

web_scraping_wiki

May 17, 2019

Table of Contents

1 Web scraping com Python

1.1 Introdução

2 Conceitos básicos

2.1 HTML

2.2 CSS

2.3 JS

2.4 Hands On

2.4.1 Explicando como extraímos a explicação da Wikipedia

Autor: Matheus de Vasconcellos Barroso

1 Web scraping com Python

1.1 Introdução

Aprenderemos alguns conceitos básicos de **web scraping** e como utilizar o [Python](#) para essa tarefa. Mais precisamente esse material foi preparado utilizando um [Jupyter Notebook](#) e algumas extensões úteis ([nbextensions](#)). O exemplo que utilizaremos será a definição de **web scraping** pela [Wikipedia](#):

```
In [48]: #from IPython.display import display, HTML
         from IPython.display import IFrame
```

```
IFrame(src = 'https://en.wikipedia.org/wiki/Web_scraping', width = 900, height=650)
```

```
Out[48]: <IPython.lib.display.IFrame at 0x902e0f0>
```

Imagine como seria útil desenvolver uma técnica para ler o primeiro ou **n** primeiros parágrafos de uma página da [Wikipedia](#) para testar algoritmos de sumarização? Utilizando as técnicas de raspagem de dados podemos obter de forma simples os três primeiros parágrafos da url:

```
In [50]: import requests
         from bs4 import BeautifulSoup

         page = requests.get("https://en.wikipedia.org/wiki/Web_scraping")
         soup = BeautifulSoup(page.content, 'html.parser')
         for item in soup.find_all('p')[:3]: print(item.get_text())
```

Web scraping, web harvesting, or web data extraction is data scraping used for extracting data

Web scraping a web page involves fetching it and extracting from it.[1][2] Fetching is the down

Web scraping is used for contact scraping, and as a component of applications used for web ind

Para alcançar esse objetivo final serão abordados alguns tópicos necessários para compreender como funciona esse processo: + [HTML](#) + [CSS](#) (Não será abordado) + [JS](#) (Não será abordado)

Após um melhor entedimento teórico utilizaremos algumas bibliotecas do Python para realizar essa tarefa como [requests](#) e [bs4](#). Algumas importante, mas que não serão utilizadas: [scrapy](#) e [selenium](#)

2 Conceitos básicos

2.1 HTML

HyperText Markup Language (HTML) ou Linguagem de Marcação de Hipertexto é uma linguagem que server para a construção de páginas web. Ela serve para informar ao navegador como exibir o conteúdo da página.

O HTML é formado por tags, é importante compreendê-las e reconhecê-las já que fascilitam o trabalho de raspagem de dados na web. Algumas tags principais: - **html**: informar ao navegador aonde temos código em HTML - **head**: contém informações como o título da página - **body**: é aonde o conteúdo principal da página está inserido, usualmente é onde o scraping ocorre. - **p**: delimita um parágrafo - **a**: para links - **div**: aponta uma região na página, útil para dividir o conteúdo - **b**: texto em negrito - **i**: texto em itálico - **table**: cria uma tablea - **form**: cria um formulário

2.2 CSS

Cascading Style Sheets (CSS) é um mecanismo para adicionar estilos (cores, fontes, espaçamento, etc.) a um documento web

2.3 JS

JavaScript (JS) é uma linguagem de programação que adiciona interatividade às paginas web.

2.4 Hands On

2.4.1 Explicando como extraímos a explicação da Wikipedia

Primeiro precisamos importas as bibliotecas:

```
In [ ]: import requests
        from bs4 import BeautifulSoup
```

Posteriormente, precisamos obter o conteúdo da página (HTML, CSS, JS). Podemos utilizar o método getdo módulo **requests** para a *url* desejada e atribuí-lo à variável 'page':

```
In [51]: page = requests.get("https://en.wikipedia.org/wiki/Web_scraping")
```

Se a conexão for bem sucedida teremos um status 200 para a página:

```
In [52]: page.status_code
```

Out [52]: 200

Podemos obter o conteúdo da página:

In [53]: `page.content`

Chit [52] arranges the schedule PE into two dimensions of classes' sizes and stages' changes according to "data" distribution to construct

```
In [59]: from bs4 import BeautifulSoup
```

```
soup = BeautifulSoup(page.content, 'html.parser')
soup
```

Out [59]: <!DOCTYPE html>

```
<html class="client-nojs" dir="ltr" lang="en">
<head>
<meta charset="utf-8"/>
<title>Web scraping - Wikipedia</title>
<script>document.documentElement.className = document.documentElement.className.replace(/no-js/i, 'js');
<script>(window.RLQ=window.RLQ||[]).push(function(){mw.config.set({"wgCanonicalNamespaces":{}
});RLPAGEMODULES=["ext.cite.ux-enhancements","site","mediawiki.page.startup","mediawiki
<link href="/w/load.php?lang=en&modules=ext.3d.styles%7Cext.cite.styles%7Cext.uls
<script async="" src="/w/load.php?lang=en&modules=startup&only=scripts&skin=vector
<meta content="" name="ResourceLoaderDynamicStyles"/>
<link href="/w/load.php?lang=en&modules=ext.gadget.charinsert-styles&only=styles&skin=vector
<link href="/w/load.php?lang=en&modules=site.styles&only=styles&skin=vector
<meta content="MediaWiki 1.34.0-wmf.4" name="generator"/>
<meta content="origin" name="referrer"/>
<meta content="origin-when-crossorigin" name="referrer"/>
<meta content="origin-when-cross-origin" name="referrer"/>
<link href="android-app://org.wikipedia/http/en.m.wikipedia.org/wiki/Web_scraping" rel="android-app"
<link href="/w/index.php?title=Web_scraping&action=edit" rel="alternate" title="Edit this page"
<link href="/w/index.php?title=Web_scraping&action=edit" rel="edit" title="Edit this page"
<link href="/static/apple-touch/wikipedia.png" rel="apple-touch-icon"/>
<link href="/static/favicon/wikipedia.ico" rel="shortcut icon"/>
<link href="/w/opensearch_desc.php" rel="search" title="Wikipedia (en)" type="application/opensearchdescription+xml"
<link href="//en.wikipedia.org/w/api.php?action=rdsd" rel="EditURI" type="application/javascript"
<link href="//creativecommons.org/licenses/by-sa/3.0/" rel="license"/>
<link href="https://en.wikipedia.org/wiki/Web_scraping" rel="canonical"/>
<link href="//login.wikimedia.org" rel="dns-prefetch"/>
<link href="//meta.wikimedia.org" rel="dns-prefetch"/>
```

```

<!--[if lt IE 9]><script src="/w/load.php?lang=qqx&modules=html5shiv&only=scr
</head>
<body class="mediawiki ltr sitedir-ltr mw-hide-empty-elt ns-0 ns-subject mw-editable p
<div class="noprint" id="mw-page-base"></div>
<div class="noprint" id="mw-head-base"></div>
<div class="mw-body" id="content" role="main">
<a id="top"></a>
<div class="mw-body-content" id="siteNotice"><!-- CentralNotice --></div>
<div class="mw-indicators mw-body-content">
</div>
<h1 class="firstHeading" id="firstHeading" lang="en">Web scraping</h1>
<div class="mw-body-content" id="bodyContent">
<div class="noprint" id="siteSub">From Wikipedia, the free encyclopedia</div>
<div id="contentSub"></div>
<div id="jump-to-nav"></div>
<a class="mw-jump-link" href="#mw-head">Jump to navigation</a>
<a class="mw-jump-link" href="#p-search">Jump to search</a>
<div class="mw-content-ltr" dir="ltr" id="mw-content-text" lang="en"><div class="mw-p
<div class="hatnote navigation-not-searchable" role="note">For broader coverage of th
<p><b>Web scraping</b>, <b>web harvesting</b>, or <b>web data extraction</b> is <a href
</p><p>Web scraping a web page involves fetching it and extracting from it.<sup class
</p><p>Web scraping is used for <a href="/wiki/Contact_scraping" title="Contact scrap
</p><p><a href="/wiki/Web_page" title="Web page">Web pages</a> are built using text-ba
</p><p>Newer forms of web scraping involve listening to data feeds from web servers.
</p><p>There are methods that some websites use to prevent web scraping, such as dete
</p>
<div class="toc" id="toc"><input class="toctogglecheckbox" id="toctogglecheckbox" rol
<ul>
<li class="toclevel-1 tocsection-1"><a href="#History"><span class="tocnumber">1</span>
<li class="toclevel-1 tocsection-2"><a href="#Techniques"><span class="tocnumber">2</span>
<ul>
<li class="toclevel-2 tocsection-3"><a href="#Human_copy-and-paste"><span class="tocn
<li class="toclevel-2 tocsection-4"><a href="#Text_pattern_matching"><span class="toc
<li class="toclevel-2 tocsection-5"><a href="#HTTP_programming"><span class="tocnumbe
<li class="toclevel-2 tocsection-6"><a href="#HTML_parsing"><span class="tocnumber">2
<li class="toclevel-2 tocsection-7"><a href="#DOM_parsing"><span class="tocnumber">2.
<li class="toclevel-2 tocsection-8"><a href="#Vertical_aggregation"><span class="tocn
<li class="toclevel-2 tocsection-9"><a href="#Semantic_annotation_recognizing"><span
<li class="toclevel-2 tocsection-10"><a href="#Computer_vision_web-page_analysis"><sp
</ul>
</li>
<li class="toclevel-1 tocsection-11"><a href="#Software"><span class="tocnumber">3</span>
<ul>
<li class="toclevel-2 tocsection-12"><a href="#Example_tools"><span class="tocnumber">
<ul>
<li class="toclevel-3 tocsection-13"><a href="#Javascript_tools"><span class="tocnumb
<li class="toclevel-3 tocsection-14"><a href="#Web_crawling_frameworks"><span class="
</ul>

```

```

</li>
</ul>
</li>
<li class="toclevel-1 tocsection-15"><a href="#Legal_issues"><span class="tocnumber">4.1</span><span class="tocsection">Legal issues</span></a>
<ul>
<li class="toclevel-2 tocsection-16"><a href="#United_States"><span class="tocnumber">4.2</span><span class="tocsection">United States</span></a>
<li class="toclevel-2 tocsection-17"><a href="#The_EU"><span class="tocnumber">4.2</span><span class="tocsection">The EU</span></a>
<li class="toclevel-2 tocsection-18"><a href="#Australia"><span class="tocnumber">4.3</span><span class="tocsection">Australia</span></a>
</ul>
</li>
<li class="toclevel-1 tocsection-19"><a href="#Methods_to_prevent_web_scraping"><span class="tocnumber">4.4</span><span class="tocsection">Methods to prevent web scraping</span></a>
<li class="toclevel-1 tocsection-20"><a href="#See_also"><span class="tocnumber">6</span><span class="tocsection">See also</span></a>
<li class="toclevel-1 tocsection-21"><a href="#References"><span class="tocnumber">7</span><span class="tocsection">References</span></a>
</ul>
</div>
<h2><span class="mw-headline" id="History">History</span><span class="mw-editsection">edit</span></h2>
<table class="box-Unreferenced_section plainlinks metadata ambox ambox-content ambox-1">
<tr>
<td>
<p>The history of the web scraping is actually much longer, dating back significantly<span> </span></p>
<ul>
<li>After the birth of <a href="/wiki/History_of_the_World_Wide_Web" title="History of the World Wide Web">the World Wide Web</a>
<li>In 1993,December, the First <b>crawler-based web search engine</b> - <a href="/wiki/Excite">Excite</a>
<li>In 2000, the <b>first Web API and API crawler</b> came. <a href="/wiki/Application Programming Interface">Application Programming Interface</a>
<li>In 2004, <a href="/wiki/Beautiful_Soup_(HTML_parser)" title="Beautiful Soup (HTML parsing module)">Beautiful Soup (HTML parsing module)</a>
<h2><span class="mw-headline" id="Techniques">Techniques</span><span class="mw-editsection">edit</span></h2>
<p>Web scraping is the process of automatically mining data or collecting information from a website<span> </span></p>
<h3><span class="mw-headline" id="Human_copy-and-paste">Human copy-and-paste</span><span class="mw-editsection">edit</span></h3>
<p>Sometimes even the best web-scraping technology cannot replace a humans manual examination<span> </span></p>
<h3><span class="mw-headline" id="Text_pattern_matching">Text pattern matching</span><span class="mw-editsection">edit</span></h3>
<p>A simple yet powerful approach to extract information from web pages can be based on regular expressions<span> </span></p>
<h3><span class="mw-headline" id="HTTP_programming">HTTP programming</span><span class="mw-editsection">edit</span></h3>
<p><a href="/wiki/Static_web_page" title="Static web page">Static</a> and <a href="/wiki/Server-side_scripting">server-side</a> programming<span> </span></p>
<h3><span class="mw-headline" id="HTML_parsing">HTML parsing</span><span class="mw-editsection">edit</span></h3>
<p>Many websites have large collections of pages generated dynamically from an underlying database<span> </span></p>
<h3><span class="mw-headline" id="DOM_parsing">DOM parsing</span><span class="mw-editsection">edit</span></h3>
<div class="hatnote navigation-not-searchable" role="note">Further information: <a href="/wiki/Document_Object_Model">Document Object Model</a>
<p>By embedding a full-fledged web browser, such as the <a href="/wiki/Internet_Explorer">Internet Explorer</a> or <a href="/wiki/Firefox">Firefox</a>
</p>
<h3><span class="mw-headline" id="Vertical_aggregation">Vertical aggregation</span><span class="mw-editsection">edit</span></h3>
<p>There are several companies that have developed vertical specific harvesting platforms<span> </span></p>
<h3><span class="mw-headline" id="Semantic_annotation_recognizing">Semantic annotation recognizing</span><span class="mw-editsection">edit</span></h3>
<p>The pages being scraped may embrace <a href="/wiki/Metadata" title="Metadata">metadata</a> or <a href="/wiki/Structured_data">structured data</a>
</p>

```

<h3>Computer vision v

<p>There are efforts using

</p>

<h2>Software

<p>There are many software tools available that can be used to customize web-scraping

</p>

<h3>Example tools

- cURL command line tool and library for
- Data Toolbar web scraping
- Diffbot uses computer vision and ma
- Heritrix gets pages (lots of them)
- HtmlUnit headless browser that can
- HTTrack free and open source Web cra
- iMacros a browser extension to recor
- Selenium (software
- Jaxer
- Mozenda is a WYSIWYG software that
- nokogiri
- OutWit Hub Web scraping applica
- >watir
- Wget computer program that retrieves conten
- WSO2 Mashup Server
- Yahoo! Query L

<h4>Javascript tools

- Greasemonkey
- Node.js
- PhantomJS scripted,
- jQuery

<h4>Web crawling frameworks

<p>These can be used to build web scrapers.

</p>

- Scrapy

<h2>Legal issues

--

<p>The legality of web scraping varies across the world. In general, web scraping may

</p>

<h3>United States

--

<p>In the United States, website owners can use three major

</p><p>U.S. courts have acknowledged that users of "scrapers" or "robots" may be held

</p><p>One of the first major tests of

</p><p>Southwest Airlines

In 2012, a startup called 3Taps scraped classified housing ads from Craigslist. Craig

</p><p>Although these are early scraping decisions, and the theories of liability are

</p><p>While the law in this area becomes more settled, entities contemplating using

</p><p>In the plaintiff's web site during the period of this trial the terms of use l

</p><p>In <i><a href="/wiki/Facebook,_Inc._v._Power_Ventures,_Inc." title="Facebook, I

</p><p>Internet Archive collects and distributes significant number of publicly avail

The EU

In February 2006, the Danish Maritime and Commercial Court (Copenhagen) ruled that

In a February 2010 case complicated by matters of jurisdiction, Ireland's High Court^[22] The decision

Australia

In Australia, the [Spam Act 2003](/wiki/Spam_Act_2003 "Spam Act 2003")

Methods to prevent web scraping

The administrator of a website can use various measures to stop or slow a bot. Some

- Blocking an [IP address](/wiki/IP_address "IP address") either
- Disabling any [web service](/wiki/Web_service "Web service") [a](#)
- Bots sometimes declare who they are (using [User agent](/wiki/User_agent "User agent"))
- Bots can be blocked by monitoring excess traffic
- Bots can sometimes be blocked with tools to verify that it is a real person accessing
- Commercial anti-bot services: Companies offer anti-bot and anti-scraping services
- Locating bots with a [Honeypot \(computing\)](/wiki/Honeypot_(computing) "Honeypot (computing)")
- [Obfuscation](/wiki/Obfuscation "Obfuscation") using [a](#)
- Because bots rely on consistency in the front-end code of a target website, adding
- Websites can declare if crawling is allowed or not in the [Robots exclusion protocol](/wiki/Robots_exclusion_protocol "Robots exclusion protocol")

See also

- [Archive.is](/wiki/Archive.is "Archive.is")
- [Comparison of feed aggregators](/wiki/Comparison_of_feed_aggregators "Comparison of feed aggregators")
- [Data scraping](/wiki/Data_scraping "Data scraping")
- [Data wrangling](/wiki/Data_wrangling "Data wrangling")
- [Importer \(computing\)](/wiki/Importer_(computing) "Importer (computing)")
- [Job wrapping](/wiki/Job_wrapping "Job wrapping")
- [Knowledge extraction](/wiki/Knowledge_extraction "Knowledge extraction")
- [OpenSocial](/wiki/OpenSocial "OpenSocial")
- [Scraper site](/wiki/Scraper_site "Scraper site")
- [Fake news website](/wiki/Fake_news_website "Fake news website")
- [Blog scraping](/wiki/Blog_scraping "Blog scraping")
- [Spamdexing](/wiki/Spamdexing "Spamdexing")
- [Domain name drop list](/wiki/Domain_name_drop_list "Domain name drop list")
- [Text corpus](/wiki/Text_corpus "Text corpus")
- [Web archiving](/wiki/Web_archiving "Web archiving")
- [Blog network](/wiki/Blog_network "Blog network")
- [Search Engine Scraping](/wiki/Search_Engine_Scraping "Search Engine Scraping")
- [Web crawlers](/wiki/Category:Web_crawlers "Category:Web crawlers")

References

- ↑ [a](#)


```

<li id="cite_note-26"><span class="mw-cite-backlink"><b><a href="#cite_ref-26">^</a></b>
</li>
</ol></div>
<!--
NewPP limit report
Parsed by mw1264
Cached time: 20190516030922
Cache expiry: 2592000
Dynamic content: false
CPU time usage: 0.492 seconds
Real time usage: 0.660 seconds
Preprocessor visited node count: 1805/1000000
Preprocessor generated node count: 0/1500000
Postexpand include size: 65153/2097152 bytes
Template argument size: 772/2097152 bytes
Highest expansion depth: 15/40
Expensive parser function count: 8/500
Unstrip recursion depth: 1/20
Unstrip postexpand size: 70826/5000000 bytes
Number of Wikibase entities loaded: 3/400
Lua time usage: 0.279/10.000 seconds
Lua memory usage: 6.63 MB/50 MB
-->
<!--
Transclusion expansion time report (%,ms,calls,template)
100.00% 548.078      1 -total
 59.50% 326.090      1 Template:Reflist
 38.08% 208.730      5 Template:Cite_journal
 19.94% 109.301      1 Template:More_citations_needed
 18.46% 101.199      4 Template:Ambox
 16.32%  89.451     18 Template:Cite_web
  6.60%  36.163      1 Template:Find_sources_mainspace
  5.52%  30.244      1 Template:Split_section
  4.82%  26.437      1 Template:Split_portions
  2.82%  15.468      1 Template:US-centric
-->
<!-- Saved in parser cache with key enwiki:pcache:idhash:2696619-0!canonical and times
-->
</div><noscript>Retrieved from "<a dir="ltr" href="https://en.wikipedia.org/v
<div class="catlinks" data-mw="interface" id="catlinks"><div class="mw-normal-catlinks
<div class="visualClear"></div>
</div>
</div>
<div id="mw-navigation">
<h2>Navigation menu</h2>
<div id="mw-head">
<div aria-labelledby="p-personal-label" id="p-personal" role="navigation">

```

```

<h3 id="p-personal-label">Personal tools</h3>
<ul>
<li id="pt-anonuserpage">Not logged in</li><li id="pt-anontalk"><a accesskey="n" href=
</div>
<div id="left-navigation">
<div aria-labelledby="p-namespaces-label" class="vectorTabs" id="p-namespaces" role="n
<h3 id="p-namespaces-label">Namespaces</h3>
<ul>
<li class="selected" id="ca-nstab-main"><span><a accesskey="c" href="/wiki/Web_scrapin
</div>
<div aria-labelledby="p-variants-label" class="vectorMenu emptyPortlet" id="p-variants
<input aria-labelledby="p-variants-label" class="vectorMenuCheckbox" type="checkbox"/>
<h3 id="p-variants-label">
<span>Variants</span>
</h3>
<ul class="menu">
</ul>
</div>
</div>
<div id="right-navigation">
<div aria-labelledby="p-views-label" class="vectorTabs" id="p-views" role="navigation
<h3 id="p-views-label">Views</h3>
<ul>
<li class="collapsible selected" id="ca-view"><span><a href="/wiki/Web_scraping">Read
</div>
<div aria-labelledby="p-cactions-label" class="vectorMenu emptyPortlet" id="p-cactions
<input aria-labelledby="p-cactions-label" class="vectorMenuCheckbox" type="checkbox"/>
<h3 id="p-cactions-label"><span>More</span></h3>
<ul class="menu">
</ul>
</div>
<div id="p-search" role="search">
<h3>
<label for="searchInput">Search</label>
</h3>
<form action="/w/index.php" id="searchform">
<div id="simpleSearch">
<input accesskey="f" id="searchInput" name="search" placeholder="Search Wikipedia" ti
</form>
</div>
</div>
</div>
<div id="mw-panel">
<div id="p-logo" role="banner"><a class="mw-wiki-logo" href="/wiki/Main_Page" title="
<div aria-labelledby="p-navigation-label" class="portal" id="p-navigation" role="navig
<h3 id="p-navigation-label">Navigation</h3>
<div class="body">
<ul>

```

```

<li id="n-mainpage-description"><a accesskey="z" href="/wiki/Main_Page" title="Visit t
</div>
</div>
<div aria-labelledby="p-interaction-label" class="portal" id="p-interaction" role="na
<h3 id="p-interaction-label">Interaction</h3>
<div class="body">
<ul>
<li id="n-help"><a href="/wiki/Help:Contents" title="Guidance on how to use and edit W
</div>
</div>
<div aria-labelledby="p-tb-label" class="portal" id="p-tb" role="navigation">
<h3 id="p-tb-label">Tools</h3>
<div class="body">
<ul>
<li id="t-whatlinkshere"><a accesskey="j" href="/wiki/Special:WhatLinksHere/Web_scrap
</div>
</div>
<div aria-labelledby="p-coll-print_export-label" class="portal" id="p-coll-print_expor
<h3 id="p-coll-print_export-label">Print/export</h3>
<div class="body">
<ul>
<li id="coll-create_a_book"><a href="/w/index.php?title=Special:Book&bookcmd=book
</div>
</div>
<div aria-labelledby="p-lang-label" class="portal" id="p-lang" role="navigation">
<h3 id="p-lang-label">Languages</h3>
<div class="body">
<ul>
<li class="interlanguage-link interwiki-ar"><a class="interlanguage-link-target" href=
<div class="after-portlet after-portlet-lang"><span class="wb-langlinks-edit wb-langl
</div>
</div>
</div>
<div id="footer" role="contentinfo">
<ul id="footer-info">
<li id="footer-info-lastmod"> This page was last edited on 15 May 2019, at 09:54<span
<li id="footer-info-copyright">Text is available under the <a href="//en.wikipedia.org
additional terms may apply. By using this site, you agree to the <a href="//foundati
</ul>
<ul id="footer-places">
<li id="footer-places-privacy"><a class="extiw" href="https://foundation.wikimedia.org
<li id="footer-places-about"><a href="/wiki/Wikipedia:About" title="Wikipedia:About">
<li id="footer-places-disclaimer"><a href="/wiki/Wikipedia:General_disclaimer" title=
<li id="footer-places-contact"><a href="//en.wikipedia.org/wiki/Wikipedia:Contact_us">
<li id="footer-places-developers"><a href="https://www.mediawiki.org/wiki/Special:MyL
<li id="footer-places-cookiestatment"><a href="https://foundation.wikimedia.org/wiki
<li id="footer-places-mobileview"><a class="noprint stopMobileRedirectToggle" href="//
</ul>

```

```

<ul class="noprint" id="footer-icons">
<li id="footer-copyrightico">
<a href="https://wikimediafoundation.org/"><img alt="Wikimedia Foundation" height="31"
<li id="footer-poweredbyico">
<a href="//www.mediawiki.org/"></div>
</div>
<script>(window.RLQ=window.RLQ||[]).push(function(){mw.config.set({"wgPageParseReport
<script type="application/ld+json">{"@context":"https://schema.org","@type":"Article
<script>(window.RLQ=window.RLQ||[]).push(function(){mw.config.set({"wgBackendResponse
</body>
</html>

```

É possível melhorar ainda mais usando o método prettify

In [57]: `print(soup.prettify())`

```

<!DOCTYPE html>
<html class="client-nojs" dir="ltr" lang="en">
<head>
<meta charset="utf-8"/>
<title>
Web scraping - Wikipedia
</title>
<script>
document.documentElement.className = document.documentElement.className.replace( /(^\s)cli
</script>
<script>
(window.RLQ=window.RLQ||[]).push(function(){mw.config.set({"wgCanonicalNamespace":"","wgCan
});RLPAGEMODULES=["ext.cite.ux-enhancements","site","mediawiki.page.startup","mediawiki.page.r
</script>
<link href="/w/load.php?lang=en&modules=ext.3d.styles%7Cext.cite.styles%7Cext.uls.interl
<script async="" src="/w/load.php?lang=en&modules=startup&only=scripts&skin=vector
</script>
<meta content="" name="ResourceLoaderDynamicStyles"/>
<link href="/w/load.php?lang=en&modules=ext.gadget.charinsert-styles&only=styles&
<link href="/w/load.php?lang=en&modules=site.styles&only=styles&skin=vector" rel=
<meta content="MediaWiki 1.34.0-wmf.4" name="generator"/>
<meta content="origin" name="referrer"/>
<meta content="origin-when-crossorigin" name="referrer"/>
<meta content="origin-when-cross-origin" name="referrer"/>
<link href="android-app://org.wikipedia/http/en.m.wikipedia.org/wiki/Web_scraping" rel="alter
<link href="/w/index.php?title=Web_scraping&action=edit" rel="alternate" title="Edit this
<link href="/w/index.php?title=Web_scraping&action=edit" rel="edit" title="Edit this page
<link href="/static/apple-touch/wikipedia.png" rel="apple-touch-icon"/>
<link href="/static/favicon/wikipedia.ico" rel="shortcut icon"/>
<link href="/w/openserach_desc.php" rel="search" title="Wikipedia (en)" type="application/op

```

```

<link href="//en.wikipedia.org/w/api.php?action=rsd" rel="EditURI" type="application/rsd+xml" />
<link href="//creativecommons.org/licenses/by-sa/3.0/" rel="license"/>
<link href="https://en.wikipedia.org/wiki/Web_scraping" rel="canonical"/>
<link href="//login.wikimedia.org" rel="dns-prefetch"/>
<link href="//meta.wikimedia.org" rel="dns-prefetch"/>
<!--[if lt IE 9]><script src="/w/load.php?lang=qqx&modules=html5shiv&only=scripts&am
</head>
<body class="mediawiki ltr sitedir-ltr mw-hide-empty-elt ns-0 ns-subject mw-editable page-Web.
  <div class="noprint" id="mw-page-base">
  </div>
  <div class="noprint" id="mw-head-base">
  </div>
  <div class="mw-body" id="content" role="main">
    <a id="top">
    </a>
    <div class="mw-body-content" id="siteNotice">
      <!-- CentralNotice -->
    </div>
    <div class="mw-indicators mw-body-content">
    </div>
    <h1 class="firstHeading" id="firstHeading" lang="en">
      Web scraping
    </h1>
    <div class="mw-body-content" id="bodyContent">
      <div class="noprint" id="siteSub">
        From Wikipedia, the free encyclopedia
      </div>
      <div id="contentSub">
      </div>
      <div id="jump-to-nav">
      </div>
      <a class="mw-jump-link" href="#mw-head">
        Jump to navigation
      </a>
      <a class="mw-jump-link" href="#p-search">
        Jump to search
      </a>
      <div class="mw-content-ltr" dir="ltr" id="mw-content-text" lang="en">
        <div class="mw-parser-output">
          <table class="box-More_citations_needed plainlinks metadata ambox ambox-content ambox-Re
          <tbody>
            <tr>
              <td class="mbox-image">
                <div style="width:52px">
                  <a class="image" href="/wiki/File:Question_book-new.svg">
                    <img alt="" data-file-height="399" data-file-width="512" decoding="async" height="3
                  </a>
                </div>
              </td>
            </tr>
          </tbody>
        </div>
      </div>
    </div>
  </div>
</body>
</html>

```

```

</td>
<td class="mbox-text">
<div class="mbox-text-span">
  This article
  <b>
    needs additional citations for
    <a href="/wiki/Wikipedia:Verifiability" title="Wikipedia:Verifiability">
      verification
    </a>
  </b>
  .
  <span class="hide-when-compact">
    Please help
    <a class="external text" href="//en.wikipedia.org/w/index.php?title=Web_scraping&ar
      improve this article
    </a>
    by
    <a href="/wiki/Help:Introduction_to_referencing_with_Wiki_Markup/1" title="Help:In
      adding citations to reliable sources
    </a>
    . Unsourced material may be challenged and removed.
  <br/>
  <small>
    <span class="plainlinks">
      <i>
        Find sources:
      </i>
      <a class="external text" href="//www.google.com/search?as_eq=wikipedia&amp;q=%22
        "Web scraping"
      </a>

      <a class="external text" href="//www.google.com/search?tbm=nws&amp;q=%22Web+scrap
        news
      </a>
      <b>
        ù
      </b>
      <a class="external text" href="//www.google.com/search?&amp;q=%22Web+scraping%22
        newspapers
      </a>
      <b>
        ù
      </b>
      <a class="external text" href="//www.google.com/search?tbs=bks:1&amp;q=%22Web+sc
        books
      </a>
      <b>
        ù

```

```

        </b>
        <a class="external text" href="//scholar.google.com/scholar?q=%22Web+scraping%22"
          scholar
        </a>
        <b>
          ũ
        </b>
        <a class="external text" href="https://www.jstor.org/action/doBasicSearch?Query=JSTOR"
          JSTOR
        </a>
      </span>
    </small>
  </span>
  <small class="date-container">
    <i>
      (
      <span class="date">
        June 2017
      </span>
      )
    </i>
  </small>
  <small class="hide-when-compact">
    <i>
      (
      <a href="/wiki/Help:Maintenance_template_removal" title="Help:Maintenance template removal"
        Learn how and when to remove this template message
      </a>
      )
    </i>
  </small>
</div>
</td>
</tr>
</tbody>
</table>
<div class="hatnote navigation-not-searchable" role="note">
  For broader coverage of this topic, see
  <a href="/wiki/Data_scraping" title="Data scraping">
    Data scraping
  </a>
  .
</div>
<p>
  <b>
    Web scraping
  </b>
  ,

```

web harvesting
 , or
web data extraction
 is
[data scraping](/wiki/Data_scraping "Data scraping")
 used for
[extracting data](/wiki/Data_extraction "Data extraction")
 from
[websites](/wiki/Website "Website")
 .
<sup>class="reference" id="cite_ref-Boeing2016JPER_1-0">

 [1]</sup>

Web scraping software may access the World Wide Web directly using the
[Hypertext Transfer Protocol](/wiki/Hypertext_Transfer_Protocol "Hypertext Transfer Protocol")
 , or through a web browser. While web scraping can be done manually by a software user,
[bot](/wiki/Internet_bot "Internet bot")
 or
[web crawler](/wiki/Web_crawler "Web crawler")
 . It is a form of copying, in which specific data is gathered and copied from the web,
[database](/wiki/Database "Database")
 or spreadsheet, for later
[retrieval](/wiki/Data_retrieval "Data retrieval")
 or
[analysis](/wiki/Data_analysis "Data analysis")

 .
 </p>
 <p>
 Web scraping a web page involves fetching it and extracting from it.
 <sup class="reference" id="cite_ref-Boeing2016JPER_1-1">

 [1]

 </sup>
 <sup class="reference" id="cite_ref-2">

 [2]

 </sup>
 Fetching is the downloading of a page (which a browser does when you view the page). The

 parsed

 , searched, reformatted, its data copied into a spreadsheet, and so on. Web scrapers typ
 </p>
 <p>
 Web scraping is used for

 contact scraping

 , and as a component of applications used for

 web indexing

 ,

 web mining

 and

 data mining

 , online price change monitoring and

 price comparison

 , product review scraping (to watch the competition), gathering real estate listings, w
 <a href="/wiki/Change_detection_and_notification" title="Change detection and notificat
 website change detection

 , research, tracking online presence and reputation,


```

    web mashup
  </a>
  and,
  <a href="/wiki/Web_data_integration" title="Web data integration">
    web data integration
  </a>
  .
</p>
<p>
  <a href="/wiki/Web_page" title="Web page">
    Web pages
  </a>
  are built using text-based mark-up languages (
  <a href="/wiki/HTML" title="HTML">
    HTML
  </a>
  and
  <a href="/wiki/XHTML" title="XHTML">
    XHTML
  </a>
  ), and frequently contain a wealth of useful data in text form. However, most web pages
  <a class="mw-redirect" href="/wiki/End-user_(computer_science)" title="End-user (computer science)">
    end-users
  </a>
  and not for ease of automated use. Because of this, tool kits that scrape web content w
  <a class="mw-redirect" href="/wiki/Application_Programming_Interface" title="Application Programming Interface">
    Application Programming Interface
  </a>
  (API) to extract data from a web site. Companies like
  <a class="mw-redirect" href="/wiki/Amazon_AWS" title="Amazon AWS">
    Amazon AWS
  </a>
  and
  <a href="/wiki/Google" title="Google">
    Google
  </a>
  provide web scraping tools, services and public data available free of cost to end users
</p>
<p>
  Newer forms of web scraping involve listening to data feeds from web servers. For example
  <a href="/wiki/JSON" title="JSON">
    JSON
  </a>
  is commonly used as a transport storage mechanism between the client and the web server
</p>
<p>
  There are methods that some websites use to prevent web scraping, such as detecting and
  <a href="/wiki/Document_Object_Model" title="Document Object Model">

```

```

    DOM
  </a>
  parsing,
  <a href="/wiki/Computer_vision" title="Computer vision">
    computer vision
  </a>
  and
  <a href="/wiki/Natural_language_processing" title="Natural language processing">
    natural language processing
  </a>
  to simulate human browsing to enable gathering web page content for offline parsing.
</p>
<div class="toc" id="toc">
  <input class="toctogglecheckbox" id="toctogglecheckbox" role="button" style="display:none" type="checkbox"/>
  <div class="toctitle" dir="ltr" lang="en">
    <h2>
      Contents
    </h2>
    <span class="toctogglespan">
      <label class="toctogglelabel" for="toctogglecheckbox">
        </label>
      </span>
    </div>
    <ul>
      <li class="toclevel-1 tocsection-1">
        <a href="#History">
          <span class="tocnumber">
            1
          </span>
          <span class="toctext">
            History
          </span>
        </a>
      </li>
      <li class="toclevel-1 tocsection-2">
        <a href="#Techniques">
          <span class="tocnumber">
            2
          </span>
          <span class="toctext">
            Techniques
          </span>
        </a>
      </li>
      <li class="toclevel-2 tocsection-3">
        <a href="#Human_copy-and-paste">
          <span class="tocnumber">
            2.1

```

```

    </span>
    <span class="toctext">
      Human copy-and-paste
    </span>
  </a>
</li>
<li class="toclevel-2 tocsection-4">
  <a href="#Text_pattern_matching">
    <span class="tocnumber">
      2.2
    </span>
    <span class="toctext">
      Text pattern matching
    </span>
  </a>
</li>
<li class="toclevel-2 tocsection-5">
  <a href="#HTTP_programming">
    <span class="tocnumber">
      2.3
    </span>
    <span class="toctext">
      HTTP programming
    </span>
  </a>
</li>
<li class="toclevel-2 tocsection-6">
  <a href="#HTML_parsing">
    <span class="tocnumber">
      2.4
    </span>
    <span class="toctext">
      HTML parsing
    </span>
  </a>
</li>
<li class="toclevel-2 tocsection-7">
  <a href="#DOM_parsing">
    <span class="tocnumber">
      2.5
    </span>
    <span class="toctext">
      DOM parsing
    </span>
  </a>
</li>
<li class="toclevel-2 tocsection-8">
  <a href="#Vertical_aggregation">

```

```

        <span class="tocnumber">
          2.6
        </span>
        <span class="toctext">
          Vertical aggregation
        </span>
      </a>
    </li>
    <li class="toclevel-2 tocsection-9">
      <a href="#Semantic_annotation_recognizing">
        <span class="tocnumber">
          2.7
        </span>
        <span class="toctext">
          Semantic annotation recognizing
        </span>
      </a>
    </li>
    <li class="toclevel-2 tocsection-10">
      <a href="#Computer_vision_web-page_analysis">
        <span class="tocnumber">
          2.8
        </span>
        <span class="toctext">
          Computer vision web-page analysis
        </span>
      </a>
    </li>
  </ul>
</li>
<li class="toclevel-1 tocsection-11">
  <a href="#Software">
    <span class="tocnumber">
      3
    </span>
    <span class="toctext">
      Software
    </span>
  </a>
  <ul>
    <li class="toclevel-2 tocsection-12">
      <a href="#Example_tools">
        <span class="tocnumber">
          3.1
        </span>
        <span class="toctext">
          Example tools
        </span>
      </a>
    </li>
  </ul>
</li>

```

```

</a>
<ul>
  <li class="toclevel-3 tocsection-13">
    <a href="#Javascript_tools">
      <span class="tocnumber">
        3.1.1
      </span>
      <span class="toctext">
        Javascript tools
      </span>
    </a>
  </li>
  <li class="toclevel-3 tocsection-14">
    <a href="#Web_crawling_frameworks">
      <span class="tocnumber">
        3.1.2
      </span>
      <span class="toctext">
        Web crawling frameworks
      </span>
    </a>
  </li>
</ul>
</li>
<ul>
  <li class="toclevel-1 tocsection-15">
    <a href="#Legal_issues">
      <span class="tocnumber">
        4
      </span>
      <span class="toctext">
        Legal issues
      </span>
    </a>
    <ul>
      <li class="toclevel-2 tocsection-16">
        <a href="#United_States">
          <span class="tocnumber">
            4.1
          </span>
          <span class="toctext">
            United States
          </span>
        </a>
      </li>
      <li class="toclevel-2 tocsection-17">
        <a href="#The_EU">

```

```

        <span class="tocnumber">
            4.2
        </span>
        <span class="toctext">
            The EU
        </span>
    </a>
</li>
<li class="toclevel-2 tocsection-18">
    <a href="#Australia">
        <span class="tocnumber">
            4.3
        </span>
        <span class="toctext">
            Australia
        </span>
    </a>
</li>
</ul>
</li>
<li class="toclevel-1 tocsection-19">
    <a href="#Methods_to_prevent_web_scraping">
        <span class="tocnumber">
            5
        </span>
        <span class="toctext">
            Methods to prevent web scraping
        </span>
    </a>
</li>
<li class="toclevel-1 tocsection-20">
    <a href="#See_also">
        <span class="tocnumber">
            6
        </span>
        <span class="toctext">
            See also
        </span>
    </a>
</li>
<li class="toclevel-1 tocsection-21">
    <a href="#References">
        <span class="tocnumber">
            7
        </span>
        <span class="toctext">
            References
        </span>
    </a>
</li>

```

```

    </a>
  </li>
</ul>
</div>
<h2>
  <span class="mw-headline" id="History">
    History
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=1" title="Edit se
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h2>
<table class="box-Unreferenced_section plainlinks metadata ambox ambox-content ambox-Unr
  <tbody>
    <tr>
      <td class="mbox-image">
        <div style="width:52px">
          <a class="image" href="/wiki/File:Question_book-new.svg">
            <img alt="" data-file-height="399" data-file-width="512" decoding="async" height="3
          </a>
        </div>
      </td>
      <td class="mbox-text">
        <div class="mbox-text-span">
          This section
          <b>
            does not
            <a href="/wiki/Wikipedia:Citing_sources" title="Wikipedia:Citing sources">
              cite
            </a>
            any
            <a href="/wiki/Wikipedia:Verifiability" title="Wikipedia:Verifiability">
              sources
            </a>
          </b>
          .
          <span class="hide-when-compact">
            Please help
            <a class="external text" href="//en.wikipedia.org/w/index.php?title=Web_scraping&a
              improve this section

```



```

    </a>
    by
    <a href="/wiki/Help:Introduction_to_referencing_with_Wiki_Markup/1" title="Help:In
    adding citations to reliable sources
    </a>
    . Unsourced material may be challenged and
    <a href="/wiki/Wikipedia:Verifiability#Burden_of_evidence" title="Wikipedia:Verifi
    removed
    </a>
    .
</span>
<small class="date-container">
    <i>
    (
    <span class="date">
    October 2018
    </span>
    )
    </i>
</small>
<small class="hide-when-compact">
    <i>
    (
    <a href="/wiki/Help:Maintenance_template_removal" title="Help:Maintenance templat
    Learn how and when to remove this template message
    </a>
    )
    </i>
</small>
</div>
</td>
</tr>
</tbody>
</table>
<p>
    The history of the web scraping is actually much longer, dating back significantly to t
</p>
<ul>
<li>
    After the birth of
    <a href="/wiki/History_of_the_World_Wide_Web" title="History of the World Wide Web">
    <b>
    World Wide Web
    </b>
</a>
    in 1989, the First web robot -
    <b>
    <a href="/wiki/World_Wide_Web_Wanderer" title="World Wide Web Wanderer">

```

World Wide Web Wanderer

 was created in 1993, June, which was intended only to measure the size of the web.

 In 1993, December, the First

 crawler-based web search engine

 -

 JumpStation

 . As there were not so many websites available on the web, search engines at that time

 In 2000, the

 first Web API and API crawler

 came.
 <a href="/wiki/Application_programming_interface" title="Application programming inter
API">
 API

 stands for

 Application Programming Interface

 . It is an interface that makes it much easier to develop a program by providing the b

 Salesforce

 and

 eBay

 launched their own API, with which programmers were enabled to access and download some

 In 2004,

 Beautiful Soup

 was released. It is a library designed for Python. As not all websites offer APIs, prog


```

</ul>
<h2>
  <span class="mw-headline" id="Techniques">
    Techniques
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=2" title="Edit section: Techniques">
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h2>
<p>
  Web scraping is the process of automatically mining data or collecting information from the web. It is a technique used to extract data from websites and store it in a structured format. It is a common technique used in data mining, machine learning, and artificial intelligence. It is a technique used to extract data from websites and store it in a structured format. It is a common technique used in data mining, machine learning, and artificial intelligence.
  <a class="mw-redirect" href="/wiki/Semantic_web" title="Semantic web">
    semantic web
  </a>
  vision, an ambitious initiative that still requires breakthroughs in text processing, speech recognition, and natural language processing. It is a common technique used in data mining, machine learning, and artificial intelligence.
  <a class="mw-redirect" href="/wiki/Human-computer_interaction" title="Human-computer interaction">
    human-computer interactions
  </a>
  . Current web scraping solutions range from the ad-hoc, requiring human effort, to fully automated solutions that can scrape data from a large number of websites.
</p>
<h3>
  <span class="mw-headline" id="Human_copy-and-paste">
    Human copy-and-paste
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=3" title="Edit section: Human copy-and-paste">
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h3>
<p>
  Sometimes even the best web-scraping technology cannot replace a humans manual examination of the data.
</p>
<h3>

```

```

<span class="mw-headline" id="Text_pattern_matching">
  Text pattern matching
</span>
<span class="mw-editsection">
  <span class="mw-editsection-bracket">
    [
  </span>
  <a href="/w/index.php?title=Web_scraping&action=edit&section=4" title="Edit se
    edit
  </a>
  <span class="mw-editsection-bracket">
    ]
  </span>
</span>
</h3>
<p>
  A simple yet powerful approach to extract information from web pages can be based on the
  <a href="/wiki/Grep" title="Grep">
    grep
  </a>
  command or
  <a href="/wiki/Regular_expression" title="Regular expression">
    regular expression
  </a>
  -matching facilities of programming languages (for instance
  <a href="/wiki/Perl" title="Perl">
    Perl
  </a>
  or
  <a href="/wiki/Python_(programming_language)" title="Python (programming language)">
    Python
  </a>
  ).
</p>
<h3>
  <span class="mw-headline" id="HTTP_programming">
    HTTP programming
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=5" title="Edit se
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>

```

```

    </span>
</h3>
<p>
  <a href="/wiki/Static_web_page" title="Static web page">
    Static
  </a>
  and
  <a href="/wiki/Dynamic_web_page" title="Dynamic web page">
    dynamic web pages
  </a>
  can be retrieved by posting HTTP requests to the remote web server using
  <a class="mw-redirect" href="/wiki/Socket_programming" title="Socket programming">
    socket programming
  </a>
  .
</p>
<h3>
  <span class="mw-headline" id="HTML_parsing">
    HTML parsing
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=6" title="Edit section: Web scraping">
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h3>
<p>
  Many websites have large collections of pages generated dynamically from an underlying database.
  <a href="/wiki/Wrapper_(data_mining)" title="Wrapper (data mining)">
    wrapper
  </a>
  . Wrapper generation algorithms assume that input pages of a wrapper induction system contain
  <sup class="reference" id="cite_ref-3">
    <a href="#cite_note-3">
      [3]
    </a>
  </sup>
  Moreover, some
  <a href="/wiki/Semi-structured_data" title="Semi-structured data">
    semi-structured data
  </a>
  query languages, such as

```

```

<a href="/wiki/XQuery" title="XQuery">
  XQuery
</a>
and the HTQL, can be used to parse HTML pages and to retrieve and transform page content
</p>
<h3>
  <span class="mw-headline" id="DOM_parsing">
    DOM parsing
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=7" title="Edit section: DOM parsing" >
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h3>
<div class="hatnote navigation-not-searchable" role="note">
  Further information:
  <a href="/wiki/Document_Object_Model" title="Document Object Model">
    Document Object Model
  </a>
</div>
<p>
  By embedding a full-fledged web browser, such as the
  <a href="/wiki/Internet_Explorer" title="Internet Explorer">
    Internet Explorer
  </a>
  or the
  <a href="/wiki/Mozilla" title="Mozilla">
    Mozilla
  </a>
  browser control, programs can retrieve the dynamic content generated by client-side scripts
</p>
<h3>
  <span class="mw-headline" id="Vertical_aggregation">
    Vertical aggregation
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=8" title="Edit section: Vertical aggregation" >
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h3>

```

```

    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h3>
<p>
  There are several companies that have developed vertical specific harvesting platforms.
  <a class="mw-redirect" href="/wiki/Long_Tail" title="Long Tail">
    Long Tail
  </a>
  of sites that common aggregators find complicated or too labor-intensive to harvest con
</p>
<h3>
  <span class="mw-headline" id="Semantic_annotation_recognizing">
    Semantic annotation recognizing
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=9" title="Edit se
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h3>
<p>
  The pages being scraped may embrace
  <a href="/wiki/Metadata" title="Metadata">
    metadata
  </a>
  or semantic markups and annotations, which can be used to locate specific data snippets
  <a href="/wiki/Microformat" title="Microformat">
    Microformat
  </a>
  does, this technique can be viewed as a special case of DOM parsing. In another case, tl
  <sup class="reference" id="cite_ref-4">
    <a href="#cite_note-4">
      [4]
    </a>
  </sup>
  are stored and managed separately from the web pages, so the scrapers can retrieve data
</p>
<h3>
  <span class="mw-headline" id="Computer_vision_web-page_analysis">

```

```

    Computer vision web-page analysis
</span>
<span class="mw-editsection">
  <span class="mw-editsection-bracket">
    [
  </span>
  <a href="/w/index.php?title=Web_scraping&action=edit&section=10" title="Edit s
    edit
  </a>
  <span class="mw-editsection-bracket">
    ]
  </span>
</span>
</h3>
<p>
  There are efforts using
  <a href="/wiki/Machine_learning" title="Machine learning">
    machine learning
  </a>
  and
  <a href="/wiki/Computer_vision" title="Computer vision">
    computer vision
  </a>
  that attempt to identify and extract information from web pages by interpreting pages v
  <sup class="reference" id="cite_ref-5">
    <a href="#cite_note-5">
      [5]
    </a>
  </sup>
</p>
<h2>
  <span class="mw-headline" id="Software">
    Software
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=11" title="Edit s
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h2>
<p>
  There are many software tools available that can be used to customize web-scraping solu

```



```

</p>
<h3>
  <span class="mw-headline" id="Example_tools">
    Example tools
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=12" title="Edit s
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h3>
<ul>
  <li>
    <a href="/wiki/CURL" title="CURL">
      cURL
    </a>
    command line tool and library for transferring (including getting) data with URLs sup
  </li>
  <li>
    <a href="/wiki/Data_Toolbar" title="Data Toolbar">
      Data Toolbar
    </a>
    web scraping add-on for Internet Explorer, Mozilla Firefox, and Google Chrome Web br
  </li>
  <li>
    <a href="/wiki/Diffbot" title="Diffbot">
      Diffbot
    </a>
    uses computer vision and machine learning to automatically extract data from web pag
  </li>
  <li>
    <a href="/wiki/Heritrix" title="Heritrix">
      Heritrix
    </a>
    gets pages (lots of them). It is a web crawler designed for web archiving, written by
    <a href="/wiki/Wayback_Machine" title="Wayback Machine">
      Wayback Machine
    </a>
    ).
  </li>
  <li>
    <a href="/wiki/HtmlUnit" title="HtmlUnit">

```

```

    HtmlUnit
  </a>
  headless browser that can be used for retrieving web pages, web scraping, and more.
</li>
<li>
  <a href="/wiki/HTTrack" title="HTTrack">
    HTTrack
  </a>
  free and open source Web crawler and offline browser, designed to download websites.
</li>
<li>
  <a href="/wiki/IMacros" title="IMacros">
    iMacros
  </a>
  a browser extension to record, code, share and replay browser automation (javascript)
</li>
<li>
  <a href="/wiki/Selenium_(software)" title="Selenium (software)">
    Selenium (software)
  </a>
  a portable software-testing framework for web applications
</li>
<li>
  <a href="/wiki/Aptana#Aptana_Jaxer" title="Aptana">
    Jaxer
  </a>
</li>
<li>
  <a href="/wiki/Mozenda" title="Mozenda">
    Mozenda
  </a>
  is a WYSIWYG software that offers cloud, onsite, and data wrangling services.
</li>
<li>
  <a href="/wiki/Nokogiri_(software)" title="Nokogiri (software)">
    nokogiri
  </a>
</li>
<li>
  <a href="/wiki/OutWit_Hub" title="OutWit Hub">
    OutWit Hub
  </a>
  Web scraping application including built-in data, image, document extractors and edit
</li>
<li>
  <a href="/wiki/Watir" title="Watir">
    watir
  </a>

```

```

</li>
<li>
  <a href="/wiki/Wget" title="Wget">
    Wget
  </a>
  computer program that retrieves content from web servers. It is part of the GNU Project
</li>
<li>
  <a href="/wiki/WSO2_Mashup_Server" title="WSO2 Mashup Server">
    WSO2 Mashup Server
  </a>

</li>
<li>
  <a href="/wiki/Yahoo!_Query_Language" title="Yahoo! Query Language">
    Yahoo! Query Language
  </a>
  (YQL)
</li>
</ul>
<h4>
  <span class="mw-headline" id="Javascript_tools">
    Javascript tools
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=13" title="Edit section: Javascript tools">
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h4>
<ul>
  <li>
    <a href="/wiki/Greasemonkey" title="Greasemonkey">
      Greasemonkey
    </a>
  </li>
  <li>
    <a href="/wiki/Node.js" title="Node.js">
      Node.js
    </a>
  </li>

```

```

<li>
  <a href="/wiki/PhantomJS" title="PhantomJS">
    PhantomJS
  </a>
  scripted,
  <a href="/wiki/Headless_browser" title="Headless browser">
    headless browser
  </a>
  used for automating web page interaction.
</li>
<li>
  <a href="/wiki/JQuery" title="jQuery">
    jQuery
  </a>
</li>
</ul>
<h4>
  <span class="mw-headline" id="Web_crawling_frameworks">
    Web crawling frameworks
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=14" title="Edit s
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h4>
<p>
  These can be used to build web scrapers.
</p>
<ul>
  <li>
    <a href="/wiki/Scrapy" title="Scrapy">
      Scrapy
    </a>
  </li>
</ul>
<h2>
  <span class="mw-headline" id="Legal_issues">
    Legal issues
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">

```

```

[
</span>
<a href="/w/index.php?title=Web_scraping&action=edit&section=15" title="Edit s
edit
</a>
<span class="mw-editsection-bracket">
]
</span>
</span>
</h2>
<table class="box-Globalize plainlinks metadata ambox ambox-content ambox-globalize" rol
<tbody>
<tr>
<td class="mbox-image">
<div style="width:52px">
<img alt="Globe icon." data-file-height="290" data-file-width="350" decoding="async
</div>
</td>
<td class="mbox-text">
<div class="mbox-text-span">
The examples and perspective in this section
<b>
deal primarily with the United States and do not represent a
<a href="/wiki/Wikipedia:WikiProject_Countering_systemic_bias" title="Wikipedia:Wil
worldwide view
</a>
of the subject
</b>
.
<span class="hide-when-compact">
You may
<a class="external text" href="//en.wikipedia.org/w/index.php?title=Web_scraping&a
improve this section
</a>
, discuss the issue on the
<a href="/wiki/Talk:Web_scraping" title="Talk:Web scraping">
talk page
</a>
, or
<a href="/wiki/Wikipedia:Article_wizard" title="Wikipedia:Article wizard">
create a new article
</a>
, as appropriate.
</span>
<small class="date-container">
<i>
(
<span class="date">

```

```

        October 2015
    </span>
    )
</i>
</small>
<small class="hide-when-compact">
    <i>
    (
        <a href="/wiki/Help:Maintenance_template_removal" title="Help:Maintenance template
        Learn how and when to remove this template message
    </a>
    )
    </i>
</small>
</div>
</td>
</tr>
</tbody>
</table>
<p>
    The legality of web scraping varies across the world. In general, web scraping may be a
    <a class="mw-redirect" href="/wiki/Terms_of_use" title="Terms of use">
        terms of use
    </a>
    of some websites, but the enforceability of these terms is unclear.
    <sup class="reference" id="cite_ref-6">
        <a href="#cite_note-6">
            [6]
        </a>
    </sup>
</p>
<h3>
    <span class="mw-headline" id="United_States">
        United States
    </span>
    <span class="mw-editsection">
        <span class="mw-editsection-bracket">
            [
        </span>
        <a href="/w/index.php?title=Web_scraping&action=edit&section=16" title="Edit s
        edit
        </a>
        <span class="mw-editsection-bracket">
            ]
        </span>
    </span>
</h3>
<table class="plainlinks metadata ambox ambox-move" role="presentation">

```

```

<tbody>
<tr>
<td class="mbox-image">
<div style="width:52px">
<img alt="" data-file-height="100" data-file-width="300" decoding="async" height="100" data-bbox="199 128 436 191"/>
</div>
</td>
<td class="mbox-text">
<div class="mbox-text-span">
It has been suggested that this section be
<a href="/wiki/Wikipedia:Splitting" title="Wikipedia:Splitting">
split
</a>
out into another article titled
<i>
<a class="new" href="/w/index.php?title=Web_scraping_in_the_United_States&action=edit" title="Edit page">
Web scraping in the United States
</a>
</i>
. (
<a href="/wiki/Talk:Web_scraping#Split_section" title="Talk:Web scraping">
Discuss
</a>
)
<small>
<i>
(July 2018)
</i>
</small>
</div>
</td>
</tr>
</tbody>
</table>
<p>
In the United States, website owners can use three major
<a href="/wiki/Cause_of_action" title="Cause of action">
legal claims
</a>
to prevent undesired web scraping: (1) copyright infringement (compilation), (2) violation of
<a href="/wiki/Computer_Fraud_and_Abuse_Act" title="Computer Fraud and Abuse Act">
Computer Fraud and Abuse Act
</a>
(CFAA), and (3)
<a href="/wiki/Trespass_to_chattels" title="Trespass to chattels">
trespass to chattel
</a>
.

```

^{[\[7\]](#cite_note-7)}

However, the effectiveness of these claims relies upon meeting various criteria, and the
[*Feist Publications v. Rural Telephone Service*](/wiki/Feist_Publications,_Inc.,_v._Rural_Telephone_Service_Co. "Feist Publications v. Rural Telephone Service")
that duplication of facts is allowable.

U.S. courts have acknowledged that users of "scrapers" or "robots" may be held liable for
[trespass to chattels](/wiki/Trespass_to_chattels "Trespass to chattels")

^{[\[8\]](#cite_note-8)}

^{[\[9\]](#cite_note-9)}

which involves a computer system itself being considered personal property upon which the
[eBay v. Bidder's Edge](/wiki/EBay_v._Bidder%27s_Edge "EBay v. Bidder's Edge")
, resulted in an injunction ordering Bidder's Edge to stop accessing, collecting, and in
[auction sniping](/wiki/Auction_sniping "Auction sniping")
. However, in order to succeed on a claim of trespass to
[chattels](/wiki/Personal_property "Personal property")
, the
[plaintiff](/wiki/Plaintiff "Plaintiff")
must demonstrate that the

[defendant](/wiki/Defendant "Defendant")
intentionally and without authorization interfered with the plaintiff's possessory interest.
<sup>class="reference" id="cite_ref-10">
[\[10\]](#cite_note-10)</sup>

One of the first major tests of
[screen scraping](/wiki/Screen_scraping "Screen scraping")
involved
[American Airlines](/wiki/American_Airlines "American Airlines")
(AA), and a firm called FareChase.
<sup>class="reference" id="cite_ref-11">
[\[11\]](#cite_note-11)</sup>

AA successfully obtained an
[injunction](/wiki/Injunction "Injunction")
from a Texas trial court, stopping FareChase from selling software that enables users to
<sup>class="reference" id="cite_ref-12">
[\[12\]](#cite_note-12)</sup>

[Southwest Airlines](/wiki/Southwest_Airlines "Southwest Airlines")
has also challenged screen-scraping practices, and has involved both FareChase and another firm.
[US Copyright law](/wiki/US_Copyright_law "US Copyright law")
, and that under copyright, the pieces of information being scraped would not be subject to copyright.
[Supreme Court of the United States](/wiki/Supreme_Court_of_the_United_States "Supreme Court of the United States")

, FareChase was eventually shuttered by parent company
[Yahoo!](/wiki/Yahoo! "Yahoo!")
 , and Outtask was purchased by travel expense company Concur.
<sup>class="reference" id="cite_ref-impervawp2011_13-0">

 [13]

</sup>
 In 2012, a startup called 3Taps scraped classified housing ads from Craigslist. *Craigslist v. 3Taps*
[Craigslist v. 3Taps](/wiki/Craigslist_v._3Taps "Craigslist v. 3Taps")
 . The court held that the cease-and-desist letter and IP blocking was sufficient for Cr
[Computer Fraud and Abuse Act](/wiki/Computer_Fraud_and_Abuse_Act "Computer Fraud and Abuse Act")
 .
</p>
<p>
 Although these are early scraping decisions, and the theories of liability are not unif
<sup>class="reference" id="cite_ref-14">

 [14]

</sup>
</p>
<p>
 While the law in this area becomes more settled, entities contemplating using scraping
Cvent, Inc. v. Eventbrite, Inc.
[Cvent, Inc.](/wiki/Cvent,_Inc. "Cvent, Inc.")
[Eventbrite, Inc.](/wiki/Eventbrite "Eventbrite")
 In the United States district court for the eastern district of Virginia, the court rule
[browse wrap](/wiki/Browse_wrap "Browse wrap")
 contract or license to be enforced.
^{class="reference" id="cite_ref-15">}</sup></sup>

[15]

</sup>

In a 2014, filed in the

<a href="/wiki/United_States_District_Court_for_the_Eastern_District_of_Pennsylvania" t

United States District Court for the Eastern District of Pennsylvania

,

<sup class="reference" id="cite_ref-16">

[16]

</sup>

e-commerce site

QVC

objected to the Pinterest-like shopping aggregator Resultlys `scraping of QVCs site for

<sup class="reference" id="cite_ref-17">

[17]

</sup>

QVC's complaint alleges that the defendant disguised its web crawler to mask its source

</p>

<p>

In the plaintiff's web site during the period of this trial the terms of use link is di

<sup class="reference" id="cite_ref-18">

[18]

</sup>

</p>

<p>

In

<i>

<a href="/wiki/Facebook,_Inc._v._Power_Ventures,_Inc." title="Facebook, Inc. v. Power V

Facebook, Inc. v. Power Ventures, Inc.

</i>

, a district court ruled in 2012 that Power Ventures could not scrape Facebook pages on

Electronic Frontier Foundation

filed a brief in 2015 asking that it be overturned.

<sup class="reference" id="cite_ref-19">

[19]

 </sup>
 <sup class="reference" id="cite_ref-20">

 [20]

 </sup>
 In
 <i>

 Associated Press v. Meltwater U.S. Holdings, Inc.

 </i>
 , a court in the US held Meltwater liable for scraping and republishing news information
 </p>
 <p>
 Internet Archive collects and distributes significant number of publicly available webpages
 </p>
 <h3>

 The EU

 [

 edit

]

 </h3>
 <p>
 In February 2006, the Danish Maritime and Commercial Court (Copenhagen) ruled that systematic scraping of news information is illegal.
 <sup class="reference" id="cite_ref-21">

 [21]

 </sup>
 </p>
 <p>
 In a February 2010 case complicated by matters of jurisdiction, Ireland's High Court decided that the scraping of news information is illegal.

 inchoate

 state of developing case law. In the case of

<i>
 Ryanair Ltd v Billigfluege.de GmbH
 </i>
 , Ireland's High Court ruled

 Ryanair's

 "click-wrap" agreement to be legally binding. In contrast to the findings of the United
 <sup class="reference" id="cite_ref-22">

 [22]

 </sup>
 The decision is under appeal in Ireland's Supreme Court.
 <sup class="reference" id="cite_ref-23">

 [23]

 </sup>
 </p>
 <h3>

 Australia

 [

 <a href="/w/index.php?title=Web_scraping&action=edit§ion=18" title="Edit s
 edit

]

 </h3>
 <p>
 In Australia, the

 Spam Act 2003

 outlaws some forms of web harvesting, although this only applies to email addresses.
 <sup class="reference" id="cite_ref-24">

 [24]

 </sup>
 <sup class="reference" id="cite_ref-25">

```

    <a href="#cite_note-25">
      [25]
    </a>
  </sup>
</p>
<h2>
  <span class="mw-headline" id="Methods_to_prevent_web_scraping">
    Methods to prevent web scraping
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=19" title="Edit s
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</h2>
<p>
  The administrator of a website can use various measures to stop or slow a bot. Some tech
</p>
<ul>
  <li>
    Blocking an
    <a href="/wiki/IP_address" title="IP address">
      IP address
    </a>
    either manually or based on criteria such as
    <a href="/wiki/Geolocation" title="Geolocation">
      geolocation
    </a>
    and
    <a href="/wiki/DNSBL" title="DNSBL">
      DNSRBL
    </a>
    . This will also block all browsing from that address.
  </li>
  <li>
    Disabling any
    <a href="/wiki/Web_service" title="Web service">
      web service
    </a>
    <a href="/wiki/Application_programming_interface" title="Application programming inter
      API
    </a>

```

that the website's system might expose.

- Bots sometimes declare who they are (using

[user agent](/wiki/User_agent "User agent")

[strings](/wiki/String_(computer_science) "String (computer science)")

) and can be blocked on that basis using

[robots.txt](/wiki/Robots_exclusion_standard "Robots exclusion standard")

; '

[googlebot](/wiki/Googlebot "Googlebot")

' is an example. Other bots make no distinction between themselves and a human using a
- Bots can be blocked by monitoring excess traffic
- Bots can sometimes be blocked with tools to verify that it is a real person accessing t

[CAPTCHA](/wiki/CAPTCHA "CAPTCHA")

. Bots are sometimes coded to explicitly break specific CAPTCHA patterns or may employ
- Commercial anti-bot services: Companies offer anti-bot and anti-scraping services for v

[application firewalls](/wiki/Application_firewall "Application firewall")

have limited bot detection capabilities as well. However, many such solutions are not v

^{class="reference" id="cite_ref-26">}

[26]

</sup>

.

- Locating bots with a

[>](/wiki/Honeypot_(computing) "Honeypot (computing)")

```

    honeypot
  </a>
  or other method to identify the IP addresses of automated crawlers.
</li>
<li>
  <a href="/wiki/Obfuscation" title="Obfuscation">
    Obfuscation
  </a>
  using
  <a class="mw-redirect" href="/wiki/CSS_sprite" title="CSS sprite">
    CSS sprites
  </a>
  to display such data as phone numbers or email addresses, at the cost of
  <a href="/wiki/Web_accessibility" title="Web accessibility">
    accessibility
  </a>
  to
  <a href="/wiki/Screen_reader" title="Screen reader">
    screen reader
  </a>
  users.
</li>
<li>
  Because bots rely on consistency in the front-end code of a target website, adding sma
</li>
<li>
  Websites can declare if crawling is allowed or not in the
  <a href="/wiki/Robots_exclusion_standard" title="Robots exclusion standard">
    robots.txt
  </a>
  file and allow partial access, limit the crawl rate, specify the optimal time to crawl
</li>
</ul>
<h2>
  <span class="mw-headline" id="See_also">
    See also
  </span>
  <span class="mw-editsection">
    <span class="mw-editsection-bracket">
      [
    </span>
    <a href="/w/index.php?title=Web_scraping&action=edit&section=20" title="Edit s
      edit
    </a>
    <span class="mw-editsection-bracket">
      ]
    </span>
  </span>
</span>

```



```

</h2>
<div class="div-col columns column-width" style="-moz-column-width: 22em; -webkit-column:
<ul>
  <li>
    <a class="mw-redirect" href="/wiki/Archive.is" title="Archive.is">
      Archive.is
    </a>
  </li>
  <li>
    <a href="/wiki/Comparison_of_feed_aggregators" title="Comparison of feed aggregators">
      Comparison of feed aggregators
    </a>
  </li>
  <li>
    <a href="/wiki/Data_scraping" title="Data scraping">
      Data scraping
    </a>
  </li>
  <li>
    <a href="/wiki/Data_wrangling" title="Data wrangling">
      Data wrangling
    </a>
  </li>
  <li>
    <a href="/wiki/Importer_(computing)" title="Importer (computing)">
      Importer
    </a>
  </li>
  <li>
    <a href="/wiki/Job_wrapping" title="Job wrapping">
      Job wrapping
    </a>
  </li>
  <li>
    <a href="/wiki/Knowledge_extraction" title="Knowledge extraction">
      Knowledge extraction
    </a>
  </li>
  <li>
    <a href="/wiki/OpenSocial" title="OpenSocial">
      OpenSocial
    </a>
  </li>
  <li>
    <a href="/wiki/Scraper_site" title="Scraper site">
      Scraper site
    </a>
  </li>

```

```

<li>
  <a href="/wiki/Fake_news_website" title="Fake news website">
    Fake news website
  </a>
</li>
<li>
  <a href="/wiki/Blog_scraping" title="Blog scraping">
    Blog scraping
  </a>
</li>
<li>
  <a href="/wiki/Spamdexing" title="Spamdexing">
    Spamdexing
  </a>
</li>
<li>
  <a href="/wiki/Domain_name_drop_list" title="Domain name drop list">
    Domain name drop list
  </a>
</li>
<li>
  <a href="/wiki/Text_corpus" title="Text corpus">
    Text corpus
  </a>
</li>
<li>
  <a href="/wiki/Web_archiving" title="Web archiving">
    Web archiving
  </a>
</li>
<li>
  <a class="mw-redirect" href="/wiki/Blog_network" title="Blog network">
    Blog network
  </a>
</li>
<li>
  <a class="mw-redirect" href="/wiki/Search_Engine_Scraping" title="Search Engine Scraping">
    Search Engine Scraping
  </a>
</li>
<li>
  <a href="/wiki/Category:Web_crawlers" title="Category:Web crawlers">
    Web crawlers
  </a>
</li>
</ul>
</div>
<h2>

```

```

<span class="mw-headline" id="References">
  References
</span>
<span class="mw-editsection">
  <span class="mw-editsection-bracket">
    [
  </span>
  <a href="/w/index.php?title=Web_scraping&action=edit&section=21" title="Edit s
    edit
  </a>
  <span class="mw-editsection-bracket">
    ]
  </span>
</span>
</h2>
<div class="reflist columns references-column-width" style="-moz-column-width: 20em; -we
<ol class="references">
  <li id="cite_note-Boeing2016JPER-1">
    <span class="mw-cite-backlink">
      ^
    <a href="#cite_ref-Boeing2016JPER_1-0">
      <sup>
        <i>
          <b>
            a
          </b>
        </i>
      </sup>
    </a>
    <a href="#cite_ref-Boeing2016JPER_1-1">
      <sup>
        <i>
          <b>
            b
          </b>
        </i>
      </sup>
    </a>
    </span>
    <span class="reference-text">
      <cite class="citation journal">
        Boeing, G.; Waddell, P. (2016). "New Insights into Rental Housing Markets across the
        <i>
          Journal of Planning Education and Research
        </i>
        (0739456X16664789).
        <a href="/wiki/ArXiv" title="ArXiv">
          arXiv

```

```

</a>
:
<span class="cs1-lock-free" title="Freely accessible">
  <a class="external text" href="//arxiv.org/abs/1605.05397" rel="nofollow">
    1605.05397
  </a>
</span>
.
<a href="/wiki/Digital_object_identifier" title="Digital object identifier">
  doi
</a>
:
<a class="external text" href="//doi.org/10.1177%2F0739456X16664789" rel="nofollow">
  10.1177/0739456X16664789
</a>
.
</cite>
<span class="Z3988" title="ctx_ver=Z39.88-2004&amp;rft_val_fmt=info%3Aofi%2Ffmt%3Ake
</span>
<style data-mw-deduplicate="TemplateStyles:r886058088">
  .mw-parser-output cite.citation{font-style:inherit}.mw-parser-output .citation q{qu
</style>
</span>
</li>
<li id="cite_note-2">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-2">
        ^
      </a>
    </b>
  </span>
  <span class="reference-text">
    <cite class="citation journal">
      Vargiu &amp; Urru (2013). "Exploiting web scraping in a collaborative filtering- ba
    <i>
      Artificial Intelligence Research
    </i>
    .
    <b>
      2
    </b>
    (1).
    <a href="/wiki/Digital_object_identifier" title="Digital object identifier">
      doi
    </a>
    :
    <a class="external text" href="//doi.org/10.5430%2Fair.v2n1p44" rel="nofollow">

```

```

10.5430/air.v2n1p44
</a>
.
</cite>
<span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake
</span>
<link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
</span>
</li>
<li id="cite_note-3">
<span class="mw-cite-backlink">
<b>
<a href="#cite_ref-3">
^
</a>
</b>
</span>
<span class="reference-text">
<cite class="citation journal">
Song, Ruihua; Microsoft Research (Sep 14, 2007).
<a class="external text" href="https://pdfs.semanticscholar.org/4fb4/3c5a212df751e8
"Joint Optimization of Wrapper Generation and Template Detection"
</a>
<span class="cs1-format">
(PDF)
</span>
.
<i>
The 13th International Conference on Knowledge Discovery and Data Mining
</i>
.
</cite>
<span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake
</span>
<link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
</span>
</li>
<li id="cite_note-4">
<span class="mw-cite-backlink">
<b>
<a href="#cite_ref-4">
^
</a>
</b>
</span>
<span class="reference-text">
<a class="external text" href="http://www.gooseeker.com/en/node/knowledgebase/freefo
Semantic annotation based web scraping

```

```

    </a>
  </span>
</li>
<li id="cite_note-5">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-5">
        ^
      </a>
    </b>
  </span>
  <span class="reference-text">
    <cite class="citation web">
      Roush, Wade (2012-07-25).
      <a class="external text" href="http://www.xconomy.com/san-francisco/2012/07/25/diffbot-is-using-computer-vision-to-reinvent-the-semantic-web/">
        "Diffbot Is Using Computer Vision to Reinvent the Semantic Web"
      </a>
      . www.xconomy.com
    <span class="reference-accessdate">
      . Retrieved
      <span class="nowrap">
        2013-03-15
      </span>
    </span>
    .
  </cite>
  <span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Akev%3AURI%2Fhttp%3A%2Fwww.xconomy.com%2Fs%2Fsan-francisco%2F2012-07-25%2Fdiffbot-is-using-computer-vision-to-reinvent-the-semantic-web/">
  </span>
  <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
</span>
</li>
<li id="cite_note-6">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-6">
        ^
      </a>
    </b>
  </span>
  <span class="reference-text">
    <cite class="citation web">
      <a class="external text" href="https://web.archive.org/web/20020308222536/http://www.chillingeffects.org/linking/faq.cgi#QID590">
        "FAQ about linking Are website terms of use binding contracts?"
      </a>
      . www.chillingeffects.org. 2007-08-20. Archived from
      <a class="external text" href="http://www.chillingeffects.org/linking/faq.cgi#QID590">
        the original
      </a>
    </cite>
  </span>
</li>

```

```

on 2002-03-08
<span class="reference-accessdate">
  . Retrieved
  <span class="nowrap">
    2007-08-20
  </span>
</span>
.
</cite>
<span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake
</span>
<link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
</span>
</li>
<li id="cite_note-7">
<span class="mw-cite-backlink">
<b>
  <a href="#cite_ref-7">
    ^
  </a>
</b>
</span>
<span class="reference-text">
<cite class="citation journal">
  Kenneth, Hirschey, Jeffrey (2014-01-01).
  <a class="external text" href="http://scholarship.law.berkeley.edu/btlj/vol29/iss4/"
    "Symbiotic Relationships: Pragmatic Acceptance of Data Scraping"
  </a>
  .
  <i>
    Berkeley Technology Law Journal
  </i>
  .
  <b>
    29
  </b>
  (4).
  <a href="/wiki/Digital_object_identifier" title="Digital object identifier">
    doi
  </a>
  :
  <a class="external text" href="//doi.org/10.15779/2FZ38B39B" rel="nofollow">
    10.15779/Z38B39B
  </a>
  .
  <a href="/wiki/International_Standard_Serial_Number" title="International Standard
    ISSN
  </a>

```

```

    <a class="external text" href="//www.worldcat.org/issn/1086-3818" rel="nofollow">
      1086-3818
    </a>
    .
  </cite>
  <span class="Z3988" title="ctx_ver=Z39.88-2004&rt_val_fmt=info%3Aofi%2Ffmt%3Ake
  </span>
  <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
</span>
</li>
<li id="cite_note-8">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-8">
        ^
      </a>
    </b>
  </span>
  <span class="reference-text">
    <cite class="citation web">
      <a class="external text" href="http://www.tomwbell.com/NetLaw/Ch06.html" rel="nofol
        "Internet Law, Ch. 06: Trespass to Chattels"
      </a>
      . www.tomwbell.com. 2007-08-20
      <span class="reference-accessdate">
        . Retrieved
        <span class="nowrap">
          2007-08-20
        </span>
      </span>
    </cite>
    .
    <span class="Z3988" title="ctx_ver=Z39.88-2004&rt_val_fmt=info%3Aofi%2Ffmt%3Ake
    </span>
    <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
  </span>
</li>
<li id="cite_note-9">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-9">
        ^
      </a>
    </b>
  </span>
  <span class="reference-text">
    <cite class="citation web">
      <a class="external text" href="https://web.archive.org/web/20020308222536/http://ww

```


"What are the "trespass to chattels" claims some companies or website owners have l

 . www.chillingeffects.org. 2007-08-20. Archived from
 <a class="external text" href="http://www.chillingeffects.org/linking/faq.cgi#QID46
 the original

 on 2002-03-08

 . Retrieved

 2007-08-20

 .
 </cite>
 <span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake

 <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>

 <li id="cite_note-10">

 ^

 <cite class="citation web">
 <a class="external text" href="http://www.tomwbell.com/NetLaw/Ch07/Ticketmaster.htm
 "Ticketmaster Corp. v. Tickets.com, Inc"

 . 2007-08-20

 . Retrieved

 2007-08-20

 .
 </cite>
 <span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake

 <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>

 <li id="cite_note-11">


```

</a>
. The Free Library. 2003-06-13
<span class="reference-accessdate">
. Retrieved
<span class="nowrap">
2012-02-26
</span>
</span>
.
</cite>
<span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake
</span>
<link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
</span>
</li>
<li id="cite_note-impervawp2011-13">
<span class="mw-cite-backlink">
<b>
<a href="#cite_ref-impervawp2011_13-0">
^
</a>
</b>
</span>
<span class="reference-text">
Imperva (2011).
<a class="external text" href="http://www.imperva.com/docs/WP_Detecting_and_Blocking
Detecting and Blocking Site Scraping Attacks
</a>
. Imperva white paper..
</span>
</li>
<li id="cite_note-14">
<span class="mw-cite-backlink">
<b>
<a href="#cite_ref-14">
^
</a>
</b>
</span>
<span class="reference-text">
<cite class="citation web">
Adler, Kenneth A. (2003-07-29).
<a class="external text" href="http://library.findlaw.com/2003/Jul/29/132944.html"
"Controversy Surrounds 'Screen Scrapers': Software Helps Users Access Web Sites Bu
</a>
<span class="reference-accessdate">
. Retrieved
<span class="nowrap">

```

2010-10-27

.

</cite>

<link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>

<li id="cite_note-15">

^

<cite class="citation web">

"QVC Inc. v. Resultly LLC, No. 14-06714 (E.D. Pa. filed Nov. 24, 2014)"

(PDF)

. 2014-11-24

. Retrieved

2015-11-05

.

</cite>

<link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>

<li id="cite_note-16">

^

QVC Inc. v. Resultly LLC, No. 14-06714 (E.D. Pa. filed Nov. 24, 2014).

United States District Court for the Eastern District of Pennsylvania

Retrieved
5 November
2015

Neuburger, Jeffrey D (5 December 2014).

The National Law Review

Proskauer Rose LLP

Retrieved
5 November
2015


```

<a class="external text" href="https://www.techdirt.com/articles/20090605/222820514"
  "Can Scraping Non-Infringing Content Become Copyright Infringement... Because Of H
</a>
.
<i>
  Techdirt
</i>
. 2009-06-10
<span class="reference-accessdate">
  . Retrieved
  <span class="nowrap">
    2016-05-24
  </span>
</span>
.
</cite>
<span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake
</span>
<link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
</span>
</li>
<li id="cite_note-20">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-20">
        ^
      </a>
    </b>
  </span>
  <span class="reference-text">
    <cite class="citation web">
      <a class="external text" href="https://www.eff.org/cases/facebook-v-power-ventures"
        "Facebook v. Power Ventures"
      </a>
      .
    <i>
      Electronic Frontier Foundation
    </i>
    <span class="reference-accessdate">
      . Retrieved
      <span class="nowrap">
        2016-05-24
      </span>
    </span>
    .
  </cite>
  <span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake
</span>

```

```

    <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
  </span>
</li>
<li id="cite_note-21">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-21">
        ^
      </a>
    </b>
  </span>
  <span class="reference-text">
    <cite class="citation web">
      <a class="external text" href="https://web.archive.org/web/20071012005033/http://www.UDSKRIFT AF SØ- & HANDELSRETTENS DOMBOG">
        </a>
      <span class="cs1-format">
        (PDF)
      </span>
      (in Danish). bvhd.dk. 2006-02-24. Archived from
      <a class="external text" href="http://www.bvhd.dk/uploads/tx_mocarticles/S_-_og_Han
        the original
      </a>
      <span class="cs1-format">
        (PDF)
      </span>
      on 2007-10-12
      <span class="reference-accessdate">
        . Retrieved
        <span class="nowrap">
          2007-05-30
        </span>
      </span>
      .
    </cite>
    <span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake
  </span>
  <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
</span>
</li>
<li id="cite_note-22">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-22">
        ^
      </a>
    </b>
  </span>

```



```

    <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
  </span>
</li>
<li id="cite_note-24">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-24">
        ^
      </a>
    </b>
  </span>
  <span class="reference-text">
    <cite class="citation web">
      National Office for the Information Economy (February 2004).
      <a class="external text" href="https://www.lloyds.com/~media/5880dae185914b2487bed"
        "Spam Act 2003: An overview for business"
      </a>
      . Australian Communications Authority. p.ã6
    <span class="reference-accessdate">
      . Retrieved
    <span class="nowrap">
      2017-12-07
    </span>
  </span>
  .
    </cite>
  <span class="Z3988" title="ctx_ver=Z39.88-2004&amp;rft_val_fmt=info%3Aofi%2Ffmt%3Ake"
  </span>
  <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
</span>
</li>
<li id="cite_note-25">
  <span class="mw-cite-backlink">
    <b>
      <a href="#cite_ref-25">
        ^
      </a>
    </b>
  </span>
  <span class="reference-text">
    <cite class="citation web">
      National Office for the Information Economy (February 2004).
      <a class="external text" href="http://www.webstartdesign.com.au/spam_business_pract"
        "Spam Act 2003: A practical guide for business"
      </a>
      <span class="cs1-format">
        (PDF)
      </span>
    </cite>
  </span>

```

```

        . Australian Communications Authority. p.20
        <span class="reference-accessdate">
            . Retrieved
            <span class="nowrap">
                2017-12-07
            </span>
        </span>
        .
    </cite>
    <span class="Z3988" title="ctx_ver=Z39.88-2004&rft_val_fmt=info%3Aofi%2Ffmt%3Ake
    </span>
    <link href="mw-data:TemplateStyles:r886058088" rel="mw-deduplicated-inline-style"/>
    </span>
</li>
<li id="cite_note-26">
    <span class="mw-cite-backlink">
        <b>
            <a href="#cite_ref-26">
                ^
            </a>
        </b>
    </span>
    <span class="reference-text">
        Mayank Dhiman
        <a class="external text" href="https://s3.us-west-2.amazonaws.com/research-papers-my
        Breaking Fraud & Bot Detection Solutions
    </a>
    <i>
        OWASP AppSec Cali' 2018
    </i>
    Retrieved February 10, 2018.
    </span>
</li>
</ol>
</div>
<!--

```

```

NewPP limit report
Parsed by mw1264
Cached time: 20190516030922
Cache expiry: 2592000
Dynamic content: false
CPU time usage: 0.492 seconds
Real time usage: 0.660 seconds
Preprocessor visited node count: 1805/1000000
Preprocessor generated node count: 0/1500000
Postexpand include size: 65153/2097152 bytes
Template argument size: 772/2097152 bytes
Highest expansion depth: 15/40

```

```

Expensive parser function count: 8/500
Unstrip recursion depth: 1/20
Unstrip postexpand size: 70826/5000000 bytes
Number of Wikibase entities loaded: 3/400
Lua time usage: 0.279/10.000 seconds
Lua memory usage: 6.63 MB/50 MB
-->
    <!--
Transclusion expansion time report (%,ms,calls,template)
100.00%  548.078      1 -total
 59.50%  326.090      1 Template:Reflist
38.08%  208.730      5 Template:Cite_journal
19.94%  109.301      1 Template:More_citations_needed
18.46%  101.199      4 Template:Ambox
16.32%   89.451     18 Template:Cite_web
  6.60%   36.163      1 Template:Find_sources_mainspace
  5.52%   30.244      1 Template:Split_section
  4.82%   26.437      1 Template:Split_portions
  2.82%   15.468      1 Template:US-centric
-->
    <!-- Saved in parser cache with key enwiki:pcache:idhash:2696619-0!canonical and timestar
-->
    </div>
    <noscript>
        
    Retrieved from "
    <a dir="ltr" href="https://en.wikipedia.org/w/index.php?title=Web_scraping&oldid=8971
        https://en.wikipedia.org/w/index.php?title=Web_scraping&oldid=897184947
    </a>
    "
</div>
<div class="catlinks" data-mw="interface" id="catlinks">
    <div class="mw-normal-catlinks" id="mw-normal-catlinks">
        <a href="/wiki/Help:Category" title="Help:Category">
            Categories
        </a>
        :
        <ul>
            <li>
                <a href="/wiki/Category:Web_scraping" title="Category:Web scraping">
                    Web scraping
                </a>
            </li>
        </ul>
    </div>
</div>

```

<div class="mw-hidden-catlinks mw-hidden-cats-hidden" id="mw-hidden-catlinks">

Hidden categories:

CS1 Danish-language sources (da)

Articles needing additional references from June 2017

All articles needing additional references

Articles needing additional references from October 2018

Articles with limited geographic scope from October 2015

USA-centric

USA-centric

Articles to be split from July 2018

All articles to be split

</div>

</div>

<div class="visualClear">

</div>

```

    </div>
</div>
<div id="mw-navigation">
  <h2>
    Navigation menu
  </h2>
  <div id="mw-head">
    <div aria-labelledby="p-personal-label" id="p-personal" role="navigation">
      <h3 id="p-personal-label">
        Personal tools
      </h3>
      <ul>
        <li id="pt-anonuserpage">
          Not logged in
        </li>
        <li id="pt-anontalk">
          <a accesskey="n" href="/wiki/Special:MyTalk" title="Discussion about edits from this IP address" id="pt-anontalk-link">Talk</a>
        </li>
        <li id="pt-anoncontribs">
          <a accesskey="y" href="/wiki/Special:MyContributions" title="A list of edits made from this IP address" id="pt-anoncontribs-link">Contributions</a>
        </li>
        <li id="pt-createaccount">
          <a href="/w/index.php?title=Special:CreateAccount&returnto=Web+scraping" title="You are encouraged to create an account and log in; you can view and post your own messages" id="pt-createaccount-link">Create account</a>
        </li>
        <li id="pt-login">
          <a accesskey="o" href="/w/index.php?title=Special:UserLogin&returnto=Web+scraping" title="Please log in if you have an account" id="pt-login-link">Log in</a>
        </li>
      </ul>
    </div>
    <div id="left-navigation">
      <div aria-labelledby="p-namespaces-label" class="vectorTabs" id="p-namespaces" role="navigation">
        <h3 id="p-namespaces-label">Namespaces</h3>
        <ul>
          <li class="selected" id="ca-nstab-main">
            <span>Article</span>
            <a accesskey="c" href="/wiki/Web_scraping" title="View the content page [c]">View the content page [c]</a>
          </li>
        </ul>
      </div>
    </div>
  </div>

```

```

        </span>
    </li>
    <li id="ca-talk">
        <span>
            <a accesskey="t" href="/wiki/Talk:Web_scraping" rel="discussion" title="Discussion ab
            Talk
            </a>
        </span>
    </li>
</ul>
</div>
<div aria-labelledby="p-variants-label" class="vectorMenu emptyPortlet" id="p-variants" r
    <input aria-labelledby="p-variants-label" class="vectorMenuCheckbox" type="checkbox"/>
    <h3 id="p-variants-label">
        <span>
            Variants
        </span>
    </h3>
    <ul class="menu">
    </ul>
</div>
</div>
<div id="right-navigation">
    <div aria-labelledby="p-views-label" class="vectorTabs" id="p-views" role="navigation">
        <h3 id="p-views-label">
            Views
        </h3>
        <ul>
            <li class="collapsible selected" id="ca-view">
                <span>
                    <a href="/wiki/Web_scraping">
                        Read
                    </a>
                </span>
            </li>
            <li class="collapsible" id="ca-edit">
                <span>
                    <a accesskey="e" href="/w/index.php?title=Web_scraping&action=edit" title="Edit t
                    Edit
                    </a>
                </span>
            </li>
            <li class="collapsible" id="ca-history">
                <span>
                    <a accesskey="h" href="/w/index.php?title=Web_scraping&action=history" title="Pas
                    View history
                    </a>
                </span>
            </li>
        </ul>
    </div>
</div>

```

```

        </li>
    </ul>
</div>
<div aria-labelledby="p-cactions-label" class="vectorMenu emptyPortlet" id="p-cactions" r
    <input aria-labelledby="p-cactions-label" class="vectorMenuCheckbox" type="checkbox"/>
    <h3 id="p-cactions-label">
        <span>
            More
        </span>
    </h3>
    <ul class="menu">
    </ul>
</div>
<div id="p-search" role="search">
    <h3>
        <label for="searchInput">
            Search
        </label>
    </h3>
    <form action="/w/index.php" id="searchform">
        <div id="simpleSearch">
            <input accesskey="f" id="searchInput" name="search" placeholder="Search Wikipedia" tit
            <input name="title" type="hidden" value="Special:Search"/>
            <input class="searchButton mw-fallbackSearchButton" id="mw-searchButton" name="fulltex
            <input class="searchButton" id="searchButton" name="go" title="Go to a page with this c
        </div>
    </form>
</div>
</div>
</div>
<div id="mw-panel">
    <div id="p-logo" role="banner">
        <a class="mw-wiki-logo" href="/wiki/Main_Page" title="Visit the main page">
        </a>
    </div>
    <div aria-labelledby="p-navigation-label" class="portal" id="p-navigation" role="navigation
    <h3 id="p-navigation-label">
        Navigation
    </h3>
    <div class="body">
        <ul>
            <li id="n-mainpage-description">
                <a accesskey="z" href="/wiki/Main_Page" title="Visit the main page [z]">
                    Main page
                </a>
            </li>
            <li id="n-contents">
                <a href="/wiki/Portal:Contents" title="Guides to browsing Wikipedia">

```



```

        Contents
    </a>
</li>
<li id="n-featuredcontent">
    <a href="/wiki/Portal:Featured_content" title="Featured content the best of Wikipedia">
        Featured content
    </a>
</li>
<li id="n-currentevents">
    <a href="/wiki/Portal:Current_events" title="Find background information on current events">
        Current events
    </a>
</li>
<li id="n-randompage">
    <a accesskey="x" href="/wiki/Special:Random" title="Load a random article [x]">
        Random article
    </a>
</li>
<li id="n-sitesupport">
    <a href="https://donate.wikimedia.org/wiki/Special:FundraiserRedirector?utm_source=donatewikimediafoundation">
        Donate to Wikipedia
    </a>
</li>
<li id="n-shoplink">
    <a href="//shop.wikimedia.org" title="Visit the Wikipedia store">
        Wikipedia store
    </a>
</li>
</ul>
</div>
</div>
<div aria-labelledby="p-interaction-label" class="portal" id="p-interaction" role="navigation">
    <h3 id="p-interaction-label">
        Interaction
    </h3>
    <div class="body">
        <ul>
            <li id="n-help">
                <a href="/wiki/Help:Contents" title="Guidance on how to use and edit Wikipedia">
                    Help
                </a>
            </li>
            <li id="n-aboutsite">
                <a href="/wiki/Wikipedia:About" title="Find out about Wikipedia">
                    About Wikipedia
                </a>
            </li>
            <li id="n-portal">

```

```

        <a href="/wiki/Wikipedia:Community_portal" title="About the project, what you can do, v
        Community portal
    </a>
</li>
<li id="n-recentchanges">
    <a accesskey="r" href="/wiki/Special:RecentChanges" title="A list of recent changes in
    Recent changes
    </a>
</li>
<li id="n-contactpage">
    <a href="//en.wikipedia.org/wiki/Wikipedia:Contact_us" title="How to contact Wikipedia
    Contact page
    </a>
</li>
</ul>
</div>
</div>
<div aria-labelledby="p-tb-label" class="portal" id="p-tb" role="navigation">
    <h3 id="p-tb-label">
        Tools
    </h3>
    <div class="body">
        <ul>
            <li id="t-whatlinkshere">
                <a accesskey="j" href="/wiki/Special:WhatLinksHere/Web_scraping" title="List of all Eng
                What links here
                </a>
            </li>
            <li id="t-recentchangeslinked">
                <a accesskey="k" href="/wiki/Special:RecentChangesLinked/Web_scraping" rel="nofollow" t
                Related changes
                </a>
            </li>
            <li id="t-upload">
                <a accesskey="u" href="/wiki/Wikipedia:File_Upload_Wizard" title="Upload files [u]">
                    Upload file
                </a>
            </li>
            <li id="t-specialpages">
                <a accesskey="q" href="/wiki/Special:SpecialPages" title="A list of all special pages
                Special pages
                </a>
            </li>
            <li id="t-permalink">
                <a href="/w/index.php?title=Web_scraping&oldid=897184947" title="Permanent link to
                Permanent link
                </a>
            </li>

```

```

<li id="t-info">
  <a href="/w/index.php?title=Web_scraping&action=info" title="More information about
    Page information
  </a>
</li>
<li id="t-wikibase">
  <a accesskey="g" href="https://www.wikidata.org/wiki/Special:EntityPage/Q665452" title=
    Wikidata item
  </a>
</li>
<li id="t-cite">
  <a href="/w/index.php?title=Special:CiteThisPage&page=Web_scraping&id=89718494"
    Cite this page
  </a>
</li>
</ul>
</div>
</div>
<div aria-labelledby="p-coll-print_export-label" class="portal" id="p-coll-print_export" r
<h3 id="p-coll-print_export-label">
  Print/export
</h3>
<div class="body">
<ul>
<li id="coll-create_a_book">
  <a href="/w/index.php?title=Special:Book&bookcmd=book_creator&referer=Web+scrap
    Create a book
  </a>
</li>
<li id="coll-download-as-rdf2latex">
  <a href="/w/index.php?title=Special:ElectronPdf&page=Web+scraping&action=show-c
    Download as PDF
  </a>
</li>
<li id="t-print">
  <a accesskey="p" href="/w/index.php?title=Web_scraping&printable=yes" title="Print
    Printable version
  </a>
</li>
</ul>
</div>
</div>
<div aria-labelledby="p-lang-label" class="portal" id="p-lang" role="navigation">
<h3 id="p-lang-label">
  Languages
</h3>
<div class="body">
<ul>

```

```

<li class="interlanguage-link interwiki-ar">
  <a class="interlanguage-link-target" href="https://ar.wikipedia.org/wiki/%D8%A7%D8%B3%
    </a>
</li>
<li class="interlanguage-link interwiki-ca">
  <a class="interlanguage-link-target" href="https://ca.wikipedia.org/wiki/Web_scraping"
    Català
  </a>
</li>
<li class="interlanguage-link interwiki-de">
  <a class="interlanguage-link-target" href="https://de.wikipedia.org/wiki/Screen_Scraping"
    Deutsch
  </a>
</li>
<li class="interlanguage-link interwiki-es">
  <a class="interlanguage-link-target" href="https://es.wikipedia.org/wiki/Web_scraping"
    Español
  </a>
</li>
<li class="interlanguage-link interwiki-eu">
  <a class="interlanguage-link-target" href="https://eu.wikipedia.org/wiki/Web_scraping"
    Euskara
  </a>
</li>
<li class="interlanguage-link interwiki-fr">
  <a class="interlanguage-link-target" href="https://fr.wikipedia.org/wiki/Web_scraping"
    Français
  </a>
</li>
<li class="interlanguage-link interwiki-id">
  <a class="interlanguage-link-target" href="https://id.wikipedia.org/wiki/Web_scraping"
    Bahasa Indonesia
  </a>
</li>
<li class="interlanguage-link interwiki-is">
  <a class="interlanguage-link-target" href="https://is.wikipedia.org/wiki/Vefs%C3%B6fnun"
    Íslenska
  </a>
</li>
<li class="interlanguage-link interwiki-it">
  <a class="interlanguage-link-target" href="https://it.wikipedia.org/wiki/Web_scraping"
    Italiano
  </a>
</li>
<li class="interlanguage-link interwiki-lv">
  <a class="interlanguage-link-target" href="https://lv.wikipedia.org/wiki/Rasmo%C5%A1lana"
    Latviešu

```

```

    </a>
</li>
<li class="interlanguage-link interwiki-nl">
    <a class="interlanguage-link-target" href="https://nl.wikipedia.org/wiki/Scrapen" href=
    Nederlands
    </a>
</li>
<li class="interlanguage-link interwiki-ja">
    <a class="interlanguage-link-target" href="https://ja.wikipedia.org/wiki/%E3%82%A6%E3%

    </a>
</li>
<li class="interlanguage-link interwiki-sr">
    <a class="interlanguage-link-target" href="https://sr.wikipedia.org/wiki/Web_scraping"
    / srpski
    </a>
</li>
<li class="interlanguage-link interwiki-tr">
    <a class="interlanguage-link-target" href="https://tr.wikipedia.org/wiki/Web_kaz%C4%B1
    Türkçe
    </a>
</li>
<li class="interlanguage-link interwiki-uk">
    <a class="interlanguage-link-target" href="https://uk.wikipedia.org/wiki/Web_scraping"

    </a>
</li>
<li class="interlanguage-link interwiki-zh">
    <a class="interlanguage-link-target" href="https://zh.wikipedia.org/wiki/%E7%BD%91%E9%

    </a>
</li>
</ul>
<div class="after-portlet after-portlet-lang">
    <span class="wb-langlinks-edit wb-langlinks-link">
        <a class="wbc-editpage" href="https://www.wikidata.org/wiki/Special:EntityPage/Q665452"
        Edit links
        </a>
    </span>
</div>
</div>
</div>
<div id="footer" role="contentinfo">
    <ul id="footer-info">
        <li id="footer-info-lastmod">
            This page was last edited on 15 May 2019, at 09:54

```

```

    <span class="anonymous-show">
      (UTC)
    </span>
  .
</li>
<li id="footer-info-copyright">
  Text is available under the
  <a href="//en.wikipedia.org/wiki/Wikipedia:Text_of_Creative_Commons_Attribution-ShareAlike_License">
    Creative Commons Attribution-ShareAlike License
  </a>
  <a href="//creativecommons.org/licenses/by-sa/3.0/" rel="license" style="display:none;">
  </a>
  ;
  additional terms may apply. By using this site, you agree to the
  <a href="//foundation.wikimedia.org/wiki/Terms_of_Use">
    Terms of Use
  </a>
  and
  <a href="//foundation.wikimedia.org/wiki/Privacy_policy">
    Privacy Policy
  </a>
  . Wikipedia is a registered trademark of the
  <a href="//www.wikimediafoundation.org/">
    Wikimedia Foundation, Inc.
  </a>
  , a non-profit organization.
</li>
</ul>
<ul id="footer-places">
  <li id="footer-places-privacy">
    <a class="extiw" href="https://foundation.wikimedia.org/wiki/Privacy_policy" title="wmf:Privacy policy">
      Privacy policy
    </a>
  </li>
  <li id="footer-places-about">
    <a href="/wiki/Wikipedia:About" title="Wikipedia:About">
      About Wikipedia
    </a>
  </li>
  <li id="footer-places-disclaimer">
    <a href="/wiki/Wikipedia:General_disclaimer" title="Wikipedia:General disclaimer">
      Disclaimers
    </a>
  </li>
  <li id="footer-places-contact">
    <a href="//en.wikipedia.org/wiki/Wikipedia:Contact_us">
      Contact Wikipedia
    </a>
  </li>
</ul>

```

```

</li>
<li id="footer-places-developers">
  <a href="https://www.mediawiki.org/wiki/Special:MyLanguage/How_to_contribute">
    Developers
  </a>
</li>
<li id="footer-places-cookistatement">
  <a href="https://foundation.wikimedia.org/wiki/Cookie_statement">
    Cookie statement
  </a>
</li>
<li id="footer-places-mobileview">
  <a class="noprint stopMobileRedirectToggle" href="//en.m.wikipedia.org/w/index.php?title=
    Mobile view
  </a>
</li>
</ul>
<ul class="noprint" id="footer-icons">
<li id="footer-copyrightico">
  <a href="https://wikimediafoundation.org/">
    
</li>
<li id="footer-poweredbyico">
  <a href="//www.mediawiki.org/">
    
</div>
</div>
<script>
  (window.RLQ=window.RLQ||[]).push(function(){mw.config.set({"wgPageParseReport":{"limitreport
</script>
<script type="application/ld+json">
  {"@context":"https://schema.org","@type":"Article","name":"Web scraping","url":"https://\
</script>
<script>
  (window.RLQ=window.RLQ||[]).push(function(){mw.config.set({"wgBackendResponseTime":807,"wgH
</script>
</body>
</html>

```

Podemos utilizar um seletor CSS para extrair somente o corpo do HTML, veja que agora já está mais fácil encontrar o primeiro parágrafo, basta localizar o < p> a esquerda:

```
In [68]: #print(soup.find('body').prettify()) # é uma outra opção via find
soup.select("html body")
```

```
Out [68]: [<body class="mediawiki ltr sitedir-ltr mw-hide-empty-elt ns-0 ns-subject mw-editable"
  <div class="noprint" id="mw-page-base"></div>
  <div class="noprint" id="mw-head-base"></div>
  <div class="mw-body" id="content" role="main">
  <a id="top"></a>
  <div class="mw-body-content" id="siteNotice"><!-- CentralNotice --></div>
  <div class="mw-indicators mw-body-content">
  </div>
  <h1 class="firstHeading" id="firstHeading" lang="en">Web scraping</h1>
  <div class="mw-body-content" id="bodyContent">
  <div class="noprint" id="siteSub">From Wikipedia, the free encyclopedia</div>
  <div id="contentSub"></div>
  <div id="jump-to-nav"></div>
  <a class="mw-jump-link" href="#mw-head">Jump to navigation</a>
  <a class="mw-jump-link" href="#p-search">Jump to search</a>
  <div class="mw-content-ltr" dir="ltr" id="mw-content-text" lang="en"><div class="mw-jump-link">
  <div class="hatnote navigation-not-searchable" role="note">For broader coverage of this topic, see
  <p><b>Web scraping</b>, <b>web harvesting</b>, or <b>web data extraction</b> is a computer
  </p><p>Web scraping a web page involves fetching it and extracting from it.<sup class="wikitext">
  </p><p>Web scraping is used for Contact scraping
  </p><p><a href="/wiki/Web_page" title="Web page">Web pages</a> are built using text-based
  </p><p>Newer forms of web scraping involve listening to data feeds from web servers.
  </p><p>There are methods that some websites use to prevent web scraping, such as detecting
  </p>
  <div class="toc" id="toc"><input class="toctogglecheckbox" id="toctogglecheckbox" type="checkbox">
  <ul>
  <li class="toclevel-1 tocsection-1"><a href="#History"><span class="tocnumber">1</span><span>History</span></a>
  <li class="toclevel-1 tocsection-2"><a href="#Techniques"><span class="tocnumber">2</span><span>Techniques</span></a>
  <ul>
  <li class="toclevel-2 tocsection-3"><a href="#Human_copy-and-paste"><span class="tocnumber">2.1</span><span>Human copy-and-paste</span></a>
  <li class="toclevel-2 tocsection-4"><a href="#Text_pattern_matching"><span class="tocnumber">2.2</span><span>Text pattern matching</span></a>
  <li class="toclevel-2 tocsection-5"><a href="#HTTP_programming"><span class="tocnumber">2.3</span><span>HTTP programming</span></a>
  <li class="toclevel-2 tocsection-6"><a href="#HTML_parsing"><span class="tocnumber">2.4</span><span>HTML parsing</span></a>
  <li class="toclevel-2 tocsection-7"><a href="#DOM_parsing"><span class="tocnumber">2.5</span><span>DOM parsing</span></a>
  <li class="toclevel-2 tocsection-8"><a href="#Vertical_aggregation"><span class="tocnumber">2.6</span><span>Vertical aggregation</span></a>
  <li class="toclevel-2 tocsection-9"><a href="#Semantic_annotation_recognizing"><span class="tocnumber">2.7</span><span>Semantic annotation recognizing</span></a>
  <li class="toclevel-2 tocsection-10"><a href="#Computer_vision_web-page_analysis"><span class="tocnumber">2.8</span><span>Computer vision web-page analysis</span></a>
  </ul>
  </li>
  <li class="toclevel-1 tocsection-11"><a href="#Software"><span class="tocnumber">3</span><span>Software</span></a>
  <ul>
  <li class="toclevel-2 tocsection-12"><a href="#Example_tools"><span class="tocnumber">3.1</span><span>Example tools</span></a>
  <ul>
  <li class="toclevel-3 tocsection-13"><a href="#Javascript_tools"><span class="tocnumber">3.1.1</span><span>Javascript tools</span></a>
  <li class="toclevel-3 tocsection-14"><a href="#Web_crawling_frameworks"><span class="tocnumber">3.1.2</span><span>Web crawling frameworks</span></a>
```



```

</ul>
</li>
</ul>
</li>
<li class="toclevel-1 tocsection-15"><a href="#Legal_issues"><span class="tocnumber">
<ul>
<li class="toclevel-2 tocsection-16"><a href="#United_States"><span class="tocnumber">
<li class="toclevel-2 tocsection-17"><a href="#The_EU"><span class="tocnumber">4.2</span>
<li class="toclevel-2 tocsection-18"><a href="#Australia"><span class="tocnumber">4.3
</ul>
</li>
<li class="toclevel-1 tocsection-19"><a href="#Methods_to_prevent_web_scraping"><span
<li class="toclevel-1 tocsection-20"><a href="#See_also"><span class="tocnumber">6</span>
<li class="toclevel-1 tocsection-21"><a href="#References"><span class="tocnumber">7
</ul>
</div>
<h2><span class="mw-headline" id="History">History</span><span class="mw-editsection">
<table class="box-Unreferenced_section plainlinks metadata ambox ambox-content ambox">
<p>The history of the web scraping is actually much longer, dating back significantly
</p>
<ul><li>After the birth of <a href="/wiki/History_of_the_World_Wide_Web" title="Hist
<li>In 1993,December, the First <b>crawler-based web search engine</b> - <a href="/w
<li>In 2000, the <b>first Web API and API crawler</b> came. <a href="/wiki/Applicati
<li>In 2004, <a href="/wiki/Beautiful_Soup_(HTML_parser)" title="Beautiful Soup (HTML
<h2><span class="mw-headline" id="Techniques">Techniques</span><span class="mw-edits
<p>Web scraping is the process of automatically mining data or collecting information
</p>
<h3><span class="mw-headline" id="Human_copy-and-paste">Human copy-and-paste</span><
<p>Sometimes even the best web-scraping technology cannot replace a humans manual ex
</p>
<h3><span class="mw-headline" id="Text_pattern_matching">Text pattern matching</span>
<p>A simple yet powerful approach to extract information from web pages can be based
</p>
<h3><span class="mw-headline" id="HTTP_programming">HTTP programming</span><span clas
<p><a href="/wiki/Static_web_page" title="Static web page">Static</a> and <a href="/v
</p>
<h3><span class="mw-headline" id="HTML_parsing">HTML parsing</span><span class="mw-e
<p>Many websites have large collections of pages generated dynamically from an under
</p>
<h3><span class="mw-headline" id="DOM_parsing">DOM parsing</span><span class="mw-ed
<div class="hatnote navigation-not-searchable" role="note">Further information: <a href="/wiki/Internet_Expl
<p>By embedding a full-fledged web browser, such as the <a href="/wiki/Internet_Expl
</p>
<h3><span class="mw-headline" id="Vertical_aggregation">Vertical aggregation</span><
<p>There are several companies that have developed vertical specific harvesting plat
</p>
<h3><span class="mw-headline" id="Semantic_annotation_recognizing">Semantic annotati
<p>The pages being scraped may embrace <a href="/wiki/Metadata" title="Metadata">meta

```

Computer vision

There are efforts using [Machine learning](/wiki/Machine_learning "Machine learning").

Software

There are many software tools available that can be used to customize web-scraping.

Example tools

- [CURL](/wiki/CURL "CURL") command line tool and library for
- [Data Toolbar](/wiki/Data_Toolbar "Data Toolbar") web scraping
- [Diffbot](/wiki/Diffbot "Diffbot") uses computer vision and machine learning
- [Heritrix](/wiki/Heritrix "Heritrix") gets pages (lots of them)
- [HtmlUnit](/wiki/HtmlUnit "HtmlUnit") headless browser that can
- [HTTrack](/wiki/HTTrack "HTTrack") free and open source Web crawler
- [iMacros](/wiki/IMacros "IMacros") a browser extension to record
- [Selenium \(software\)](/wiki/Selenium_(software) "Selenium (software)")
- [Aptana](/wiki/Aptana#Aptana_Jaxer "Aptana") Jaxer
- [Mozenda](/wiki/Mozenda "Mozenda") is a WYSIWYG software that
- [nokogiri](/wiki/Nokogiri_(software) "Nokogiri (software)")
- [OutWit Hub](/wiki/OutWit_Hub "OutWit Hub") Web scraping application
- [watir](/wiki/Watir "Watir")
- [Wget](/wiki/Wget "Wget") computer program that retrieves content
- [WSO2 Mashup Server](/wiki/WSO2_Mashup_Server "WSO2 Mashup Server")
- [Yahoo! Query Language](/wiki/Yahoo!_Query_Language "Yahoo! Query Language")

Javascript tools

- [Greasemonkey](/wiki/Greasemonkey "Greasemonkey")
- [Node.js](/wiki/Node.js "Node.js")
- [PhantomJS](/wiki/PhantomJS "PhantomJS") scripted, [Puppeteer](/wiki/Puppeteer "Puppeteer")
- [jQuery](/wiki/JQuery "jQuery")

Web crawling frameworks

These can be used to build web scrapers.

- [Scrapy](/wiki/Scrapy "Scrapy")

Legal issues

--

The legality of web scraping varies across the world. In general, web scraping may be illegal.

United States

--

In the United States, website owners can use three major [Cause of action](/wiki/Cause_of_action "Cause of action") to sue web scrapers.

U.S. courts have acknowledged that users of "scrapers" or "robots" may be held liable for copyright infringement.

One of the first major tests of [Screen scraping](/wiki/Screen_scraping "Screen scraping") was in [Southwest Airlines](/wiki/Southwest_Airlines "Southwest Airlines") v. [Flightradar24](/wiki/Flightradar24 "Flightradar24").

In 2012, a startup called 3Taps scraped classified housing ads from Craigslist. Craigslist sued 3Taps for copyright infringement.

Although these are early scraping decisions, and the theories of liability are still developing, they have set a precedent.

While the law in this area becomes more settled, entities contemplating using web scrapers should consult with legal counsel.

In the plaintiff's web site during the period of this trial the terms of use of the defendant's website stated:

In *Facebook, Inc. v. Power Ventures, Inc.*, the Ninth Circuit held that

Internet Archive collects and distributes significant number of publicly available documents.

The EU

In February 2006, the Danish Maritime and Commercial Court (Copenhagen) ruled that the Danish Copyright Act of 1997 was not compatible with the EU Directive on Rental Right and Lending Right.

In a February 2010 case complicated by matters of jurisdiction, Ireland's High Court ruled that the Copyright Act of 1962 was not compatible with the EU Directive on Rental Right and Lending Right.

Australia

In Australia, the [Spam Act 2003](/wiki/Spam_Act_2003 "Spam Act 2003") provides a legal framework for dealing with spam.

Methods to prevent web scraping

The administrator of a website can use various measures to stop or slow a bot. Some of the measures are:

- Blocking an [IP address](/wiki/IP_address "IP address") either temporarily or permanently.
- Disabling any [web service](/wiki/Web_service "Web service") that is used by the bot.
- Bots sometimes declare who they are (using [User agent](/wiki/User_agent "User agent")).
- Bots can be blocked by monitoring excess traffic.
- Bots can sometimes be blocked with tools to verify that it is a real person accessing the website.
- Commercial anti-bot services: Companies offer anti-bot and anti-scraping services.

- Locating bots with a [Honeypot \(computing\)](/wiki/Honeypot_(computing) "Honeypot (computing)").
- [Obfuscation](/wiki/Obfuscation "Obfuscation") using [JavaScript](/wiki/JavaScript "JavaScript").
- Because bots rely on consistency in the front-end code of a target website, adding random elements can help.
- Websites can declare if crawling is allowed or not in the [Robots.txt](/wiki/Robots_txt "Robots.txt") file.

See also

- [Archive.is](/wiki/Archive.is "Archive.is")
- [Comparison of feed aggregators](/wiki/Comparison_of_feed_aggregators "Comparison of feed aggregators")
- [Data scraping](/wiki/Data_scraping "Data scraping")
- [Data wrangling](/wiki/Data_wrangling "Data wrangling")
- [Importer \(computing\)](/wiki/Importer_(computing) "Importer (computing)")
- [Job wrapping](/wiki/Job_wrapping "Job wrapping")
- [Knowledge extraction](/wiki/Knowledge_extraction "Knowledge extraction")
- [OpenSocial](/wiki/OpenSocial "OpenSocial")
- [Scraper site](/wiki/Scraper_site "Scraper site")
- [Fake news website](/wiki/Fake_news_website "Fake news website")
- [Blog scraping](/wiki/Blog_scraping "Blog scraping")
- [Spamdexing](/wiki/Spamdexing "Spamdexing")
- [Domain name drop list](/wiki/Domain_name_drop_list "Domain name drop list")
- [Text corpus](/wiki/Text_corpus "Text corpus")
- [Web archiving](/wiki/Web_archiving "Web archiving")
- [Blog network](/wiki/Blog_network "Blog network")
- [Search Engine Scraping](/wiki/Search_Engine_Scraping "Search Engine Scraping")
- [Web crawlers](/wiki/Category:Web_crawlers "Category:Web crawlers")

References

- ↑ [^](#cite_note-Boeing2016JPER-1) [^](#cite_note-Boeing2016JPER-1)


```

</li>
<li id="cite_note-26"><span class="mw-cite-backlink"><b><a href="#cite_ref-26">^</a>
</li>
</ol></div>
<!--
NewPP limit report
Parsed by mw1264
Cached time: 20190516030922
Cache expiry: 2592000
Dynamic content: false
CPU time usage: 0.492 seconds
Real time usage: 0.660 seconds
Preprocessor visited node count: 1805/1000000
Preprocessor generated node count: 0/1500000
Postexpand include size: 65153/2097152 bytes
Template argument size: 772/2097152 bytes
Highest expansion depth: 15/40
Expensive parser function count: 8/500
Unstrip recursion depth: 1/20
Unstrip postexpand size: 70826/5000000 bytes
Number of Wikibase entities loaded: 3/400
Lua time usage: 0.279/10.000 seconds
Lua memory usage: 6.63 MB/50 MB
-->
<!--
Transclusion expansion time report (%,ms,calls,template)
100.00% 548.078      1 -total
 59.50% 326.090      1 Template:Reflist
 38.08% 208.730      5 Template:Cite_journal
 19.94% 109.301      1 Template:More_citations_needed
 18.46% 101.199      4 Template:Ambox
 16.32%  89.451     18 Template:Cite_web
  6.60%  36.163      1 Template:Find_sources_mainspace
  5.52%  30.244      1 Template:Split_section
  4.82%  26.437      1 Template:Split_portions
  2.82%  15.468      1 Template:US-centric
-->
<!-- Saved in parser cache with key enwiki:pcache:idhash:2696619-0!canonical and time
-->
</div><noscript>Retrieved from "<a dir="ltr" href="https://en.wikipedia.org
<div class="catlinks" data-mw="interface" id="catlinks"><div class="mw-normal-catlin
<div class="visualClear"></div>
</div>
</div>
<div id="mw-navigation">
<h2>Navigation menu</h2>
<div id="mw-head">

```

```

<div aria-labelledby="p-personal-label" id="p-personal" role="navigation">
<h3 id="p-personal-label">Personal tools</h3>
<ul>
<li id="pt-anonuserpage">Not logged in</li><li id="pt-anontalk"><a accesskey="n" href=
</div>
<div id="left-navigation">
<div aria-labelledby="p-namespaces-label" class="vectorTabs" id="p-namespaces" role=
<h3 id="p-namespaces-label">Namespaces</h3>
<ul>
<li class="selected" id="ca-nstab-main"><span><a accesskey="c" href="/wiki/Web_scrap
</div>
<div aria-labelledby="p-variants-label" class="vectorMenu emptyPortlet" id="p-variant
<input aria-labelledby="p-variants-label" class="vectorMenuCheckbox" type="checkbox"/
<h3 id="p-variants-label">
<span>Variants</span>
</h3>
<ul class="menu">
</ul>
</div>
</div>
<div id="right-navigation">
<div aria-labelledby="p-views-label" class="vectorTabs" id="p-views" role="navigation
<h3 id="p-views-label">Views</h3>
<ul>
<li class="collapsible selected" id="ca-view"><span><a href="/wiki/Web_scraping">Rea
</div>
<div aria-labelledby="p-cactions-label" class="vectorMenu emptyPortlet" id="p-caction
<input aria-labelledby="p-cactions-label" class="vectorMenuCheckbox" type="checkbox"/
<h3 id="p-cactions-label"><span>More</span></h3>
<ul class="menu">
</ul>
</div>
<div id="p-search" role="search">
<h3>
<label for="searchInput">Search</label>
</h3>
<form action="/w/index.php" id="searchform">
<div id="simpleSearch">
<input accesskey="f" id="searchInput" name="search" placeholder="Search Wikipedia" t
</form>
</div>
</div>
</div>
<div id="mw-panel">
<div id="p-logo" role="banner"><a class="mw-wiki-logo" href="/wiki/Main_Page" title=
<div aria-labelledby="p-navigation-label" class="portal" id="p-navigation" role="nav
<h3 id="p-navigation-label">Navigation</h3>
<div class="body">

```

```

<ul>
<li id="n-mainpage-description"><a accesskey="z" href="/wiki/Main_Page" title="Visit
</div>
</div>
<div aria-labelledby="p-interaction-label" class="portal" id="p-interaction" role="n
<h3 id="p-interaction-label">Interaction</h3>
<div class="body">
<ul>
<li id="n-help"><a href="/wiki/Help:Contents" title="Guidance on how to use and edit
</div>
</div>
<div aria-labelledby="p-tb-label" class="portal" id="p-tb" role="navigation">
<h3 id="p-tb-label">Tools</h3>
<div class="body">
<ul>
<li id="t-whatlinkshere"><a accesskey="j" href="/wiki/Special:WhatLinksHere/Web_scrap
</div>
</div>
<div aria-labelledby="p-coll-print_export-label" class="portal" id="p-coll-print_exp
<h3 id="p-coll-print_export-label">Print/export</h3>
<div class="body">
<ul>
<li id="coll-create_a_book"><a href="/w/index.php?title=Special:Book&bookcmd=boo
</div>
</div>
<div aria-labelledby="p-lang-label" class="portal" id="p-lang" role="navigation">
<h3 id="p-lang-label">Languages</h3>
<div class="body">
<ul>
<li class="interlanguage-link interwiki-ar"><a class="interlanguage-link-target" hre
<div class="after-portlet after-portlet-lang"><span class="wb-langlinks-edit wb-lang
</div>
</div>
</div>
<div id="footer" role="contentinfo">
<ul id="footer-info">
<li id="footer-info-lastmod"> This page was last edited on 15 May 2019, at 09:54<span
<li id="footer-info-copyright">Text is available under the <a href="//en.wikipedia.o
additional terms may apply. By using this site, you agree to the <a href="//foundat
</ul>
<ul id="footer-places">
<li id="footer-places-privacy"><a class="extiw" href="https://foundation.wikimedia.o
<li id="footer-places-about"><a href="/wiki/Wikipedia:About" title="Wikipedia:About"
<li id="footer-places-disclaimer"><a href="/wiki/Wikipedia:General_disclaimer" title
<li id="footer-places-contact"><a href="//en.wikipedia.org/wiki/Wikipedia:Contact_us
<li id="footer-places-developers"><a href="https://www.mediawiki.org/wiki/Special:My
<li id="footer-places-cookiestatement"><a href="https://foundation.wikimedia.org/wik
<li id="footer-places-mobileview"><a class="noprint stopMobileRedirectToggle" href="

```

```

</ul>
<ul class="noprint" id="footer-icons">
<li id="footer-copyrightico">
<a href="https://wikimediafoundation.org/">
<li id="footer-poweredbyico">
<a href="//www.mediawiki.org/">
</ul>
<div style="clear: both;"></div>
</div>
<script>(window.RLQ=window.RLQ||[]).push(function(){mw.config.set({"wgPageParseReportUrl":"/w/index.php?title=Special:PageParseReport&from=Web_scraping&to=Web_scraping"});});
<script type="application/ld+json">{"@context":"https://schema.org","@type":"Article","mainEntity":{},"mainEntityList":[]};
<script>(window.RLQ=window.RLQ||[]).push(function(){mw.config.set({"wgBackendResponseTime":0});});
</body>

```

Podemos selecionar o primeiro parágrafo utilizando o método `find` para localizar uma *tag* HTML e o `get_text` para extrair somente o texto:

```
In [71]: print(soup.find('p').get_text())
```

Web scraping, web harvesting, or web data extraction is data scraping used for extracting data

Veja que estamos próximos do nosso objetivo final, basta usar o método `find_all` que conseguiremos encontrar todas as tags `p`'s. Atente que os resultados serão retornados em uma lista:

```
In [73]: for item in soup.find_all('p')[:3]: print(item.get_text())
```

Web scraping, web harvesting, or web data extraction is data scraping used for extracting data

Web scraping a web page involves fetching it and extracting from it.[1][2] Fetching is the downloading of a web page from a web server.

Web scraping is used for contact scraping, and as a component of applications used for web indexing.

Resumindo, poderíamos ter utilizado somente os seguintes comandos:

```
In [74]: import requests
        from bs4 import BeautifulSoup

        page = requests.get("https://en.wikipedia.org/wiki/Web_scraping")
        soup = BeautifulSoup(page.content, 'html.parser')
        for item in soup.find_all('p')[:3]: print(item.get_text())
```

Web scraping, web harvesting, or web data extraction is data scraping used for extracting data

Web scraping a web page involves fetching it and extracting from it.[1][2] Fetching is the downloading of a web page from a web server.

Web scraping is used for contact scraping, and as a component of applications used for web indexing.

Para fixar as ideias veja como seria simples retornar os dois primeiros parágrafo na wikipédia sobre HTML, CSS e JS em três passos: + Lista de *url*'s + Função para retornar o parágrafo + Loop para aplicar a função a cada item na lista

Passo 1:

```
In [105]: paginas = [  
    'https://pt.wikipedia.org/wiki/HTML',  
    'https://pt.wikipedia.org/wiki/Cascading_Style_Sheets',  
    'https://pt.wikipedia.org/wiki/JavaScript'  
]
```

paginas

```
Out[105]: ['https://pt.wikipedia.org/wiki/HTML',  
    'https://pt.wikipedia.org/wiki/Cascading_Style_Sheets',  
    'https://pt.wikipedia.org/wiki/JavaScript']
```

Passo 2:

```
In [106]: def retorna_n_paragrafos(url, n = 1):  
    '''  
    Essa função retorna os n primeiros parágrafos da url selecionada.  
    Possui os seguintes argumentos:  
    @url: a url desejada  
    @n: número de parágrafos para retornar, default = 1  
    '''  
  
    page = requests.get(url)  
    soup = BeautifulSoup(page.content, 'html.parser')  
    paragrafos = soup.find_all('p')[:n]  
    paragrafos_texto = []  
    for item in paragrafos: paragrafos_texto.append(item.get_text())  
  
    return '/n'.join(paragrafos_texto) #concatenar elementos e separar por nova linha  
  
    print(retorna_n_paragrafos('https://pt.wikipedia.org/wiki/HTML', n = 2))
```

HTML (abreviação para a expressão inglesa HyperText Markup Language, que significa Linguagem de Marcação de Hipertexto) é um padrão para a representação estruturada de hipermídia e conteúdo baseado em tempo.

Vemos que a função funciona, o primeiro /n ocorre porque o primeiro parágrafo não tinha conteúdo.

Passo 3 (criando o loop):

```
In [107]: for pagina in paginas:  
    out = retorna_n_paragrafos(pagina, n = 2)  
    print(pagina + '\n' + out)
```

<https://pt.wikipedia.org/wiki/HTML>

HTML (abreviação para a expressão inglesa HyperText Markup Language, que significa Linguagem de Marcação de Hipertexto) é uma linguagem de marcação para a representação estruturada de hipermídia e conteúdo baseado em tempo.

https://pt.wikipedia.org/wiki/Cascading_Style_Sheets

Cascading Style Sheets (CSS) é um mecanismo para adicionar estilo (cores, fontes, espaçamento, etc.) a documentos HTML. O código CSS pode ser aplicado diretamente nas tags ou ficar contido dentro das tags <style>.

<https://pt.wikipedia.org/wiki/JavaScript>

JavaScript, frequentemente abreviado como JS, é uma linguagem de programação interpretada de script. É atualmente a principal linguagem para programação client-side em navegadores web. É também

In []: Fim.