

Отчет по проекту Веб скрейпинг и парсинг сайтов

МАХНЕВ КИРИЛЛ

КОМАНДА: IFTISHACKERS

Содержание

- ▶ 1. Введение
- ▶ 2. Цели и задачи проекта
- ▶ 3. Основные понятия
- ▶ 4. Методы и инструменты веб скрейпинга
- ▶ 5. Исследование целевых веб-сайтов
- ▶ 6. Разработка скриптов для сбора данных
- ▶ 7. Обход препятствий и обеспечение безопасности
- ▶ 8. Тестирование и оптимизация скриптов
- ▶ 9. Результаты
- ▶ 10. Рекомендации по дальнейшему развитию проекта

Введение

- ▶ В нашем современном цифровом мире приобретает все большее значение умение извлекать ценные данные из веб-ресурсов для принятия обоснованных решений. В рамках данного проекта было проведено исследование и разработка скриптов для веб скрейпинга с целью сбора и анализа информации с различных интернет-ресурсов. В данном отчете будут представлены результаты работы, методы, используемые инструменты. Проект нацелен на эффективное получение и использование данных из сети Интернет с целью улучшения бизнес-процессов и принятия обоснованных решений.

Цели и задачи

▶ Цели проекта:

- ▶ 1. Получение актуальных данных с веб-сайтов для использования в аналитике.
- ▶ 2. Автоматизация процесса сбора информации с различных онлайн-ресурсов.
- ▶ 3. Создание дата-фрейма на основе полученных данных для дальнейшего анализа.

▶ Задачи проекта:

- ▶ 1. Исследование целевых веб-сайтов и определение структуры данных для скрейпинга.
- ▶ 2. Выбор подходящих инструментов и технологий для проведения веб-скрейпинга.
- ▶ 3. Написание скриптов для сбора, обработки и сохранения данных с веб-страниц.
- ▶ 4. Разработка алгоритмов для обхода препятствий, таких как капчи или ограничения скорости запросов.
- ▶ 5. Создание удобного интерфейса для визуализации и анализа полученных данных.

ОСНОВНЫЕ ПОНЯТИЯ

► Парсинг — это процесс синтаксического анализа текста с целью извлечения определенной информации. Скрапинг — это процесс загрузки веб-страницы и попытка извлечения информации из неё, обычно в обход API, ограничений и правил пользования сайтом

► Зачем нужен веб-скрапинг?

► Веб-скрапинг может использоваться для различных целей, таких как:

- сбор данных для анализа (например, статистика посещаемости сайтов)
- мониторинг цен на товары и услуги
- создание баз данных контента для исследований
- автоматизация рутинных задач, связанных с работой в интернете

HTML

- HTML (от английского HyperText Markup Language) — это язык гипертекстовой разметки текста. Он нужен, чтобы размещать на веб-странице элементы: текст, картинки, таблицы и видео.
- Теги (с англ. tag) — это метки, которые классифицируют контент и облегчают его поиск для пользователей.

ТЕГИ ДЛЯ ФОРМАТИРОВАНИЯ ТЕКСТА	
<code><pre></pre></code>	Обрамляет предварительно отформатированный текст.
<code><h1></h1></code>	Создает САМЫЙ БОЛЬШОЙ заголовок
<code><h6></h6></code>	Создает самый маленький заголовок
<code></code>	Создает жирный текст
<code><i></i></code>	Создает наклонный текст
<code><tt></tt></code>	Создает текст - имитирующий стиль печатной машинки.
<code><cite></cite></code>	Используется для цитат, обычно наклонный текст.
<code></code>	Используется для выделения из текста слова (наклонный или жирный текст)
<code></code>	Устанавливает размер текста в пределах от 1 до 7.
<code></code>	Устанавливает цвет текста, используя значение цвета в виде RRGGBB.
ГИПЕРССЫЛКИ	
<code></code>	Создает гиперссылку на другие документы или часть текущего документа.
<code></code>	Создает гиперссылку вызова почтовой программы для написания письма по указанному адресу.
<code></code>	Отмечает часть текста как цель для гиперссылок в документе.
<code></code>	Создает гиперссылку на часть текущего документа.

Методы и инструменты веб скрейпинга

- ▶ Библиотека `Requests`:
- ▶ Библиотека `Requests` для Python позволяет работать с HTTP-запросами любого уровня сложности, используя простой синтаксис. Это помогает не тратить время на написание кода, а быстро взаимодействовать с серверами.
- ▶ Библиотека `BeautifulSoup4`
- ▶ `BeautifulSoup` – это библиотека для парсинга HTML и XML документов в Python. Она предоставляет простой и удобный способ извлекать данные из веб-страниц, а также облегчает работу с этими данными. У библиотеки `BeautifulSoup` удобный интерфейс для взаимодействия с HTML-кодом, который позволяет легко находить нужные элементы и извлекать из них информацию. Эта библиотека является одной из наиболее популярных и широко используется для работы с парсерами и при анализе данных.



Исследование целевых веб сайтов

- ▶ Для успешной реализации проекта по веб скрапингу и парсингу было проведено исследование различных веб-ресурсов с целью выбора оптимального для скрапинга. Одним из ключевых сайтов, который был выбран в качестве объекта скрапинга, стал Avito.
- ▶ Почему Avito был выбран в качестве сайта для скрапинга:
 - ▶ - Широкий ассортимент данных: Avito предоставляет широкий спектр информации о товарах и услугах, что позволяет получить разнообразные данные для анализа.
 - ▶ - Популярность и актуальность: Avito является одним из самых популярных онлайн-площадок для объявлений, что обеспечивает доступ к актуальным данным.
 - ▶ - Структурированность данных: Сайт Avito хорошо структурирован, что упрощает процесс парсинга и извлечения требуемой информации.
- ▶ Выбор Avito в качестве целевого сайта для скрапинга обусловлен возможностью получения высококачественных данных, актуальности информации и удобством работы с веб-страницами данного ресурса.



Разработка скриптов для сбора данных

- ▶ Принцип работы кода:
- ▶ Код осуществляет скрапинг данных с страниц Avito, используя прокси и случайные задержки во избежание блокировки. Каждая страница парсится с помощью BeautifulSoup, чтобы извлечь информацию о цене, названии и описании объявлений об автомобилях. Полученные данные структурируются в DataFrame для последующего анализа или использования.

Рабочий код

```
1 import requests
2 import pandas as pd
3 from bs4 import BeautifulSoup
4 import time
5 import random
6
7 auto = []
8 prox = {'https': 'http://87.267.186.66:8080', 'http': 'http://88.222.245.88:80'}
9 head = {'Accept': '*/*',
10         'User-Agent': 'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/126.0.0.0 Safari/537.36'}
11 url = 'https://www.belta.by/moskva/avtomobili/baza-3808491CAU7utpiklyg?cf-id=1&radius=1000&searchRadius=1000&...'
12 page = 1
13 while page < 30:
14     get = requests.get(url + str(page), headers=head, proxies=prox)
15     soup = BeautifulSoup(get.text, 'html')
16     # Поиск пометки для парсинга (что она помечает себя помечу с атрибутом code=display: none в атрибуте style, что является одной из особенностей списка hidden)
17     for link in soup.select('a'):
18         if 'display' in link.get('style', 'default: ''') and 'none' in link['style']:
19             continue
20     auto_info = soup.find_all('div', class_='div-item-body-MlBuy')
21     if auto_info is not []:
22         auto.extend(auto_info)
23         # random задержку на выходе между записями чтобы уменьшить нагрузку на сайт
24         time_delay = random.random()
25         scaled_time_delay = 1 + (time_delay * (9 - 0))
26         time.sleep(scaled_time_delay)
27     else:
28         print("Больше данных нет")
29         break
30     page += 1
31     print(auto)
32     count = 0
33     lat_info = []
34     lat_price = []
35     lat_description = []
36     while count < len(auto):
37         price = elem.find('strong', class_='styles-module-root-BXMRa')
38         if price is not None:
39             price = price.text
40         else:
```

Рабочий код

```
1 let_info = []
2 let_price = []
3 let_description = []
4 while count < let(auto):
5     price = elem.find("strong", class="style-module-root-134K style-module-size_1_10_1K style-module-size_1_compensated-2WFF style-module-size_1_17n1 style-module-size_1_17n1")
6     if price is not None:
7         price = price.text
8     else:
9         price = "No info"
10    info = elem.find("p", class="name")
11    if info is not None:
12        info = info.text
13    else:
14        info = "No info"
15    description = elem.find("p", class="style-module-root-134K style-module-size_1_10_1K style-module-size_1_compensated-2WFF style-module-size_1_17n1 style-module-size_1_17n1")
16    if description is not None:
17        description = description.text
18    else:
19        description = "No info"
20    let_price.append(price)
21    let_info.append(info)
22    let_description.append(description)
23    count += 1
24 #F1 = pd.DataFrame({'Crosshair': let_price, 'Obama response': let_info})
25 #F2 = pd.DataFrame({'Rugby season': let_description})
26 print(F1)
27 print(F2)
```

1

	Стоимость	Общие сведения
0	7 420 000 Р	BMW X7 3.0 AT, 2021, 31 000 км
1	3 750 000 Р	BMW 4 серия 2.0 AT, 2017, 95 000 км
2	1 750 000 Р	BMW X6 M 4.4 AT, 2011, 171 000 км
3	13 500 000 Р	BMW X7 3.0 AT, 2022, 22 000 км
4	2 430 000 Р	BMW X1 1.5 AMT, 2020, 118 485 км
5	10 500 000 Р	BMW X7 3.0 AT, 2021, 64 000 км
6	7 830 000 Р	BMW X5 3.0 AT, 2019, 50 000 км
7	8 150 000 Р	BMW X7 3.0 AT, 2019, 48 000 км
8	8 590 000 Р	BMW X3 2.0 AT, 2024
9	10 430 000 Р	BMW 5 серия 2.0 AT, 2023
10	от 6 390 000 Р	BMW 3 серия 2.0 AT, 2024
11	6 750 000 Р	BMW X7 3.0 AT, 2019, 82 000 км
12	4 580 000 Р	BMW X3 3.0 AT, 2020, 126 000 км
13	1 750 000 Р	BMW X1 2.0 AT, 2014, 162 000 км
14	3 750 000 Р	BMW X6 3.0 AT, 2012, 210 829 км
15	5 700 000 Р	BMW X3 2.0 AT, 2021, 57 000 км
16	2 300 000 Р	BMW 5 серия 2.0 AT, 2012, 175 000 км
17	2 650 000 Р	BMW X2 2.0 AT, 2019, 165 000 км
18	1 345 000 Р	BMW X1 2.0 AT, 2012, 214 484 км
19	4 650 000 Р	BMW X3 2.0 AT, 2018, 48 000 км
20	3 990 000 Р	BMW X6 3.0 AT, 2015, 167 100 км
21	4 800 000 Р	BMW 6 серия GT 2.0 AT, 2021, 49 000 км
22	6 950 000 Р	BMW X5 3.0 AT, 2019, 112 648 км

Подробное описание

- 0 Автомобиль пригнан из Америки, с минимальными ...
- 1 Продан невероятно крутой авто.\nПродажа только ...
- 2 Продам свой автомобиль. В отличном состоянии, ...
- 3 Продаю X7 в идеальнейшем состоянии. Самая прих... ..
- 4 100% гарантия кристической чистоты.\nбез скрыти...
- 5 Продаю в отличном состоянии свой x7 был на кам... ..
- 6 Мораль x5 xdrive 30d Sport.\nпозвоните из Германи...
- 7 Дилерский автомобиль. Куплен в Рольф Премиум Х... ..
- 8 No info
- 9 No info
- 10 No info
- 11 Продаю отличный авто Бmw X-7, был привезен из ...
- 12 Bmw X3 gbl полный привод 3.0 tdi. Чистый 2020 ...
- 13 В продаже Bmw X1 xdrive28d!\n2014 год ресталин...
- 14 Bmw X6 M 50 D 3.0 AT (500 л. С).\n\nM Perform...
- 15 X3 200 в цвете black sapphire metallic.\nM-Spo...
- 16 Продан Bmw 525D Xdrive (218 л\с полный привод)...
- 17 В продаже X2 с оргл пробегом, пригнан из Герма...
- 18 Продаётся официальным дилером Bmw, Jаesso и
- 19 Не пригнан. Куплен в Дилерском центре Краснодар...
- 20 Продаётся Bmw X6,\nгод выпуска – 2015, из сало...
- 21 Продаю свой автомобиль Bmw 6 серии 2021 года н...
- 22 Bmw X5 3.0d 60S белый перламутр.\n\nлиб собствен...
- 23 В продаже Bmw 118i 2021 года выпуска. Я явлес...

Разбор кода

- ▶ 1. Импортирование необходимых библиотек:
- ▶ 2. Определение списка `auto`, который будет содержать информацию об автомобилях.
- ▶ 3. Создание словаря `proxies` с прокси-серверами для отправки запросов.
- ▶ 4. Создание словаря `head` с заголовками HTTP, включая `user-agent`.
- ▶ 5. Установка URL, который будет использоваться для парсинга объявлений на сайте Avito. Параметры URL включают страницу и радиус поиска.
- ▶ 6. Начало цикла `while`, который будет проходить по нескольким страницам объявлений на сайте Avito.
- ▶ 7. Отправка GET запроса к URL страницы с помощью `requests.get()` и передача заголовков и прокси.
- ▶ 8. Парсинг HTML страницы с помощью `BeautifulSoup` и поиск информации об объявлениях.
- ▶ 9. Получение информации об автомобилях и добавление её в список `auto`.
- ▶ 10. Добавление случайной задержки между запросами с помощью `time.sleep()`.
- ▶ 11. Проверка наличия данных об автомобилях на странице. Если данных нет, происходит остановка цикла.
- ▶ 12. Увеличение номера страницы для следующего запроса.
- ▶ 13. Создание списков `lst_info`, `lst_price` и `lst_description` для хранения данных из объявлений.
- ▶ 14. Цикл `while` для обработки каждого элемента из списка `auto`.
- ▶ 15. Для каждого элемента из списка `auto` происходит поиск информации о цене, общих сведениях и подробном описании.
- ▶ 16. Если информация есть, она добавляется в соответствующие списки, иначе добавляется "No info".
- ▶ 17. Создание двух `DataFrame` из списков `lst_info`, `lst_price` и `lst_description` с использованием библиотеки `pandas`.
- ▶ 18. Вывод на экран двух `DataFrame` с информацией об объявлениях.
- ▶ Этот код работает следующим образом: он отправляет запросы на сайт Avito, получает информацию об автомобилях с нескольких страниц, парсит данные и сохраняет их в таблицы, используя библиотеку `pandas`.

Обход препятствий и обеспечение безопасности

- ▶ Навигация по сайту может быть сложной из-за таких препятствий, как honeypot (ловушки для хакеров). Honeypot — это скрытые ссылки, специально разработанные для того, чтобы оставаться незамеченными обычными пользователями, однако быть обнаруженными парсерами и ботами. Эти ссылки часто скрываются с помощью HTML-элементов, для которых установлено значение hidden или none, или маскируются под кнопки, чей цвет соответствует фону страницы. Основная цель внедрения “медовых точек” — выявление и занесение ботов в черный список.
- ▶ Простой фрагмент кода, который можно использовать, чтобы попытаться избежать ловушек в Python:

```
for link in soup.select('a'):
    if 'display' in link.get('key: 'style', default: '') and 'none' in link['style']:
        continue
auto_info = soup.find_all(name="div", class_="iva-item-body-KLUuy")
if auto_info is not []:
    auto.extend(auto_info)
```

Обход препятствий и обеспечение безопасности

- ▶ Задержка по времени для снижения нагрузки на сайт

```
# сделаем задержку по времени между запросами чтобы уменьшить нагрузку на сайт
time_delay = random.random()
scaled_time_delay = 1 + (time_delay * (9 - 5))
time.sleep(scaled_time_delay)
else:
    print("Больше данных нет")
    break
file_in.write(res)
```

Тестирование и оптимизация скриптов

- ▶ Тестирование программы производилось с помощью парсинга первой страницы сайта и записи ее содержимого в файл. Далее производилась работа с содержимым этого файла и отладка скрипта именно с этим набором входных данных.

Тестирование и оптимизация скриптов

```
1 > Import ...
2
3 # Получение данных с первой страницы сайта и сохранения их в файл
4
5 #auto = []
6 #url = ["https://178.218.44.79:3128", "http://172.18.164.178:3125"]
7 #head = {"Accept": "*/*"}
8 #
9 # "User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/126.0.0.0 Safari/537.36"
10 #url = "https://www.avito.ru/moskva/avtomobili/bmw-A5g8AgITAN7gtg3klp4cd-14e-2d-radius-38866aaaarchadius-10886e-1"
11 #url = requests.get(url, headers=head, proxies=proxies)
12 #res = get_text
13 #with open("Users/user/Desktop/book3.html", "w") as file_in:
14 #    file_in.write(res)
15
16 auto = []
17 lst_info = []
18 lst_price = []
19 lst_description = []
20
21 with open("Users/user/Desktop/book3.html", "r") as file_read:
22     stran = file_read.read()
23
24 soup = BeautifulSoup(stran, parser="lxml")
25 auto_info = soup.find_all(name="div", class_="iva-item-body-MLBuy")
26 auto.extend(auto_info)
27
28 for elem in auto:
29     price = elem.find("strong", class_="style-module-root-biknd")
30     if price is not None:
31         price = price.text
32     else:
33         price = "No info"
34     info = elem.find("p", itemprop="name")
35     if info is not None:
36         info = info.text
37     else:
38         info = "No info"
39     description = elem.find("p", class_="style-module-root-VarKl style-module-size_s-xb_uk style-module-size_s_compensated-QWfF style-module-size_s_27nf style")
40     if description is not None:
41         description = description.text
42     else:
43         description = "No info"
44     lst_price.append(price)
45     lst_info.append(info)
46     lst_description.append(description)
47
48 df1 = pandas.DataFrame({"Сτοιимость": lst_price, "Общие сведения": lst_info})
49 df2 = pandas.DataFrame({"Подробное описание": lst_description})
```


Результаты

- Проект веб-скрапинга сайта Avito был успешно завершен. Были разработаны скрипты для автоматического сбора данных о объявлениях, таких как цена, модель автомобиля, год выпуска, пробег и дополнительная информация. Собранные данные были обработаны и организованы для анализа. Результатом работы стала подробная база данных с информацией о различных автомобилях, представленных на Avito. Полученные данные позволили провести анализ рынка автомобилей, выявить тенденции и сделать выводы о предложениях на платформе. В целом, проект по веб-скрапину сайта Avito был полезным для извлечения ценной информации и исследования рынка автомобилей.

Рекомендации по дальнейшему развитию проекта

- ▶ 1. Улучшить точность скрапинга: Необходимо постоянно мониторить и обновлять скрипты скрапинга, чтобы обеспечить высокую точность сбора данных.
- ▶ 2. Добавить функциональность: Рассмотрите возможность добавления функций, таких как автоматическое обновление данных, уведомления о новых объявлениях или фильтрация информации для более удобного анализа.
- ▶ 3. Исследование конкурентов: Проводить анализ других площадок с объявлениями о продаже автомобилей для сравнения предложений и цен, что поможет определить конкурентное положение.
- ▶ 4. Визуализация данных: Рассмотреть создание дашбордов или отчетов на основе собранных данных для более наглядного и удобного представления информации.
- ▶ 5. Расширение сферы применения: Подумать о возможностях использования аналитических данных не только для внутренних целей, но и для предоставления информации сторонним пользователям или компаниям по запросу.