

Python Workshop: Gathering data from the web

Ties de Kok

Tilburg University



Agenda

What are we going to discuss today?

1. How does the web work?
2. Terminology
3. Ethics
4. Tools
5. Specific topics:
 - Interacting with an API
 - Web scrape a page
 - Reverse-engineer HTTP requests
 - Dealing with Javascript elements

Agenda

The web

How does the web work?

What we see:

TIES DE KOK

[Home](#) [Publications](#) [Projects](#) [Teaching](#) [C.V.](#)



Ties de Kok

PhD in Accounting
Tilburg University



Biography

Ties de Kok is a PhD researcher at Tilburg University that specializes in combining computer science with empirical Accounting research. His research interest is in financial accounting, capital markets, empirical management accounting, computer science, and natural language processing.

My university page is located here:

<https://www.tilburguniversity.edu/webwijs/show/t.c.j.dekok.htm>

Interests

- Financial Accounting
- Management Accounting
- Computational Linguistics
- Big Data

Education

- 🎓 Research Master in Accounting
Tilburg University (2013–2015)
- 🎓 Bachelor Business Administration
Tilburg University (2010–2013)

Agenda

The web

How does the web work?

What computers see:

```
<h1 id="biography">Biography</h1>
<p></p>
<p>Ties de Kok is a PhD researcher at Tilburg University that specializes in combining computer science with empirical Accounting research. His research interest is in financial accounting, capital markets, empirical management accounting, computer science, and natural language processing.</p>
<p>My university page is located here: <a href="https://www.tilburguniversity.edu/webwijs/show/t.c.j.dekok.htm" target="_blank">https://www.tilburguniversity.edu/webwijs/show/t.c.j.dekok.htm</a></p>

<div class="row">
  <div class="col-sm-5">
    <h3>Interests</h3>
    <ul class="ul-interests">
      <li>Financial Accounting</li>
      <li>Management Accounting</li>
      <li>Computational Linguistics</li>
      <li>Big Data</li>
    </ul>
  </div>
  <div class="col-sm-7">
    <h3>Education</h3>
    <ul class="ul-edu fa-ul">
      <li>
        <i class="fa-li fa fa-graduation-cap"></i>
        <div class="description">
          <p class="course">Research Master in Accounting</p>
          <p class="institution">Tilburg University (2013-2015)</p>
        </div>
      </li>
      <li>
        <i class="fa-li fa fa-graduation-cap"></i>
        <div class="description">
          <p class="course">Bachelor Business Administration</p>
          <p class="institution">Tilburg University (2010-2013)</p>
        </div>
      </li>
    </ul>
  </div>
</div>
```

Agenda

The web

How does the web work?

Web browsers are awesome!

```
<h1 id="biography">Biography</h1>
<p></p>
<p>Ties de Kok is a PhD researcher at Tilburg University that specializes in c
research interest is in financial accounting, capital markets, empirical manag
</p>
<p>My university page is located here: <a href="https://www.tilburguniversity.
target="_blank">https://www.tilburguniversity.edu/webwijs/show/t.c.j.dekok.htm
```

```
<div class="row">
  <div class="col-sm-5">
    <h3>Interests</h3>
    <ul class="ul-interests">
      <li>Financial Accounting</li>
      <li>Management Accounting</li>
      <li>Computational Linguistics</li>
      <li>Big Data</li>
    </ul>
  </div>
</div>
```

```
<div class="col-sm-7">
  <h3>Education</h3>
  <ul class="ul-edu fa-ul">
    <li>
      <i class="fa-li fa fa-graduation-cap"></i>
      <div class="description">
        <p class="course">Research Master in Accounting</p>
        <p class="institution">Tilburg University (2013-2015)</p>
      </div>
    </li>
    <li>
      <i class="fa-li fa fa-graduation-cap"></i>
      <div class="description">
        <p class="course">Bachelor Business Administration</p>
        <p class="institution">Tilburg University (2010-2013)</p>
      </div>
    </li>
  </ul>
</div>
```



TIES DE KOK



Ties de Kok
PhD in Accounting
Tilburg University



Biography

Ties de Kok is a PhD researcher at Tilburg University that specializes in combining computer science with empirical Accounting research. His research interest is in financial accounting, capital markets, empirical management accounting, computer science, and natural language processing.

My university page is located here:

<https://www.tilburguniversity.edu/webwijs/show/t.c.j.dekok.htm>

Interests

- Financial Accounting
- Management Accounting
- Computational Linguistics
- Big Data

Education

- Research Master in Accounting
Tilburg University (2013-2015)
- Bachelor Business Administration
Tilburg University (2010-2013)

Agenda

The web

Terminology

Web-development terms

- HTML
 - ▶ The structure + static content
- CSS
 - ▶ Determines the way things look
- Javascript
 - ▶ Code to make stuff happen

A screenshot of a JSFiddle interface. At the top, there are tabs for 'Result', 'HTML', 'CSS', and 'JavaScript'. To the right of these tabs is a link 'Edit in JSFiddle' and a blue cloud icon. The main area shows a dark grey background with a white rectangular box in the center. Inside the box, the text 'Hello World' is displayed in a large, bold, black font. Below it is a blue button with the text 'Change color' in white. The overall layout is clean and professional.

Agenda

The web

Terminology

Gathering data from the web

- **Web scraping**
 - ▶ Extract data from a webpage
- Web crawling
 - ▶ Automatically move across webpages

Agenda

The web

Terminology

Gathering data from the web

- **Web scraping**
 - ▶ Extract data from a webpage
- Web crawling
 - ▶ Automatically move across webpages
- API
 - ▶ "Webpages" for computers, not for humans.

Let's say we want to know the current price of Bitcoin:

Humans: Google "Bitcoin price", open website, click on some menu.

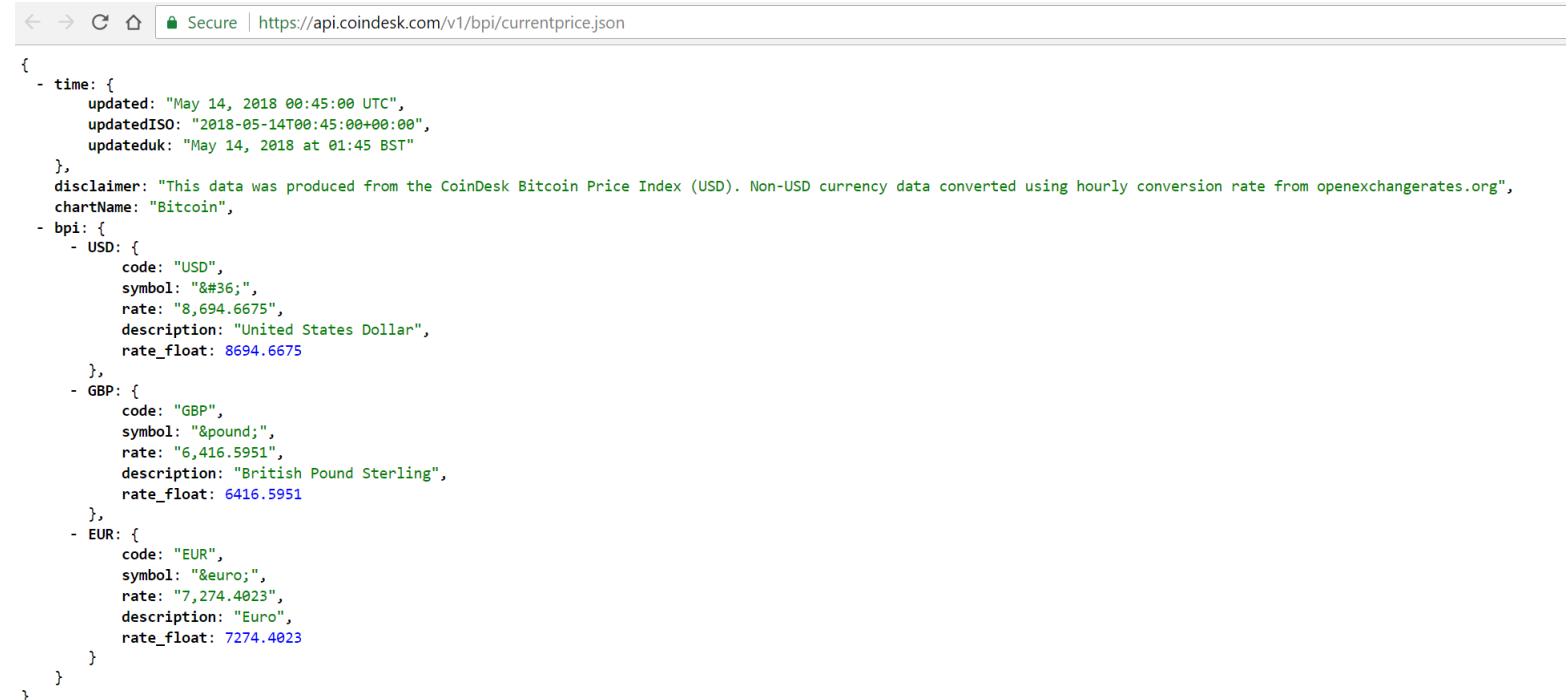
Computers: ??

Agenda

The web

Terminology

Computers: <https://api.coindesk.com/v1/bpi/currentprice.json>



A screenshot of a web browser displaying a JSON response. The URL in the address bar is <https://api.coindesk.com/v1/bpi/currentprice.json>. The page content shows a JSON object representing current Bitcoin prices in USD, GBP, and EUR.

```
{
  - time: {
    updated: "May 14, 2018 00:45:00 UTC",
    updatedISO: "2018-05-14T00:45:00+00:00",
    updateduk: "May 14, 2018 at 01:45 BST"
  },
  disclaimer: "This data was produced from the CoinDesk Bitcoin Price Index (USD). Non-USD currency data converted using hourly conversion rate from openexchangerates.org",
  chartName: "Bitcoin",
  - bpi: {
    - USD: {
      code: "USD",
      symbol: "$",
      rate: "8,694.6675",
      description: "United States Dollar",
      rate_float: 8694.6675
    },
    - GBP: {
      code: "GBP",
      symbol: "&pound;",
      rate: "6,416.5951",
      description: "British Pound Sterling",
      rate_float: 6416.5951
    },
    - EUR: {
      code: "EUR",
      symbol: "&euro;",
      rate: "7,274.4023",
      description: "Euro",
      rate_float: 7274.4023
    }
  }
}
```

Agenda

The web

Terminology

Ethics

Ethics of gathering data from the web

General rules:

1. Always use the API if there is an API available
2. Never allow your scraper to turn into an unintentional DoS attack
3. Read the ToS, explicitly prohibited? Probably don't do it

There is a lot of conflicting information and advice...

My golden principles:

1. Make sure nothing is harmed by gathering the data
2. Never distribute gathered data without permission

Agenda

The web

Terminology

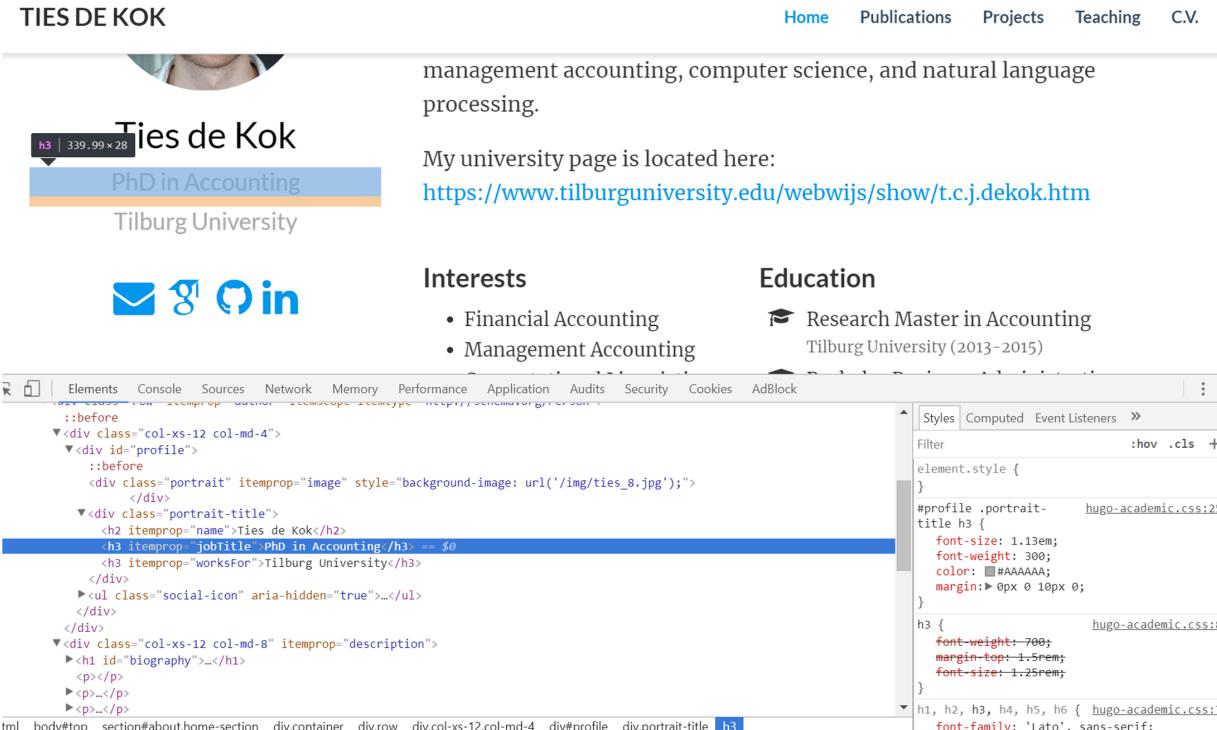
Ethics

Tools

Basic tools

- Chrome + Chrome DevTools

► Open DevTools: **CTRL + SHIFT + J**



The screenshot shows a portion of the Tilburg University website for Ties de Kok. At the top, there's a navigation bar with links to Home, Publications, Projects, Teaching, and C.V. Below the navigation, there's a portrait of Ties de Kok, his name, and his PhD in Accounting from Tilburg University. On the right side of the page, there are sections for Interests (Financial Accounting, Management Accounting) and Education (Research Master in Accounting at Tilburg University from 2013-2015). In the bottom right corner of the page, the developer tools of Google Chrome are open. The Elements tab is selected, showing the HTML structure of the page. A specific element, a blue button containing the text "PhD in Accounting", is highlighted with a blue selection bar. The right panel of the developer tools shows the Styles tab, displaying CSS rules for the highlighted element and other parts of the page. One rule is explicitly named "hugo-academic.css:2". The overall layout is clean and professional, typical of an academic profile page.

Agenda

The web

Terminology

Ethics

Tools

Basic tools

- SelectorGadget Chrome extension

► Download + Install SelectorGadget



Ties de Kok

PhD in Accounting
Tilburg University



Biography

Ties de Kok is a PhD researcher at Tilburg University that specializes in combining computer science with empirical Accounting research. His research interest is in financial accounting, capital markets, empirical management accounting, computer science, and natural language processing.

My university page is located here:

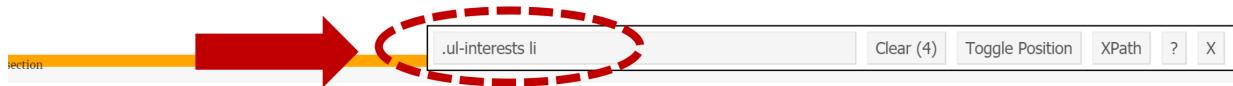
<https://www.tilburguniversity.edu/webwijs/show/t.c.j.dekok.htm>

Interests

- Financial Accounting
- Management Accounting
- Computational Linguistics
- Big Data

Education

- Research Master in Accounting
Tilburg University (2013- 2015)
- Bachelor Business Administration
Tilburg University (2010- 2013)



Agenda

The web

Terminology

Ethics

Tools

Basic tools

- Network Sniffer Chrome extension

► Download + Install Network Sniffer



requestId	statusCode	ip	url	type	timeStamp	tabId	fromCache	method
30797	200	185.199.110.153	https://www.tiesdekok.com/img/icon.png	image	1526261694081.187	338	false	GET
30794	200	74.125.135.156	https://stats.g.doubleclick.net/r/collect?v=1&aip=1&t=dc&_r=3&tid=UA-9085...	image	1526261693842.674	338	false	GET
30791	200	185.199.110.153	https://www.tiesdekok.com/fonts/academicons.ttf?v=1.5.0	font	1526261693608.339	338	false	GET
30790	200	185.199.110.153	https://www.tiesdekok.com/fonts/fontawesome-webfont.woff2?v=4.7.0	font	1526261693580.492	338	false	GET
30787	200	185.199.110.153	https://www.tiesdekok.com/img/ties_8.jpg	image	1526261693957.175	338	false	GET
30785	200	185.199.111.153	https://www.tiesdekok.com/js/hugo-academic.js	script	1526261693166.877	338	true	GET
30783	200	185.199.111.153	https://www.tiesdekok.com/js/isotope.pkgd.min.js	script	1526261693167.4111	338	true	GET
30782	200	185.199.111.153	https://www.tiesdekok.com/js/bootstrap.min.js	script	1526261693162.49	338	true	GET
30781	200	185.199.111.153	https://www.tiesdekok.com/js/jquery-1.12.3.min.js	script	1526261693339.983	338	true	GET
30778	200	185.199.111.153	https://www.tiesdekok.com/css/hugo-academic.css	stylesheet	1526261693146.528	338	true	GET
30776	200	185.199.111.153	https://www.tiesdekok.com/css/academicons.min.css	stylesheet	1526261693146.997	338	true	GET
30775	200	185.199.111.153	https://www.tiesdekok.com/css/font-awesome.min.css	stylesheet	1526261693145.237	338	true	GET
30774	200	185.199.111.153	https://www.tiesdekok.com/css/bootstrap.min.css	stylesheet	1526261693149.069	338	true	GET
30773	200	185.199.111.153	https://www.tiesdekok.com/css/highlight.min.css	stylesheet	1526261693142.9329	338	true	GET
30771	200	185.199.111.153	https://www.tiesdekok.com/	main_frame	1526261693044.897	338	true	GET

Agenda

The web

Terminology

Ethics

Tools

Python packages

To retrieve a webpage:

1) `Requests`

To process a webpage:

2) `LXML` and the higher-level wrapper `Requests-HTML`

To deal with JavaScript heavy-webpages:

3) `Selenium with Python`

(Note: Selenium is not covered in this workshop)

Agenda

The web

Terminology

Ethics

Tools

Topics

Specific topics:

- Interacting with an API
- Web scrape a page
- Reverse-engineer HTTP requests
- Dealing with Javascript elements

API Requests

Basics of using an API

An API, in a simplified sense, has two characteristics:

1. A request is made using a URL + parameters
2. A response is returned in a machine-readable format.

The machine-readable formats are usually either:

- JSON
- XML
- (sometimes plain text)

API Requests

Making API requests using requests



Basics of web scraping

Scraping a page consists of 4 steps:

1. Construct / determine the URL
2. Retrieve the webpage data (usually HTML)
3. Parse the HTML
4. Extract information using HTML structure

Basics of web scraping

Scraping a page consists of 4 steps:

1. Construct / determine the URL
2. Retrieve the webpage data (usually HTML)
3. Parse the HTML
4. Extract information using HTML structure

Step 1 and Step 2:

- ▶ Mostly the same compared to interacting with an API.

Step 3 and Step 4:

- ▶ Different, a HTML page is meant for humans!

API Requests

Web Scraping

Retrieve webpage data (step 1 & 2):

In [8]:

```
import requests
url = 'https://www.tiesdekok.com/'
res = requests.get(url)
```

API Requests

Web Scraping

Retrieve webpage data (step 1 & 2):

In [8]:

```
import requests
url = 'https://www.tiesdekok.com/'
res = requests.get(url)
```

Retrieve data (step 3 & 4):

In [10]:

res.text

???

API Requests

Web Scraping

Parse HTML (step 3):

HTML has a clear structure, but we need to parse (i.e. interpret) it first!

The `lxml.html` library (or `requests-html` wrapper) make this very easy.

Using LXML.HTML

In [17]:

```
import lxml.html
import requests

res = requests.get('https://www.tiesdekok.com')
tree = lxml.html.fromstring(res.text)
```

Using Requests-HTML

In [15]:

```
from requests_html import HTMLSession
session = HTMLSession()
res = session.get('https://www.tiesdekok.com')
```

Note: as the names implies `requests-html` combines `requests` with `lxml.html`

We can now use the HTML structure to extract information!

API
Requests

Web
Scraping

CSS Selectors (step 4, extract information)

Result HTML CSS [Edit in JSFiddle](#) 

Click me!

CSS Selectors (step 4, extract information)

Result HTML CSS [Edit in JSFiddle](#) 

Click me!

CSS Selector make it very easy to select elements to extract!

In [27]:

```
from requests_html import HTMLSession
session = HTMLSession()
res = session.get('https://www.tiesdekok.com')

for element in res.html.find('.ul-interests > li'):
    print(element.text)
```

Financial Accounting
Management Accounting
Computational Linguistics
Big Data



This is the CSS selector

Syntax of a css selector

Most frequent options:

1. Use a dot to select based on **class**: `.classname`
2. Use a hash to select based on **id**: `#idname`
3. Selected directly on the **element**: `p, span, h1`

You can chain multiple conditions together using: `>`, `+`, and `~`.

Example: get `<p>` elements with the `title` class and `<div>` parent:

► CSS Selector: `div > p.title`

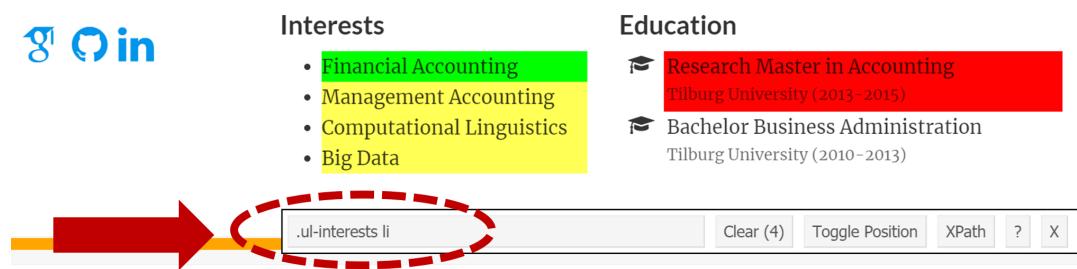
For a full overview I recommend checking this page:

https://www.w3schools.com/cssref/css_selectors.asp

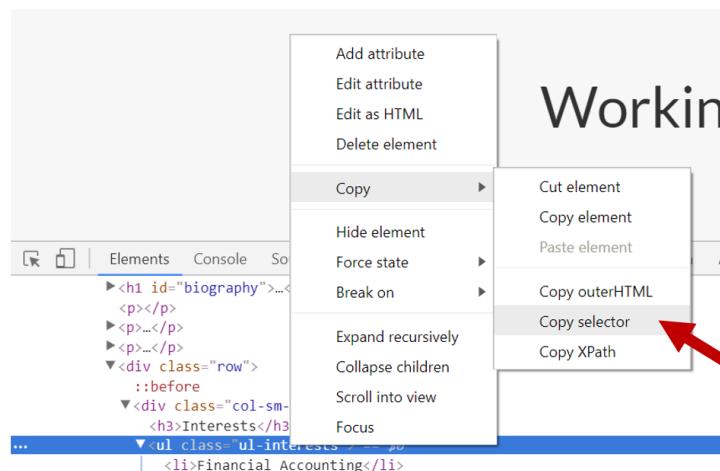
API Requests

Web Scraping

SelectorGadget Chrome extension:



Chrome DevTools:



API Requests

Web Scraping

Recap: web-scraping a page

Step 1: determine the URL of the page you need

URL = <https://www.tiesdekok.com>

Step 2 and Step 3: download and parse the HTML of the webpage

```
from requests_html import HTMLSession
session = HTMLSession()
res = session.get('https://www.tiesdekok.com')
```

Note: the Requests-HTML does the HTML parsing automatically for you.

Step 4: use CSS Selectors to extract information

```
for element in res.html.find('.ul-interests > li'):
    print(element.text)
```

Financial Accounting
Management Accounting
Computational Linguistics
Big Data

API
Requests

Web
Scraping

HTTP
Requests

HTTP Requests

Modern webpages often "load" data to the page using HTTP Requests.

Tip: reverse-engineer the APIs that are used and mimic them!

API
Requests

Web
Scraping

HTTP
Requests

HTTP Requests

Modern webpages often "load" data to the page using HTTP Requests.

Tip: reverse-engineer the APIs that are used and mimic them!

Example:

Let's say we want to get data on the approval rating for Jeff Bezos:



API Requests

Web Scraping

HTTP Requests

What do we see in the Network Sniffer Chrome extension?



The screenshot shows the Network Sniffer extension interface. At the top, there's a search bar with the filter set to "api". Below it are buttons for "on capture", "apply", "clear filter", and "clear result". The main area displays a table of network requests:

requestId	statusCode	ip	url	type	timeStamp	tabId	fromCache	method
11571	200	104.16.232.151	https://www.glassdoor.com/api/employer/6036-rating.htm?dataType=distr...	xmlhttprequest	1526321133053.8372	368	false	GET
11570	200	104.16.232.151	https://www.glassdoor.com/api/employer/6036-rating.htm?dataType=tren...	xmlhttprequest	1526321135769.052	368	false	GET

Two red arrows point to specific rows in the table. One arrow points to the first row with the text "API!" below it, and another arrow points to the second row with the text "HTTP Requests!" below it.

API Requests

Web Scraping

HTTP Requests

What do we see in the Network Sniffer Chrome extension?

The screenshot shows the Network Sniffer extension interface. At the top, there is a search bar with the filter set to "api". Below the search bar are three buttons: "on capture", "apply", "clear filter", and "clear result". The main area displays a table of network requests:

requestId	statusCode	ip	url	type	timeStamp	tabId	fromCache	method
11571	200	104.16.232.151	https://www.glassdoor.com/api/employer/6036-rating.htm?dataType=distr...	xmlhttprequest	1526321133053.8372	368	false	GET
11570	200	104.16.232.151	https://www.glassdoor.com/api/employer/6036-rating.htm?dataType=tren...	xmlhttprequest	1526321135769.052	368	false	GET

Two red arrows point from the text "API!" and "HTTP Requests!" to the respective rows in the table.

The image shows a comparison between the raw API response and its visual representation on a website.

Left (Raw API Response):

```
{  
    employerId: 6036,  
    - labels: [  
        "Approve",  
        "Disapprove"  
    ],  
    - values: [  
        7216,  
        1222  
    ]  
}
```

Middle (VS):

Right (Visualized Form):

Amazon Ratings and Trends

Category	Value	Score
Overall	4.0	4.0
Culture & Values	3.9	3.9
Work/Life Balance	3.3	3.3
Senior Management	3.7	3.7
Comp & Benefits	3.9	3.9
Career Opportunities	4.2	4.2

CEO Approval Trend

CEO Approval Distribution

Category	Percentage
Approve	82%
Disapprove	18%

API
Requests

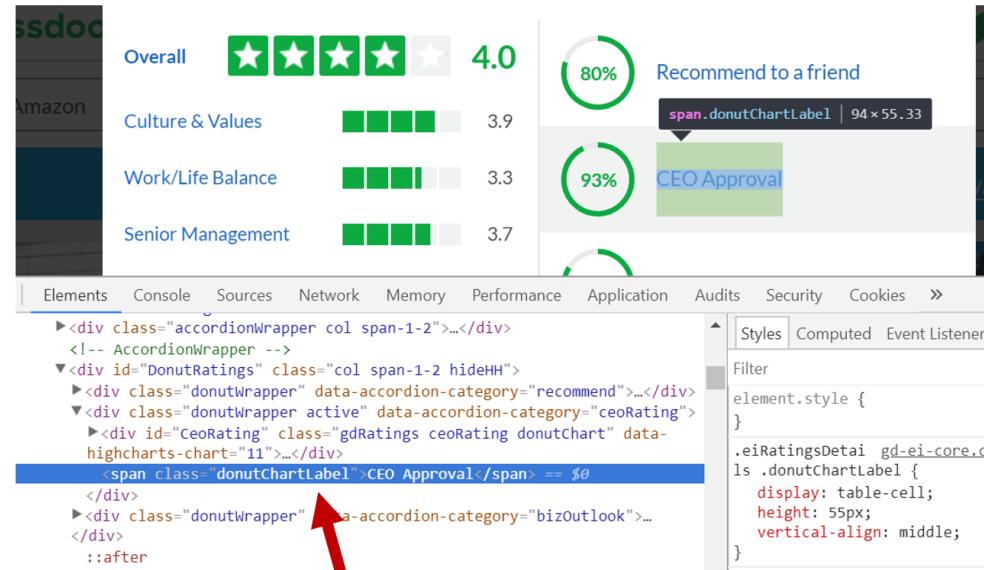
Web
Scraping

HTTP
Requests

JavaScript
Pages

JavaScript heavy webpages

Some webpages rely heavily on JavaScript to load in data-elements:



We can see it

But Python
cannot!?

```
from requests_html import HTMLSession
session = HTMLSession()
res = session.get('https://www.glassdoor.com/Re
res.html.find('.donutChartLabel')
```

Out[5]: []

API
Requests

Web
Scraping

HTTP
Requests

JavaScript
Pages

JavaScript heavy webpages

Can we still scrape them?

Sure, but with a different approach:

Option 1: try to reverse-engineer the HTTP Requests

- ▶ As demonstrated on previous slide

Option 2: use browser automation tools

- ▶ Two primary tools:

1. Use a headless browser (`requests-html` can do this)
2. Use `Selenium` with Chrome bindings

API
Requests

Web
Scraping

HTTP
Requests

JavaScript
Pages

Option 2: use browser automation tools

- ▶ Use a headless browser (`requests-html` can do this)

**Webpage opened by a
headless (i.e. hidden) browser**

In [7]:

```
from requests_html import HTMLSession
session = HTMLSession()
res = session.get('https://www.glassdoor.com/Reviews/Ama:
res.html.render()
res.html.find('.donutChartLabel')
```

Out[7]: [<Element 'span' class='donutChartLabel',>,<Element 'span' class='donutChartLabel',>,<Element 'span' class='donutChartLabel',>,<Element 'span' class='donutChartLabel',>,<Element 'span' class='donutChartLabel',>,<Element 'span' class='donutChartLabel',>]

**Now Python
can see it!**

Note: first time you run `html.render()` it will download some dependencies.

API
Requests

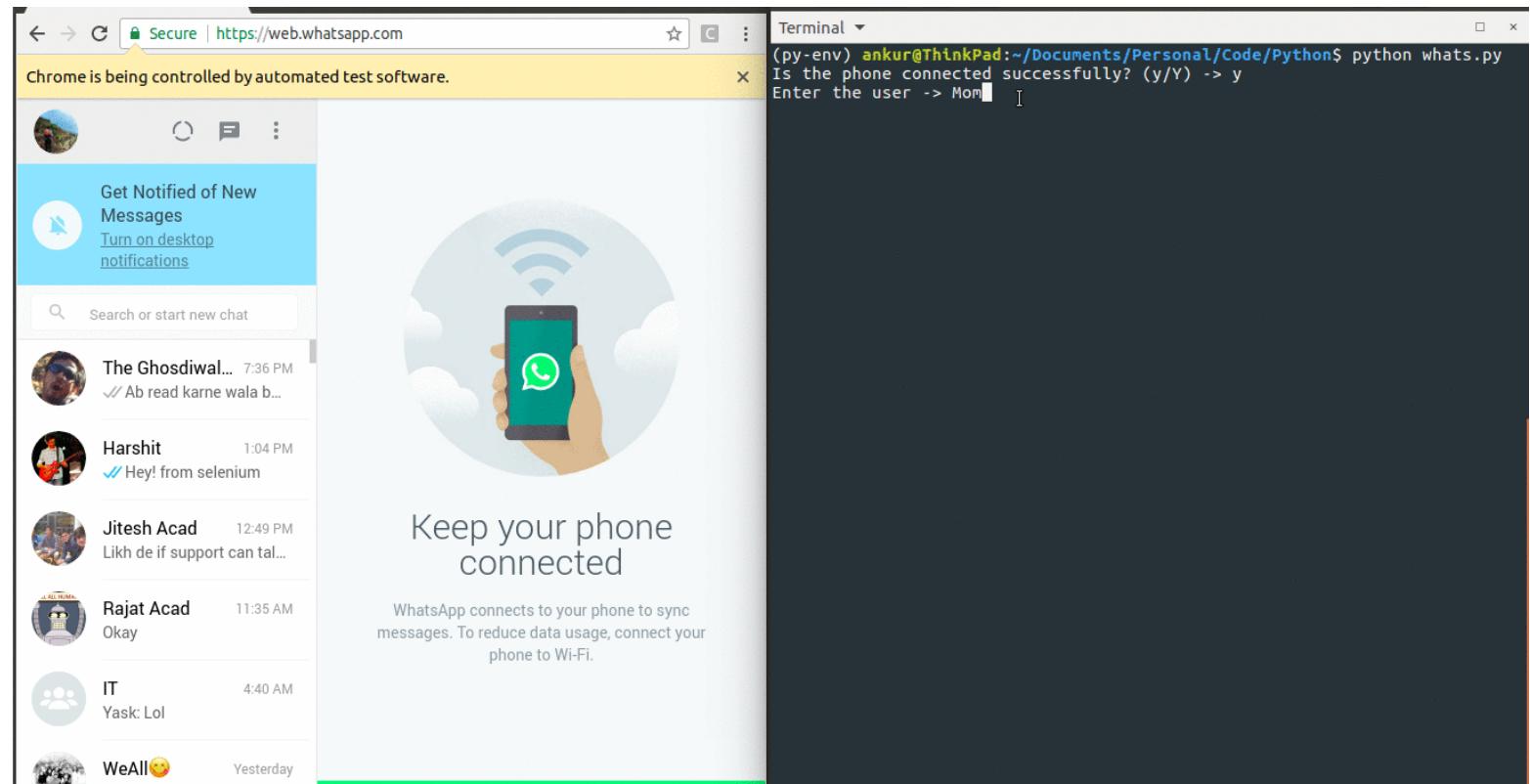
Web
Scraping

HTTP
Requests

JavaScript
Pages

Option 2: use browser automation tools

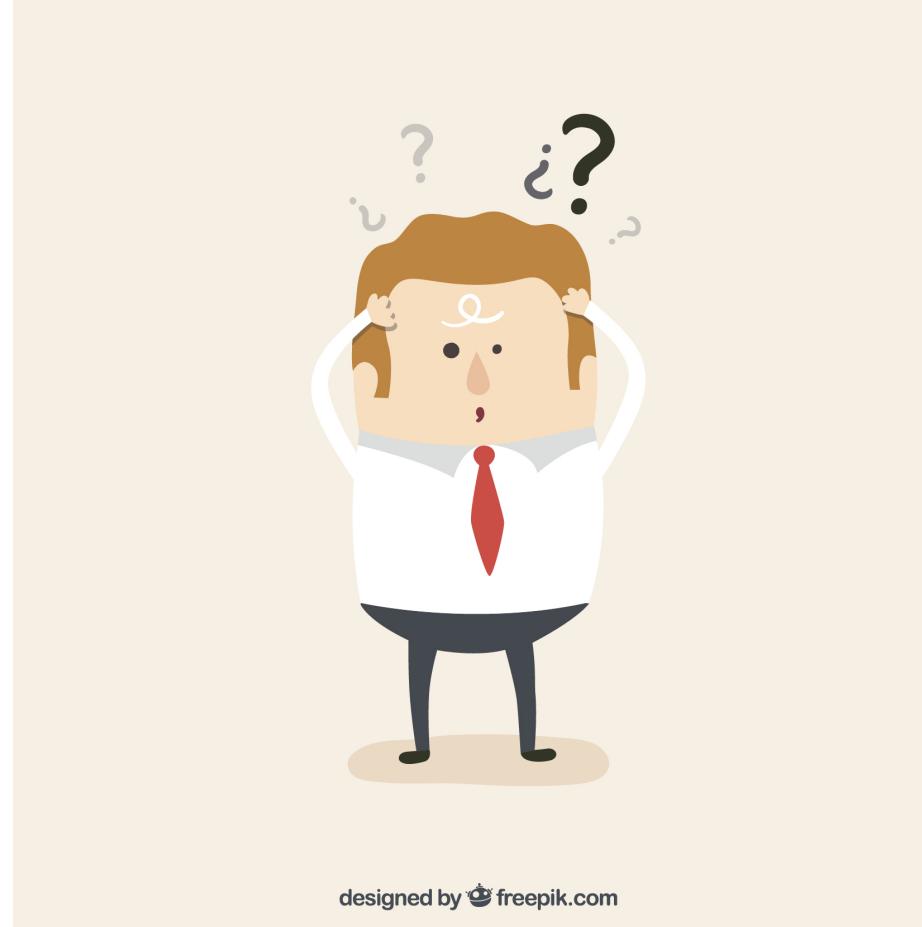
- ▶ Use **Selenium** with Chrome bindings



GIF courtesy of the [PyWhatsApp GitHub page](#)

Closing remarks

Questions?



designed by  freepik.com

Closing remarks

Setup + Get Started

Setup:

1. Make sure you have Chrome installed
2. Install `SelectorGadget` extension
3. Install `networksniffer` extension

Get Started:

Goal: retrieve my name, job title, and institution from my personal page:
<http://www.tiesdekok.com/>

1. Open up a notebook and import `requests_html`
2. Retrieve website data + parse HTML using `requests_html`
3. Use Chrome DevTools + `SelectorGadget` to find the right CSS Selectors
4. Retrieve name, job title, and institution

For help: [see notebook on webscraping](#)

Closing
remarks

Setup +
Get Started

Mini-Task
Instructions

Mini Task

Goal: Get hands-on experience with extracting real data from the web.

Instructions

1. Open (start) a Jupyter Notebook in the `Materials` folder
2. Solve tasks in: `day_2 > mini_task > webscraping_mini_task.ipynb`

For help:

- Python tutorial
- Python Basics Notebook
- Data handling with Pandas
- Web scraping with Python