

## DYNAMIC WEB SCRAPING THROUGH SELENIUM WEBDRIVER IN PYTHON

**Aneesha Gudavalli**

M.Tech Scholar, Department of Information Technology, V R Siddhartha Engineering College,  
Vijayawada, Andhra Pradesh, India, gudavallianeesha@gmail.com

**G. JayaLakshmi**

Assistant Professor, Department of Information Technology, V R Siddhartha Engineering College,  
Vijayawada, Andhra Pradesh, India,  
jaya1123@vrsiddahrtha.ac.in

### **Abstract:**

Machine learning is used to power today's technological marvels, such as self-driving cars, picture and speech recognition, and space transportation. To build a stable & dependable machine learning model for such business challenges, Data Science experts would require a large amount of data. The first step in the data science life cycle is data mining or data collection. Depending on the business needs, data might be collected through SAP servers, logs, databases, APIs, online repositories, or the web. Web scraping, often known as "crawling" or "spidering," is a method of extracting information from an online source, most typically a webpage. Crawling the web is one of the methods for extracting large volumes of data in a short amount of time. This paper proposes a way of scraping data from numerous dynamic websites or web pages via the internet using Selenium WebDriver and Python. Selenium is capable of scraping large amounts of data in a short amount of time, such as text and images.

**Keywords:** data scraping, web scraping, selenium, selenium WebDriver, dynamic webpage scraping.

### **I. INTRODUCTION**

On the internet, there is a significant amount of data that can be used to meet business objectives. To get this information from the web, one requires a tool or technique. This is where the concept of web-scraping comes into play. Web scraping [7] can help us collect large amounts of data about customers, products, people, stock markets, and other topics. Customer buying trends, employee churn, customer attitudes, and other factors are all evaluated, data acquired from websites like as e-commerce portals, job portals, and social media channels can be used. Sites are a huge expanse of information that everybody can get to. Innovation's recent fad has constrained us to adjust our organization rehearses. Data crawling, also known as online scraping or data harvesting, has existed since the beginning of the internet. Although it is now synonymous with web content extraction, it wasn't originally used for that. Selenium is a web application testing framework that is portable. It's free software that runs on Windows, Linux, and Mac OS X and is licensed under the Apache License 2.0. Selenium is also used as a web scraping tool, in addition to its primary purpose. We won't go over all of Selenium's components; instead, we'll concentrate on WebDriver shown in fig1, a single component essential for web scraping. Selenium WebDriver allows us to develop and execute test cases by controlling a web browser through a programming interface. When websites display dynamic content, Selenium comes very helpful. Despite the fact that Scrapy is a sophisticated online scraping technology, it is rendered useless when dealing with dynamic websites.

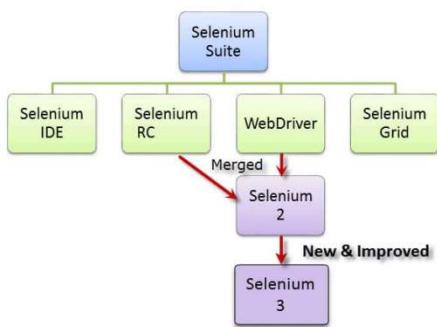


**Fig1: Selenium Webdriver**

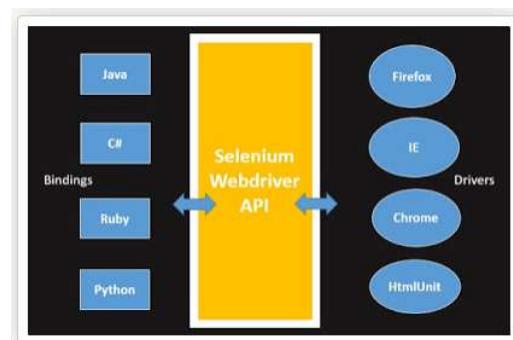
To get rendered HTML content from dynamic pages, Selenium web drivers are used. Selenium is a Python library that allows you to run a web driver from a remote computer. To browse to the given data source, the selenium package's driver.get method is used. The drivers (such as ChromeDriver and FirefoxDriver) web browsers are directed to deliver commands in question and receive the response. To receive reliable information, user input is sometimes required, such as clicking buttons and picking relevant alternatives from dropdown menus, which Selenium supports effectively.

## SELENIUM SUITE AND WEBDRIVER ARCHITECTURE

Selenium is an open source program automation tool that is commonly used for web application testing. It automates cyclic tasks by computerizing the control of an internet browser. Selenium is a collection of testing tools that may be used to test a variety of programs, frameworks, and languages such as C#, Java, Ruby, and Python. Figure 2 depicts the four-part selenium system. Selenium IDE, for example, is a Firefox add-on that enables you record and replay tests. Selenium Webdriver, the third component, offers APIs in a range of languages, giving you greater control and allowing you to follow best practices in software development.



**Fig2: Selenium Suite**



**Fig3: selenium WebDriver architecture**

WebDriver was designed by Simon Stewart in 2006, at a time when JavaScript tools like Selenium Core were becoming increasingly powerful and restrictive in browsers and web apps. It was the first browser-controlling cross-platform testing framework. To build a programming interface that is more intuitive and concise. It supports [2] dynamic web pages, which allow elements of a page to change without requiring the entire page to be refreshed. The Web Driver (shown in fig.3) .

## II. LITERATURE SURVEY

Web scraping is a technique for obtaining large volumes of data from a variety of websites. This information can then be saved in a file or database. Web scraping is the process of extracting a web page's HTML code. Selenium [1] is utilized to achieve this goal. Selenium is an automated testing tool that is commonly used for web browser automation. It can, however, be used to scrape the web as well. The URLs of websites connected to the search query are acquired using the Selenium web driver.

Following that, each URL is automatically opened, and the data on each webpage is scraped one by one.

Selenium WebDriver settles on direct decisions to the program, exploiting every program's intrinsic robotization usefulness. Since there are such countless programs and programming dialects, WebDriver API[3] gives a standard determination. The expression "far off Webdriver" alludes to an API that should be carried out by every program. The language ties will provide orders to the normal driver API, and on the opposite end, a driver will tune in for those orders and execute them in the program utilizing far off WebDriver, returning the outcome/reaction by means of API to the code/Binding.

An inquiry demand is made to the Twitter Search[4] endpoint, which then, at that point makes a URL, utilizing any legitimate arrangement of HTTP-HEADERS, a variety of words (otherwise called looking through terms), a variety of dates (date range), and discretionary contentions.

A web bug measures a delivered HTTP reaction, and the HTML payload is sent to a download layer. At long last, natural information including tweets is provided into the Scrapy motor, which uses label selectors to extricate hypertext labels; each tweet is dealt with as an individual design comprised of plain content, the date it was made, and geographic information (if exists). Reactions are scratched utilizing a greatest position class quality from a div > HTML tag; this data is a pagination identifier from the most as of late appended tweets.

Web scratching of plans and fixings to acquire plans dependent on a specific fixing, which helps the culinary specialist in managing an element of interest, picking dishes dependent on fixings, or adhering to an eating routine arrangement dependent on the components utilized in plans. The data assembled is saved in the MongoDB data set. The interaction of web scratching [5] is utilized to move web information into a predefined storable arrangement. The culinary specialist can plan an eating regimen plan for the customers dependent on the fixing data or prescribe whether to utilize a specific thing in the formula dependent on the clients' inclinations or wellbeing concerns.

A toolset using cloud-based web scratching [6] to remove, refine, unite, and store COVID-19 cases data at different scales for each and every available country all through the planet normally. Changing the format crawler unit to meet the specific prerequisites of every information source. Testing and guaranteeing that simply the information expected from the objective information source is gathered.

The python-based web crawler [8] scrapy may likewise help us in recovering the ideal outcome, as we dissect the cycle utilizing explicit code and give the required url to the emphasis to scratch the information from the source url. The undertaking's essential web slithering content, which shows the information scratched and put away in the data set of items from an interpersonal interaction webpage, for this situation.

The initial step in this module is to scrape [9] Google results for a given string. Selenium, a portable web application testing framework, is used to achieve this goal. On the basis of a list compiled by us, stratifying the connections is in charge of identifying the links as reliable or unreliable. Extracting the links, on the other hand, entails scraping the content of each reliable link. We do it by utilizing the Python-based Scrappy framework.

Web scraping technology was utilized to obtain weather data [10] from websites in the South Sumatera area. Online scraping is a technique for retrieving individual elements from a web page. The information acquired through this web scraping technique will be saved in a database or data warehouse, where it can be used for future research into weather forecast data mining in South Sumatra, as well as the creation of weather-based decision-making applications.

## II. PROPOSED SYSTEM

The proposed a method which can scrap information from dynamic website pages utilizing selenium WebDriver in python. During the scratching interaction, any client activity on a program window can interfere with the stream and can cause a sudden conduct. Thus, for scratching applications, it is essential to stay away from any outer reliance while making applications, for example, browser. To overcome such issues Headless programs is used, which can work without showing any graphical UI and furthermore permits applications to be a solitary wellspring of cooperation for clients and gives a smooth client experience.

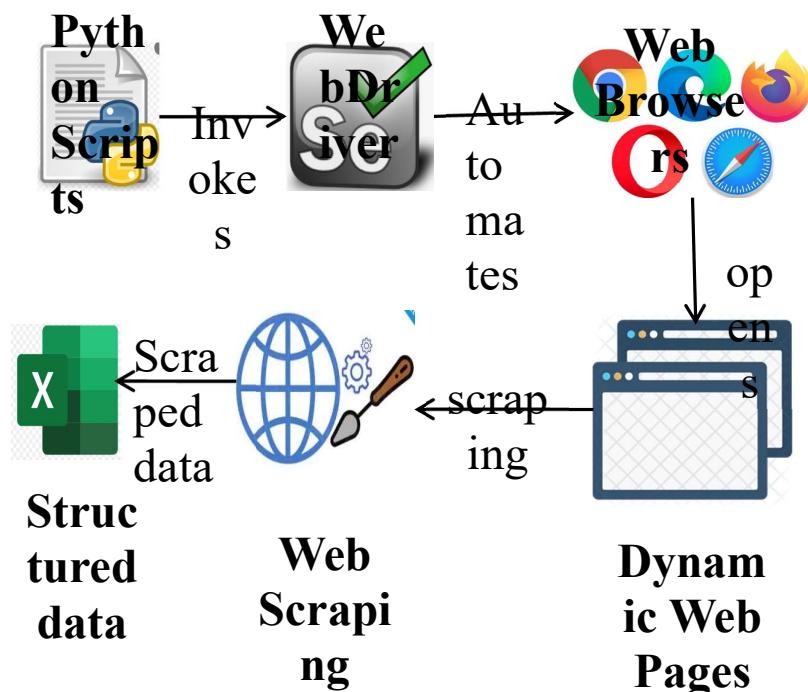


Fig4: Architecture diagram

As shown in the fig4, The webdriver will be executed by the scripts. The webdriver will automate the web browsers and open dynamic web pages that are directed to open, then scrape the data from those web pages and save it in an excel sheet or database. Selenium WebDriver is a popular technology for automating web user interfaces. It enables the automatic execution of online browser window tasks such as travelling to a website, filling out forms (including dealing with text boxes, radio buttons, and drop-downs), submitting the forms, browsing around web pages, and handling pop-ups, among others.

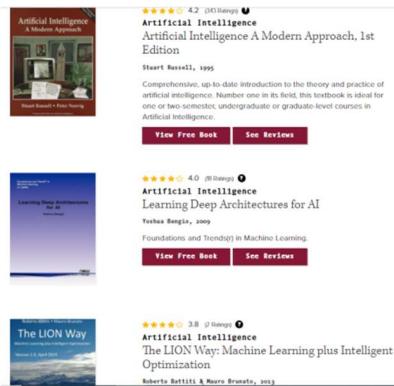


Fig5: Image of website that is been scraped

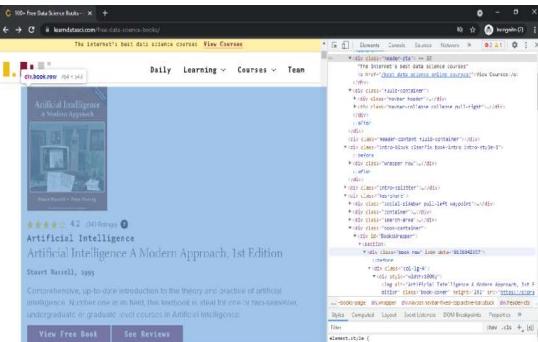


Fig6: Image of inspected code of one book.

**Locating the Elements:** The most crucial aspect of web scraping is finding the elements. A class, id, tag, name, or xpath can be used to identify elements on a web page. Classes are not unique, but ids are. This means that a class can identify several web elements, whereas an id can only identify one element. Any of the methods listed below can be used to identify an HTML element.  
**driver.find\_element\_by\_id, driver.find\_element\_by\_name, driver.find\_element\_by\_xpath, driver.find\_element\_by\_tag\_name, driver.find\_element\_by\_class\_name**

Any of the following can be used to identify multiple HTML

**elements.driver.find\_elements\_by\_name, driver.find\_elements\_by\_xpath, driver.find\_elements\_by\_tag\_name, driver.find\_elements\_by\_class\_name**

Scratch subtleties for each book on the page. Each page has 20 books. The subtleties of each book can be found by utilizing the URL on each card. Along these lines, to get the book subtleties we need this links. Scrape books in every single page. This implies that we will have a circle to scratch each book in the page and another to repeat through pages. Moving starting with one page then onto the next includes a change of the URL such that it is unimportant to anticipate a connection to any page. Fig6 shows the HTML code for the highlighted region after inspecting the site. To inspect the code right click on the particular element of the webpage. To access to href, we must go via the following hierarchy: class= “product prod” > h3 tag > a tag and the get value of the href property. In fact, all books on all pages are of the same product prod class and include the article tag.

### III. Results and Conclusion

A1	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1. title	book_cat	author	year	rating	total_rating	descriptive	book_link	review_link												
2. Artificial I Artificial I Stuart Russel			1995	4.2	343	Comprehere	http://www.amazon.com/gp/product/0136042597/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=0136042597&linkCode													
3. Learning I Artificial I Yoshua Bengio			2009	4	18	Foundatice	http://www.amazon.com/Learning-Architectures-Findings-Trends-Machine/dp/1601982941/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=1601982941&linkCode													
4. The LION Way I Yoshua Bengio			2013	3.8	2	Learning I	http://www.amazon.com/The-LION-Way-Intelligent-Optimization-Findings-Trends-Machine/dp/169604023/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=169604023&linkCode													
5. Disruptive Big Data I Jeffrey Na			2013	3.5	114	This book	http://www.amazon.com/gp/product/184823877N/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=184823877N&linkCode													
6. Computer Computer I Richard N			2010	4.2	101	Challenge	http://www.amazon.com/gp/product/1848239345/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=1848239345&linkCode													
7. Natural La Computer Steven Brin			2009	4.1	461	This book	http://www.amazon.com/gp/product/0596514495/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=0596514495&linkCode													
8. Programmin Computer Jan Erik S			2012	4	50	If you var J	http://ocrhttps://www.amazon.com/Programming-Computer-Vision-Python-algorithms/dp/1449316549/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=B00JUD80YY&linkCode													
9. The Eleme Data Anali Jeff Leek			2014	3.6	171	Data anal	https://leefilehttps://www.amazon.com/gp/product/B00JUD80YY/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=B00JUD80YY&linkCode													
10. A Course I Data Mini Hal Daum			2014	0	0	None	https://GinNone													
11. First End Data Mini Max Welli			2011	0	0	None	https://WnNone													
12. Algorithm Data Mini Cabo Stell			2009	4.2	4	This book	http://www.amazon.com/gp/product/1804549241/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=1804549241&linkCode													
13. A Program Data Mini Ron Zacha			2015	0	0	A Guide tc	http://guiNone													
14. Bayesian I Data Mini David Barl			2014	4.1	166	For finaly	http://ivehttps://www.amazon.com/gp/product/0521518148/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=0521518148&linkCode													
15. Data Mini Data Mini Wikibook			2014	0	0	None	https://erNone													
16. Data Mini Data Mini Mohammad			2014	4.1	11	The main	http://ivnhttps://www.amazon.com/gp/product/0596514495/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=0596514495&linkCode													
17. Data Mini Data Mini H. Wit			2005	3.8	158	Offers a t	http://tpihttps://www.amazon.com/gp/product/0521766328/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=0521766328&linkCode													
18. Data Mini Data Mini Graham V			2011	4.1	36	This book	http://mrihttps://www.amazon.com/gp/product/0521766328/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=0521766328&linkCode													
19. Deep Lear Data Mini Yoshua Be			2015	0	0	The Deep	http://vniNone													
20. Gaussian I Data Mini C. E. Raspi			2006	4.2	82	A competi	http://vnihttps://www.amazon.com/gp/product/026321825X/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=026321825X&linkCode													
21. Informatic I Data Mini David J.C.			2005	4.5	399	Essential	http://vnihttps://www.amazon.com/gp/product/026321825X/refas_li_t1ie=UTF8&camp=1788&creative=9325&creativeASIN=026321825X&linkCode													
22. Introducti Product I Data Mini Almon S-			2008	0	0	None	http://vniNone													
23. Introducti Product I Data Mini Alex Smol			2008	0	0	None	http://vniNone													
24. KB 4C I Ne Data Mini Roberto B			2013	0	0	None	http://vniNone													
25. Machine L Data Mini Abdchahid Melloul			0	0	None	http://vniNone														

Fig7: Scrapped data saved in an excel sheet.

Exploratory Data Analysis is the crucial process of using summary statistics and graphical representations to undertake preliminary investigations on data in order to uncover patterns, spot anomalies, test hypotheses, and verify assumptions. Results obtained after EDA are as follows

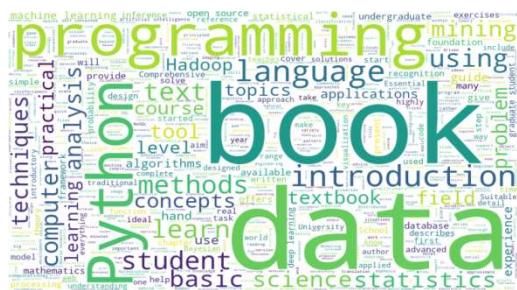


Fig8: Word cloud for book descriptions



Fig9: Word cloud for book titles

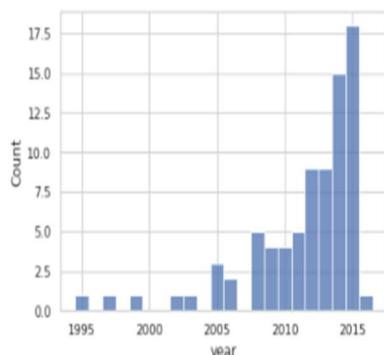


Fig10: Histogram of year.

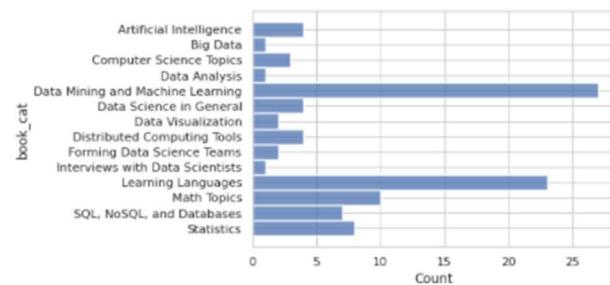


Fig11: Histogram for book category.

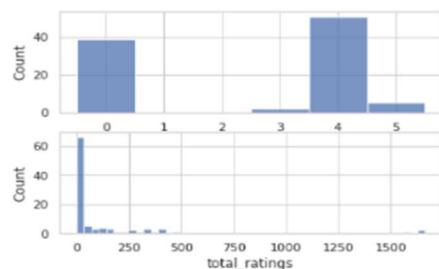


Fig12: plot for rating and total rating.

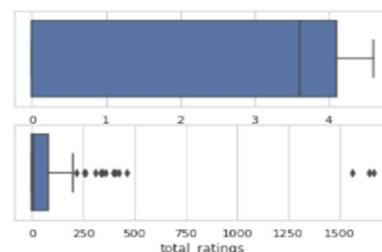


Fig13: Box plot for rating and total rating and outliers.

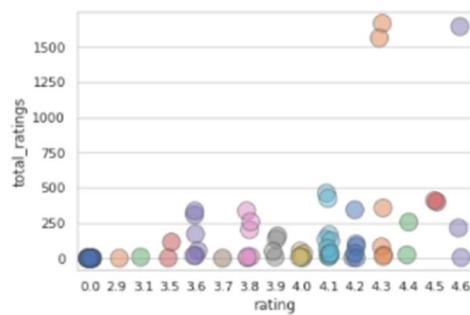


Fig14: strip plot of rating and total ratings

Web scraping is a common method of gathering information from the internet. Because many websites have distinct layout, there isn't a single scraper that can be used on all of them. The most important

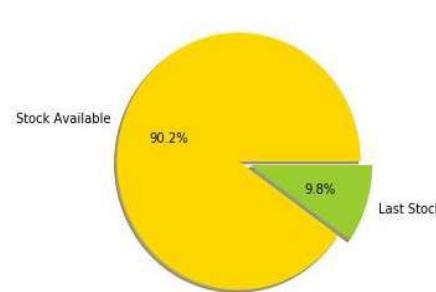


Fig15: Pie chart show stock availability.

skill is a thorough understanding of online scraping, which includes knowing how to locate web items and recognizing and dealing with issues, This work has covered all those techniques to scrape the dynamic web pages easily and fastly.

## REFERENCES

- [1] K Usha Manjari, Syed Rousha, Dasi Sumanth, Dr. J Sirisha Devi," Extractive Text Summarization from Web pages using Selenium and TF-IDF algorithm", Proceedings of the Fourth International Conference on Trends in Electronics and Informatics (ICOEI 2020), IEEE Xplore Part Number: CFP20J32-ART; ISBN: 978-1-7281-5518-0.
- [2] Revathi.K1 , Prof.V.Janani," SELENIUM TEST AUTOMATION FRAMEWORK IN ON-LINE BASED APPLICATION", International Journal of Advance Research In Science And Engineering, IJARSE, Vol. No.4, Special Issue (02), February 2015.
- [3] Paruchuri Ramya, Vemuri Sindhura, P Vidya Sagar," Testing using Selenium Web Driver", IEEE,2017.
- [4] A. Hernandez-Suarez,G. Sanchez-Perez, K. Toscano-Medina, V. Martinez-Hernandez, V. Sanchez and H. Perez-Meana ,” A Web Scraping Methodology for Bypassing Twitter API Restrictions”, arXiv:1803.09875v1 [cs.IR] 27 Mar 2018.
- [5] Shilpa Chaudhari, Aparna R., Vinay G Tekkur, Pavan G L., and Shreekanth R Karki,” Ingredient/Recipe Algorithm using Web Mining and Web Scraping for Smart Chef”, 978-1-7281-6828-9/20/\$31.00 ©2020 IEEE.
- [6] HAI LAN, DEXUAN SHA,ANUSHA SRIRENGANATHAN MALARVIZHI,YI LIU, YUN LI, NADINE MEISTER, QIAN LIU , ZIFU WANG, JINGCHAO YANG, AND CHAOWEI PHIL YANG ,” COVID-Scraper: An Open-Source Toolset for Automatically Scraping and Processing Global Multi-Scale Spatiotemporal COVID-19 Records”, VOLUME 9,IEEE Access, 2021.
- [7] Kasereka Henrys,” Importance of web scraping in e-commerce and e-marketing”, HKCorporation IT Official Journal January 2021.
- [8] David Mathew Thomas, Sandeep Mathur,” Data Analysis by Web Scraping using Python”, Proceedings of the Third International Conference on Electronics Communication and Aerospace Technology [ICECA 2019] IEEE Conference Record # 45616; IEEE Xplore ISBN: 978-1-7281-0167-5.
- [9] Dinesh Kumar Vishwakarma , Deepika Varshney, Ashima Yadav,” Detection and veracity analysis of fake news via scrapping and authenticating the web search”, 1389-0417/ 2019 Elsevier B.V.
- [10] Fatmasari, Yesi Novaria Kunang, Susan Dian Purnamasari,” Web Scraping Techniques to Collect Weather Data in South Sumatera”,international conference on electrical engineering and computer science(ICECOS) 2018.