

Scrapping

Internet rvest selenium
web technologies

Чувакин Сергей

R meetup

7 декабря 2019 г.

Internet

WWW vs. Internet ?

Internet

- Серверы
- Клиенты
- Третья сторона...

HTML - язык разметки, который отвечает за содержание страницы.

CSS - язык стилей, который отвечает, за то, как будет выглядеть страница.

Пример: ...



Интернет запрос - структурированная конструкция позволяющая получать содержимое страниц с удаленных серверов.

Сценарий: request - response.

- headers
- body
- cookies
- data
- parameters (params)
- cache

◀ ▶ ⏪ ⏩ 🔍 ↺

body

Методы:

- GET
- HEAD
- POST
- PUT
- PATCH
- DELETE
- etc

Request URL: <https://stepik.azureedge.net/static/frontend/cli-build/critical.css?v=1575568415>

Request method: GET

Status code: 200 OK ?

Version: HTTP/2.0

Edit and Resend



cookies

▼ Response cookies

▼ csrftoken:

expires: 2020-12-04T22:35:26.000Z

path: /

value: UcWluP9uAaxsRQdpQq3pVSxGalDNvew49bpTFRk8CSXh4byAykHHgjc096yUCWU2

▼ Request cookies

csrftoken: UcWluP9uAaxsRQdpQq3pVSxGalDNvew49bpTFRk8CSXh4byAykHHgjc096yUCWU2

sessionId: pbcwdu9jblbcw92wn60iyqaj4z6dzp7

params

▼ Query string

client: ubuntu

channel: fs

biw: 1920

bih: 970

sxsrfr: ACYBGNSCJNhbWK9LR0Twtm2E1UUTH3DhQ:1575671808262

ei: ANjqXd_ZD_Glk74Pi-uz0A0

q: web scraping memes

oq: web scraping memes

gs_l: psy-ab.12..35i39.0.0..977...0.0..1.175.312.0j2.....0.....gws-wiz.RM5vqfMy6to

ved: 0ahUKEwjf2afziqLmAhhXxxcQBHYv1DNoQ4dUDCAo

Cache

▼ Cache

Last Fetched 12/7/2019 1:36:28 AM

Fetch Count 3

Data Size 2325

Last Modified 12/7/2019 1:36:29 AM

Expires 12/7/2019 4:04:22 AM

Device Not Available



Что еще?

Что еще?

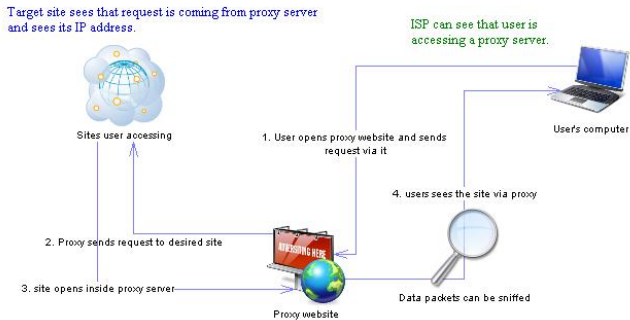
- Протоколы
- Архитектура DOM
- MVC фреймворки
- frontend-backend
- Sessions
- Работа с СУБД
- etc

CANT SCRAPE WEBSITE

WEBSITE USES AJAX

Navigation icons: back, forward, search, etc.

1



vpn

Target site sees that request is coming from VPN computer and sees its IP address.

