

ХЕРСОНСЬКИЙ НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ  
(повне найменування вищого навчального закладу)  
ФАКУЛЬТЕТ ІНФОРМАЦІЙНИХ ТЕХНОЛОГІЙ ТА ДИЗАЙНУ  
(повне найменування інституту, назва факультету (відділення))  
КАФЕДРА ПРОГРАМНИХ ЗАСОБІВ І ТЕХНОЛОГІЙ  
(повна назва кафедри (предметної, циклової комісії))

## Пояснювальна записка

до атестаційної роботи  
бакалавра  
(освітньо-кваліфікаційний рівень)

на тему: «Розробка веб-додатку для парсингу інтернет-магазинів»

Виконав: студент 4 курсу, групи 4ПР2  
напряму підготовки (спеціальності)  
121 -«Інженерія програмного забезпечення»  
(шифр і назва напряму підготовки, спеціальності)

Дембовський А.В.

(прізвище та ініціали)

Керівник \_д.т.н., професор Жарікова М.В.

(прізвище та ініціали)

Рецензент \_\_\_\_\_

(прізвище та ініціали)

Херсонський національний технічний університет \_\_\_\_\_  
( повне найменування вищого навчального закладу )

Факультет Інформаційних технологій та дизайну \_\_\_\_\_  
Кафедра \_\_\_\_\_ Програмних засобів і технологій \_\_\_\_\_

Освітньо-кваліфікаційний рівень \_\_\_\_\_ бакалавр \_\_\_\_\_

Галузь знань \_\_\_\_\_ 12 – Інформаційні технології \_\_\_\_\_

Спеціальність \_\_\_\_\_ 121 «Інженерія програмного забезпечення» \_\_\_\_\_

ЗАТВЕРДЖУЮ

Завідувач кафедри ПЗ і Т

д.т.н. проф. В.Г. Шерстюк \_

“ \_\_\_ ” \_\_\_\_\_ 2021 р.

## З А В Д А Н Н Я

### НА ДИПЛОМНИЙ ПРОЕКТ (РОБОТУ) СТУДЕНТУ

\_\_\_\_\_ Дембовському А.В. \_\_\_\_\_

1. Тема проекту (роботи) \_\_\_\_\_ Розробка веб-додатку для парсингу інтернет-магазинів \_\_\_\_\_

керівник проекту (роботи) \_\_\_\_\_ д.т.н. професор Жарікова М.В. \_\_\_\_\_,  
(прізвище, ім'я, по батькові, науковий ступінь, вчене звання)

затверджена наказом вищого навчального закладу від “ \_\_\_ ” \_\_\_\_\_ № \_\_\_\_\_

2. Строк подання студентом проекту (роботи) \_\_\_\_\_ 25.05.2021 \_\_\_\_\_

3. Вихідні дані до проекту (роботи) \_\_\_\_\_ ДСТУ з обробки інформації, літературні та періодичні джерела, матеріали практики \_\_\_\_\_

4. Зміст розрахунково-пояснювальної записки (перелік питань, які потрібно розробити) \_\_\_\_\_ 1. Аналіз предметної області та постановка завдання; 2. Алгоритми машинного навчання та методи дослідження; 3. Проектування та розробка програмного застосунку; 4. Хід дослідження ефективності алгоритмів \_\_\_\_\_

5. Перелік графічного матеріалу (з точним зазначенням обов'язкових креслень)

1. Мета, об'єкт, предмет та задачі дослідження;

## 6. Консультанти розділів проекту (роботи)

Розділ	Прізвище, ініціали та посада консультанта	Підпис, дата	
		завдання видав	завдання прийняв

7. Дата видачі завдання \_\_\_\_\_

## КАЛЕНДАРНИЙ ПЛАН

№ з/п	Назва етапів дипломного проекту (роботи)	Строк виконання етапів проекту (роботи)	Примітка
1	Відбір та вивчення літературних джерел	25.09.2021	
2	Аналіз стану вирішення завдання на сучасному етапі	30.09.2021	
3	Побудова концептуальної моделі	3.10.2021	
4	Розробка математичної моделі	10.10.2021	
5	Побудова алгоритму функціонування програмного продукту	15.10.2021	
6	Написання вихідного коду програми	01.11.2021	
7	Налагодження програмного коду	15.11.2021	
8	Оформлення пояснювальної записки	02.12.2021	

Студент \_\_\_\_\_ Дембовський А.В.  
( підпис ) (прізвище та ініціали)Керівник проекту (роботи) \_\_\_\_\_ Жарікова М.В.  
( підпис ) (прізвище та ініціали)

## РЕФЕРАТ

Об'єкт дослідження – веб-додаток для парсингу інтернет-магазинів.

В першому розділі проаналізовано предметну область за поставлено задачу дослідження.

В другому розділі описано використані алгоритми і технології парсингу, такі як SQL, CSV, RSS, XLS, JSON, Document Object Model.

В третьому розділі описано проектування та розробку веб-додатку.

В червертому розділі розділі описано практичну реалізацію веб-додатку.

Повний обсяг роботи 72 сторінки. Містить 21 рисунок, 1 таблицю, 2 додаток, 12 літературних джерел.

Ключові слова: парсинг, веб-сайт, посилання, контент, інформація, алгоритм, структура, веб-ресурс.

## ЗМІСТ

ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ ТА СКОРОЧЕНЬ	7
ВСТУП	8
РОЗДІЛ 1 АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ ТА ПОСТАНОВКА ЗАВДАННЯ	9
1.1 Поняття парсингу	9
1.2 Можливості, які надає парсинг	12
1.3 Парсинг для моніторингу цін інтернет-магазинів	14
1.4 Постановка завдання	21
РОЗДІЛ 2 ОПИС ВИКОРИСТАНИХ ТЕХНОЛОГІЙ ТА АЛГОРИТМІВ ПАРСИНГУ	22
2.1 Інструменти для написання парсера	22
2.2 Етапи парсинга	23
2.3 Імпорт контенту	25
2.4 Синтаксичний аналіз	26
2.5 Експорт даних	26
2.5.1 SQL	27
2.5.2 CSV	27
2.5.3 RSS	28
2.5.4 XLS	29
2.5.5 JSON	29
2.5.6 Document Object Model	30
2.5.7 Захист від парсингу та методи його обходу	31
РОЗДІЛ 3 ПРОЕКТУВАННЯ ТА РОЗРОБКА ВЕБ-ДОДАТКУ	36
3.1 Вибір програмних засобів	36
3.2 Завантаження веб-сторінок	37
3.2.1 Повторна спроба завантаження	38

3.2.2 Налаштування агента користувача	39
3.2.3 Підтримка проксі-сервера	41
3.2.4 Веб-сторінки з динамічним контентом	43
3.2.5 Взаємодія з формами	45
3.2.6 Обхід CAPTCHA	46
3.3 Аналіз та обробка даних	50
3.3.1 Аналіз структури веб-сторінок	50
3.3.2 Три підходи обробки веб-сторінок	51
3.4 Експорт даних у базу	55
РОЗДІЛ 4 ПРАКТИЧНА РЕЛІЗАЦІЯ	50
4.1 Проектування розробки програми	58
4.2 Розробка функціоналу програми	62
4.3 Розробка графічного інтерфейсу програми	63
ВИСНОВОК	64
ПЕРЕЛІК ВИКОРИСТАНИХ ДЖЕРЕЛ	66
ДОДАТОК А	68
ДОДАТОК Б	72

## ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ ТА СКОРОЧЕНЬ

Web – система доступу до пов'язаних між собою документів на різних комп'ютерах, підключених до Інтернету

HTML – HyperText Markup Language

URL – Uniform Resource Locator

HTTP – HyperText Transfer Protocol

XML – eXtensible Markup Language (розширювана мова розмітки)

БД – база даних

NoSQL – not only Structured query language

CSV - Comma-Separated Values

JSON - JavaScript Object Notation

DOM - Document Object Model

URL - Uniform Resource Locator

CAPTCHA - Completely Automated Public Turing test to tell Computers and Humans Apart

DFD - Data Flow Diagrams

## ВСТУП

Парсинг даних – це автоматичний збір інформації розміщеної у відкритому доступі в Інтернет та її систематизація за допомогою спеціальних програм. В цілому, парсер виконує ті ж функції, що й ручний збір інформації, тільки в рази швидше та зводячи до 0 можливість помилки.

Парсинг – це лінійне зіставлення послідовності слів з правилами мови. Поняття «мова» розглядається в найширшому контексті. Це може бути звичайна людська мова (наприклад, українська), яка використовується для комунікації людей. А може бути і формалізована мова, зокрема, будь-яка мова програмування.

Парсинг сайтів – послідовний синтаксичний аналіз інформації, розміщеної на Інтернет-сторінках. Текст Інтернет-сторінок – це ієрархічний набір даних, структурований за допомогою людських і комп'ютерних мов. Людською мовою надана інформація, знання, заради яких, власне, люди і користуються мережею Інтернет. Комп'ютерні мови (HTML, JavaScript, CSS) визначають, як інформація виглядає на моніторі, розмітку сторінки, її стиль.