Вивчення базових операцій обробки ХМС-документів

Метою роботи є здобуття практичних навичок створення програм, орієнтованих на обробку XML-документів засобами мови Python.

Завдання роботи полягає у наступному:

- 1. Виконати збір інформації зі сторінок Web-сайту за варіантом.
- 2. Виконати аналіз сторінок Web-сайту для подальшої обробки текстової та графічної інформації, розміщеної на ньому.
- 3. Реалізувати функціональні можливості згідно вимог, наведених нижче.

Функціональні вимоги

1. На основі базової адреси Web-сайту виконати обхід наявних сторінок сайту, відокремлюючи текстову та графічну інформацію від тегів HTML. Пошук вузлів виконувати засобами XPath. Наступну сторінку для аналізу **цього ж сайту** обрати як одне із гіперпосилань на даній сторінці (тег). Обмежитись аналізом 20 сторінок сайту. Зберегти XML у вигляді файлу. Формат XML-документу:

```
<page url="wwww.server.com/index.hml">
  <fragment type="text">
.... знайдений текст
 </fragment>
  <fragment type="image">
.... url зображення
 </fragment>
 </page>
 <page url="wwww.server.com/index1.hml">
  <fragment type="text">
.... знайдений текст
</fragment>
 <fragment type="image">
.... url зображення
</fragment>
 </page>
</data>
```

2. Виконати аналіз отриманих даних засобами XML згідно варіанту та вивести результати у консольне вікно. Відбір вузлів та розрахунки за варіантом виконувати засобами XPath.

Дисципліна «Бази даних. Частина 2» весна 2021 року

- 3. Проаналізувати вміст Web-сторінок інтернет-магазину (див. варіант). Отримати ціну, опис та зображення для 20 товарів з нього за допомогою DOM-парсеру та мови XPath для пошуку відповідних вузлів. Результат записати в XML-файл.
- 4. Перетворити отриманий XML-файл у XHTML-сторінку за допомогою мови XSLT. Дані подати у вигляді XHTML-таблиці та записати його у файл.

Вимоги до інтерфейсу користувача

Використовувати консольний (текстовий) інтерфейс користувача.

Вимоги до інструментарію

- 1. Мова програмування Python 3
- 2. Бібліотека сканування та отримання даних пошуковий робот Scrapy.
- 3. Бібліотека для генерації та перетворення XML libxml2 або (xml.dom, xml.xpath.. вбудовані пакети), lxml на вибір студента.
- 4. Середовище розробки програмного забезпечення PyCharm Community Edition (*опціонально*)

Вибір варіанту

Робота виконується індивідуально. Варіант обирається шляхом вибору останніх двох цифр номеру залікової книжки студента.

Варіанти

№ вар	Базова сторінка (завдання 1)	Зміст завдання 2	Адреса інтернет-магазину (завдання 3)
1.	www.kpi.ua	Максимальна кількість текстових фрагментів	www.rozetka.ua
2.	www.ukr.net	Середня кількість текстових фрагментів	www.repka.ua
3.	www.bigmir.net	Мінімальна кількість графічних фрагментів	www.sokol.ua

Дисципліна «Бази даних. Частина 2» весна 2021 року

4.	www.korrespondent.net	Кількість текстових	www.hotline.ua
		фрагментів по кожному	
		документу	
5.	www.football.ua	Кількість графічних	www.moyo.ua
		фрагментів по кожному	
		документу	
6.	www.isport.ua	Вивести список	www.portativ.ua
		гіперпосилань	
7.	www.tsn.ua	Мінімальна кількість	<u>www.wallet.ua</u>
		графічних фрагментів	
8.	www.golos.ua	Середня кількість	www.petmarket.ua
		текстових фрагментів	
9.	www.xsport.ua	Вивести список	www.meblium.com.ua
		гіперпосилань	
10.	www.ridna.ua	Кількість графічних	www.veliki.com.ua
		фрагментів по кожному	
		документу	
11.	www.ukraine-is.com	Вивести список	www.instrument.in.ua
11.		гіперпосилань	
12.	www.uartlib.org	Середня кількість	www.hozmart.com.ua
12.		текстових фрагментів	
13.	www.doroga.ua	Мінімальна кількість	www.freedelivery.in.ua
13.		графічних фрагментів	
14.	www.stejka.com	Вивести список	www.mebli-lviv.com.ua
14.		гіперпосилань	
15.	www.ua.igotoworld.co	Середня кількість	https://allo.ua
	<u>m</u>	графічних фрагментів	
16.	www.posolstva.org.ua	Максимальна кількість	www.odissey.kiev.ua
10.		текстових фрагментів	
17.	www.uahotels.info	Середня кількість	www.zvetsad.com.ua
		текстових фрагментів	
18.	www.osvita.ua	Мінімальна кількість	www.auto-store.kiev.ua
		графічних фрагментів	
19.	www.shkola.ua	Кількість текстових	www.tennismag.com.ua
		фрагментів по кожному	
		документу	

Дисципліна «Бази даних. Частина 2» весна 2021 року

20.	www.ostriv.in.ua	Вивести список	www.fishing-mart.com.
		гіперпосилань	<u>ua</u>

Вимоги до оформлення лабораторної роботи у електронному вигляді

Опис лабораторної роботи у репозиторії включає: назву лабораторної роботи, варіант студента, 2-3 копії екранних форм (screenshots).

Контрольні запитання

- 1. Дати визначення well-formed та коректному XML-документу.
- 2. Назвати способи обробки ХМL-документів.
- 3. Охарактеризувати мову XPath.
- 4. Охарактеризувати мову XSLT.