# Web Data Mining
## Lecture 3: Data Access and Acquisition Methods 2

**Jaroslav Kuchař & Milan Dojčinovski**
jaroslav.kuchar@fit.cvut.cz, milan.dojchinovski@fit.cvut.cz

Czech Technical University in Prague - Faculty of Information Technologies - Software and Web Engineering

OPPA EVROPSKÁ UNIE

Summer semester 2019/2020
Humla v0.3

---

## Overview

- Parsing and Extracting
- Web Scraping
- Accessing Specific Content

## Data Types

- Unstructured
  - *Lacks any structure*
  - *e.g. free texts, documents, multimedia, ...*
- Semistructured
  - *Inconsistent structure*
  - *Often self-describing (e.g. key-value pairs), flexibility*
  - *e.g. CSV, XML, ...*
- Structured
  - *Conform to the predefined model*
  - *Predefined schema*
  - *Highly structured*
  - *e.g. numbers, dates, structured entities, ...*

## Recall: Parsing

- Parsing the content of the HTTP payload
  - *extracting content for indexing*
  - *extracting links to be added to the frontier*
  - *extracting additional crawling and indexing directives*
  - *headers* Cache-Control, Content-Type, X-Robots-Tag, ...
- HTML code very often contains invalid markup
  - *unclosed elements, unencoded special characters, missing required attributes, improperly nested tags, missing quotes, ...*
- Bad HTML markup should be fixed
  - *a preprocessing step is required co clean up the HTML*
  - *many tools available,* tidy *- a tool provided by W3C*

## HTML markup fix

- Tidy:

```
1  from tidylib import tidy_document
2  document, errors = tidy_document('''<p>f&otilde;o <img src="bar.jpg">''',
3      options={'numeric-entities':1})
4  print(document)
5  print(errors)
```
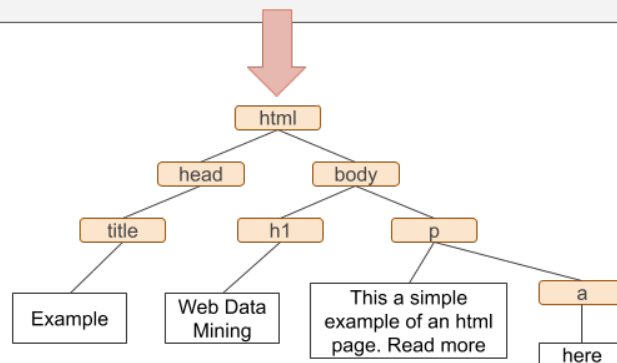
- Output:

```
1   <!DOCTYPE html PUBLIC "-//W3C//DTD HTML 4.01//EN">
2   <html>
3     <head>
4       <title></title>
5     </head>
6     <body>
7       <p>
8         fõo <img src="bar.jpg" alt="">
9       </p>
10    </body>
11  </html>
```

```
1   line 1 column 1 - Warning: missing <!DOCTYPE> declaration
2   line 1 column 1 - Warning: inserting missing 'title' element
3   line 1 column 18 - Warning: <img> lacks "alt" attribute
```

## Illustration of the DOM tree model

- Parsing of a simple HTML page

# Extraction of Relevant Information

- Web Information Extraction
  - *Natural language text processing*
  - *Extracting structured data from Web pages*

- Basic approaches
  - *Text fields identification*
    - → *Title, meta information, text blocks*
    - → *Each can have different importance (e.g. HTML title, h1, h3, …)*
  - *Anchors' texts extraction*
    - → *They are acting as a short description of the target.*
    - → *Important for search engines.*

- Detection of relevant structures/information
  - *Recognizers*
  - *Wrappers*
  - *Automatic approaches*

- Other approaches
  - *Machine readable annotations, JS variables, API calls, …*

# Recognizers

- Follow a procedure to find a piece of information based on its appearance
  - *e.g. e-mail addresses, phone numbers or street addresses*

- Most of nowadays crawlers and search engines are able to automatically collect such information.

- Can be executed as a set of regular expression patterns:

```python
import re
phonePattern = re.compile(r'^(\d{3})-(\d{3})-(\d{4})$')
phonePattern.search('800-555-1212').groups()
# ('800', '555', '1212')
```
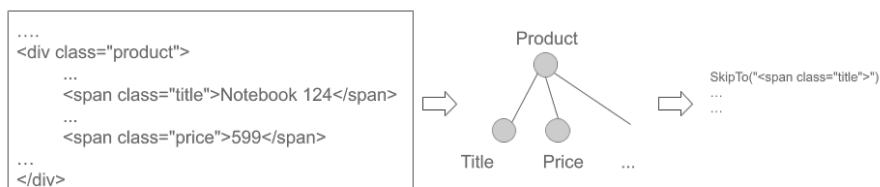
```python
import re
emailPattern = re.compile(r'([a-zA-Z0-9._%+-]+)@([a-zA-Z0-9.-]+\.[a-zA-Z]{2,
emailPattern.search('user@sub.example.com').groups()
# ('user', 'sub.example.com')
```

# Wrappers

- Allows to a semi-structured web data source be consulted as if it was a common database.

- Approaches
  - *Manual*
    - → *Observing the source code, finding patterns and writing the crawler*
    - → *Not scalable*
  - *Induction*
    - → *Supervised learning approach, semiautomatic*
    - → *Set of manually labeled pages*
  - *Automatic*
    - → *Unsupervised approach*
    - → *Automatically finds patterns (e.g. repetitive structures, visual aspects, ...)*
    - → *Can scale*

# Wrapper Induction Approaches

- Learning from a set of labeled instances
  - *Manual annotations of informations that are required for extraction.*
  - *e.g. product name, price and associated attributes*

- Approach
  - *Tree based representations of annotated information*
    - → *Possibilities for generalizations/pruning*
  - *Learning of extraction rules:*
    - → *e.g.* SkipTo("<span class="title">")



- Limitations
  - *Manual annotation is not possible in large scale.*
  - *Web is dynamic. Cost of maintenance.*

## Automatic Extraction

- Based on mining repeated patterns in multiple structures.
- Uses matching approaches
  - *String matching*
    - → *e.g. edit distances - minimum number of changes to change one string to another one ("Page 1" vs "Page 2", "Notebook" vs "CPU")*
  - *Tree matching*
    - → *minimum set of operation to transform trees*
- Building and processing trees
  - *Pure DOM based tree and tree pruning*
  - *Visual aspects can influence the importance of parts*
    - → *e.g. size, colors, visibility, ...*
- Patterns
  - *Repetitive tree structures with the same content can represent headers, footers, menus or ads*
  - *Varying subtrees across multiple pages identify main content*
    - → *High density of textual nodes can represent the main textual content*
  - *Structures close to each other but different colors represents different things*
  - *...*
- Tools
  - *boilerpipe, readability,*
  - *diffbot - analyze the visual layout, dragnet - rule-based models, octoparse*
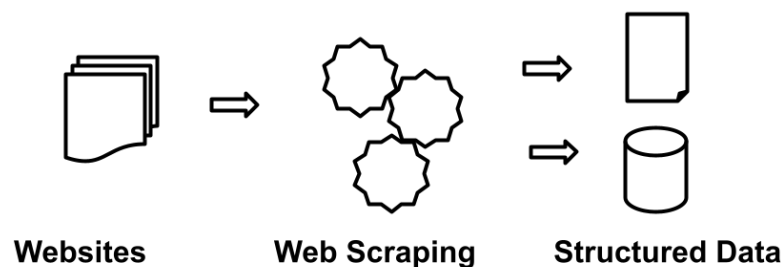  - *mercury, fathom*

## Overview

- Parsing and Extracting
- Web Scraping
- Accessing Specific Content

# Web Scraping

- Web Scraping vs Crawling
  - *Web Crawling*
    - → *usually provides data for indexing, search engines etc.*
    - → *following links, finding new pages and collecting content*
  - *Web Scraping*
    - → *for pages without API*
    - → *extracting structured information from a web page*
    - → *usually specific for the target website and structured data*
    - → *more focused process*
  - *Many overlapping actions*
    - → *execute JavaScript, emulate human user behavior, submit forms, log in to a website etc.*

# Web Scraping (cont.)

- Main process
  1. *Pick a URL*
     - *e.g. one from the predefined or collected list of URLs*
  2. *Collect the URL content*
     - *HTTP request/response*
  3. *Parse and extract desired information from the structure*
     - *using predefined rules or automatic detection mechanisms*



**Websites**      **Web Scraping**      **Structured Data**

# Web Scraping Example - Fetching Data

- Sending HTTP requests

```python
1  import requests
2
3  url="https://en.wikipedia.org/wiki/Web_scraping"
4
5  # Make a GET request to fetch the raw HTML content
6  html_content = requests.get(url).text
7
8  print(html_content)
```

- Output:

```html
1  <!DOCTYPE html>
2  <html class="client-nojs" lang="en" dir="ltr">
3  <head>
4  <meta charset="UTF-8"/>
5  <title>Web scraping - Wikipedia</title>
6  ...
```

# Parsing Data - Beautifulsoup

- Beautifulsoup
  - *Python library for pulling data out of HTML and XML files*

- Parsers
  - lxml
    - → *fast*
    - → *parses broken HTML*
    - → *external dependency*
  - html.parser
    - → *internal dependency in Python*
  - html5lib
    - → *parses broken HTML*
    - → *external dependency*

- Parsing
  - *navigating HTML tree*
    - → children, next_siblings, previous_sibling, parent
  - *selecting elements*
    - → find, find_all
  - *using CSS selectors, ...*

# Web Scraping Example - Parsing Data

```
1  # Parse the html content
2  soup = BeautifulSoup(html_content, "lxml")
3  # soup = BeautifulSoup(html_content, "html.parser")
4  print(soup.prettify()) # print the parsed data of html
5
6  # accessing DOM tree elements
7  print(soup.title)
8  print(soup.title.text)
9
10 # finding elements
11 print(soup.find_all('section', class_='categories'))
12 print(soup.find(id='product-price').children)
13
14 # using CSS selectors
15 print([e.text for e in soup.select("h2:has(> span#See_also) + div > ul > li
```

```
1  ....
2
3  ['Archive.is', 'Comparison of feed aggregators', 'Data scraping',
4  'Data wrangling', 'Importer', 'Job wrapping', 'Knowledge extraction',
5  'OpenSocial', 'Scraper site', 'Fake news website', 'Blog scraping',
6  'Spamdexing', 'Domain name drop list', 'Text corpus',
7  'Web archiving', 'Blog network', 'Search Engine Scraping',
8  'Web crawlers']
```

# Web Scraping Example - Parsing Data 2

```
1  import re
2  # Parse the html content
3  soup = BeautifulSoup(html_content, "lxml")
4  # soup = BeautifulSoup(html_content, "html.parser")
5
6  soup.find_all('a', {'href':re.compile('*.html') }).attrs['href']
7
8  soup.find_all(lambda t: len(t.attrs)==3)
```

```
1  ....
2
3  [ ... ]
```

## Overview

- Parsing and Extracting
- Web Scraping
- Accessing Specific Content

## Crawling Deep Web

- Deep Web
  - *The content hidden for crawling engines - e.g. behind HTML forms.*
  - *Usually user databases, registration-required web forums, web mail pages, online banking, pages behind paywalls (video on demand), ...*

https://www.cisoplatform.com/profiles/blogs/surface-web-deep-web-and-dark-web-are-they-different

## Crawling Deep Web (cont.)

- Categories
  - *Contextual web - content dependent on the context (e.g. location)*
  - *Dynamic pages - response to a query (e.g. search)*
  - *Restricted access - technical restrictions (e.g. robots.txt, CAPTCHA)*
  - *Non-HTML content - multimedia or other file formats*
  - *Private Web - login required pages*
  - *Scripted content - AJAX, Flash, ...*
  - *Unlinked - no backlinks*
  - *Hidden - using only specific software, protocols (e.g. Tor)*

## Crawling Dynamic pages

- Detection of forms
  - *e.g. search form*

```
1  <input type="search">
```

- Query generation
  - *pre-computing submissions for each HTML form*
    → *Using known entities from a knowledge base or search engine query logs*

- URL generation
  - *replacing parameters in URLs*
    → *e.g. https://www.google.com/search?q=deep+web*

- Empty page filtering
  - *for detection of irrelevant generations*

- Deduplications of results

## Selenium

- Selenium
  - *primarily for automating web applications for testing purposes*
  - *can be used for scraping/crawling as well*

```python
# Example for forms
from selenium.webdriver import Firefox
from selenium.webdriver.firefox.options import Options
opts = Options()
opts.set_headless()
browser = Firefox(options=opts)
browser.get('https://duckduckgo.com')
search_form = browser.find_element_by_id('search_form_input_homepage')
search_form.send_keys('python')
search_form.submit()
results = browser.find_elements_by_class_name('result__a')
print(results[0].text)
browser.close()
```

## Crawling Restricted/Private pages

- Limited possibility to access the data
  - *Login forms with email confirmations*
  - *CAPTCHA, reCAPTCHA, ...*
    - → *Completely Automated Public Turing test to tell Computers and Humans Apart (CAPTCHA)*
    - → *Visuals, sounds, ...*
    - → *Secondary goals - image annotations, digitizing texts, ...*

- Approaches
  - *Special agents for session/cookie based approaches*
    - → *Automatic or manual registrations/logins*
  - *CAPTCHA solvers, ...*
  - *Simulations of real human behavior, ...*

# Web Scraping Example - "Real Browser"

- ## Sending HTTP requests

```
1   import requests
2
3   headers = {'user-agent': 'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_1
4   url="https://en.wikipedia.org/wiki/Web_scraping"
5
6   # Make a GET request to fetch the raw HTML content
7   html_content = requests.get(url, headers=headers).text
8
9   print(html_content)
```

- ## Output:

```
1   <!DOCTYPE html>
2   <html class="client-nojs" lang="en" dir="ltr">
3   <head>
4   <meta charset="UTF-8"/>
5   <title>Web scraping - Wikipedia</title>
6   ...
```

# Web Scraping Example - Session

- ## Sending HTTP requests

```
1    import requests
2
3    # requests.get(url, cookies={...}})
4
5    session = requests.Session()
6    login_data = {'user': 'user', 'password': 'password'}
7    session.post('https://example.com/login', login_data)
8
9    result = s.get('https://example.com/data')
10
11   print(result.text)
```

- ## Output:

```
1   <!DOCTYPE html>
2   <html class="client-nojs" lang="en" dir="ltr">
3   <head>
4   <meta charset="UTF-8"/>
5   <title>Web scraping - Wikipedia</title>
6   ...
```

# Crawling AJAX applications

- Modern applications produce content dynamically, faster and richer

```
1  <div id="main">
2  </div>
3  <script>
4  for(var i=0; i<=5; i++){
5  var data = "<a href=\"#\">Link"+i+"</a><br />";
6      document.getElementById('main').innerHTML += data;
7  }
8  </script>
```

Link0

Link1

Link2

Link3

Link4

Link5

# Crawling AJAX applications

- AJAX application's content is produced dynamically by the browser
  - *dynamically fetch data from the server using JavaScript and display to the user*
  - *e.g., by clicking on button or link*
  - *this dynamically created content is not visible to crawlers*

- Typical scenario
  - *user opens a web page, e.g., http://example.com/*
  - *user clicks on link products*
  - *a JavaScript code is executed, which dynamically changes the URL to a pretty URL http://example.com/#products and content of the current page*

- Problem: a browser can execute JavaScript and produce content on-the-fly, the crawler generally can have limitations

- Headless browsers solutions
  - *No GUI, emulated DOM, ...*
  - *can be expensive, resource-wise, ...*

# Crawling AJAX applications

- Another problem: only the URL part before the hashtag (#) is processed from the server
  - *hash fragments are not part of the HTTP requests*
  - *hash fragments are not sent to the server*

- HTTP request containing hash fragment

```
1  curl -v https://mail.google.com/mail/u/0/#inbox
2
3  > GET /mail/u/0/ HTTP/1.1
4  > Host: mail.google.com
5  > User-Agent: curl/7.30.0
6  > Accept: */*
7
8  < ...HTTP response...
```

- Note the path /mail/u/0/ in the HTTP request
  - *it does not contain the hash fragment*

---

# Making AJAX Applications Crawlable

- Solution:
  - *The crawler detects a pretty URL*
    - → *a URL containing hash fragment beginning with* !
    - → *e.g., https://mail.google.com/mail/u/0/#!inbox*
  - *The crawler transforms the pretty URL to a ugly URL*
    - → *the #! is replaced with _escaped_fragment_*
    - → *the _escaped_fragment_ becomes part of the query parameters*
    - → *e.g., https://mail.google.com/mail/u/0/#!inbox becomes https://mail.google.com/mail/u/0/?_escaped_fragment_=inbox*
  - *The crawler requests the ugly URL*
  - *The server interprets the ugly URL and serves the content*
  - *The crawler processes the content*

- Reference
  - *see Google AJAX Crawling for more details (officially deprecated as of October 2015)*

## Spider Traps

- Spider Trap - a set of web pages that create an infinite number of URLs for crawling
  - *finding new previously unvisited URLs for infinite amount of time*
- Harmful for crawlers
  - *waste of bandwidth and storage to download and store duplicate and useless data*
  - *waste of server's bandwidth and space with fake information*
- Creation of spider's traps for good causes
  - *to catch spambots, which send emails on behalf of others*
  - *to catch crawlers, which waste a website's bandwidth*
- Example - Calendar with "next day/month/week/year" link
  - *creation of dynamic pages with links, which point to next year or month*
  - *http://cal.org/01/2014/ has "next" link pointing to http://cal.org/01/2015/*
  - *again http://cal.org/01/2015/ has "next" link pointing to http://cal.org/01/2016/*

## Honeypot Traps

- Honeypot Traps
  - *In the form of links which are not visible to the typical user on the browser*
    - → *e.g. setting the CSS as* display: none
  - *hidden form inputs*
    - → *e.g. predefined random/empty value expected (e.g. GitHub login page)*
    - → *specific name, that the crawler fill in e.g.* username, password
- Crawler
  - *can try to scrape the information from the link*
  - *can be detected and owner blocks the source IP address*
- Detection is not easy

# Web Scraping Example - Proxies

- Sending HTTP requests
  - *https://free-proxy-list.net/*

```python
import requests
url = 'https://httpbin.org/ip'
proxies = {
    "http": 'http://62.244.49.202:52323',
    "https": 'http://62.244.49.202:52323'
}
response = requests.get(url,proxies=proxies)
print(response.json())
```

- Output:

```
{'origin': '62.244.49.202'}
```

# Other Tools, Demos, Examples

- Headless browsers tools
  - *https://github.com/GoogleChrome/puppeteer*
  - *https://github.com/emadehsan/thal*
  - *https://nickjs.org/*

```javascript
const nick = new Nick()
;(async () => {
    const tab = await nick.newTab()
    await tab.open("news.ycombinator.com")
    await tab.untilVisible("#hnmain")
    await tab.inject("../injectables/jquery-3.0.0.min.js")
    const hackerNewsLinks = await tab.evaluate((arg, callback) => {
        const data = []
        $(".athing").each((index, element) => {
            data.push({
                title: $(element).find(".storylink").text(),
                url: $(element).find(".storylink").attr("href")
            })
        })
        callback(null, data)
    })
    console.log(JSON.stringify(hackerNewsLinks, null, 2))
})()
.then(() => {
    console.log("Job done!")
    nick.exit()
})
.catch((err) => {
    console.log(`Something went wrong: ${err}`)
    nick.exit(1)
})
```

# Other Tools, Demos, Examples

- https://blog.phantombuster.com/web-scraping-in-2017-headless-chrome-tips-tricks-4d6521d695e8

```
1   if (await tab.isVisible(".captchaImage")) {
2      // Get the URL of the generated CAPTCHA image
3      // Note that we could also get its base64-encoded value and solve it too
4      const captchaImageLink = await tab.evaluate((arg, callback) => {
5        callback(null, $(".captchaImage").attr("src"))
6      })
7
8      // Make a call to a CAPTCHA solving service
9      const captchaAnswer = await buster.solveCaptchaImage(captchaImageLink)
10
11     // Fill the form with our solution
12     await tab.fill(".captchaForm",
13       { "captcha-answer": captchaAnswer },
14       { submit: true })
15  }
```

# Machine Readable Annotations

- Significantly reduce the effort to properly extract several entities
  - *Standard based*
  - *Limited amount of entities*
  - *Limited amount of relevant information*

- Content annotated by providers in any existing format
  - *Content-level*
    - → *Microformats*
    - → *Microdata*
    - → *RDFa*
  - *Page-level*
    - → *OpenGraph*
    - → *JSON-LD*

# Content-Level Annotations

- Microformats

```
1  <p class="h-card">
2    <img class="u-photo" src="http://example.org/photo.png" alt="" />
3    <a class="p-name u-url" href="http://example.org">Joe Bloggs</a>
4    <a class="u-email" href="mailto:joebloggs@example.com">joebloggs@example.c
5    <span class="p-street-address">17 Austerstræti</span>
6    <span class="p-locality">Reykjavík</span>
7    <span class="p-country-name">Iceland</span>
8  </p>
```

- Microdata

```
1  <p itemscope itemprop="Person" itemtype="http://schema.org/Person">
2    <span itemprop="name">Christopher Froome</span> was sponsored by
3    <span itemprop="sponsor" itemtype="http://schema.org/Organization">
4      <a itemprop="url" href="http://www.skysports.com/">Sky</a></span> in the
5  </p>
```

- RDFa

```
1  <p vocab="http://schema.org/" typeof="Person">
2    <span property="name">Christopher Froome</span> was sponsored by
3    <span property="sponsor" typeof="http://schema.org/Organization">
4      <a property="url" href="http://www.skysports.com/">Sky</a></span> in the
5  </p>
```

# Page-Level Annotations

- OpenGraph

```
1  <head prefix="og: http://ogp.me/ns# fb: http://ogp.me/ns/fb# ">
2    <meta property="fb:app_id"      content="45646845135" />
3    <meta property="og:url"         content="http://example.com/profile/tom_h
4    <meta property="og:type"        content="profile" />
5    <meta property="og:title"       content="Tom Hanks" />
6    <meta property="og:description" content="Tom Hanks profile" />
7    <meta property="og:image"       content="http://example.com/profile/hanks
```

- JSON-LD

```
1  <script type="application/ld+json">
2  {
3    "@context": "http://schema.org/",
4    "@type": "Person",
5    "name": "Christopher Froome",
6    "sponsor":
7    {
8      "@type": "Organization",
9      "name": "Sky",
10     "url": "http://www.skysports.com/"
11   }
12 }
13 </script>
```

# Machine Readable Annotations - Example

- BeautifulSoup

```
1  url="https://en.wikipedia.org/wiki/Web_scraping"
2  html_content = requests.get(url).text
3  soup = BeautifulSoup(html_content, "html.parser")
4  p = soup.find('script', {'type':'application/ld+json'})
5  print(p.contents)
```

```
1  ['{"@context":"https:\\/\\/schema.org","@type":"Article","name":"Web scrapin
```

- extruct

```
1  import extruct
2  import requests
3  import pprint
4  pp = pprint.PrettyPrinter(indent=2)
5  data = extruct.extract(html_content)
6  pp.pprint(data)
```

```
1  { 'json-ld': [ { '@context': 'https://schema.org',
2                   '@type': 'Article',
3                   'author': { '@type': 'Organization',
4                               'name': 'Contributors to Wikimedia projects'},
5                   'dateModified': '2019-12-16T19:58:14Z',
6                   'datePublished': '2005-09-17T18:57:30Z',
7                   'headline': 'data scraping used for extracting data from '
8                               'websites',
```

# Other "Sources" for Scraping

- JavaScript variables
  - <script> *tag on the page*

```
1  from selenium.webdriver import Firefox
2  from selenium.webdriver.firefox.options import Options
3  opts = Options()
4  browser = Firefox(options=opts)
5  browser.get('https://www.bbc.com/news/world-51235555')
6
7  browser.execute_script('return config;')
8
9  browser.close()
```

```
1  {
2    'headline': 'Space cookies: First food baked in space by astronauts',
3    'iStats_counter_name': 'news.world.story.51235555.page',
4    'language': 'en-gb',
5    'last_updated': {'date': '2020-01-24 12:26:02',
6     'timezone': 'Europe/London',
7     'timezone_type': 3},
8    'length': 2994,
9    'mediaType': 'video',
10   ...
11   'section': {'id': '99115',
12    'name': 'World',
13    'uri': '/news/world',
14    'urlIdentifier': '/news/world'},
```

## Other "Sources" for Scraping (cont.)

- API calls, XHRs
  - *website loads data dynamically*
  - *it typically uses XMLHttpRequests (XHRs)*

```
1  url = 'https://aukro.cz/backend/api/offers/search?page=0&size=60&sort='
2  data = {'text': 'notebook', "splitGroupKey": "listing", "splitGroupValue": "
3  headers = { 'Content-Type': 'application/json'}
4  requests.post(url, json = data, headers = headers).json()
```

```
1  ...
2  'content': [{'itemId': 121212121212,
3     'itemName': 'Notebook',
4     'catalog': {'id': 1, 'name': 'Elektronika'},
5  ...
```

## Scraping Strategy

- Scrape Websites Without Being Blocked
  - *Slow down the scraping*
    - → *too many request*
    - → *too fast form submit etc.*
  - *Use proxy servers*
  - *Apply different scraping patterns*
  - *Look like a human*
    - → *Use (switch) correct user-agents*
    - → *Handle cookies*
  - *Be careful of honeypot traps*
  - *Use headless or real browser*