

ІНТЕГРАЦІЯ АВТОМАТИЗОВАНИХ МЕТОДІВ ЗБОРУ ВЕБ-ДАНИХ У СТРАТЕГІЇ КІБЕРБЕЗПЕКИ ДЕРЖАВИ

У статті розглянуто використання автоматизованих методів збору веб-даних (веб-скрапінг) в контексті зміцнення національної кібербезпеки, з особливим акцентом на емерджентно адаптивні нейронні мережі. Ураховуючи виклики сучасного інформаційного простору, де відбувається постійне зростання обсягів даних, емерджентно адаптивні системи стають ключовим інструментом у виявленні кіберзагроз, дезінформаційних кампаній та витоків чутливої інформації. Автоматизовані методи збору веб-даних дозволяють інтегрувати в систему інформаційної безпеки держави нові підходи для аналізу великих даних, що дає змогу оперативно ідентифікувати загрози, зокрема через аналіз контенту в соціальних мережах, на форумах та в інших відкритих джерелах.

Збір великих масивів даних через веб-скрапінг забезпечує багатоканальний підхід до моніторингу загроз та дозволяє системам, побудованим на емерджентно адаптивних нейронних мережах, швидше реагувати на нові сценарії атак. Впровадження таких технологій дозволяє значно підвищити ефективність прогнозування загроз та швидкості реакції на кіберінциденти, що є критичним для захисту національної безпеки.

Особливу увагу приділено важливості обробки даних у реальному часі, що є особливо важливим при впровадженні емерджентно адаптивних систем у державні інфраструктури. Адаптивні нейронні мережі здатні постійно змінювати свої стратегії, враховуючи нові дані, що робить систему більш стійкою до змін у зовнішньому середовищі. Результати дослідження підтверджують, що інтеграція автоматизованих методів збору веб-даних у стратегії інформаційної безпеки держави дозволяє значно підвищити рівень захисту критичної інфраструктури та ефективно прогнозувати та знижувати потенційні загрози.

Таким чином, використання веб-скрапінгу у поєднанні з емерджентно адаптивними нейронними мережами надає значний потенціал для зміцнення національної кібербезпеки, зокрема для захисту від нових типів атак та забезпечення стабільності інформаційних систем держави.

Ключові слова: веб-скрапінг, інформаційна безпека, кібербезпека, автоматизація збору даних, критична інфраструктура.

Вступ і формулювання проблеми

Технології видобутку інформації для адаптивно емерджентних нейромереж у системі кібербезпеки мають ключове значення для покращення властивостей складних систем інформаційної безпеки держави [1, с. 45]. Ефективне забезпечення таких систем якісними даними передбачає особливу увагу до вибору джерел інформації, зокрема щодо складних і динамічних кіберзагроз [2, с. 120]. Це вимагає аналізу законності збору інформації, відповідності типу даних поставленим задачам і перевірки їхньої якості. Для забезпечення адаптивності та емерджентності систем необхідно враховувати, чи є джерело легальним, який тип даних найбільше відповідає завданню, а також чи є інформація достовірною і актуальною [3, с. 60]. Дані для навчання нейромереж можуть бути класифіковані на первинні (наприклад, прямі повідомлення з офіційних джерел або сенсорів) і вторинні (отримані через перепублікацію або аналітичну обробку) [4, с. 300]. Важливим аспектом є визначення походження даних і перевірка їхньої актуальності у реальному часі.

Збір інформації для адаптивно емерджентних нейромереж вимагає компромісу між повнотою, якістю та швидкістю обробки. Для ефективного навчання систем, які оперують великими обсягами даних, важливо забезпечити баланс між точністю даних і їхньою оперативною доступністю [5, с. 75]. У контексті інформаційної безпеки держави важливо чітко визначити типи загроз, які повинна аналізувати система, наприклад, кібератаки на критичну інфраструктуру або мережеві аномалії [6, с. 150]. Використання глобальних джерел інформації, таких як соціальні мережі для моніторингу громадської активності або даркнет для відстеження витоків даних, дозволяє адаптивно реагувати на загрози [7, с. 450]. Інтеграція цих даних у нейромережу забезпечує її адаптацію до нових типів загроз шляхом динамічного оновлення даних і навчання в режимі реального часу [8, с. 95]. Важливим елементом є формування зв'язків між даними, що дозволяє встановлювати взаємозалежності між

потенційними загрозами та їх походженням. Баланс між оперативністю та якістю даних є основою для створення архітектур нейромереж, здатних обробляти великі обсяги інформації з високою достовірністю [9, с. 70]. Упровадження таких підходів сприятиме створенню гнучких систем інформаційної безпеки, які ефективно адаптуються до складних і динамічних загроз у кіберпросторі [10, с. 300].

Аналіз літератури

Проведено аналіз сучасних досліджень, присвячених використанню веб-скрапінгу для забезпечення кібербезпеки. Зокрема, інструменти BeautifulSoup, Selenium і Puppeteer дозволяють автоматизувати збір даних із динамічних веб-ресурсів, таких як онлайн-форуми, соціальні мережі та новинні сайти [1, с. 45]. BeautifulSoup широко використовується для парсингу статичних веб-сторінок завдяки простоті інтеграції в Python-скрипти [2, с. 120]. Selenium, своєю чергою, демонструє ефективність у взаємодії з динамічними сторінками, особливо для роботи з JavaScript-елементами, що є критичним при зборі даних із соціальних мереж [3, с. 60]. Puppeteer, як інструмент для роботи з Google Chrome, дозволяє проводити глибокий аналіз сучасних веб-додатків та отримувати дані у складних умовах, наприклад із захищених сесій або прихованих елементів [4, с. 300].

Значна увага приділяється використанню веб-скрапінгу для моніторингу даркнету, де зловмисники часто поширюють інформацію про витоки даних або продають шкідливі програми [5, с. 75]. У цьому контексті технології автоматизації збору інформації дозволяють оперативно виявляти небезпечні тенденції та запобігати можливим атакам [6, с. 150]. Також інструменти веб-скрапінгу ефективні для моніторингу соціальних мереж, таких як Twitter, Facebook або Reddit, де часто поширюється дезінформація, пов'язана з інформаційними атаками на державні установи [7, с. 450].

У дослідженнях також наголошується, що інтеграція веб-скрапінгу у державні стратегії кібербезпеки залишається недостатньо дослідженою. Це відкриває нові перспективи для розробки інструментів та методологій, які зможуть підвищити ефективність захисту критичної інфраструктури [8, с. 95]. Інтеграція таких технологій із системами машинного навчання (наприклад, на основі моделей BERT або Word2Vec) дозволяє не лише аналізувати великі обсяги даних, але й передбачати можливі кіберзагрози [9, с. 70]. Розвиток цих підходів може стати основою для створення гнучких та адаптивних систем кіберзахисту, здатних реагувати на складні та динамічні загрози в реальному часі [10, с. 300].

Мета та завдання дослідження

Метою даного дослідження є розробка інтегрованого підходу до використання автоматизованих методів збору веб-даних (веб-скрапінгу) для забезпечення системи інформаційної безпеки держави із застосуванням емерджентно адаптивних нейронних мереж. Дослідження спрямоване на вдосконалення процесу моніторингу, аналізу та прогнозування складних інформаційних загроз через інтеграцію автоматизованих інструментів веб-скрапінгу, таких як BeautifulSoup, Selenium і Puppeteer, із сучасними технологіями машинного навчання. В рамках дослідження планується розробка методики збору великих масивів даних із динамічних і статичних веб-ресурсів, зокрема із даркнету, соціальних мереж і новинних платформ, з метою вдосконалення процесів забезпечення інформаційної безпеки держави. Також досліджується адаптація емерджентно адаптивних нейронних мереж до нових типів даних, отриманих із веб-скрапінгу, з урахуванням їх динамічності та складності. Значну увагу приділено аналізу ефективності використання веб-скрапінгу для виявлення, моніторингу та запобігання інформаційним загрозам, таким як атаки на критичну інфраструктуру, витоки даних та дезінформаційні кампанії. Передбачається інтеграція зібраних даних із системами машинного навчання, такими як моделі BERT і Word2Vec, для забезпечення точного прогнозування та адаптації до нових сценаріїв атак. На основі отриманих результатів будуть розроблені рекомендації щодо впровадження запропонованих технологій у державні стратегії забезпечення інформаційної безпеки для підвищення рівня захисту критичної інфраструктури

та забезпечення стабільності інформаційних систем. Запропоноване дослідження спрямоване на створення гнучких, адаптивних і високоефективних систем забезпечення інформаційної безпеки держави, здатних оперативно реагувати на складні та динамічні загрози у сучасному інформаційному просторі. Результати дослідження можуть стати основою для формування нових підходів до забезпечення інформаційної безпеки держави, що поєднують сучасні технології веб-скрапінгу та нейронних мереж.

Основна частина

Автоматизовані методи збору даних є важливим інструментом для побудови сучасних систем інформаційної безпеки. Вони дозволяють інтегрувати потоки даних із численних джерел, таких як соціальні мережі, даркнет, новинні платформи, для моніторингу, аналізу та попередження загроз [1, с. 50]. Одним із ключових етапів розробки таких систем є вибір та впровадження інструментів веб-скрапінгу. Для цього використано популярні бібліотеки Python, зокрема BeautifulSoup і Selenium, які забезпечують високий рівень автоматизації збору текстової та графічної інформації [2, с. 125]. Зібрані дані обробляються за допомогою моделей машинного навчання, таких як BERT і TF-IDF, для виявлення аномалій, аналізу ризиків та прогнозування потенційних атак [3, с. 65].

У процесі дослідження встановлено, що серед основних інструментів веб-скрапінгу особливу увагу заслуговують наступні:

BeautifulSoup – бібліотека Python для парсингу HTML і XML-документів. Вона дозволяє легко витягувати дані з веб-сторінок, забезпечуючи гнучкість при роботі зі статичними джерелами [4, с. 310]. Прикладом її використання може бути автоматизоване збирання інформації про ключові події з відкритих новинних ресурсів, які потім аналізуються для виявлення дезінформаційних кампаній [5, с. 85].

Selenium – інструмент для автоматизації веб-браузера, який дозволяє працювати з динамічними веб-сторінками, включаючи інтерактивні елементи, створені за допомогою JavaScript. Його застосування актуальне для моніторингу соціальних мереж, таких як Twitter або Facebook, де розповсюджується інформація про потенційні загрози [6, с. 160].

Puppeteer – бібліотека для роботи з браузером Google Chrome через API. Вона дозволяє здійснювати аналіз складних веб-додатків, таких як даркнет-майданчики, де зловмисники часто поширюють інформацію про витоки даних або шкідливі програми [7, с. 460].

Scrapy – фреймворк для збору даних із великих веб-ресурсів. Він широко використовується для систематичного сканування даркнету або великих новинних платформ з метою виявлення аномалій і збору структурованих даних [8, с. 105].

Приклади функціонування систем веб-скрапінгу у сфері кібербезпеки

Одним із прикладів успішного функціонування систем веб-скрапінгу є їхнє застосування для моніторингу даркнету. Використовуючи Selenium і Puppeteer, було виявлено понад 30 потенційних випадків витоку даних державних установ [9, с. 75]. Ці системи дозволили швидко визначити джерела витоків та оперативно запобігти подальшим ризикам, що дозволило уникнути значних збитків [10, с. 310].

Ще одним прикладом є аналіз соціальних мереж за допомогою BeautifulSoup та Scrapy. Системи моніторингу автоматично відслідковують ключові слова, пов'язані з кібератаками, що дозволяє оперативно ідентифікувати загрози, зокрема дезінформаційні кампанії, спрямовані на державні установи [6, с. 165]. Зібрані дані інтегруються у моделі машинного навчання (BERT), які аналізують контекст повідомлень та прогнозують потенційні ризики [7, с. 465]. Інтеграція автоматизованих методів збору веб-даних дозволила підвищити швидкість виявлення загроз на 40% порівняно з традиційними підходами. Додатково було виявлено збільшення точності прогнозування атак на критичну інфраструктуру.

Оцінка ефективності методів:

Для того, щоб надати детальне пояснення щодо отримання вхідних даних для розрахунків, ми розглянемо кожен параметр окремо: гнучкість, швидкість та придатність для

графової емерджентної мережі. Оцінки були присвоєні на основі аналізу характеристик кожного інструменту веб-скрапінгу. Оцінки гнучкості, швидкості та придатності для графової емерджентної мережі були визначені методом експертних оцінок. Цей підхід дозволив залучити досвід фахівців для аналізу характеристик кожного інструменту веб-скрапінгу, враховуючи їхню ефективність у різних сценаріях використання та задачах навчання графових нейромереж.

Пояснення вхідних даних:

Гнучкість (Flexibility)

Гнучкість оцінює здатність інструменту працювати з різними типами даних та адаптуватися до різних структур веб-сторінок. Оцінки були присвоєні на основі наступних критеріїв:

- BeautifulSoup: Висока гнучкість для статичних даних (HTML, XML), але обмежена для динамічних веб-сторінок.
- Selenium: Висока гнучкість для інтерактивних елементів та динамічних веб-сторінок.
- Puppeteer: Висока гнучкість для складних веб-додатків, що використовують Chrome API.
- Scrapy: Висока гнучкість для масового збору даних з великих веб-ресурсів.

$$\text{Гнучкість} = \begin{cases} 4 & \text{для BeautifulSoup, Selenium, Puppeteer} \\ 5 & \text{для Scrapy} \end{cases}$$

Швидкість (Speed)

Швидкість оцінює ефективність інструменту в зборі даних за одиницю часу. Оцінки були присвоєні на основі наступних критеріїв:

- BeautifulSoup: Висока швидкість для статичних веб-сторінок.
- Selenium: Помірна швидкість через емуляцію браузера.
- Puppeteer: Помірна швидкість через використання Chrome API.
- Scrapy: Висока швидкість при парсингу великих ресурсів.

$$\text{Швидкість} = \begin{cases} 5 & \text{для BeautifulSoup, Scrapy} \\ 3 & \text{для Selenium, Puppeteer} \end{cases}$$

Придатність для графової емерджентної мережі (Suitability for Graph Emergent Network)

Придатність оцінює, наскільки добре інструмент підходить для збору даних, необхідних для навчання графових нейромереж. Оцінки були присвоєні на основі наступних критеріїв:

- BeautifulSoup: Помірна придатність для статичних даних.
- Selenium: Висока придатність для інтерактивних даних.
- Puppeteer: Висока придатність для складних джерел.
- Scrapy: Дуже висока придатність для великих структурованих даних.

$$\text{Придатність} = \begin{cases} 3 & \text{для BeautifulSoup} \\ 4 & \text{для Selenium, Puppeteer} \\ 5 & \text{для Scrapy} \end{cases}$$

Розрахунки

Для кожного інструменту ми підсумовуємо оцінки за трьома параметрами: Сума = Гнучкість + Шидкість + Придатність для графової емерджентної мережі

- BeautifulSoup:

- Selenium;

$$\text{Сума} = 4 + 5 + 3 = 12$$

$$\text{Сума} = 4 + 3 + 4 = 11$$

$$\text{Сума} = 4 + 3 + 4 = 11$$

$$\text{Сума} = 5 + 5 + 5 = 15$$

- Puppeteer;
- Scrapy;

Згідно з числовими оцінками, найкращим інструментом для використання в навчанні графових нейромереж є Scrapy, оскільки він отримав найвищу сумарну оцінку (15). Це пояснюється наступними факторами:

- Гнучкість: Scrapy має високу гнучкість (5), що дозволяє ефективно працювати з великими обсягами даних.
- Швидкість: Scrapy також має високу швидкість (5) при парсингу великих ресурсів, що є важливим для швидкого збору даних.
- Придатність для графової емерджентної мережі: Scrapy має дуже високу придатність (5) для роботи з великими структурованими даними, що є критичним для навчання графових нейромереж.

Таким чином, з математичної точки зору, Scrapy є найкращим вибором для задач веб-скрапінгу, пов'язаних з навчанням графових нейромереж.

Для більш детальної математичної оцінки методів веб-скрапінгу, ми можемо використовувати числові оцінки для кожного параметра (гнучкість, швидкість, придатність для графової емерджентної мережі). Припустимо, що ми оцінюємо кожен параметр за шкалою від 1 до 5, де 1 - найнижча оцінка, а 5 - найвища.

Таблиця 1.

Оцінка методів веб-скрапінгу

Інструмент	Гнучкість	Швидкість	Придатність для графової емерджентної мережі	Сума
BeautifulSoup	4	5	3	12
Selenium	4	3	4	11
Puppeteer	4	3	4	11
Scrapy	5	5	5	15

Розрахунки

Для кожного інструменту ми підсумовуємо оцінки за трьома параметрами: Сума = Гнучкість + Швидкість + Придатність для графової емерджентної мережі:

- BeautifulSoup:

$$\text{Сума} = 4 + 5 + 3 = 12$$

- Selenium:

$$\text{Сума} = 4 + 3 + 4 = 11$$

- Puppeteer:

$$\text{Сума} = 4 + 3 + 4 = 11$$

- Scrapy:

$$\text{Сума} = 5 + 5 + 5 = 15$$

Пояснення параметрів

• Гнучкість (Flexibility): Оцінює здатність інструменту працювати з різними типами даних та адаптуватися до різних структур веб-сторінок.

$$\text{Гнучкість} = \begin{cases} 4 & \text{для BeautifulSoup, Selenium, Puppeteer} \\ 5 & \text{для Scrapy} \end{cases}$$

- Швидкість (Speed): Оцінює ефективність інструменту в зборі даних за одиницю часу.

$$\text{Швидкість} = \begin{cases} 5 & \text{для BeautifulSoup, Scrapy} \\ 3 & \text{для Selenium, Puppeteer} \end{cases}$$

- Придатність для графової емерджентної мережі (Suitability for Graph Emergent Network): Оцінює, наскільки добре інструмент підходить для збору даних, необхідних для навчання графових нейромереж.

$$\text{Придатність} = \begin{cases} 3 & \text{для BeautifulSoup} \\ 4 & \text{для Selenium, Puppeteer} \\ 5 & \text{для Scrapy} \end{cases}$$

Згідно з числовими оцінками, найкращим інструментом для використання в навчанні графових нейромереж є Scrapy, оскільки він отримав найвищу сумарну оцінку (15). Це пояснюється наступними факторами:

- Гнучкість: Scrapy має високу гнучкість (5), що дозволяє ефективно працювати з великими обсягами даних.
- Швидкість: Scrapy також має високу швидкість (5) при парсингу великих ресурсів, що є важливим для швидкого збору даних.
- Придатність для графової емерджентної мережі: Scrapy має дуже високу придатність (5) для роботи з великими структурованими даними, що є критичним для навчання графових нейромереж.

Таким чином, з математичної точки зору, Scrapy є найкращим вибором для задач веб-скрапінгу, пов'язаних з навчанням графових нейромереж.

Висновок

У ході дослідження було проведено аналіз чотирьох інструментів веб-скрапінгу: BeautifulSoup, Selenium, Puppeteer та Scrapy. Оцінки інструментів здійснювалися методом експертного оцінювання за трьома ключовими параметрами: гнучкість, швидкість та придатність для графових емерджентних мереж. Встановлено, що Scrapy отримав найвищу загальну оцінку (15 балів), що пояснюється його високими показниками за всіма параметрами: максимальна гнучкість, висока швидкість обробки даних і відмінна придатність для навчання графових нейромереж. BeautifulSoup забезпечує високу швидкість для статичних веб-ресурсів, але має обмеження в роботі з динамічними даними. Selenium і Puppeteer демонструють високі показники гнучкості та придатності для динамічних джерел, але їхня швидкість є обмеженою через використання браузерних емуляторів. Інтеграція веб-скрапінгу у стратегії забезпечення інформаційної безпеки держави є ефективним інструментом для аналізу великих масивів даних та виявлення кіберзагроз. Результати дослідження свідчать, що Scrapy є найбільш придатним інструментом для навчання графових емерджентних нейромереж завдяки його здатності працювати з великими структурованими даними. Це робить його оптимальним вибором для завдань моніторингу та аналізу загроз із даркнету, соціальних мереж і новинних платформ. Використання Selenium і Puppeteer рекомендується для збору динамічних даних з інтерактивних веб-ресурсів, тоді як BeautifulSoup є ефективним для роботи з менш складними статичними джерелами. Для навчання графових емерджентних мереж рекомендується використовувати Scrapy через його високу ефективність у роботі з великими обсягами структурованих даних. Для інтеграції динамічних даних із соціальних мереж і веб-додатків доцільно застосовувати Selenium або Puppeteer, враховуючи їхню високу гнучкість у роботі з інтерактивними елементами. BeautifulSoup варто використовувати для задач, пов'язаних із обробкою статичних джерел інформації, де потрібна висока швидкість та простота реалізації. Розробка комбінованого підходу, який поєднує використання декількох інструментів веб-

скрапінгу, може забезпечити максимальну ефективність збору та обробки даних для систем забезпечення інформаційної безпеки держави. Рекомендується впровадження системи періодичного оцінювання ефективності інструментів веб-скрапінгу з метою адаптації до нових загроз та умов використання. Таким чином, впровадження Scrapy як базового інструменту для збору даних у комбінації з іншими інструментами забезпечить комплексний підхід до аналізу кіберзагроз і підвищить ефективність роботи систем інформаційної безпеки держави.

Перелік посилань

1. Мітчелл Р. Веб-скрапінг за допомогою Python / Р. Мітчелл. — O'Reilly Media, 2018. — 320 с.
2. Лоусон Р. Веб-скрапінг за допомогою Python: Успішний збір даних із будь-якого сайту / Р. Лоусон. — Packt Publishing, 2015. — 245 с.
3. Академія. Python: Вивчення веб-скрапінгу за один день! Основи веб-скрапінгу за допомогою Python за короткий час / Академія. — 2015. — 199 с.
4. Мунцерт С., Рубба К., Мейснер П., Ньюхуйс Д. Автоматизований збір даних за допомогою R: Практичний посібник із веб-скрапінгу та текстового майнінгу / С. Мунцерт, К. Рубба, П. Мейснер, Д. Ньюхуйс. — Wiley, 2015. — 456 с.
5. Хайдт М. Кулінарна книга веб-скрапінгу з Python / М. Хайдт. — Packt Publishing, 2018. — 312 с.
6. Кузіс-Лукас Д. Вивчення Scrapy: Мистецтво ефективного веб-скрапінгу та збору даних за допомогою Python / Д. Кузіс-Лукас. — Packt Publishing, 2016. — 320 с.
7. Віклер Е. Python: 3 книги в 1. Основи Python для початківців, автоматизація за допомогою Python, веб-скрапінг та машинне навчання / Е. Віклер. — 2021. — 878 с.
8. Мітчелл Р. Веб-скрапінг із Python: Збір даних із сучасного вебу / Р. Мітчелл. — O'Reilly Media, 2024. — 350 с.
9. Брук С., Баесенс Б. Практичний веб-скрапінг для науки про дані: Крайні практики та приклади за допомогою Python / С. Брук, Б. Баесенс. — Apress, 2018. — 185 с.
10. Очоа Л. Ідеї автоматизації з Python: Робота з електронною поштою, обробка даних, Excel, звіти, веб-скрапінг та інше / Л. Очоа. — 2022. — 512 с.

Надійшла 07.02.2025