

# **Wykop.pl – Webscrapping using BS4, Selenium, Scrapy**

## 1. Introduction

The site **wykop.pl** is regarded (and rightly so) as a forum where an atmosphere of heckling and name-calling prevails. Nevertheless, the main page of wykop contains articles on current news from Poland and the whole World. The site is updated practically every minute, so by spending 15-20 minutes at evening, you can get reliable news of the day. In our work, we decided to use popular webscrapping libraries, so that we could check what topics currently dominate on the site, what hashtags are the most used and what people think about each article based on their reactions.

Unfortunately, due to the restrictive policy of the wykop.pl website, scraping a very large number of articles is impossible. **Therefore, the scope of our analysis will be limited to scraping 120 pages of the site, where from each subpage we will only select the first article that is on the page.** Once we are inside the article, we try to scrape the following things from it:

- title
- the nickname of the user who posted the article
- the number of likes under the article
- the number of dislikes
- number of views
- all the hashtags that have been placed under the article

We would like to **warn you in advance** that the final number of records you receive may vary slightly from 120, but will certainly be higher than the lower limit required to pass the course (i.e. 100). We have discovered two reasons for this:

- the first article may be an advertisement, which when clicked on, automatically leaves wykop.pl and goes to the advertiser's website
- the first article, may also contain content allowed only for persons over 18 years of age, which requires a prior account creation and proof of age.

Furthermore, the site is **updated on a regular basis**, so running through the codes at different time intervals will **undoubtedly give us different results**. Nevertheless, below is a brief roadmap of how our tool should work in general.

wykop.pl Wykopaliszko 446 Hity Mikroblog

Ciekawostki Polska Covid19 Nauka Historia Rozrywka Sport Motoryzacja Świat Giełda Programowanie Technologia

STRONA GŁÓWNA: najnowsze aktywne ulubione

55 Piosenka z ruin Azovstału....  
Poleca maro-grzechotnik wykopowe uniwsko

138 wykop SzeF Volkswagena chce żeby wojna się skończyła bo interes mu źle idzie .  
@125procent forsal.pl #niemcy #ukraina +1 inny  
Herbert Diess wezwał Ukrainę do wynegocjowania porozumienia z Moskwą, dzięki czemu będzie można znieść sankcje na Rosję, a niemiecka gospodarka uniknie szkód. Kontrowersyjne słowa szefa Volkswagena spotkały się z natychmiastową reakcją ukraińskich władz.  
18 komentarzy opublikowany 3 min. temu

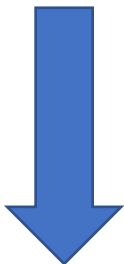
136 wykop Zelenski: Wojna się skończy, gdy Ukraina odzyska wszystkie swoje terytoria!  
@Bryzoll twitter.com #swiat #wojna +2 inne  
Prezydent Wołodymyr Zełenski powiedział, że jego chęć kontynuowania negocjacji maleje z dnia na dzień po zobaczeniu dowodów masakr i okrucieństw popełnionych przez Rosjan w całym kraju.  
28 komentarzy opublikowany 9 min. temu

139 wykop Prywatny właściciel ogrodził plotem jezioro na Kaszubach.  
@xallemorph onet.pl #polska #ekologia +2 inne  
Nowy właściciel jeziora Wielkie Oczko ogrodził je ponaddwumetrowym plotem. Jezioro to znajduje się na terenie Wdzydzkiego Parku Krajobrazowego i obszaru Natura 2000. Nadzór budowlany nie widzi jednak problemu, gdyż plot nie przekracza wysokości 2,2 m, a więc nie wymaga zgody.  
57 komentarzy opublikowany 14 min. temu

117 wykop Płonie niemal 800 hektarów. Ukraińcy mają problem z gaszeniem pożarów lasów  
@szyna352 miz4.pl #europa #ukraina +3 inne  
Prawie 800 hektarów lasów płonie w obwodzie chersońskim na południu Ukrainy. Walkę z żywiołem utrudnia okupacja tych terenów oraz działania bojowe – poinformował w środę doradca ministra ds. ochrony środowiska i zasobów naturalnych Serhij Własenko.  
7 komentarzy opublikowany 18 min. temu

Select only the first article from each subpage.

From both red boxes, scrape those things that will be useful for further analysis.



wykop.pl Wykopaliszko 447 Hity Mikroblog Szukaj Dodaj Zaloguj się

Wojna za naszą wschodnią granicą toczy się również w sieci. Atak Rosji na Ukrainę wywołał lawinę artykułów, informacji i komentarzy. Nie wszystkie doniesienia są jednak prawdziwe, część z nich ma za zadanie wywołać chaos i niepewność. Nie pozwólmy na to. Starajmy się wspólnie o wysoki poziom społecznej weryfikacji treści na naszym portalu. Sprawdzajmy i weryfikujmy informacje w rzetelnych oraz oficjalnych źródłach. Wszelkie próby rozpowszechniania informacji nieprawdziwych raportujemy moderacji.

182 wykop SzeF Volkswagena chce żeby wojna się skończyła bo interes mu źle idzie .  
@125procent forsal.pl #niemcy #ukraina #rosja  
Herbert Diess wezwał Ukrainę do wynegocjowania porozumienia z Moskwą, dzięki czemu będzie można znieść sankcje na Rosję, a niemiecka gospodarka uniknie szkód. Kontrowersyjne słowa szefa Volkswagena spotkały się z natychmiastową reakcją ukraińskich władz.

Dodany 1 godz. temu przez:  
125procent dołączył 12 lat 8 mies. temu  
182 wykopał 6 zakopał <1 tys. wyświetleń

OSTATNIO POPULARNE  
Warszawa: ambasador Rosji obłany farbą. Bójka przed cmentarzem 4234  
Mam dość monopolu Orlenu i wysokich cen paliw! 4158

## 2. Description of coding parts

### 2.1. BeautifulSoup

[Script name BS\_final.py.]

At the very beginning we put two parameters to change. The first one *START\_100\_SCRAPE* is boolean and setting it to True, will make our scraper start retrieving information from the next 120 pages. In case of parameter *SET\_PAGES\_TO\_SCRAPE\_PER\_ONE\_LOOP*, we can control how many links to articles our scraper should scrape. Default parameter is equal to 1, but it can be set to a maximum of 35. However, setting this parameter so high would result in 4200 requests, which would cause our connection to be broken by the anti-scraping program.

Our scraper starts by creating an empty list where we will store our results. Next, we create a loop which goes through the following subpages of the *wykop.pl/strona/{number}*. Inside the function we save to the *links* variable all articles found in a given tab. After that, we are looping over the values of the *links* variable which opens a link to each article from the subpage. Once we are inside the article, we try to scrape our 6 variables of interest.

Finally, the results are stored in the *post\_summary* dictionary, which is then added to the *list\_of\_result*, containing all the results scraped so far. We also include a 10 second pause between subpage changes.

Below, we've included an example of the result for our scraper. In order to demonstrate how our artificial log works, we have changed the range of scraped pages to 2 on one occasion and the range of scraped articles per page to 10.

Figure 1. Example output from tool based on beautifulsoup library.

```
ARTICLE 19 FROM PAGE 1 CONTAINS CONTENT ONLY FOR ADULTS. CANNOT BE SCRAPPED.
CURRENT NON-SCRAPPED ARTICLES: 1

Article: 1 from subpage: 1 has been successfully scraped !
Article: 2, SUBPAGE: 1 - ACCESS DENIED: THIS IS ADVERTISEMENT !!!
Article: 3 from subpage: 1 has been successfully scraped !
Article: 4 from subpage: 1 has been successfully scraped !
Article: 5 from subpage: 1 has been successfully scraped !
Article: 6 from subpage: 1 has been successfully scraped !
Article: 7 from subpage: 1 has been successfully scraped !
Article: 8 from subpage: 1 has been successfully scraped !
Article: 9 from subpage: 1 has been successfully scraped !
Article: 10 from subpage: 1 has been successfully scraped !
The subpage will be changed in 2 seconds.

Article: 1 from subpage: 2 has been successfully scraped !
Article: 2, SUBPAGE: 2 - ACCESS DENIED: THIS IS ADVERTISEMENT !!!
Article: 3 from subpage: 2 has been successfully scraped !
Article: 4 from subpage: 2 has been successfully scraped !
Article: 5 from subpage: 2 has been successfully scraped !
Article: 6 from subpage: 2 has been successfully scraped !
Article: 7 from subpage: 2 has been successfully scraped !
Article: 8 from subpage: 2 has been successfully scraped !
Article: 9 from subpage: 2 has been successfully scraped !
Article: 10 from subpage: 2 has been successfully scraped !
The subpage will be changed in 2 seconds.

Below, you can find your final output, number of scraped articles is equal to 18
[{'title': 'Prezes Głapiński: od polskiego "cudu gospodarczego" do galopującej inflacji', 'usern
```

## 2.2. Selenium

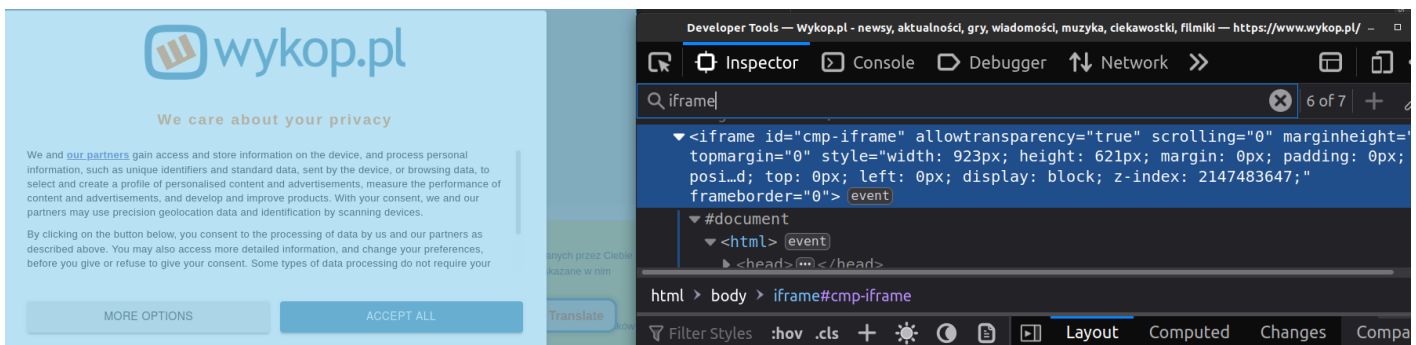
[Script name selenium\_wykop.py.]

As with scrapping through the beautiful soup, we have entered the same two launch parameters (boolean: *START\_100\_SCRAPE*, integer: *SET\_PAGES\_TO\_SCRAPE\_PER\_ONE\_LOOP*) that define the number of scrapped pages and articles.

Initially, when we start firefox geckodriver, we get a pop-up about accepting cookies. Acceptance of this, is important because the window is in a special "iframe". In a nutshell, the Iframe HTML element (<iframe>) represents a nested browsing context (en-US) by embedding another HTML page in the current one.

In our case, if you don't click on anything, the page doesn't allow you to continue browsing the site. This required us to change the format of the page to the "frame" structure, accept cookies and then return to the default character of the page (driver.switch\_to.default\_content()). After this step, we additionally accept privacy politics, located at the bottom of the page.

Figure 2. Wykop's iframe element.



Moving on to the scraping process itself, we have created a loop within a loop that goes through each article from a given tab and initially checks if the article is an advertisement (variable *check\_adv\_txt*). If not, the scraping process begins and the same things already mentioned are retrieved, otherwise the article will be skipped. The results obtained, are written to the *post\_summary* dictionary and then added to the *list\_of\_result*. After a 2-second pause, our tool returns to the previous page. If the limit of scraped articles for a tab is reached, after 5 seconds, the program automatically clicks the "next" button and the user is redirected to the next page.

## 2.3. Scrapy

[Script name scrapy\_wykop.py.]

Unlike the previous libraries, with scrapy we leave only one boolean control parameter which controls the number of scraped pages.

Firstly, we instantiate the *Wykop* class, which defines which parameters from the page we expect to be scraped. Next, we create another class - *WykopSpider*, which will send the requests to subsequent pages of the wykop.pl/strona/ website.

Again, when we open each subpage, the program will select only the FIRST ARTICLE that appears on the page and then extract the information that we are interested in. Obviously, we are aware that this first article may be an advertisement or an 18+ link (which will not be scraped). That's why, the whole loop has been encapsulated in the try function that will return an empty field if the data cannot be scraped. After calling the appropriate command (scrapy crawl wykop -o wykop.csv), we should get a csv file in our working directory, which in a later stage of our work was used for a short data analysis.

Below, we show the first 5 records of our exemplary desired output.

Figure 3. Scrapy csv output.

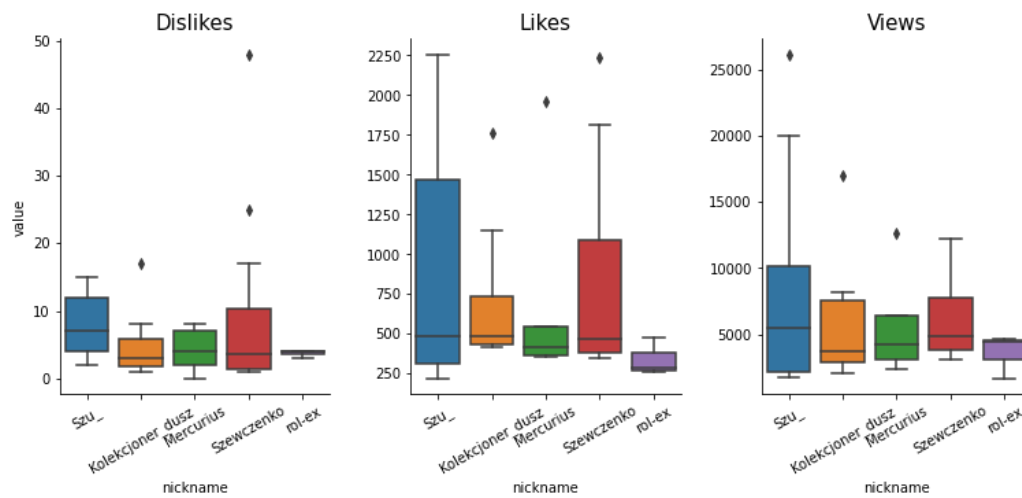
	dislikes	hashtags	likes	nickname	title	views
0	0	technologia,nauka,zainteresowania,kosmos,mars,...	186	elektryk91	NASA ma problem z Ingenuity, marsjańskim wirol...	2.4 tys.
1	0	rosja,historia,polityka,zbrodniekomunistyczne	86	Qtasus2Wielgus	Badacz zbrodni NKWD został wysłany do łagru	<1 tys.
2	14	polska,chlewoobsranygownem	269	BennyHarvey	Łódzkie. Kradł znicze wrzucając do automatu st...	12.5 tys.
3	23	polska,kultura,rozrywka	359	Wirtualnedia_pl	Polsat zrealizuje polską wersję „Sługi Narodu”	7.9 tys.
4	39	polska,technologia,amazon,afery	476	Mirxar	Strach coś zamawiać na Amazonie. Polacy płacą ...	13.4 tys.

## 3. Exploratory Data Analysis

[Script name wykop\_analysis.ipynb]

For our analysis we used 120 articles which are the first articles from 120 pages (as of 11.05.2022, 21.30). First step, aims to show the distribution of dislikes, likes and views per top 5 most frequent users.

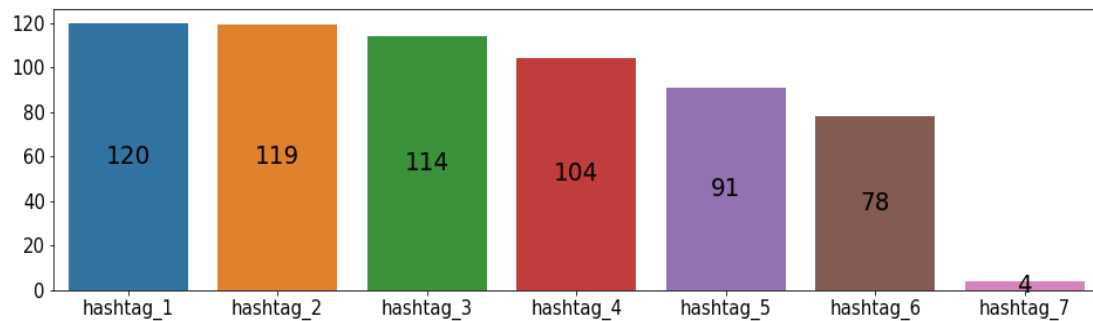
Figure 4. Boxplots for top 5 users.



From the charts above, we can see that the posts of two users, stand out from the others in terms of the number of responses they received. However, for all users, the median scores across all variables are similar.

The second step was to check the frequency of hashtag usage per article.

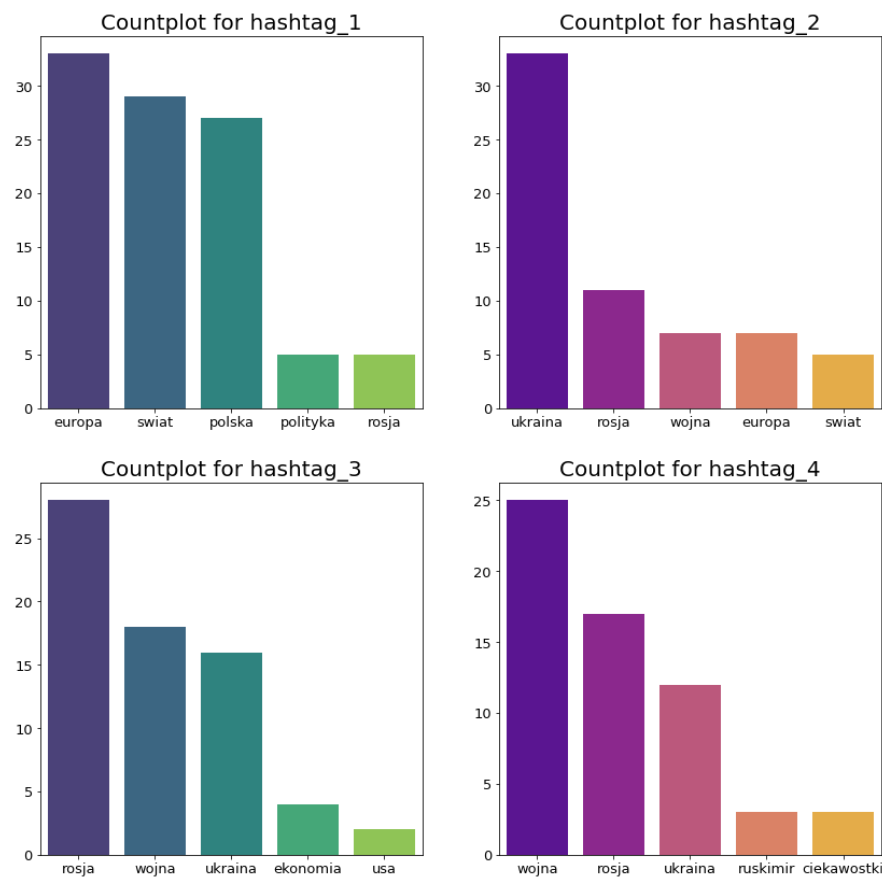
Figure 5. Distribution of hashtags used.



From the graph above, we can see that at least one hashtag was used for all articles. We can therefore assume that this is a restriction of the portal. Interestingly, for the scraped articles, the average hashtag used is around 5.

In order to find out what topics are currently on trend, we checked the relevance of the top 5 words for 4 hashtags.

Figure 6. Top five hashtags.



From the graph above we can see that the articles were mostly about Europe and in particular the war between Russia and Ukraine.

#### 4. Conclusions

In our work, we collected data on articles that are found on wykop.pl. For this purpose, we used three different libraries. We believe that the aim of the work, i.e. to retrieve information from at least 100 consecutive pages, has been achieved. The quality of the data is good enough that we were able to carry out a short analysis of the data, which allowed us to find out what news from the world are currently on top.

Should you wish to continue the work, we suggest that you also download descriptions for each article, which could be used to build an NLP model. Another step worth considering is to try to circumvent the anti-scrap system of the wykop.pl website (in the good spirit of the law, of course). This would make it possible, without any restrictions, to scrape every article from the site.