
Web Scraping

A brief overview by Adam Stapleton

What is Web Scraping?

“Web scraping, web harvesting, or web data extraction is data scraping used for extracting data from websites.”

- Wikipedia



Is web scraping ILLEGAL?

<https://en.wikipedia.org/robots.txt>

<https://www.google.com/robots.txt>

<https://www.facebook.com/robots.txt>

<http://www.daft.ie/robots.txt>



Tip

Ask Google.

They've pretty much perfected it.

Add "/robots.txt" to the end of any homepage url and you'll see

—

So, is it illegal?

Well, is it illegal to rearrange the furniture in a restaurant that you're eating at?



Packages

Necessary

- requests
- bs4 (BeautifulSoup)
- Selenium
 - GeckoDriver (Firefox)

Useful

- Pandas
- random



Approach

In general, this is how to go about extracting data from a website.

→ **Inspect html**

Find what container(s) your information is in

→ **requests/selenium**

Extract website body (which will contain your information)

→ **BeautifulSoup**

Parse the body to find the containers you want, and extract the information you need

Example

data science jobs in Dublin

Recommended Jobs - 10 new

SORT BY:
Relevance - Date

Distance:
within 25 kilometers

Job Type

Full-time (344)
Permanent (231)
Contract (58)
Internship (12)
Part-time (6)

More »

What
data science
job title, keywords or company

Where
Dublin
city or county

Find Jobs

Advanced Job Search

Submit Your CV Sign In

New! Join Indeed Prime - Get offers from great tech companies

Jobs 1 to 10 of 1,142

Show: All Jobs - 307 New Jobs

Senior Director - Data Science, Advanced Analytics Lab - Dub...

UnitedHealth Group - ★★★★★ 11,464 reviews
Computer Science, Science, Engineering, Mathematics, etc.). Be a thought leader and expounder of data methods and applications...

Sponsored Save Job

Sr Software Engineer - Cloud Data

Veritas Technologies - ★★★★★ 244 reviews
Data is an organization's digital currency; Did you know? International Data Corporation (IDC) predicts that every 2 years, reaching 44...

Sponsored Save Job

Get new jobs for this search by email

My email:

email with jobs recommended just for

alert or receiving recommended jobs, you can change your consent settings or as detailed in our terms.

Company with data science jobs

Example

The screenshot displays the Indeed website interface for searching jobs. The URL in the browser is [https://ie.indeed.com/jobs?q=data science&l=Dublin](https://ie.indeed.com/jobs?q=data%20science&l=Dublin). The search results show 'data science jobs in Dublin' with 10 new recommended jobs. The search filters on the left include 'SORT BY: Relevance - Date' and 'Distance: within 25 kilometers'. The job listings on the right include a 'New! Join Indeed Prime' banner and several job postings. The first job listing is for 'Senior Director - Data Science, Advanced Analytics Lab' at 'UnitedHealth Group', which is highlighted. The job description mentions 'Computer Science, Science, Engineering, Mathematics, Statistics etc.' and 'Be a thought leader and expounder of data science methods and applications...'. The second job listing is for 'Sr Software Engineer - Cloud Data Protection' at 'Veritas Technologies', and the third is for 'Cloud Support Associate (Security) - Amazon Web Services' at 'Amazon.com'. The browser's developer tools are open on the right, showing the HTML structure of the highlighted job listing. The HTML code includes a script tag, a style tag, and a div element with the job title and company name.

Find Jobs Find CVs Employers / Post Job Submit Your CV Sign In

indeed

data science jobs in Dublin

Recommended Jobs - 10 new

SORT BY: Relevance - Date

Distance: within 25 kilometers

Job Type

Full-time (344)

Permanent (231)

Contract (58)

Internship (12)

Part-time (6)

What: data science

Where: Dublin

job title, keywords or company

city or county

New! Join Indeed Prime - Get offers from great tech companies

Jobs 1 to 10

Show: All Jobs

div#pj_3a98e971e416346a.row.result.clickcard | 529.617 x 113.933

Senior Director - Data Science, Advanced Analytics Lab - UnitedHealth Group - ★★★★★ 11,464 reviews - Dublin

Computer Science, Science, Engineering, Mathematics, Statistics etc.). Be a thought leader and expounder of data science methods and applications....

Sponsored Save Job

Sr Software Engineer - Cloud Data Protection Veritas Technologies - ★★★★★ 244 reviews - Dublin

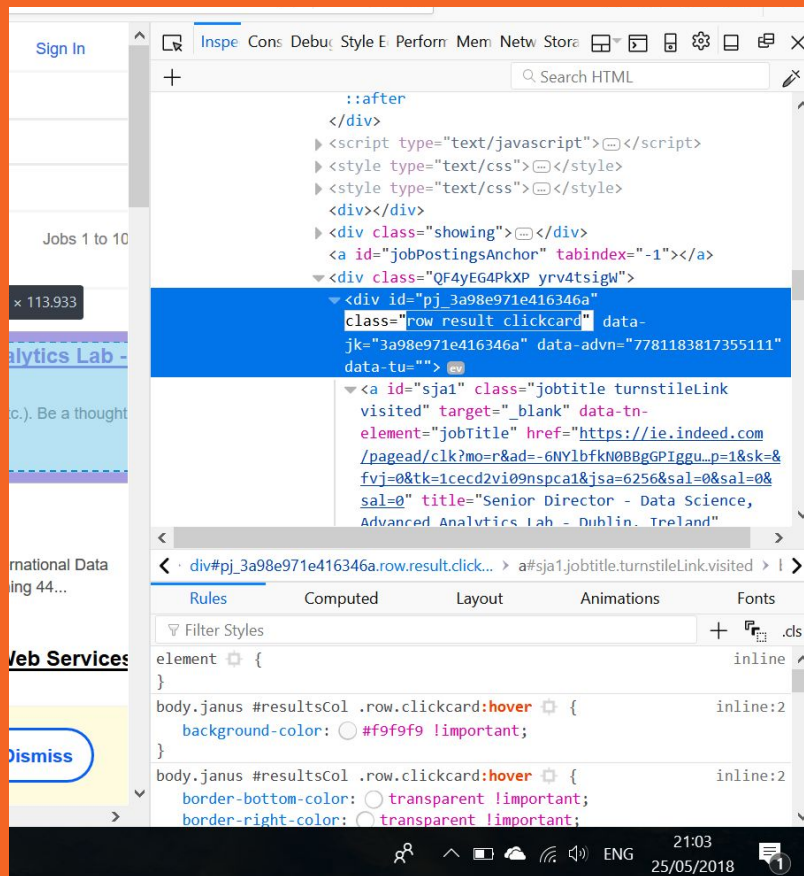
Data is an organization's digital currency; Did you know that the International Data Corporation (IDC) predicts that data will double every 2 years, reaching 44...

Sponsored Save Job

Cloud Support Associate (Security) - Amazon Web Services Amazon.com - ★★★★★ 23,651 reviews - Dublin

```
<script type="text/javascript"></script>
<style type="text/css"></style>
<style type="text/css"></style>
<div></div>
<div class="showing"></div>
<a id="jobPostingsAnchor" tabindex="-1"></a>
<div class="QF4yEG4PkXP yrv4tsigw">
  <div id="pj_3a98e971e416346a" class="row result clickcard" data-jk="3a98e971e416346a" data-advn="7781183817355111" data-tu="">
    <a id="sja1" class="jobtitle turnstileLink visited" target="_blank" data-tn-element="jobTitle" href="https://ie.indeed.com/pagead/clk?mo=r&ad=6NY1bfkN0BBgGPtguu.p-1&sk-fvj=0&tk=1cccd2vi09nspca1&jsa=6256&sal=0&sal=0" title="Senior Director - Data Science, Advanced Analytics Lab - Dublin, Ireland" rel="noopener nofollow" onmouseover="sjomd('sja1'); clk('sja1');" onclick="setRefineByCookie({}); sjoc('sja1',0);">
```


Example





Resources

- <https://automatetheboringstuff.com/chapter11/>
- https://www.seleniumhq.org/docs/03_webdriver.jsp
- <https://towardsdatascience.com/scraping-the-internets-most-popular-websites-a4c6f0be382d>
- <https://github.com/cs109/content>
 - http://nbviewer.jupyter.org/github/cs109/content/blob/master/lec_04_scraping.ipynb