



МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ
НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ
“КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ
ІМЕНІ ІГОРЯ СІКОРСЬКОГО”

Факультет прикладної математики
Кафедра програмного забезпечення комп’ютерних систем

Лабораторна робота № 1
з дисципліни “*Бази даних. Частина 2*”
на тему “*Вивчення базових операцій обробки XML-документів*”

Виконав

студент III курсу

групи КП-81

Бухаленков Дмитро Олександрович

Варіант 2

Зарахована: Петрашенко А. В.

Київ 2021

Мета роботи: 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 за варіантом.

Варіант 2

(зайт змінено на kpi.ua тому що на ukr.net немає зображень і всі гіперпосилання ведуть на інші сайти)

Перший сайт	Завдання 2	Інтернет-магазин
kpi.ua	Середня кількість текстових фрагментів	repka.ua

Результати:

Код для “павука”, що вибирає дані з kpi.ua:

```
class KPISpider(Spider):
    name = 'kpi'
    start_urls = ['https://kpi.ua/']
    allowed_domains = ['kpi.ua']
    max_pages = 20
    filename = 'output/kpi.xml'

    def __init__(self, **kwargs):
        super().__init__(**kwargs)
        self.visited_pages = []

    def parse(self, response):
        if len(self.visited_pages) >= self.max_pages:
            raise exceptions.CloseSpider('Page limit exceeded.')
        if response.url not in self.visited_pages:
            self.visited_pages.append(response.url)
            yield self.parse_page(response)
        urls = Selector(response=response).xpath('//a/@href').getall()
        for url in urls:
            yield Request(url=urljoin(response.url, url),
                           callback=self.parse)
```

```

def parse_page(self, response):
    selector = Selector(response=response)
    text_data = selector.xpath('//h1[contains(@class,
\'page-title\')]/span/text() | //div[contains(@class, \'
\'node__content\')]/div/p').getall()
    images = selector.xpath('//img/@src').getall()
    return {
        'url': response.url,
        'text': [t.strip() for t in text_data],
        'images': [response.urljoin(src) for src in images]
    }

@classmethod
def get_average_texts_count(cls):
    return
etree.parse(cls.filename).xpath("count(//page/fragment[@type='text']) div
count(//page)")

```

Код для павука інтернет магазину repka.ua:

```

class RepkaSpider(Spider):
    name = 'repka'
    start_urls = ['https://repka.ua/products/noutbuki/?brands=73,']
    max_pages = 20
    filename = 'output/repka.xml'

    def parse(self, response):
        links =
Selector(response=response).xpath("//div[@class='product-item-name']/a/@hr
ef").getall()[:RepkaSpider.max_pages]
        for link in links:
            yield Request(url=urljoin(response.url, link),
callback=self.parse_laptops)

    def parse_laptops(self, response):
        selector = Selector(response=response)
        yield {
            'name':
selector.xpath("//h1[@class='page-title']/span//text()").get(),
            'price':
selector.xpath("//span[@class='price-wrapper']/@data-price-amount").get(),
            'image':
selector.xpath("//img[@class='fotorama__img']/@src").get(),
            'description':
selector.xpath("normalize-space(//div[@id='product_description']/div[@clas
s='box'])").get()
        }

    @staticmethod
    def create_xhtml():
        parsed = etree.parse(RepkaSpider.filename)
        template = etree.parse('template.xml')
        transform = etree.XSLT(template)

```

```
table = transform(parsed)
with open('output/table.xhtml', 'w') as f:
    f.write(etree.tostring(table, pretty_print=True).decode())
```

XSLT для трансформації:

```
<?xml version="1.0" encoding="UTF-8"?>
<xsl:stylesheet version="1.0"
xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
<xsl:template match="/">

    <html xmlns="http://www.w3.org/1999/xhtml">
        <head>
            <title>Laptops</title>
            <link rel="stylesheet"
href="https://stackpath.bootstrapcdn.com/bootstrap/4.5.2/css/bootstrap.min
.css"
                crossorigin="anonymous"/>
        </head>
        <body>
            <div class="container mt-3">
                <table class="table table-bordered">
                    <tbody>
                        <tr>
                            <th>Зображення</th>
                            <th>Назва</th>
                            <th>Ціна</th>
                            <th>Опис</th>
                        </tr>
                        <xsl:for-each select="//laptop">
                            <tr>
                                <td>
                                    
                                </td>
                                <td>
                                    <strong>
                                        <xsl:value-of select="@name"/>
                                    </strong>
                                </td>
                                <td>
                                    <xsl:value-of select="price"/>
                                    грн
                                </td>
                                <td>
                                    <xsl:value-of select="description"/>
                                </td>
                            </tr>
                        </xsl:for-each>
                    </tbody>
                </table>
            </div>
        </body>
    </html>

</xsl:template>
</xsl:stylesheet>
```

Повний код доступний в репозиторію.

```
2021-02-19 18:13:25 [scrapy.core.scrapers] DEBUG: Scraped from <200 https://kpi.ua/sitemap>
None
2021-02-19 18:13:25 [scrapy.core.scrapers] DEBUG: Scraped from <200 https://kpi.ua/site>
None
2021-02-19 18:13:25 [scrapy.spidermiddlewares.offsite] DEBUG: Filtered offsite request to 'telegra.ph': <GET https://telegra.ph/FV-in-KPI-05-18>
2021-02-19 18:13:25 [scrapy.core.engine] INFO: Closing spider (Page limit exceeded.)
2021-02-19 18:13:25 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://kpi.ua/index.php/institutes> (referer: https://kpi.ua/)
2021-02-19 18:13:25 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://kpi.ua/index.php/office> (referer: https://kpi.ua/)
2021-02-19 18:13:25 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://kpi.ua/index.php/rectorate> (referer: https://kpi.ua/)
2021-02-19 18:13:25 [scrapy.spidermiddlewares.offsite] DEBUG: Filtered offsite request to 'm.youtube.com': <GET https://m.youtube.com/watch?feature=youtu.be&v=m1p25qv2hs>
2021-02-19 18:13:25 [scrapy.spidermiddlewares.offsite] DEBUG: Filtered offsite request to 'council.science': <GET https://council.science/>
2021-02-19 18:13:26 [scrapy.spidermiddlewares.offsite] DEBUG: Filtered offsite request to 'zakon.rada.gov.ua': <GET https://zakon.rada.gov.ua/laws/show/1556-18#Text>
2021-02-19 18:13:26 [scrapy.spidermiddlewares.offsite] DEBUG: Filtered offsite request to 'web.telegram.org': <GET https://web.telegram.org/#/im?p=@vpi_kpiobot>
2021-02-19 18:13:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://kpi.ua/pbf> (referer: https://kpi.ua/kpi_faculty)
2021-02-19 18:13:26 [scrapy.spidermiddlewares.offsite] DEBUG: Filtered offsite request to 'instagram.com': <GET https://instagram.com/sss.ntuu.kpi>
2021-02-19 18:13:26 [scrapy.spidermiddlewares.offsite] DEBUG: Filtered offsite request to 'bit.ly': <GET https://bit.ly/2Yx58Bp>
2021-02-19 18:13:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://kpi.ua/ixf> (referer: https://kpi.ua/kpi_faculty)
2021-02-19 18:13:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://kpi.ua/fti> (referer: https://kpi.ua/kpi_faculty)
2021-02-19 18:13:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://kpi.ua/estimate> (referer: https://kpi.ua/information)
2021-02-19 18:13:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://pk.kpi.ua/specialities/> (referer: https://kpi.ua/)
2021-02-19 18:13:26 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://def.kpi.ua/taxonomy/term/49> from <GET http://def.kpi.ua/taxonomy/term/49>
2021-02-19 18:13:26 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://pk.kpi.ua/> from <GET http://pk.kpi.ua>
2021-02-19 18:13:26 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://kpi.ua/mmi> (referer: https://kpi.ua/kpi_faculty)
2021-02-19 18:13:26 [scrapy.downloadermiddlewares.redirect] DEBUG: Redirecting (301) to <GET https://def.kpi.ua/taxonomy/term/48> from <GET http://def.kpi.ua/taxonomy/term/48>
2021-02-19 18:13:26 [scrapy.statscollectors] INFO: Dumping Scrapy stats:
{'downloader/request_bytes': 10140,
```

Рис. 1 Scrapy парсить kpi.ua

```
61 <fragment type="image">https://kpi.ua/files/images-page/mon.jpg</fragment>
62 <fragment type="image">https://kpi.ua/files/images-page/solor_0.jpg</fragment>
63 <fragment type="image">https://kpi.ua/files/images-page/prozorro.png</fragment>
64 </page>
65 <page url="https://kpi.ua/kpi_faculty">
66 <fragment type="text">Інститути та факультети</fragment>
67 <fragment type="text">&lt;p&gt;В університеті працюють 16 факультетів, 11 навчально-наукових інститутів, декілька науко
68 <fragment type="text">&lt;p&gt;&lt;strong&gt;&lt;a href="/department"&gt;Наукові підрозділи&lt;/a&gt;&lt;/strong&gt;&lt;/p&gt;
69 <fragment type="image">https://kpi.ua/files/logo.png</fragment>
70 <fragment type="image">https://kpi.ua/files/styles/medium/public/kpi_0.png?itok=f86vIBI5</fragment>
71 <fragment type="image">https://kpi.ua/images/sikorsky-distance.jpg</fragment>
72 <fragment type="image">https://kpi.ua/files/images-page/solor_0.jpg</fragment>
73 <fragment type="image">https://kpi.ua/files/images-page/vnz.jpg</fragment>
74 <fragment type="image">https://kpi.ua/files/images-page/mon.jpg</fragment>
75 <fragment type="image">https://kpi.ua/files/images-page/isc_0.png</fragment>
76 <fragment type="image">https://kpi.ua/files/images-page/prozorro.png</fragment>
77 </page>
78 <page url="https://kpi.ua/about">
79 <fragment type="text">&lt;p&gt;Національний технічний університет України «Київський політехнічний інститут імені Ігоря
80 <fragment type="text">&lt;p&gt;Національний технічний університет України «Київський політехнічний інститут імені Ігоря
81 <fragment type="text">&lt;p&gt;В університеті працюють 16 факультетів, 11 навчально-наукових інститутів, декілька науко
82 <fragment type="text">&lt;p&gt;Ідея створення технічного навчального закладу виникла у цукрозаводчиків південно-західної
83 <fragment type="text">&lt;p&gt;Кампус КНІ ім.Ігоря Сікорського займає територію близько 120 гектарів, на якій органічно
84 <fragment type="text">&lt;p&gt;Професійна спілка (профспілка) – добровільна неприбуткова громадська організація, що об'єднує
85 <fragment type="image">https://kpi.ua/files/logo.png</fragment>
86 <fragment type="image">https://kpi.ua/files/styles/thumbnail/public/images-page/zgurovskiy.jpg?itok=b1P0TynM</fragment>
87 <fragment type="image">https://kpi.ua/files/styles/thumbnail/public/accreditation.jpg?itok=vwwCXf3W</fragment>
88 <fragment type="image">https://kpi.ua/files/styles/thumbnail/public/kpi_0.png?itok=pm0I1atG</fragment>
89 <fragment type="image">https://kpi.ua/files/styles/thumbnail/public/h-02.original_0.jpg?itok=NAptVWW1</fragment>
```

Рис. 2 результуючий XML файл

```
Average text fragments: 7.9
Press enter...
```

Рис. 3 за допомогою XPath підрахуно середню кількість текстових фрагментів по сторінках

```
2 <data>
3   <laptop name="LENOVO IdeaPad S145-15API (81UT00NRRRA)">
4     <price>9999</price>
5     <image>https://m1.repka.com.ua/pub/media/catalog/product/cache/c687aa7517cf01e65c009f6943c2b1e9/
6     <description>Современные ноутбуки должны быть не только производительными, но и достаточно удоб
7   </laptop>
8   <laptop name="LENOVO V15 (82C70010RA)">
9     <price>15555</price>
10    <image>https://m1.repka.com.ua/pub/media/catalog/product/cache/c687aa7517cf01e65c009f6943c2b1e9/
11    <description>Дисплей стандарта Full HD15.6-дюймовый дисплей с разрешением матрицы 1920x1080 точ
12  </laptop>
13  <laptop name="LENOVO V14 (82C400X6RA)">
14    <price>17243</price>
15    <image>https://m1.repka.com.ua/pub/media/catalog/product/cache/c687aa7517cf01e65c009f6943c2b1e9/
16    <description>Lenovo V14 идеально сочетает в себе стиль, мощность и компактность. Работайте быст
17  </laptop>
18  <laptop name="LENOVO V15 (82C700AKRA)">
19    <price>16999</price>
20    <image>https://m1.repka.com.ua/pub/media/catalog/product/cache/c687aa7517cf01e65c009f6943c2b1e9/
21    <description>Дисплей стандарта Full HD15.6-дюймовый дисплей с разрешением матрицы 1920x1080 точ
22  </laptop>
23  <laptop name="LENOVO IdeaPad 5 (81YH00NSRA)">
24    <price>15999</price>
```

Рис. 4 XML файл з даними про ноутбуки на сайті герка.ua

	LENOVO V15 (82C7H00AKRA)	16999 грн	Дисплей стандарта Full HD15.6-дюймовый дисплей с разрешением матрицы 1920x1080 точек и LED-подсветкой обеспечивает высокую яркость и четкость изображения. ПроизводительностьНаслаждайтесь любимыми играми и приложениями благодаря высокопроизводительному процессору Intel Core 10-го поколения и быстрому накопителю. Поддержка графики высокого разрешенияДобейтесь высочайшего качества изображения, подключив устройство к телевизору или к другому монитору с помощью разъема HDMI. Высокая скорость передачи данныхНоутбук оборудован USB 3.1, что позволяет обмениваться данными с другими устройствами со скоростью, которая в 10 раз превышает скорость интерфейсов USB более ранних версий. Встроенная веб-камераВстроенная веб-камера обеспечивает четкое видеоизображение, создавая эффект полного присутствия на веб-конференциях и в интерактивных видеочатах.
	LENOVO IdeaPad 5 (81YH00NSRA)	15999 грн	Забудьте о стикерах для заклеивания камеры. Защитная шторка веб-камеры ноутбука IdeaPad 5i (14) оградит вас от вездесущих хакеров. Кроме того, IdeaPad 5i (14) оснащен устройством распознавания отпечатков пальцев на кнопке питания. Ноутбук можно настроить так, что он не будет реагировать на взаимодействия, пока не считает вашу биометрическую подпись. При разработке дисплея для ноутбука IdeaPad 5i (14) мы стремились максимально увеличить площадь экрана для удобства пользователей, уменьшая его рамки и расширяя область просмотра. Мы называем это коэффициентом активной площади, или соотношением размеров экрана и окружающей его рамки. Ноутбук IdeaPad 5i (14) имеет невероятный коэффициент активной площади — 90 % льзования.
	LENOVO IdeaPad S145-15API (81UT00MFRA)	17599 грн	Всегда готовБлагодаря процессору AMD ноутбук IdeaPad S145 поможет вам решать любые задачи. Откройте в себе страсть к путешествиямМодель IdeaPad S145 весит всего 1.85 кг и поэтому идеально подходит для путешествий. Узкая рамка делает дизайн более лаконичным и увеличивает видимую область экрана. Ноутбук представлен в разных цветах с фактурной или глянцевой отделкой и при этом продается по привлекательной цене. Превосходный звук, отличное изображениеНезависимо от того, смотрите ли вы видео, слушаете потоковую музыку или общаетесь в видеочате, вам непременно понравится звучание, которое обеспечивают динамики IdeaPad S145 с поддержкой технологии Dolby Audio. Кроме того, вас также приятно удивит изображение на 15.6-дюймовом дисплее с матовым покрытием.
	LENOVO ThinkBook 14 G2	16399 грн	Максимальная производительность Максимально реализуйте вычислительную мощь мобильного процессора AMD Ryzen™ серии 4000 для решения самых ресурсоемких задач. Независимо от того, анализируете ли вы данные или

Рис. 5 Трансформований XHTML файл

Висновки: при виконанні данії лабораторної роботи я дізнався про XPath, XSLT та XHTML, здобув навички створення програм для обробки XML документів.