="/a/news-zelenskyj-samit-myru-sproby-zirvaty/32975840.html">
<h4 class="media-block_title media-block_title--size-3" title="&#171;Є спроби зірвати&#187; – Зеленський про саміт миру у Швейцарії">
&#171;Є спроби зірвати&#187; – Зеленський про саміт миру у Швейцарії
</h4>
<p class="perex perex--mb perex--size-3">За словами президента, цей саміт є &#171;інклюзивним і справедливим форматом&#187;, на якому &#171;Україна і
```

3) Реалізуємо окрему функцію для отримання html розмітки сторінки з новинами:

```
9  def Get_soup(url,driver):
10      driver.get(url.replace(' ', ''))
11      time.sleep(5)
12      source_data = driver.page_source
13      soup = BeautifulSoup(source_data, features="html.parser")
14      return soup
```

4) Також з уже отриманої розмітки зможемо отримати інформацію, для цього теж реалізуємо окрему функцію:

```
16 def Parse_data(soup):
17     news_blocks = soup.find_all('li', class_='col-xs-12 col-sm-12 col-md-12 col-lg-12 fui-grid__inner')
18     result = []
19     for block in news_blocks:
20         title_tag = block.find('h4', class_='media-block__title')
21         perex_tag = block.find('p', class_='perex perex--mb perex--size-3')
22         date_tag = block.find('span', class_='date date--mb date--size-3')
23         link_tag = block.find('a', href=True)
24
25         title = title_tag.get_text(strip=True) if title_tag else None
26         perex = perex_tag.get_text(strip=True) if perex_tag else None
27         date = date_tag.get_text(strip=True) if date_tag else None
28         link = link_tag['href'] if link_tag else None
29         result.append({
30             'title': title,
31             'perex': perex,
32             'date': date,
33             'link': 'https://www.radiosvoboda.org'+link
34         })
35
36     return pd.DataFrame(result)
```

Будемо збирати заступні дані:

- Заголовок
- Короткий опис
- Дату публікації
- Посилання на сторінку з новиною

5) Функція яка приймає в якості аргументів місьць та рік та повертає pandas датафрейм з новинами за заданий період часу:

```
39 def Get_data(year, month):
40     result_df = pd.DataFrame()
41     for month in month:
42         for day in range(1, 31):
43             try:
44                 print('Parse date: '+year+'.'+month+'.'+str(day))
45                 soup = Get_soup(url=f"https://www.radiosvoboda.org/z/630/{year}/{month}/{day}", driver)
46                 result_df = pd.concat([result_df, Parse_data(soup)], ignore_index=True)
47             except Exception as e:
48                 print(f"Error occurred for {month}/{day}: {e}")
49     return result_df
```

Результат роботи:

```
Parse date: 2024.4.1
Parse date: 2024.4.2
Parse date: 2024.4.3
Parse date: 2024.4.4
Parse date: 2024.4.5
Parse date: 2024.4.6
Parse date: 2024.4.7
Parse date: 2024.4.8
Parse date: 2024.4.9
Parse date: 2024.4.10
```

Отриманий датафрейм збережемо в csv файл:

```
51 result=Get_data( year: '2024', month: ['4','5'])
52 result.to_csv( path_or_buf: 'news_df.csv', index=False, encoding='utf-8')
```



news_df.csv

02.06.2024 13:41

В результаті за період квітень-травень 2024р було отримано 1741 запис:

	A	B	C	D	E	F
1717	Уряд Словенії визнав Палестинську державу, рішення потребує схвалення Рішення уряду Словенії	22.09, 30 травня 202	https://www.radiosvoboda.org/a/news-uniad-slovenii-palestynska-derzhava/32972893.html			
1718	«Масмо підтвердження щодо збільшення наших спроможностей ППО» - «Буде найближчим часом	21:55, 30 травня 202	https://www.radiosvoboda.org/a/news-zelenykyi-protypovityryana-oborona/32972888.html			
1719	У Керченському порту, ймовірно, пришвартований пором, про пошкоджені Генштаб ЗСУ заявив, що	21:44, 30 травня 202	https://www.radiosvoboda.org/a/news-krum-porom/32972880.html			
1720	Міністр оборони Німеччини приїхав в Україну і оголосив про новий пакет Деталі візиту німецького	21:29, 30 травня 202	https://www.radiosvoboda.org/a/news-nimechchyna-pistorius-dopomoha/32972872.html			
1721	Представників США не буде на церемонії в пам'ять про Paici в ООН Плани провести церемоні	21:17, 30 травня 202	https://www.radiosvoboda.org/a/news-ssha-onn-raisi/32972864.html			
1722	Українські військові заявили про ураження ракетами АТАСМС поромної і У Генштабі кажуть, що ви	21:03, 30 травня 202	https://www.radiosvoboda.org/a/news-atacms-krum-porom/32972859.html			
1723	«Укренерго» вже четвертий день поспіль не планує вимкнення світла в Уї 31 травня відключення сі	20:46, 30 травня 202	https://www.radiosvoboda.org/a/news-ukrenhergo-svitlo-vidklyuchennya/32972847.html			
1724	Генштаб: на фронті за день було 84 зіткнення, половина – на двох напрямках Курхавськомому і Покр	20:34, 30 травня 202	https://www.radiosvoboda.org/a/news-henshtab-zvedennia/32972840.html			
1725	Зеленський: Україна уклала вже 12 безпекових угод на понад 23 млрд дол. Напередодні були уклад	20:25, 30 травня 202	https://www.radiosvoboda.org/a/news-zelenykyi-bezpekovi-uhody/32972834.html			
1726	Єврокомісія виступає за присутність у марафоні ЗМІ, які «виражають по В органі виконавчої влад	20:22, 30 травня 202	https://www.radiosvoboda.org/a/news-eurokomisia-telemarafon/32972833.html			
1727	Воєнжори Радіо Свобода стали лауреатами премії «Честь професії» Іхній матеріал «ЗСУ пішли	19:54, 30 травня 202	https://www.radiosvoboda.org/a/news-chest-profesii-kushnir-nuzhnenko/32972814.html			
1728	Двоє людей поранені через обстріл Руської Лозової на Харківщині – ОБВ Спочатку в ОБА повідом	19:16, 30 травня 202	https://www.radiosvoboda.org/a/news-kharkivshchyna-poranieni/32972686.html			
1729	У ГУР уточнили, що уразили чотири російські катери в Криму, з них два - «Згідно з уточненими да	19:02, 30 травня 202	https://www.radiosvoboda.org/a/news-hur-tunets-krum/32972621.html			
1730	Шрі-Ланка посилює контроль, щоб не допустити використання своїх грош Коломбо відправити дел	18:41, 30 травня 202	https://www.radiosvoboda.org/a/news-shri-lanka-rosia-viyna/32972603.html			
1731	Латвія не обговорює присутність своїх військ в Україні – голова МЗС «Я не спекулюю. Наразі і	18:22, 30 травня 202	https://www.radiosvoboda.org/a/news-latvia-ukraina-viyska/32972581.html			
1732	Литва приєдналася до коаліції з надання ППО Україні – Науседа «Україні потрібна зброя, і	17:57, 30 травня 202	https://www.radiosvoboda.org/a/news-lyuva-koalitsia-ppo/32972551.html			
1733	У МЗС Росії заявили, що готують відповідь на обмеження для російських Заходи у відпові	17:34, 30 травня 202	https://www.radiosvoboda.org/a/news-polshcha-rosia-diplomaty/32972521.html			
1734	В Ірані відкрили реєстрацію кандидатів на президентських виборах Період реєстрації триват	17:18, 30 травня 202	https://www.radiosvoboda.org/a/news-iran-reestratsiya-vybory-prezidenta/32972505.html			
1735	Неурядові організації Грузії планують оскаржувати закон про «іноагентів Кілька НУО оголосили, ш	17:06, 30 травня 202	https://www.radiosvoboda.org/a/news-hruzia-neuriadovi-orhanizatsii-inoahenty/32972369.html			
1736	Чеська ініціатива забезпечуватиме Україну десятками тисяч снарядів щ «Чехія також усеюдило	16:45, 30 травня 202	https://www.radiosvoboda.org/a/news-chekhia-ukraina-boepripsy/32972349.html			
1737	Генштаб: інтенсивні бої на Покровському напрямку тривають, найгарячіш Одну з двох атак Сили о	16:34, 30 травня 202	https://www.radiosvoboda.org/a/news-henshtab-pokrovskiy-napriamok-boii/32972339.html			
1738	Україна все ще може перемогти, але з допомогою союзників – генсекрет «Україна продовжує вою	16:27, 30 травня 202	https://www.radiosvoboda.org/a/news-stoltenberg-ukraina-viyna/32972329.html			
1739	Данія не проти застосування F-16 для ударів по військових цілях у Росії «Ми про це дуже чітко за	16:04, 30 травня 202	https://www.radiosvoboda.org/a/news-daniia-f-16-viyskovi-tsili-rosii/32972304.html			
1740	Рятувальники гасять 15 лісових пожеж на Харківщині, більшість із них сг За повідомленням ДСНС	15:30, 30 травня 202	https://www.radiosvoboda.org/a/news-lisovi-pozhezhi-kharkivshchyna-rosiiski-obstrily/32972265.html			
1741	Голова МЗС Чехії не бачить наразі перспективи надсилання чеських вій «Ми масмо захист нашог	15:03, 30 травня 202	https://www.radiosvoboda.org/a/news-mzs-chehii-perspektiva-nadsilannia-cheskyh-viiskovyh/32972222			