

A Mechanism to Manage Time and Increase Data Accuracy When Using the Selenium Library

Farnaz Taghizadeh
Kourayem 

Ph.D. Student in Information Technology Management,
Central Tehran Branch, Islamic Azad University,
Tehran, Iran

Mohammadreza
Kabaranzad
Ghadim* 

Associate Professor, Department of Industrial
Management, Central Tehran Branch, Islamic Azad
University, Tehran, Iran

Seyed Abdollah
Amin Mousavi 

Assistant Professor, Department of Information
Technology Management, Central Tehran Branch,
Islamic Azad University, Tehran, Iran

Abstract

Today, data, as one of the valuable assets of various organizations and industries, plays an important role in the development and progress of businesses. In fact, every organization uses different sources to collect its data, one of which is the web platform, where a lot of data is produced and published by different users or even robots all over the world every day. Examining, researching, studying and analyzing such data can provide useful information and knowledge for the organization. For this purpose, during the past decades, various tools have been developed that have greatly helped in extracting information from the web platform, among which we can mention Request, Selenium, Scrapy, Beautiful Soup, etc. libraries in the Python programming language. However, each of these libraries faces challenges. In this article, by studying the Selenium library and considering the existence of many challenges in it, we have presented a solution for time management and improving the challenge of its Asynchronous. Our experiments show that the use of the proposed solution increases the accuracy of the information retrieved from the

* Corresponding Author: kabaranzad@yahoo.com

This article is taken from the doctoral thesis of Information Technology Management field, Islamic Azad University, Central Tehran Branch

How to Cite: Taghizadeh Kourayem, F., Kabaranzad Ghadim, M.R., Mousavi, S.A.A. (2024). A Mechanism to Manage Time and Increase Data Accuracy When Using the Selenium Library. *Journal of knowledge retrieval and semantic systems*, 11 (41), 199-225. DOI: [10.22054/jks.2024.80235.1660](https://doi.org/10.22054/jks.2024.80235.1660).

web platform and thus improves the challenge of Asynchronous and also reduces the time to retrieve information from the web platform.

1. Introduction

The Internet platform is a very powerful source of information that can be collected with the help of various tools and techniques and used after analysis in order to make better and more efficient decisions. According to previous researchers, when it comes to automatically extracting information from the web, Selenium is always the best option, however, this library has many challenges. One of the challenges of using the Selenium library is its asynchronous and the other is the slowness of the Selenium library, which we are trying to investigate and improve in this article.

Research Question(s): How to improve the challenge of slowness and asynchronous of the Selenium library?

2. Literature Review

Selenium library, which is one of the best web scripting tools, has been used in different studies and with different purposes. This library is a free, open-source automated testing framework used to verify web applications across multiple platforms and browsers. Various programming languages such as Java, C#, and Python can be used to create Selenium test scripts (Teotia et al., 2023). But despite its many advantages, selenium also has disadvantages, including: 1. Slowness, 2. Brittleness, 3. Flakiness, 4. Maintainability, 5. Asynchronous, 6. Time-consuming, 7. Cross-browser, 8. failure analysis, 9. Infrastructure, 10. Scalability, 11. Assertability, 12. Documentation and 13. Support (Leotta et al., 2023).

3. Methodology

The first thing we examined in the Selenium library is the lack of time management. Time management actually refers to the fact that the waiting time for downloading information from the web platform in this library is not known. To solve the problem of slowness and asynchronous of the Selenium library, we have used a solution that includes three different steps:

Step 1) According to the manual checks, first we define the variable t1:

$$t1 = 0.5$$

where t1 is the value used in the sleep function (time required to open the main page in normal mode).

Step 2) We use the while loop and try, except inside it. If the page does not open after 3 seconds, or after 6 attempts with different sleep times, the desired page and products do not appear, the error of the site not being available or the internet being slow will be printed:

t1 = 0.5

try:

 while t1 <= 3:

 try:

 driver.get("The Study site") # to open the page

 time.sleep(t1) # Waiting for the page to open

 ...

 # Code related to page scrolling

 time.sleep(t1) # Waiting for information to be displayed after scrolling

 # Information collection codes

 except:

 t1 += 0.5

 Except:

 print('The internet speed is very slow or the intended site is down')

In this code, the program with time t1=0.5 first tries to display the page information in full and if it fails, it adds half a second to t1 and this repetition continues up to 6 times. If the page is displayed in full, we use the new value of t1 for the next pages.

Step 3) If the page opens, we will enter the third step, which is related to collecting the basic information of the products and we must avoid the problem of asynchronous.

At this stage, according to the existence of five different types of information (such as product name, price, amount of discount, price after discount, type of discount) of each product, we first define five different and empty lists for product information.

Then, with the help of commands related to information collection, we take the information about each product and put it in the corresponding list. Then we run the following script to prevent wrong items in the list:

```
check1 = [[product_prices], [product_off], [product_prices2],  
[Product_off_type]]  
  
for j in range(0,4):  
    if len(product_name) < check1[len[j]]:  
        product_name.append('error')  
    elif len(product_name) > check1[len[j]]:  
        check[j].append('error')
```

In other words, after taking the information of a product and adding them to the predefined main lists, the program calculates the number of items in the existing lists with the help of the len function and puts them in the checklist. Then, with the help of the for loop, the length of each list is compared with the rest of the lists, and if the number of items in a list is low, the word "error" is added to it.

4. Results

In order to evaluate the solutions presented in this research, we have reviewed the information related to supplementary drugs on the DigiKala site, which was almost 3000 different medicines, 13 times on 13 different dates, from September 23, 2023, to March 15, 2024. In this study, the codes written in order to retrieve information from the mentioned site were executed the first time with the proposed solution and the second time without the proposed solution, and in each execution, both the time of information retrieval and the number of information related to each list or the same column was recorded and compared with each other, and the error rate and its percentage were calculated based on the difference in the time of information collection in the first and second execution and the difference in the number of information collected for each column in the first and second execution. After the implementation and use of the proposed solution, the investigations show that the accuracy and correctness of the collected information have increased compared to not using the proposed solution, and the time of information collection has also improved.

5. Discussion and Conclusion

In this article, we studied and evaluated the challenges of being slow, time-consuming and asymmetric of the Selenium library. Our studies were conducted using the Python programming language. Studies show that it is very important to use the solution of checking the list and the

same length of the list at the end of the collection of each product from the web platform, so that not using it in 12 out of 13 cases of information collection from the web platform makes us encountered an error. Also, using a constant value for the sleep function significantly increases the time to retrieve information compared to using a variable value for it. In general, the findings show that the use of the proposed solution when using the Selenium library in order to extract information from the web platform helps to increase the accuracy of the information and also improves the time of complete information retrieval from the web platform.

Keywords: Web Crawler, Web Scraping, Selenium Library, Asynchronous, Data Accuracy, Data Retrieval Time

سازوکاری برای مدیریت زمان و افزایش دقت اطلاعات هنگام استفاده از کتابخانه سلینیوم

دانشجوی دکتری مدیریت فناوری اطلاعات، واحد تهران مرکزی،
دانشگاه آزاد اسلامی، تهران، ایران

فرناز تقی‌زاده کورایم 

تاریخ دریافت: ۲۰/۰۶/۱۴۰۳

دانشیار، گروه مدیریت صنعتی، واحد تهران مرکزی، دانشگاه آزاد
اسلامی، تهران، ایران

محمد رضا کاباران زاد

قدیمی * 

استادیار، گروه مدیریت فناوری اطلاعات، واحد تهران مرکزی، دانشگاه
آزاد اسلامی، تهران، ایران

سید عبداله امین موسوی 

تاریخ پذیرش: ۱۹/۰۵/۱۴۰۳

ISSN: 2980-8243

eISSN: 2783-1795

کلیدواژه‌ها: خزنده وب، وب اسکرینینگ، کتابخانه سلینیوم، نامتقارن بودن داده، دقت اطلاعات، زمان برداشت اطلاعات

* نویسنده مسئول: kabaranzad@yahoo.com

مقاله حاضر برگرفته از رساله دکتری رشته مدیریت فناوری اطلاعات دانشگاه آزاد اسلامی، واحد تهران مرکزی است.

مقدمه

امروزه داده‌ها به عنوان یکی از دارایی‌های ارزشمند سازمان‌ها، نقش مهمی را در توسعه و پیشرفت کسب‌وکارها ایفا می‌کنند. این داده‌ها در بسترهای مختلفی موجود هستند که از مهم‌ترین آن‌ها می‌توان به اینترنت و شبکه‌های اجتماعی اشاره نمود که در هر دو محیط شاهد وجود حجم بالایی از داده‌ها هستیم که در صورت پردازش و استفاده درست از آن‌ها در زمان مناسب می‌تواند کمک شایانی به سازمان‌ها و کسب‌وکارها نماید.

همان‌طور که گفته شد، بستر اینترنت منبع اطلاعاتی بسیار قدرتمندی است که می‌توان اطلاعات موجود در آن را با کمک ابزارها و فنون مختلف برداشت و پس از تحلیل در راستای تصمیم‌گیری‌های بهتر و کاراتر استفاده نمود. یکی از روش‌های موجود برای برداشت اطلاعات از بستر وب، استفاده از زبان‌های برنامه‌نویسی نظری پایتون، سی شارپ و جاوا است که در آن‌ها با استفاده از کتابخانه‌های مختلف نظری ریکوئست^۱، سلنیوم^۲، اسکرپت^۳، سوب زیبا^۴ و ... می‌توان نسبت به خزش یا برداشت اطلاعات از بستر وب اقدام نمود.

طبق صحبت‌های پژوهشگران پیشین وقتی صحبت از برداشت خودکار اطلاعات از بستر وب به میان می‌آید، سلنیوم همیشه بهترین گزینه است، با این حال این کتابخانه چالش‌های متعددی را به همراه دارد. یکی از چالش‌های استفاده از کتابخانه سلنیوم نامتقارن بودن آن است. نامتقارن به این معناست که اسکرپت‌ها یا همان کدهای مربوط به برداشت اطلاعات از بستر وب باید تعاملات ناهم‌زمان را مدیریت کنند که اغلب برطرف نمودن آن‌ها دشوار است. تعاملات ناهم‌زمان، هنگامی اتفاق می‌افتد که بخش‌های مختلف یک صفحه در وب، اطلاعات خود را از بخش‌های مختلف یک سرور، یا چند سرور مختلف و یا حتی واسطه‌های برنامه‌نویسی^۵ مختلفی دریافت می‌کنند، درنتیجه ممکن است این اطلاعات در یک‌زمان و باهم نمایش داده نشوند و بارگذاری و نمایش برخی از اطلاعات در صفحه با کندی و تأخیر مواجه شود و یا در برخی موارد حتی بخش‌هایی از صفحه وب نمایش داده نشود و نیاز به زمان بیشتر یا بارگذاری مجدد صفحه باشد. به عبارت دیگر، محتوای ناهم‌زمان در صفحه وب

¹. Requests

². Selenium

³. Scrapy

⁴. BeautifulSoup

⁵. Application Programming Interface (API)

موجب می‌شود تا بخشی از داده‌ها را از سرور بدون بارگیری مجدد کل صفحه بازیابی کند؛ و اگر اسکریپت بخواهد با محتوای صحیح صفحه تعامل داشته باشد، باید منتظر بمانید تا محتوای ناهم‌زمان بهروزرسانی شود (پاسخ از سرور دریافت شود و صفحه مطابق آن بهروزرسانی شود). چالش اصلی در رسیدگی به تعاملات ناهم‌زمان، دانستن زمان انتظار برای بارگیری محتوای صفحه وب است که گاه توسط چندین بهروزرسانی و یا افزایش زمان انتظار برای بارگذاری کامل اطلاعات، برطرف می‌گردد (Leotta et al., 2023).

اما سؤال اینجاست که بدون دخالت انسان، سیستم چه طور تشخیص دهد که چقدر باید منتظر برداشت کامل اطلاعات باشد و یا چند بار باید صفحه مجدد بارگذاری شود که اطلاعات درست را نمایش دهد؟ ما در این مقاله با مطالعه اطلاعات داروهای مکمل در سایت دیجی کالا و با استفاده از راه حل پیشنهادی در تلاش هستیم تا بهترین زمان بارگذاری را یافته و به عنوان مأذولی از آن در کتابخانه سلینیوم استفاده نماییم تا ضمن اعلام مشکل احتمالی وب، بتواند بدون دخالت انسان قادر به تشخیص بهترین زمان انتظار برای بارگذاری کامل صفحه و یا تشخیص نیاز به بارگذاری مجدد صفحه باشد.

در راستای هدف مقاله، نیازمند مرور ادبیات پژوهش هستیم که در ادامه به آن پرداخته‌ایم، سپس به تشریح راه حل پیشنهادی برای مدیریت زمان و رفع مشکل نامتقارن بودن سلینیوم پرداخته و پس از آن راه حل پیشنهادی خود را مورد ارزیابی قرار داده‌ایم و در ادامه نیز بحث و نتیجه‌گیری و پیشنهادهای آتی ذکر شده‌اند.

پیشنهاد پژوهش

امروزه، دیجیتالی کردن و مجازی‌سازی فرایندهای اجتماعی منجر به زتابایت (میلیاردها گیگابایت) داده‌های موجود در وب شده است. این داده‌ها یک نمایش بلاذرنگ از فرایندها، روابط و تعاملات متعدد در فضای مادی-اجتماعی ارائه می‌دهد؛ بنابراین، این حجم وسیع از داده‌های وب، فرصت‌های فراوانی را در اختیار پژوهشگران دانشگاهی قرار می‌دهد تا به سؤالات پژوهش‌های مربوط به پژوهش‌های جدید و قدیمی با دقت بیشتر و زمان‌بندی بهتر پاسخ دهند. همچنین شاغلین می‌توانند از این داده‌ها برای توسعه و درک بهتر مشتریان خود، تدوین راهبردها بر اساس این یافته‌ها و درنهايت بهبود عملکرد سازمانی استفاده کنند (Krotov et al., 2020).

همان‌طور که گفته شد، وب به عنوان بستری غنی از اطلاعات، دهه‌هاست که توسط پژوهشگران مختلف مورد استفاده قرار می‌گیرد؛ اما قطعاً بررسی چشمی تک‌تک اطلاعات و یا یادداشت‌برداری از حجم بالایی اطلاعات و یا به عبارت دیگر بازیابی اطلاعات به صورت دستی توسط انسان کار بسیار زمان‌بر، پرهزینه و دشواری است. به همین منظور طی دهه‌های گذشته ابزارهای مختلفی توسعه یافته‌اند که می‌توان با کمک آن‌ها اطلاعات موردنظر را به صورت خودکار و با صرف زمان بسیار کمتری بازیابی نمود که غالباً به این ابزارها و فنون برداشت اطلاعات از بستر وب، وب اسکرپینگ^۱ گفته می‌شود (Henry, 2021).

وب اسکرپینگ که به عنوان داده کاوی نیز شناخته می‌شود، فرایند جمع‌آوری مقادیر زیادی داده از وب و سپس قرار دادن آن در پایگاه‌های داده برای تجزیه و تحلیل آینده و استفاده‌های بعدی است. اسکرپینگ وب روشنی است که استفاده از آن، بینشی در مورد داده‌های قیمت، پویایی بازار، روندهای غالب، شیوه‌های به کار گرفته شده توسط رقبا و چالش‌هایی که آن‌ها با آن روبرو هستند ارائه می‌دهد (Henry, 2021).

طبق مطالعات انجام شده در سال ۲۰۲۰ که مبتنی بر جستجوی مقاله بر اساس نام یا همان عنوان مقالات بود، نشان می‌دهد که وب اسکرپینگ، یک پدیده نوظهور است؛ چراکه نخست، بیشتر مقالات حاوی این دو عبارت تنها به چند سال پیش بازمی‌گردد. دوم، تعداد کل مقالات مجلات و کنفرانس‌هایی که به صراحت به وب اسکرپینگ در چندین پایگاه داده معروف (برای نمونه IEEE Xplore, EBSCO Business Source Complete Digital Library, JSTOR وغیره) اختصاص داده شده است، کمتر از دویست است. از جمله این مقالات می‌توان به مقاله پژوهشی کروتوف و تینیسون^۲ (۲۰۱۸) اشاره نمود. این مقاله، به مطالعه داده‌های مالی و زبان‌های نشانه‌گذاری محبوب استفاده شده برای ذخیره و تبدیل این داده‌ها از یک متن ساده به یک متن قابل درک‌تر و زیباتر در بستر وب می‌پردازد و نشان می‌دهند که چگونه زبان R، همراه با کتابخانه‌هایش، می‌تواند برای جمع‌آوری، سازمان‌دهی و پیش‌پردازش داده‌های مالی مورد استفاده قرار گیرد. مقاله آموزشی مشابه دیگری توسط نیومن و همکاران^۳ (۲۰۱۷) توضیح می‌دهد که چگونه می‌توان از فناوری XPath برای برداشت فراداده از مخازن کتابخانه دیجیتال توسط پژوهشگران علاقه‌مند به

¹. Web Scraping

². Krotov & Tennyson

³. Neumann et al.

پژوهش‌های علمی استفاده نمود. همچنین مقاله بوینگ و وادل^۱ (۲۰۱۷) نشان می‌دهد که چگونه داده‌های به دست آمده از Craigslist را می‌توان بازیابی و تجزیه و تحلیل کرد تا در ک بهتری از بازار اجاره در ایالات متحده به دست آورد. در مجموع، این مقالات نشان می‌دهد که علاقه فزاینده‌ای به وب اسکرپینگ در بسیاری از رشته‌ها و صنایع علمی وجود دارد (Krotov et al., 2020).

کتابخانه سلینیوم که ابزاری برای برداشت اطلاعات از بستر وب است، در مطالعات مختلف و با اهداف مختلفی مورد استفاده قرار گرفته است، برای نمونه مقاله ثوتیا و همکاران^۲ (۲۰۲۳) از آن به منظور مطالعه و تجزیه و تحلیل داده‌های وب‌سایت‌های فیلم چینی برای درک توزیع روند ژانرهای فیلم و رتبه‌بندی آن‌ها استفاده نموده است. مقاله یوان^۳ (۲۰۲۳) روش جدیدی را برای برداشت اطلاعات از بستر وب پیشنهاد می‌کند که در آن سلینیوم (به عنوان یک ابزار مشهور برداشت اطلاعات از بستر وب) را با شبکه‌های عصبی کانولوشنال^۴ ترکیب نموده تا به طور خودکار عناصر صفحات وب را شناسایی کند. در این روش، از سلینیوم برای پیمایش صفحات وب، دستکاری عناصر و همچنین گرفتن اسکرین‌شات استفاده می‌کند. این اسکرین‌شات‌ها متعاقباً با استفاده از شبکه‌های عصبی کانولوشنال پردازش می‌شوند تا از شناسایی و طبقه‌بندی عناصر صفحه و ب اطمینان حاصل شود. مقاله نوشته شده توسط سوگانتی و وارون^۵ (۲۰۲۴) با استفاده از قابلیت‌های سلینیوم، وب‌سایت‌ها را پیمایش نموده و اطلاعات مرتبط را به طور مؤثر استخراج می‌کند. همچنین با برطرف نمودن چالش‌های ناشی از محتواهای پویا و ساختارهای پیچیده وب‌سایت، از استخراج یکپارچه داده‌ها اطمینان می‌دهد. این پژوهش به عنوان یک راه حل برای خودکار نمودن کارهای اسکرپتی و ب، توانمندسازی کاربران برای تعیین وب‌سایت‌های هدف، تعریف معیارهای استخراج و جمع‌آوری داده‌های موردنظر بدون مداخله دستی انسان عمل می‌کند و در واقع بهره‌وری و کارایی را در بازیابی داده‌ها از چشم‌انداز وسیع اینترنت به طور قابل توجهی افزایش می‌دهد (Suganthi & Varun, 2024).

¹. Beoing & Waddell

². Teotia et al.

³. Yuan, S.

⁴. Convolutional Neural Network (CNN)

⁵. Suganthi & Varun

همچنین مقاله نوشته شده توسط هنریس^۱ (۲۰۲۱) به اهمیت وب اسکرپنگ در تجارت الکترونیک و بازاریابی پرداخته و از جمله دلایل اهمیت آن را در ارائه اطلاعات مرتبط با موارد ذیل می‌داند:

نظرارت بر قیمت و تحقیقات محصول؛ مقایسه قیمت آنلاین؛ تجزیه و تحلیل بهتر مشتری؛ تحلیل بازار؛ تبلیغات بهتر؛ تأثیرگذاری بر راهبرد بازاریابی و فروش؛ نظرارت بر برنده؛ استخراج جزئیات کسب و کار از فهرست کسب و کار؛ کمک به تحلیل آینده و غیره (Henrys, 2021).

سلنیوم یک چهارچوب آزمون خودکار متن باز و رایگان است که برای تأیید برنامه‌های وب در بسیاری از سکوها و مرورگرها استفاده می‌شود. برای ساخت اسکریپت‌های تست سلنیوم، می‌توان از زبان‌های برنامه‌نویسی متنوعی از جمله جاوا، سی شارپ و پایتون استفاده کرد. سلنیوم در سال ۲۰۰۴ کشف شد. جیسون هاگینز^۲ مردی بود که سلنیوم را تولید کرد. او یک توسعه‌دهنده و مهندس سلنیوم است. او زمانی شروع به کار کرد که به طور فعال سعی کرد خود را وقف کارهای شناختی کند و متوجه شد که به جای حرکت فیزیکی در این فعالیت‌ها، ممکن است از زمان خود حداکثر استفاده را بکند. او شروع به آزمایش با جاوا اسکریپت کرد و یک چهارچوب برنامه‌نویسی جاوا ایجاد کرد که می‌تواند با یک وب‌سایت ارتباط برقرار کند و به طور طبیعی تعدادی از برنامه‌ها را آزمایش کند که درنتیجه تلاش‌های موفقیت‌آمیز او کتابخانه سلنیوم ایجاد گردید (Teotia et al., 2023).

همان‌طور که گفته شد، وقتی صحبت از برداشت خودکار اطلاعات از بستر وب به میان می‌آید، سلنیوم همیشه بهترین گزینه است، چراکه: ۱- سلنیوم به طور رایگان در دسترس است؛ ۲- این نرم‌افزار از زبان‌های مختلفی از جمله پایتون، پی‌اچ‌پی، جاوا، سی شارپ، روبي و جاوا اسکریپت و غیره پشتیبانی می‌کند؛ ۳- اسکریپت برای اجرا در هیچ سیستم عامل دیگری به جز سیستم عاملی که در آن ایجاد شده است نیازی به تغییر ندارد؛ زیرا با بسیاری از سیستم عامل‌های دیگر سازگار است؛ ۴- قابلیت سازگاری مرورگر آن به شخص اجازه می‌دهد تا اسکریپت را در هر مرورگری اجرا کند، از جمله مایکروسافت اج، گوگل کروم، فایرفاکس، سافاری، اپرا و بسیاری دیگر؛ ۵- کاربرپسند و نصب آن ساده است؛ ۶- کد به سرویس‌های وب تفسیر می‌شود و درایور راه دور آن را با استفاده از پرس‌وجوهای پروتکل

¹. Henrys, K.

². Jason Huggins

انتقال ابرمنن^۱ دریافت می‌کند، بنابراین قبل از اجرا به هیچ سروی نیاز ندارد. سلنیوم بیش از سایر فنون خراش دادن^۲ وب یا همان وب اسکرپینگ مورد علاقه است؛ زیرا امکان تعامل با صفحات وب پویا را فراهم می‌کند. سلنیوم در شرایطی که داده‌هایی که می‌خواهیم استخراج کنیم در پشت اشیای جاوا اسکرپت پنهان می‌شوند که باید روی آن‌ها کلیک کرد تا محتوا نمایان شود، به خوبی کار می‌کند (Teotia et al., 2023).

اما با وجود مزایای گفته شده، سلنیوم به ویژه در طراحی برنامه‌های آزمون خود کار نرم‌افزار معنایی را نیز داراست که از جمله آن‌ها می‌توان به ۱- کندی؛ ۲- شکنندگی^۳ (این مشکل زمانی ایجاد می‌گردد که تکامل وب‌سایت یا برنامه کاربردی و یا رفع ایراد آن، موجب ایجاد خطا در بخش‌های دیگر شده و عملکرد آن‌ها را تحت تأثیر قرار دهد که در صورت وقوع چنین مشکلی سلنیوم غالب قادر به تشخیص آن و یا حتی یافتن برخی پیوندهای و فیلد‌ها نیست)؛ ۳- پوسته‌پوسته شدن^۴ (این مشکل هنگامی رخ می‌دهد که اسکرپت‌های وب‌سایت و یا برنامه کاربردی بدون دلیل آشکاری تغییر می‌کنند، این مشکل موجب ایجاد خطا هنگام اجرای کدهای مربوط به کتابخانه سلنیوم می‌شود که گاه خطایابی آن دشوار و چالش‌برانگیز است)؛ ۴- دشواری در قابلیت نگهداری؛ ۵- نامتقارن بودن؛ ۶- زمان بربودن؛ ۷- مرورگر کراس^۵ (این مشکل هنگام استفاده از مرورگرهای متفاوت اتفاق می‌افتد، به عبارت دیگر هر یک از مرورگرهای ممکن است صفحات وب را به صورت متفاوتی ارائه دهد که این مورد به سیستم عامل مورداستفاده نیز بستگی دارد، یعنی کدهای نوشته شده برای با کمک کتابخانه سلنیوم ممکن است هنگام استفاده از یک مرور به درستی کار کند ولی هنگام استفاده از مرورگر دیگر یا حتی سیستم عامل دیگر با خطا مواجه گردد)؛ ۸- تجزیه و تحلیل شکست؛ ۹- زیرساخت؛ ۱۰- مقیاس‌پذیری؛ ۱۱- ادعای‌پذیری^۶ (این مشکل به این دلیل ایجاد می‌گردد که اغلب داده‌ها از منابع نامعتبر جمع‌آوری می‌گردد و هر لحظه امکان دارد که ظاهر و احساس صفحات وب، سطح دسترسی و یا نقشه‌های تعاملی موجود در صفحات وب تغییر کند، درنتیجه برخی ادعاهای مانند این مسئله که سلنیوم قادر به بررسی

¹. HTTP

². Scraping

³. Brittleness

⁴. Flakiness

⁵. Cross-browser

⁶. Assertability

- ویژگی‌های بصری صفحات وب است، با چالش مواجه می‌گردد؛ ۱۲ - مستندسازی؛ ۱۳ - پشتیانی کردن، اشاره نمود (Leotta et al., 2023).

روش‌شناسی پژوهش

کتابخانه سلینیوم به عنوان یکی از بهترین کتابخانه‌ها برای برداشت اطلاعات از بستر وب، مدتی است که مورد توجه پژوهشگران مختلف قرار گرفته است. با این حال، با وجود مزایای مختلف دارای معایبی نیز است که از جمله آن‌ها می‌توان به عدم وجود مدیریت زمان و چالش نامتقارن بودن در آن اشاره نمود که در ادامه ابتدا تعریف عملیاتی از این دو مورد را ارائه نموده و سپس روش پژوهش را بیان نموده و بعد از آن راه حلی مناسبی به منظور بهبود کتابخانه سلینیوم که مبتنی بر این دو مورد است را ارائه نموده‌ایم.

تعریف عملیاتی از مدیریت زمان: از جمله مشکلات کتابخانه سلینیوم کندی و زمان‌بر بودن آن است، چراکه به منظور برداشت اطلاعات از بستر وب باید از طریق مرورگر تمام لایه‌های موردنیاز را بررسی و اطلاعات موردنظر را جمع‌آوری نماید؛ اما چگونه می‌توان این کندی و زمان‌بر بودن را اندازه‌گیری نمود و آن را بهبود بخشید؟ ما به منظور افزایش سرعت برداشت اطلاعات از بستر وب و یا همان کاهش زمان برداشت اطلاعات از بستر وب و کاهش مشکلات کندی و زمان‌بر بودن کتابخانه سلینیوم که به صورت کلی آن را مدیریت زمان در کتابخانه سلینیوم نام‌گذاری نموده‌ایم، از تابع sleep استفاده نموده‌ایم. به عبارت دیگر، در این پژوهش مشخص نموده‌ایم که ایجاد چه تغییراتی با استفاده از تابع sleep می‌تواند مشکل کندی و زمان‌بر بودن کتابخانه سلینیوم را کاهش دهد.

تعریف عملیاتی از نامتقارن بودن: یکی دیگر از مشکلات کتابخانه سلینیوم نامتقارن بودن آن است. این مشکل هنگامی که تعاملات ناهم‌زمان در یک صفحه وب وجود دارد، اتفاق می‌افتد. تعاملات ناهم‌زمان، به این دلیل رخ می‌دهند که یک صفحه وب اطلاعات خود را از سرویس‌دهنده‌های مختلفی دریافت کرده و آن‌ها را نمایش می‌دهد، درنتیجه ممکن است سرعت ارتباط بین یک صفحه وب با تمام سرویس‌دهنده‌های آن یکسان نباشد و یا حتی به دلایل مختلف و در طول زمان پیوند بین سرویس‌دهنده و وب‌سایت تغییر نماید و این مسئله موجب می‌گردد تا برخی اطلاعات به نسبت اطلاعات دیگر با تأخیر نمایش داده شوند و یا اصلاً سلینیوم قادر به تشخیص وجود آن‌ها نباشد. این مسئله، اسکریپت‌های نوشته شده برای کتابخانه سلینیوم را با مشکل مواجه می‌نماید، چراکه اسکریپت‌ها قادر به برداشت بخشی از

اطلاعات هستند ولی نمی‌توانند تشخیص دهنند که اطلاعات دیگر هنوز بارگذاری نشده‌اند و یا پیوند ارتباطی آن‌ها با سرور تغییر نموده است. درنتیجه خروجی نهایی را تحت تأثیر قرار می‌دهند و موجب می‌شوند که اطلاعات در جای درست خود قرار نگیرند، به عنوان نمونه در مورد مطالعاتی ما قیمت یک محصول در مقابل نام محصولی دیگر ثبت گردد. ما به منظور برطرف نمودن این مشکل که نام آن را در عنوان مقاله افزایش دقت اطلاعات گذاشته‌ایم، از فهرست‌ها و طول آن‌ها استفاده نموده‌ایم.

ما در پژوهش خود ابتدا داده‌های کیفی را از طریق اسکریپت‌های نوشته‌شده با استفاده از کتابخانه سلینیوم از بستر وب جمع‌آوری نمودیم و سپس با محاسبه تعداد این داده‌ها و مدت زمان موردنیاز برای برداشت این اطلاعات از بستر وب، داده‌های کمی را محاسبه نمودیم. به عبارت دیگر داده‌های نخست ما کیفی و داده‌های دوم که اسکریپت‌های سلینیوم براساس آن‌ها بهبود یافت و پژوهش نیز براساس آن‌ها ارزیابی گردید، کمی است. به عبارت دیگر، داده‌های استفاده شده در این پژوهش هم شامل داده‌های کمی و هم شامل داده‌های کیفی هستند. درنتیجه داده‌های موجود در این پژوهش ترکیبی و پژوهش حاضر از نظر هدف کاربردی و از نظر روش تحلیلی است. روش جمع‌آوری اطلاعات هم به صورت مطالعه موردنی بود که در آن از مکمل‌های دارویی سایت دیجی کالا استفاده شده است. اولین موردی که آن را در کتابخانه سلینیوم موردنبررسی قرار دادیم، عدم وجود امکان مدیریت زمان در آن است. مدیریت زمان درواقع اشاره به این موضوع دارد که مدت زمان انتظار برای برداشت اطلاعات از بستر وب در این کتابخانه مشخص نیست و بستگی بسیار زیادی به سرعت اینترنت دارد. همچنین این مشکل زمانی حاد می‌شود که سیستمی تماماً خودکار برای برداشت اطلاعات از بستر وب طراحی نماییم.

اما چگونه می‌توان بر این مشکل غلبه و آن را برطرف نمود؟

قبل از پرداختن به این مشکل، ابتدا نگاهی به ساختار برخی از صفحات وب می‌اندازیم. اکثر صفحات اصلی فروشگاهی در بستر وب، خلاصه‌ای از اطلاعات محصول از قبیل نام محصول، قیمت و میزان تخفیف را در صفحه اول نمایش می‌دهند که کاربر با کلیک بر روی محصول می‌تواند وارد صفحه مربوط به آن محصول شده و اطلاعات کامل تری از قبیل مشخصات، روش استفاده، دیدگاه کاربران و کارشناسان پیرامون محصول و ... را به دست آورد.

همچنین، طبق بررسی‌های انجام شده بر روی سایت‌های فروشگاهی مختلف در بستر وب، صفحه اصلی فروشگاه‌ها بر اساس تعداد محصول موجود در آن‌ها، معمولاً سه نوع هستند؛ یکی صفحاتی که فقط نیاز به پیمایش کردن دارند تا بتوان اطلاعات اولیه (نظیر نام محصول و قیمت و میزان تخفیف آن) همه محصولات را مشاهده نمود، دوم صفحاتی که هم نیاز به پیمایش دارند و هم نیاز به رفتن به صفحات بعد دارند تا بتوان اطلاعات اولیه تمام محصولات موردنظر را مشاهده نمود و سوم صفحاتی که نه نیاز به پیمایش دارند و نه نیاز به رفتن به صفحه بعد که معمولاً در صفحات از نوع سوم، اطلاعات و محصولات کمی وجود دارد.

کتابخانه سلیوم به منظور برداشت اطلاعات از بستر وب از وب درایور استفاده می‌کند و برای هر مرورگر نیز وب درایور مناسب آن طراحی گردیده است. همچنین از جمله دستورات مهم در کتابخانه سلیوم، دستور get است. این دستور به منظور بارگذاری صفحه موردنظر با کمک وب درایور به کار گرفته می‌شود. برای نمونه به دستور ذیل که دستوری برای باز شدن یکی از صفحات دیجی کالا است، توجه نمایید:

```
driver = webdriver.Chrome()
```

```
driver.get("https://www.digikala.com/search/nutritional-supplement")
```

اما اگر فقط همین دستور را اجرا نماییم و سپس دستورات مربوط به برداشت اطلاعاتی از قبیل نام کالا، قیمت کالا و ... از صفحه موردنظر را قرار دهیم، به احتمال زیاد با خطأ مواجه خواهیم شد، زیرا احتمالاً یا صفحه اصلاً بارگذاری نشده است و یا بخشی از آن بارگذاری نشده است. برای رفع این مشکل تابع sleep در کتابخانه سلیوم طراحی گردیده است. تابع sleep که به صورت ذیل نوشته می‌شود، موجب ایجاد زمان انتظاری که توسط کاربر تعریف شده، می‌شود و احتمال خطأ را کاهش می‌دهد. برای نمونه، در کد ذیل، برنامه پس از رفتن به صفحه موردنظر، ۲ ثانیه صبر می‌کند تا صفحه بارگذاری شود:

```
Time.sleep(2)
```

اما چه میزان زمان برای بارگذاری این صفحه که برای نمونه ما در تابع بالا آن را ۲ ثانیه در نظر گرفته‌ایم، موردنیاز است؟ برای رفع این مشکل ما از راه حلی استفاده نموده‌ایم که شامل سه مرحله مختلف است و فرضیه‌های آن به شرح ذیل است:

۱. ما در پژوهش‌های خود ساعت ثابتی را برای برداشت اطلاعات در نظر گرفته‌ایم چراکه هدف ما ساخت سیستمی خودکار است که مثلاً روز پانزدهم هرماه ساعت ۸ صبح

اطلاعات موردنظر را برداشت و داده‌های موجود را به روزرسانی نماید. به همین منظور ساعت به روزرسانی اطلاعات جدول ساعت ۸ صبح است.

۱. فرض بر این است که زمان انتظار برای باز شدن هر صفحه جدید، زمان انتظار برای نمایان شدن محصول بعد از هر پیمایش و برداشت هر یک از اطلاعات مربوط به هر محصولات در شرایط با سرعت اینترنت مناسب، نیم ثانیه است.

۲. همچنین میزان پیمایش رانیز می‌دانیم و یا اصلاً نیاز به پیمایش نداریم.

۳. در صفحاتی که بررسی می‌نماییم حداقل یک کالا وجود دارد.

۴. بیشترین تعداد کالای موجود در صفحه ۲۰ است.

مرحله ۱) طبق بررسی‌هایی که به صورت دستی انجام شد، ابتدا متغیر $t1$ را به شرح ذیل تعریف می‌کنیم:

$t1 = 0.5$

که در آن $t1$ زمان موردنیاز برای باز شدن صفحه اصلی در حالت عادی است.

مرحله ۲) از حلقه while و try except خارج از آن، استفاده می‌نماییم. به این منظور که اگر بعد از ۳ ثانیه صفحه باز نشد، یا به عبارت دیگر بعد از ۶ بار تلاش با زمان‌های sleep متفاوت صفحه و کالاهای موردنظر باز نشده و نمایان نشدنند، خطای در دسترس نبودن سایت و یا کندی اینترنت چاپ می‌گردد؛ به عبارت دیگر، در بدنه اصلی try except دیگری قرار می‌دهیم، این بخش به این منظور قرار دارد که اگر صفحه در زمان $t1$ باز نشد و محصولات نمایان نشدنند، نیم ثانیه به آن اضافه می‌شود و مجدد موردنبررسی قرار گیرد که این کار تا ۵ بار ادامه پیدا می‌کند و سپس $t1$ به دست آمده به عنوان یک متغیر نگهداری می‌شود و برای صفحات بعد مورداستفاده قرار می‌گیرد.

$t1 = 0.5$

try:

 while $t1 <= 3$:

 try:

 کد باز شدن صفحه # ("")"سایت موردمطالعه""")

 انتظار برای باز شدن صفحه #

 کد مربوط به اسکرول صفحه #

 انتظار برای نمایش اطلاعات بعد اسکرول #

ما بقی بخش های کد مربوط به برداشت اطلاعات هر محصول است #

except:

t1 += 0.5

except:

```
print('The internet speed is very slow or the intended site is  
down')
```

که در آن برنامه با زمان $t=0.5$ ابتدا تلاش می کند تا اطلاعات صفحه را به صورت کامل نمایش دهد و اگر موفق نشد، نیم ثانیه به $t=1$ اضافه می نماید و این تکرار تا ۶ بار ادامه می یابد. در صورتی که صفحه به صورت کامل نمایش داده شود، از مقدار جدید $t=1$ برای صفحات بعد استفاده می نماییم. همچنین از این راه حل می توان برای صفحات مربوط به محصول نیز استفاده نمود، به طوری که بعد از کلیک بر روی صفحه محصول اگر صفحه قابل نمایش نباشد و یا خطأ مواجه شویم، می توانیم صفحه را با تابع `refresh()` به جای تابع `get()` مجدد به روزرسانی نموده و زمان تابع `sleep()` را هم بر اساس راه حل بالا افزایش دهیم تا صفحه محصول نیز به صورت کامل بارگذاری شود.

مرحله (۳) در صورت باز شدن صفحه وارد مرحله سوم می شویم که این مرحله مربوط به برداشت اطلاعات اولیه محصولات است؛ اما چه طور مانع بروز چالش نامتقارن بودن شویم؟ به عبارت دیگر چگونه اطمینان حاصل می شود که اطلاعات مربوط به تمام محصولات به صورت کامل بارگذاری شده و اطلاعات درست را برداشت خواهیم کرد.

در این مرحله با توجه به وجود پنج نوع اطلاعات مختلف (نظیر نام کالا، قیمت، میزان تخفیف، قیمت بعد از تخفیف، نوع تخفیف) از هر محصول، ابتدا پنج فهرست مختلف و خالی که نشان دهنده هر کدام از این پنج نوع است را به شکل ذیل تعریف می نماییم:

`Product_name` = [] که فهرست مربوط به نام محصول است.

`Product_price` = [] که فهرست مربوط به قیمت محصول است.

`Product_off` = [] که فهرست مربوط به میزان تخفیف محصول است.

`Product_price2` = [] که فهرست مربوط به قیمت محصول بعد از تخفیف است.

`Product_off_type` = [] که فهرست مربوط به نوع تخفیف است.

سپس با کمک دستورهای مربوط به برداشت اطلاعات، اطلاعات مربوط به هر محصول را برداشته و در فهرست مربوط به آن قرار می دهیم. به عبارت دیگر سلیم پس از خواندن اطلاعات موردنظر مانند نام کالا، آن را در متغیری ذخیره نموده و سپس آن را به فهرست

می‌افزاید. پس از چند با برداشت اطلاعات از بستر وب ممکن است اطلاعات وب‌سایت تغییر نماید. به عبارت دیگر، کدی که ما در روز اول نوشتیم و به خوبی هم کار می‌کرد، ممکن است طی هفته‌های آینده به خوبی کار نکند و اسکریپت نوشته شده را با خطأ مواجه نماید. همان‌طور که گفته شد ما در پژوهش‌های خود از برداشت نام محصول، قیمت، میزان تخفیف و نوع تخفیف از بستر وب استفاده نموده‌ایم و متوجه شدیم که بعد چند هفته تعداد آیتم موجود در فهرست‌ها باهم برابر نیست. به عبارت دیگر تعداد نام داروی برداشت شده از سایت دیجی کالا با تعداد نوع تخفیف برداشت شده برای محصول یکسان نیستند. به عنوان نمونه هنگامی که پس از اجرای کد و برداشت اطلاعات از بستر وب، فهرست مربوط به نام محصول شامل ۲۶۷۰ تعداد محصول مختلف، فهرست مربوط به قیمت شامل ۲۴۸۰ تعداد قیمت و فهرست مربوط به نوع تخفیف شامل ۱۷۸۰ تعداد نوع تخفیف باشد، ما با خطأ مواجه هستیم، چراکه تعداد اطلاعات موجود در فهرست‌ها باید باهم یکسان باشند. شاید سؤال پیش بیاید که تمام محصولات دارای تخفیف نیستند و درنتیجه تعداد اقلام موجود در فهرست تخفیف‌ها با تعداد اقلام موجود در نام کالا یکسان نیست ولی این مورد را ما از ابتدای نوشتن اسکریپت در نظر گرفته‌ایم. ولی باوجود در نظر گرفتن تمام این موارد باز هم ممکن است مواردی پیش بیاید که نمی‌توان آن را پیش‌بینی کرد. برای نمونه، در طول سال اغلب فروشگاه‌ها تخفیف‌های مختلفی را ارائه می‌دهند، مثل تخفیف‌های ویژه و یا شگفت‌انگیز ولی برخی از روزها نوع این تخفیف‌ها تغییر می‌کند، مثلاً بلک فرایدی (جمعه سیاه) که یک روز خاص در سال است، اغلب فروشگاه‌ها تخفیف‌های بیشتری را ارائه می‌دهند که اگر ما آن را در اسکریپت خود از قبل در نظر نگرفته باشیم، ممکن است کد با خطأ مواجه گردد و یا تعداد مقادیر موجود در فهرست‌ها باهم برابر نباشند و این عامل مانع از ایجاد خروجی اکسل یا SQL و درنتیجه مانع از بهروزرسانی پایگاه داده و برنامه اصلی می‌گردد. ما در این مقاله راه حلی را ارائه داده‌ایم که با کمک آن می‌توان از وقوع چنین مشکلات پیش‌بینی نشده‌ای جلوگیری نمود.

اسکریپت استفاده شده برای جلوگیری از رد شدن آیتم و عدم ثبت آن در فهرست:

```
check1 = [[product_prices], [product_off], [product_prices2],
[Product_off_type]]
```

```
for j in range(0,4):
```

```
    if len(product_name) < check1[len[j]]:
```

```
product_name.append('error')
elif len(product_name) > check1[len[j]]:
    check[j].append('error')
```

طبق بررسی های انجام شده، ممکن است پس از برداشت اطلاعات از بستر وب متوجه شویم که تعداد آیتم های موجود در فهرست های تعریف شده، مثلاً فهرست نام کالا، فهرست قیمت و ... باهم یکسان نیستند. درنتیجه کد بالا که در انتهای حلقه for مربوط به محصولات قرار دارد و بعد از افزودن مقدار به فهرست، بررسی می کند که آیا تعداد آیتم های موجود در فهرست ها باهم برابر است یا خیر، مانع از رخدادن چنین مشکلی می گردد که این کد با استفاده از ایجاد فهرستی به نام چک به این کار کمک می کند. ما در اسکریپت اصلی ابتدا کدهای مربوط به برداشت اطلاعات مربوط به یک محصول خاص را قرار داده ایم و بعد از کدهای اصلی، کد بالا را قرار دادیم. به عبارت دیگر، برنامه بعد از برداشت اطلاعات یک کالا و افزودن آنها به فهرست های اصلی از قبل تعریف شده، تعداد آیتم فهرست های موجود را با کمک تابع len محاسبه و آنها را در فهرست چک قرار می دهد. سپس با کمک حلقه for طول هر فهرست را با بقیه فهرست ها مقایسه و در صورت کم بودن تعداد آیتم های موجود در یک فهرست کلمه "error" را به آن می افزاید.

به منظور ارزیابی راه حل های ارائه شده در این پژوهش، ما اطلاعات مربوط داروهای مکمل سایت دیجی کالا که تقریباً ۳۰۰۰ داروی مختلف بودند را ۱۳ بار در ۱۳ تاریخ مختلف، از اول مهرماه ۱۴۰۲ تا ۲۵ اسفندماه ۱۴۰۲ مورد بررسی قرار داده ایم. در این بررسی، کدهای نوشته شده به منظور برداشت اطلاعات از سایت گفته شده، بار اول با راه حل پیشنهاد شده و بار دوم بدون راه حل پیشنهاد شده، اجرا شدند و در هر اجرا هم زمان برداشت اطلاعات و هم تعداد اطلاعات مربوط به هر فهرست یا همان ستون ثبت و با یکدیگر مقایسه شدند و میزان خطا و درصد آن نیز بر اساس تفاضل زمان برداشت اطلاعات در اجرای اول و دوم و تفاضل تعداد اطلاعات برداشت شده برای هر ستون در اجرای اول و دوم محاسبه گردید.

درنهایت نیز پس از اعمال دو اسکریپت بالا و بارگذاری اطلاعات در فایل موردنظر، احتمالاً ما با تعدادی سلوی یا فیلد که حاوی کلمه error هستند مواجه خواهیم بود که باید آنها را به صورت دستی بررسی و در صورت نیاز یا کد موردنظر را اصلاح نماییم و یا

اطلاعات را به صورت دستی اصلاح نماییم که این بستگی به هزینه، دشواری و زمان موردنظر برای اصلاح این موارد دارد.

یافته‌ها

پس از پیاده‌سازی و استفاده از راه حل ارائه شده، جدول ۱ حاصل گردید؛ که نشان می‌دهد دقیق و صحیح اطلاعات برداشت شده، نسبت به عدم استفاده از راه حل پیشنهادی، افزایش یافته و زمان برداشت اطلاعات نیز بهبود یافته است:

جدول ۱. تعداد اطلاعات برداشت شده از دیجی کالا و زمان موردنیاز برای برداشت آن‌ها

زمانی که از راه حل پیشنهادشده، استفاده می‌شود						هنگام عدم استفاده از راه حل پیشنهادی						تاریخ برداشت اطلاعات	ردیف
زمان برداشت اطلاعات به دقیقه	نوع تحفیف	قیمت بعد از تحفیف	تخفیف	قیمت	نام محصول	زمان برداشت اطلاعات به دقیقه	نوع تحفیف	قیمت بعد از تحفیف	تخفیف	قیمت	نام محصول		
۱۰۷	۲۶۷۰	۲۶۷۰	۲۶۷۰	۲۶۷۰	۲۶۷۰	۲۶۷	۲۶۶۳	۲۶۷۰	۲۶۷۰	۲۶۶۲	۲۶۷۰	-۰۷-۰۱ ۱۴۰۲	۱
۱۰۲	۲۵۴۲	۲۵۴۲	۲۵۴۲	۲۵۴۲	۲۵۴۲	۲۵۴	۲۴۹۳	۲۵۴۲	۲۵۴۲	۲۵۳۷	۲۵۴۲	-۰۷-۱۵ ۱۴۰۲	۲
۱۱۶	۲۸۹۸	۲۸۹۸	۲۸۹۸	۲۸۹۸	۲۸۹۸	۲۹۰	۲۸۹۸	۲۸۹۸	۲۸۹۸	۲۸۹۸	۲۸۹۸	-۰۸-۰۱ ۱۴۰۲	۳
۱۲۲	۳۰۴۲	۳۰۴۲	۳۰۴۲	۳۰۴۲	۳۰۴۲	۳۰۴	۲۹۴۳	۳۰۴۲	۳۰۴۲	۳۰۴۰	۳۰۴۲	-۰۸-۱۵ ۱۴۰۲	۴
۱۹۱	۲۸۶۰	۲۸۶۰	۲۸۶۰	۲۸۶۰	۲۸۶۰	۲۸۶	۲۵۲۴	۲۸۶۰	۲۸۶۰	۲۸۵۵	۲۸۶۰	-۰۹-۰۱ ۱۴۰۲	۵
۱۵۶	۳۱۲۳	۳۱۲۳	۳۱۲۳	۳۱۲۳	۳۱۲۳	۳۱۲	۲۸۲۳	۳۱۲۳	۳۱۲۳	۳۱۲۱	۳۱۲۳	-۰۹-۱۵ ۱۴۰۲	۶
۱۵۴	۳۰۷۰	۳۰۷۰	۳۰۷۰	۳۰۷۰	۳۰۷۰	۳۰۷	۲۹۴۳	۳۰۵۶	۳۰۵۶	۳۰۷۰	۳۰۷۹	-۱۰-۰۱ ۱۴۰۲	۷
۱۶۶	۲۹۲۶	۲۹۲۶	۲۹۲۶	۲۹۲۶	۲۹۲۶	۲۹۳	۲۸۵۶	۲۹۲۶	۲۹۲۶	۲۹۲۴	۲۹۲۶	-۱۰-۱۵ ۱۴۰۲	۸

سازوکاری برای مدیریت زمان و افزایش دقت اطلاعات هنگام استفاده...؛ تقیزاده کورایم و همکاران | ۲۱۹

۱۰۳	۲۵۶۷	۲۵۶۷	۲۵۶۷	۲۵۶۷	۲۵۶۷	۲۵۶۷	۲۵۷	۲۵۰۴	۲۵۶۷	۲۵۶۷	۲۵۶۷	۲۵۶۷	-۱۱-۰۱ ۱۴۰۲	۹
۱۴۵	۲۵۶۱	۲۵۶۱	۲۵۶۱	۲۵۶۱	۲۵۶۱	۲۵۶۱	۲۵۶	۲۵۶۱	۲۵۶۱	۲۵۶۱	۲۵۰۸	۲۵۶۱	-۱۱-۱۵ ۱۴۰۲	۱۰
۱۳۶	۲۷۲۶	۲۷۲۶	۲۷۲۶	۲۷۲۶	۲۷۲۶	۲۷۲۶	۲۷۳	۲۶۴۲	۲۷۲۶	۲۷۱۸	۲۷۲۶	۲۷۲۶	-۱۲-۰۱ ۱۴۰۲	۱۱
۲۰۷	۳۱۱۰	۳۱۱۰	۳۱۱۰	۳۱۱۰	۳۱۱۰	۳۱۱۰	۳۱۱	۲۸۶۲	۳۱۰۴	۳۱۱۰	۳۱۱۰	۳۱۱۰	-۱۲-۱۵ ۱۴۰۲	۱۲
۲۶۰	۳۱۲۴	۳۱۲۴	۳۱۲۴	۳۱۲۴	۳۱۲۴	۳۱۲۴	۳۱۲	۲۸۶۰	۳۱۲۴	۳۱۲۴	۳۰۸۹	۳۱۲۴	-۱۲-۲۵ ۱۴۰۲	۱۳
۱۹۶۴	۳۷۲۱۹	۳۷۲۱۹	۳۷۲۱۹	۳۷۲۱۹	۳۷۲۱۹	۳۷۲۱۹	۳۷۲۲	۳۵۰۷۲	۳۷۱۹۹	۳۷۱۹۷	۳۷۱۵۷	۳۷۲۱۸	مجموع	

جدول ۱ نشان می‌دهد که برنامه نوشته شده به منظور برداشت اطلاعات از بستر وب، هنگام استفاده از راه حل پیشنهادی و هنگام عدم استفاده از آن به تفکیک تاریخ، برای هر یک از فهرست‌ها یا همان ستون‌های نام محصول، قیمت، تخفیف، قیمت بعد از تخفیف و نوع تخفیف، چند برداشت موققیت‌آمیز داشته است. به عبارت دیگر، برنامه نوشته شده در تاریخ ۱۴۰۲-۰۷-۰۱ هنگام عدم استفاده از راه حل پیشنهادی، موفق شده تا ۲۶۷۰ نام محصول مختلف را برای ستون نام محصول و ۲۶۶۲ قیمت مختلف را برای ستون یا همان فهرست قیمت از سایت دیجی کالا بردارد. این در حالی است که هنگام استفاده از راه حل پیشنهادی مقادیر مربوط به ستون‌ها یا همان فهرست‌های نام محصول، قیمت و ... عددی برابر داشته و اطلاعات نام محصول، قیمت و ... ۲۶۷۰ داروی مکمل مختلف از بستر وب به صورت موققیت‌آمیز برداشت شده است. همچنین دو ستون «زمان برداشت اطلاعات به دقیقه» مدت زمانی که طول می‌کشد تا برنامه نوشته شده اطلاعات مربوط به تمام داروهای مکمل را از سایت دیجی کالا بردارد، نشان می‌دهد. این ستون‌ها نیز به تفکیک هنگام استفاده و عدم استفاده از راه حل پیشنهادی و تاریخ در جدول ۱ نشان داده شده‌اند. طبق بررسی‌های انجام شده به صورت دستی، عدد عنوان شده در ستون نام محصول هنگام استفاده از راه حل پیشنهادی، نشان‌دهنده تعداد واقعی و اصلی داروهای مکمل موجود در تاریخ‌های ذکر شده در سایت دیجی کالا بودند.

بر اساس جدول ۱ عدم استفاده از راه حل پیشنهادشده، موجب می‌گردد تا تعداد اطلاعات موجود در هر فهرست باهم یکسان نباشد و درنهایت در خروجی پایگاه داده نه تنها اطلاعات کامل نباشد، بلکه اطلاعات در جای درست خود نیز قرار نگیرند، برای نمونه تخفیف ارائه شده برای یک محصول در مقابل نام محصول دیگر قرار گیرد. همچنین زمان مربوط به تایع() نیز برای هر پیمایش و هر صفحه به صورت جداگانه عدد ثابت ۳ ثانیه در نظر گرفته شده است، درحالی که در راه حل پیشنهادشده این مقدار بر اساس ترافیک موجود در شبکه و سرعت اینترنت از نیم تا ۳ ثانیه متغیر است. اطلاعات به دست آمده پس از برداشت اطلاعات از بستر وب نشان می‌دهد که راه حل پیشنهادشده، هم دقت اطلاعات را افزایش داده و هم زمان برداشت اطلاعات از بستر وب را کاهش می‌دهد. همچنین در جدول ۲ تفاضل اطلاعات موجود در جدول شماره یک را نشان می‌دهد. برای نمونه همان‌طور که در جدول ۲ قابل مشاهده است، هنگام عدم استفاده از راه حل پیشنهادی در طول تقریباً ۶ ماه، درمجموع ۱۶۴۷ سطر اطلاعات از ستون نوع تخفیف خالی بوده و در صفحه وب بر اساس جدول یک ۳۵۵۷۲ سطر اطلاعات برای ستون نوع تخفیف یافت شده است، درحالی که در اصل ۳۷۲۱۹ سطر از اطلاعات یا همان داروی مکمل وجود داشته است. با این حال ستون نام محصول ثبات بیشتری داشته و در طول ۶ ماه، فقط یک مورد از آن توسط خزنده وب یا همان کتابخانه سلیمیوم، در هنگام عدم استفاده از راه حل پیشنهادی، یافت نشده است. همچنین زمان برداشت اطلاعات از بستر وب در هنگام استفاده از راه حل پیشنهادی ۱۷۵۸ دقیقه نسبت به عدم استفاده از آن، کمتر بوده است.

جدول ۲. قدر مطلق تفاضل هنگام استفاده از راه حل پیشنهادی و هنگام عدم استفاده از آن

قدر مطلق تفاضل هنگام استفاده از راه حل پیشنهادی و هنگام عدم استفاده از آن						تاریخ برداشت اطلاعات	ردیف
زمان برداشت اطلاعات به دقیقه	نوع تخفیف	قیمت بعد از تفاضل	تفاضل	قیمت	نام محصول		
۱۶۰	۷	۰	۰	۸	۰	۱۴۰۲-۰۷-۰۱	۱
۱۵۳	۴۹	۰	۰	۵	۰	۱۴۰۲-۰۷-۱۵	۲
۱۷۴	۰	۰	۰	۰	۰	۱۴۰۲-۰۸-۰۱	۳

سازوکاری برای مدیریت زمان و افزایش دقت اطلاعات هنگام استفاده...؛ تقیزاده کورایم و همکاران | ۲۲۱

۱۸۳	۹۹	۰	۰	۲	۰	۱۴۰۲-۰۸-۱۵	۴
۹۵	۳۳۶	۰	۰	۵	۰	۱۴۰۲-۰۹-۰۱	۵
۱۵۶	۳۰۰	۰	۰	۲	۰	۱۴۰۲-۰۹-۱۵	۶
۱۵۴	۱۲۷	۱۴	۱۴	۰	۱	۱۴۰۲-۱۰-۰۱	۷
۱۲۷	۷۰	۰	۰	۲	۰	۱۴۰۲-۱۰-۱۵	۸
۱۵۴	۶۳	۰	۰	۰	۰	۱۴۰۲-۱۱-۰۱	۹
۱۱۱	۰	۰	۰	۳	۰	۱۴۰۲-۱۱-۱۵	۱۰
۱۳۶	۸۴	۰	۸	۰	۰	۱۴۰۲-۱۲-۰۱	۱۱
۱۰۴	۲۶۸	۶	۰	۰	۰	۱۴۰۲-۱۲-۱۵	۱۲
۵۲	۲۶۴	۰-	۰	۳۵	۰	۱۴۰۲-۱۲-۲۵	۱۳
۱۷۵۸	۱۶۴۷	۲۰	۲۲	۶۲	۱	مجموع	

بحث و نتیجه‌گیری

در این مقاله ما چالش‌های کند بودن، زمان بر بودن و نامتقارن بودن کتابخانه سلنجیوم را مورد مطالعه و ارزیابی قراردادیم. مطالعات ما با استفاده از زبان برنامه‌نویسی پایتون و برداشت اطلاعات از سایت دیجی کالا در ۱۳ تاریخ مختلف از ۱۴۰۲-۰۷-۰۱ تا ۱۴۰۲-۱۲-۲۵ انجام شد. مطالعات نشان می‌دهد که استفاده از راه حل بررسی فهرست و یکسان بودن طول فهرست در پایان برداشت هر محصول از بستر وب بسیار مهم است، به طوری که عدم استفاده از آن در ۱۲ مورد از ۱۳ مورد برداشت اطلاعات از بستر وب، ما را با خطا مواجه نمود. همچنین استفاده از یک مقدار ثابت برای تابع sleep زمان برداشت اطلاعات را نسبت به زمان استفاده از مقدار متغیر برای آن به طور چشم‌گیری افزایش می‌دهد. در کل یافته‌ها نشان می‌دهد که استفاده از راه حل پیشنهادی در هنگام استفاده از کتابخانه سلنجیوم به منظور برداشت اطلاعات از بستر وب، کمک شایانی به افزایش دقت و صحت اطلاعات نموده و زمان برداشت کامل اطلاعات از بستر وب را نیز بهبود می‌دهد.

طبق بررسی‌های انجام شده و همان‌طور که در جدول ۲ نیز آمده است، عدم استفاده از کد مربوط به مدیریت زمان موجب می‌گردد تا در طول تقریباً ۶ ماه زمان یا تقریباً ۲۹ ساعت زمان بیشتری صرف برداشت اطلاعات از بستر وب گردد. به عبارت دیگر همان‌طور که در جدول ۱ آمده است، در صورت استفاده از راه حل پیشنهادی برای مدیریت زمان، در مجموع ۱۹۶۴ دقیقه زمان برای ۱۳ برداشت مختلف اطلاعات از سایت دیجی کالا

در طول تقریباً ۶ ماه موردنیاز است، این درحالی است که عدم استفاده از راه حل پیشنهادی و استفاده از یک مقدار ثابت برای تابع sleep() این زمان را به ۳۷۲۲ دقیقه افزایش می‌دهد. همچنین عدم استفاده از راه حل پیشنهادی برای مشکل نامتقارن بودن، موجب می‌گردد تا درمجموع و در طول تقریباً ۶ ماه و ۱۳ برشاشه مختلف، یک آیتم از اطلاعات مربوط به ستون محصول، ۶۲ آیتم از اطلاعات مربوط به ستون قیمت، ۲۲ آیتم از اطلاعات مربوط به ستون تخفیف، ۲۰ آیتم مربوط به ستون قیمت بعد از تخفیف و ۱۶۴۷ آیتم مربوط به ستون نوع تخفیف کمتر از میزانی که باید از بستر وب برداشت گردد، برداشت شود. درحالی که هنگام استفاده از راه حل پیشنهادی، مقادیر تمام ستون‌ها باهم یکسان بودند، این بدین معناست که اطلاعات در جای صحیح خود قرار گرفته‌اند و درصورتی که اطلاعات بخشی از سایت هنوز نمایش داده نشده و یا ارتباط بخشی از سایت تغییر یافته باشد، بهجای داده مربوط به آن آیتم، مقدار error درج گردیده است. درج کلمه error نه تنها موجب می‌گردد که اطلاعات در جای درست خود قرار گیرند، بلکه با توجه به بررسی‌های انجام‌شده اغلب کلمه error ثبت شده به دلیل تغییر ارتباط داخل وب‌سایت بوده است، درنتیجه می‌توان اسکریپت نوشته‌شده را با استفاده از آن‌ها اصلاح نمود و یا توسعه داد.

پیشنهادها

در این مقاله ما راه حلی را بهمنظور بهبود عملکرد کتابخانه سلنیوم ارائه نمودیم که بهموجب آن می‌توان زمان برداشت اطلاعات از بستر وب و چالش نامتقارن بودن را بهبود داد. با این حال، ما در مطالعات خود میزان و زمان پیمایش را نادیده گرفته‌ایم که خود یکی چالش‌های مهم در کند بودن کتابخانه سلنیوم است. بخصوص در صفحات فروشگاهی که بعد از پیمایش، صفحه مقداری به سمت بالا و پایین هدایت شده و برداشت اطلاعات و کلیک روی محصولات را با خطأ مواجه می‌نماید.

همچنین ما در پژوهش خود فقط پیرامون چالش‌های کندی سلنیوم، زمان بر بودن آن و نامتقارن بودن سلنیوم مطالعه نموده و پیشنهادهایی را بهمنظور بهبود آن‌ها ارائه نمودیم. با این حال، کتابخانه سلنیوم دارای چالش‌های دیگری نیز است که نیازمند مطالعه و بررسی بیشتر است. این چالش‌ها عبارت‌اند از ۱- شکنندگی؛ ۲- پوسته‌پوسته شدن؛ ۳- دشواری در

قابلیت نگهداری^۱؛ ۴- مرورگر کراس^۲؛ ۵- تجزیه و تحلیل شکست^۳؛ ۶- زیرساخت^۴؛ ۷- مقیاس پذیری^۵؛ ۸- ادعای پذیری^۶؛ ۹- مستندسازی^۷؛ ۱۰- پشتیبانی کردن^۸ که در مقاله لوata و همکاران (۲۰۲۳) به تفصیل پیرامون آنها صحبت شده است و می‌توان به منظور بهبود کتابخانه سلینیوم آنها را مطالعه و راههای مناسبی را جهت برطرف نمودن آنها پیشنهاد نمود. شایان ذکر است که بررسی اکثر این چالش‌ها نیازمند مورد مطالعاتی متفاوت با رویکردی متفاوت است که خارج از محدوده این مقاله است.

ORCID

Farnaz Taghizadeh Kourayem 	http://orcid.org/0000-0003-2676-2495
Mohammadreza Kabaranzad Ghadim 	http://orcid.org/0000-0002-4372-5226
Seyed Abdollah Amin Mousavi 	http://orcid.org/0009-0005-3052-5910

References

- Boeing, G., & Waddell, P. (2017). New insights into rental housing markets across the United States: Web scraping and analyzing craigslist rental listings. *Journal of Planning Education and Research*, 37(4), 457-476.
- Henry, K. (2021). Importance of web scraping in e-commerce and e-marketing. Available at SSRN 3769593.
- Krotov, V., & Tennyson, M. (2018). Scraping Financial Data from the Web Using the R Language. *Journal of Emerging Technologies in Accounting*, 15(1), 169-181.
- Krotov, V., Johnson, L., & Silva, L. (2020). Tutorial: Legality and ethics of web scraping.
- Leotta, M., García, B., Ricca, F., & Whitehead, J. (2023, April). Challenges of end-to-end testing with selenium WebDriver and how to face them: A survey. In *2023 IEEE Conference on Software Testing, Verification and Validation (ICST)* (pp. 339-350). IEEE.
- Neumann, M., Steinberg, J., & Schaer, P. (2017). Web-Scraping for non-programmers: Introducing OXPath for digital library metadata harvesting. *Code4Lib Journal*, 38.

¹. Maintainability

². Cross-browser

³. Failure analysis

⁴. Infrastructure

⁵. Scalability

⁶. Assertability

⁷. Documentation

⁸. Support

- Suganthi, V., & Varun, M. M. (2024). Automation Using Selenium. *International Journal Of Multidisciplinary Research In Science, Engineering And Technology*, e-ISSN:2582-7219.
- Teotia, H., Shishodia, G., Tyagi, E., Prakash, A., & Avasthi, S. (2023, April). Instagram Analysis and Activity Automation: Using Python and Selenium Automation Tools. In *2023 International Conference on Computational Intelligence, Communication Technology and Networking (CICTN)* (pp. 522-526). IEEE.
- Yuan, S. (2023). Design and Visualization of Python Web Scraping Based on Third-Party Libraries and Selenium Tools. *Academic Journal of Computing & Information Science*, 6(9), 25-31.

استناد به این مقاله: تقویزاده کورایم، فرناز، کابارانزاد قدیم، محمدرضا و موسوی، سید عبدالله امین. (۱۴۰۳). سازوکاری برای مدیریت زمان و افزایش دقت اطلاعات هنگام استفاده از کتابخانه سلیوم. *فصلنامه بازیابی دانش و نظام‌های معنایی*, ۱۱ (۴۱)، ۱۹۹-۲۲۵. DOI: 10.22054/jks.2024.80235.1660



Journal of Knowledge Retrieval and Semantic Systems is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License.

