| | Web scraping, also known as web harvesting or web data extraction, is the process of automatically extracting data from websites. It involves fetching a web page and extracting from it. The data can be parsed, searched, reformatted, and copied into a spreadsheet or loaded into a database. Web scraping can be done manually, but in most cases, automatically extracting data from websites. It involves fetching a web page and extracting from it. The data can be parsed, searched, reformatted, and copied into a spreadsheet or loaded into a database. Web scraping can be done manually, but in most cases, automatically extracting data from websites. It involves fetching a web page and extracting from it. The data can be parsed, searched, reformatted, and copied into a spreadsheet or loaded into a database. Web scraping can be done manually, but in most cases, automatically extracting data from websites. It involves fetching a web page and extracting from it. The data can be parsed, searched, reformatted, and copied into a spreadsheet or loaded into a database. Web scraping can be done manually, but in most cases, automatically extracting data from websites. It involves fetching a web page and extracting from it.

 |

--
--
--
--
---|
| | content aggregation. However, some websites use methods to prevent web scraping, such as detecting and disallowing bots from crawling their pages. In response, there are we scraping systems that rely on using techniques in DOM parsing, computer vision, and natural language processing to simulate human browsing to enable gathering web page co for offline parsing.

 |
| | Dependencies Actually to follow this module you only need to install beautifulsoup4 with pip install beautifulsoup4 and you are good to go. But here some libraries that needed to be installed first that I use at bis module:

 |
| | beautifulSoup4 pandas matplotlibs

 |
| | Background At this project I try to scrap Movie Name, Ratings, Duration, Votes, and Sinopsis Data IMDb data center website. IMBD, an acronym for Internet Movie Database, is an online data

 |
| | of information related to films, television series, podcasts, home videos, video games, and streaming content online. It includes information such as cast, production crew, person biographies, plot summaries, trivia, ratings, and fan and critical reviews. IMDb was first published in 1990 by Col Needham, a computer programmer, and has since grown into a comprehensive resource for movie and TV information. I will try to scrap this sites for educational purpose only. A lot of you might ask why we need to scrap this data from the sites. I want to make a report and gain some insight from that data and maybe can be useful for others. To do that

 |
| | need to have the data, and scrapping is a good way to collect the data I don't have from the public. I will scrap 7 points from this sites. That is Title, Year, Duration, Rating Category, Ratings, and Total Votes.

 |
| | What is BeautifulSoup Beautiful Soup is a Python library for pulling data out of HTML and XML files. Beautiful Soup 3 only works on Python 2.x, but Beautiful Soup 4 also works on Python 3.x. Beautiful 4 is faster, has more features, and works with third-party parsers like lxml and html5lib.

 |
| | Since beautifulsoup used to pull the data out of a HTML, so first I need to pull out the html first. How I do it? I will use default library request. So all this code is doing is sending a GET request to spesific address I give. This is the same type of request your browser sent to view this page, but the only difference is that

 |
| | Requests can't actually render the HTML, so instead you will just get the raw HTML and the other response information. I'm using the .get() function here, but Requests allows you to use other functions like .post() and .put() to send those requests as well. At this case we will going to the IMDb webs you can click here to follow what exactly that link goes to.

 |
| | Requesting the Data and Creating a BeautifulSoup Let's begin with requesting the web from the site with get method.

 |
| | Getting the HTML from the Webpage url = 'https://www.imdb.com/chart/top/?ref_=nv_mv_250'

 |
| | <pre>import requests headers = { 'User-Agent':'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/119.0.0.0 Safari/537.36'</pre>

 |
| | url_get = requests.get(url, headers=headers) To visualize what exactly you get from the request.get, we can use .content so ee what we exactly get, in here i slice it so it won't make our screen full of the html we get from page. You can delete the slicing if you want to see what we fully get.

 |
| 9]: | <pre>url_get.content[:500] b'<!DOCTYPE html> <html lang="en-US" xmlns:fb="http://www.facebook.com/2008/fbml" xmlns:og="http://opengraphprotocol.org/schema/"><head><meta t="utf-8"/><meta content="width=device-width" name="viewport"/><script>if(typeof uet === \'function\'){ uet(\'bb\', \'LoadTitle\', {wb: 1}); ipt><script>window.addEventListener(\'load\', (event) => {\n if (typeof window.csa !== \'undefined\' && typeof window.csa === \'function\')</pre></td></tr><tr><td></td><td>As we can see we get a very unstructured and complex html, which actually contains the codes needed to show the webpages on your web browser. But we as human still confu what and where we can use that piece of code, so here where we use the beautifulsoup. Beautiful soup class will result a beautifulsoup object. Beautiful Soup transforms a comp HTML document into a complex tree of Python objects. But we'll only ever have to deal with about four kinds of objects: Tag, NavigableString, BeautifulSoup, and</td></tr><tr><td></td><td>Comment . But at this project we will only use BeautifulSoup . Let's make Beautiful soup object and feel free to explore the object here.</td></tr><tr><td></td><td><pre>from bs4 import BeautifulSoup soup = BeautifulSoup(url_get.content, "html.parser") print(type(soup)) <class 'bs4.BeautifulSoup'></pre></td></tr><tr><td></td><td>Let's see how our beautifulsoup looks like. As you can see, the content is the same with our get_url object but it's tidier. Also beautifulsoup give us method to make it even more prettier, for tidyness purpouse we slice to only see first 500 character. print(soup.prettify()[:500])</td></tr><tr><td></td><td><!DOCTYPE html> <html lang="en-US" xmlns:fb="http://www.facebook.com/2008/fbml" xmlns:og="http://opengraphprotocol.org/schema/"> <head></td></tr><tr><td></td><td><pre><meta content="width=device-width" name="viewport"/> <script> if(typeof uet === 'function'){ uet('bb', 'LoadTitle', {wb: 1}); } </script> <script> window.addEventListener('load', (event) => {</pre></td></tr><tr><td></td><td><pre>if (typeof window.csa !== 'undefined' && typeof window.csa === 'function') { var csaLatencyPlugin = window.csa(</pre> Getting right key to extract right information</td></tr><tr><td></td><td>Now we already have a tidier html, now we should search the lines that we want to use. Let's back to our web page first. The information that we need are the Date. Inflation value MoM. and Inflation value YoY which contain in the table. To know which part of the code refer to that table, we can just the code of the</td></tr><tr><td></td><td>The information that we need are the Date, Inflation value MoM, and Inflation value YoY which contain in the table. To know which part of the code refer to that table, we can just our cusor there, right click, and inspect element. Then we will see something like this. From inspect element we know that we need to find the line table with class, bar is sered. We can use the find method at our beautifusoup chiest. Let's also call our chiest to</td></tr><tr><td></td><td>From inspect element we know that we need to find the line table with class baris-scroll. We can use the find method at our beautifusoup object. Let's also call our object to what we get. Finding the right key to scrap the data & Extracting the right information</td></tr><tr><td></td><td>Find the key and put the key into the .find() Put all the exploring the right key at this cell. (please change this markdown with your explanation) 1. Find Title</td></tr><tr><td></td><td><pre>title = soup.find_all("h3", attrs={'class':'ipc-titletext'})[250].text title.split('. ')[1] 'Drishyam'</pre></td></tr><tr><td>37</td><td>2. Find Year, Duration, Rating Category soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'})[749] 13+</td></tr><tr><td>37]:
]:[</td><td><pre>13+</pre> 3. Find Ratings & Total Votes</td></tr><tr><td>38 [
38]:</td><td>soup.find_all("span", attrs={'class':'ipc-rating-star ipc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating'})[1].text '9.2\xa0(2M)'</td></tr><tr><td></td><td><pre>temp = soup.find_all("span", attrs={'class':'ipc-rating-star ipc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating'})[249].te temp = temp.replace('(', '') temp = temp.replace(')','') temp '8.2\xa092K'</pre></td></tr><tr><td>40]:
41 [
41]:</td><td>temp.split('\xa0')[0]</td></tr><tr><td>42 [
42] :</td><td>temp.split('\xa0')[1] '92K'</td></tr><tr><td></td><td>As you can see get already get the necessary key to extract all needed data. To get the only text information you can add .text . Remember you need to only get one information before you use .text otherwise it will return error. Extracting the Information</td></tr><tr><td></td><td>Now all the beautiful soup part is over. All left to do is doing some programming to extract all the data automaticly, you can do this manualy at this part but if your data too lo advice you use loop. I'll show you how to make looping for extracting the data, but before that let's check how long is our data to help our looping process. Since find_all will always return data in format list, we will use len() to check how long is our list.</td></tr><tr><th></th><th>Finding row length of Date. title_row = soup.find_all("h3", attrs={'class':'ipc-titletext'}) len(title_row)</th></tr><tr><th></th><th></th></tr><tr><th>2]:</th><th>BUT the actual Title Data only exist in index list [1:251]</th></tr><tr><td>2]:
8]:</td><td></td></tr><tr><td>2]:
8]:
8]:
9]:</td><td>BUT the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'}) len(second_info)</td></tr><tr><td>2]:
8]:
8]:
9]:</td><td>BUT the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'}) 750 third_info = soup.find_all("span", attrs={'class':'ipc-rating-star ipc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating'}) len(third_info) 250 Now we know the length of our data, now here what we will do for the looping process. Title = [] Year = []</td></tr><tr><td>2]:
8]:
8]:
9]:</td><td>BUT the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'}) 750 third_info = soup.find_all("span", attrs={'class':'ipc-rating-star ipc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating'}) len(third_info) 250 Now we know the length of our data, now here what we will do for the looping process. Title = [] Year = [] Duration = [] Rating_Category = [] Ratings = [] Total_Votes = []</td></tr><tr><td>2]:
8]:
8]:
9]:</td><td>BUT the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'}) reference to third_info = soup.find_all("span", attrs={'class':'ipc-rating-star ipc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating'}) len(third_info) 250 Now we know the length of our data, now here what we will do for the looping process. Title = [] Year = [] Duration = [] Rating_Category = [] Rating_Category = [] Ratings = []</td></tr><tr><td>2]:
8]:
8]:
9]:
43</td><td>BUT the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'}) len(second_info) third_info = soup.find_all("span", attrs={'class':'ipc-rating-star ipc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating')} 250 Now we know the length of our data, now here what we will do for the looping process. Title = [] Year = [] Duration = [] Rating_Category = [] Rating_Category = [] Rating_Category = [] Total_Votes = [] Total_Votes = [] Title.append(judul) n = 0 while n < 750: tahun = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'})[n].text n = 0 the company of the</td></tr><tr><td>2]:
8]:
8]:
9]:
43</td><td>BUT the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'}) len(second_info) 750 third_info = soup.find_all("span", attrs={'class':'ipc-rating-star ipc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating'}) len(third_info) 250 Now we know the length of our data, now here what we will do for the looping process. Title = [] Year = [] Duration = [] Rating_a Letegory = [] Rating_s = [] Total_votes = [] for i in range(1,251): judul = soup.find_all("h3", attrs={'class':'ipc-title_text'})[i].text judul = judul.split('.')[i] Title.append(judul) n = 0 while n < 750: tahun = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'})[n].text</td></tr><tr><td>2]:
8]:
8]:
9]:
43</td><td>BUT the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'}) len(second_info) 750 third_info = soup.find_all("span", attrs={'class':'ipc-rating-star ipc-rating-star-base ipc-rating-star-imdb rating@roup-imdb-rating'}) len(third_info) 250 Now we know the length of our data, now here what we will do for the looping process. </td></tr><tr><td>2]:
8]:
8]:
9]:
43</td><td>BUT the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item')) third_info = soup.find_all("span", attrs=('class':'ipc-rating-star ipc-rating-starbase ipc-rating-starimdb rating@roupimdb-rating')) len(third_info) Now we know the length of our data, now here what we will do for the looping process. Title = [] vear = [] Duration = [] Rating_attegory = [] Rating_attegory = [] Rating_attegory = [] Title_append(judul) n = 0 white n < 750: tahun = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'))[n].text n = 1 durasi = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'))[n].text n = 1 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'))[n].text 7 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'))[n].text n = 1 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'))[n].text n = 1 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'))[n].text n = 1 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'))[n].text n = 1 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item'))[n].text n = 1 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item')][n].text n = 1 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item')][n].text n = 1 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item')][n].text n = 1 ratcat = soup.find_all("span", attrs=('class':'sc-479faa3c-8 bNrEFi cli-title-metadata-item')][n].text</td></tr><tr><td>2]:
8]:
9]:
43</td><td>But the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bNrFFi cli-title-metadata-item'}) len(second_info) 758 third_info = soup.find_all("span", attrs={'class':'ipc-rating-star ipc-rating-star-base ipc-rating-starimdb ratingGroupimdb-rating'}) len(third_info) 800 800 800 800 800 800 800 8</td></tr><tr><td>2]:
8]:
8]:
9]:
43</td><td>But the actual Title Data only exist in index list [1:251] second_info = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bWrEFI cli-title-metadata-item'}) inn(third_info) = soup.find_all("span", attrs={'class':'spc-rating-star ipc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating'}) inn(third_info) = soup.find_all("span", attrs={'class':'spc-rating-star ipc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating'}) inn(third_info) = soup.find_all("span", attrs={'class':'spc-rating-starbase ipc-rating-starimdb ratingGroupimdb-rating'}) inn = soup.find_all("span", attrs={'class':'spc-title_text'})]i].text inned = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bWrEFI cli-title-metadata-item'})]i].text inned = soup.find_all("span", attrs={'class':'sc-479faa3c-8 bWrEFI cli-title-metadata-item'})[i].text inned = soup.find_all("span", attrs={'class':'sc-479faa3</td></tr><tr><td>2]:
8]:
8]:
9]:
43</td><td>BUT the actual Title Data only exist in index list [1:251] second_sinfo = soup.find_ali("span", attrs=("class":'sc-479raadc-8 bbrEFs cli-title-metadata-item")) her(centum info) third_info = soup.find_ali("span", attrs=("class":'sc-479raadc-8 bbrEFs cli-title-metadata-item")) her(clitrd_info) Now we know the length of our data, now here what we will do for the looping process. Title = [] Title = [] Title = [] Title = [] Title = soup.find_ali("span", attrs=("class":'sc-479raadc-8 bbrEFs cli-title-metadata-item"))[n].text ### attrs = soup.find_ali("span", attrs=("class":'sc-479raadc-8 bbrEFs cli-title-metadata-item")][n].text ### attrs = so</td></tr><tr><td>2]:
8]:
9]:
43</td><td>But the actual Title Data only exist in index that [1.251] second_info = soop_find_all("span", attrac("class":"so-4/afact-8 directiclitie-metadata-item")) theref_info = soop_find_all("span", attrac("class":"so-4/afact-8 directiclitie-metadata-item")) theref_info = soop_find_all("span", attrac("class":"so-4/afact-8 directiclitie-metadata-item")) Now we know the length of our data. now here what we will do for the looping process. Tatle = [] **Vear = [] **Ratispa = [] **Total Voice = [] **Total</td></tr><tr><td>2]:
8]:
9]:
43
44</td><td>### BUT the actual Title Data only exist in index list [2:251] ###################################</td></tr><tr><td>2]:
8]:
8]:
9]:
43
44</td><td>But the actual Title Data only exist in index less [1:251] secund_info = soup.find_s1[(*span*, stirse(*class*:*sc-499fands-8 byteh: cli-tille-mesabas-ise*)) ter(cal_info = soup.find_s1[(*span*, attrse(*class*:*sc-499fands-8 byteh: cli-tille-mesabas-ise*)) ter(clinto = soup.find_s1[(*span*, attrse(*class*:*sc-499fands-8 byteh: cli-tille-mesabas-ise*)) ter(thre diren) The state = [] ter(thre diren) Now we know the length of our data, now hore what we will do for the looping process. **ii.la = [] **tate = [] **ta</td></tr><tr><td>2]:
8]:
9]:
43
44</td><td>Second from a chair Time Data only cost in index inst [1251] second from a soup-from all ("show", all rest("class") iso-reting-star disc return peter-base iso-reting-star-rando return) returns on control of the cont</td></tr><tr><td>2]:
8]:
9]:
43
44</td><td>Secure of the actual falls base only exist in index list [1203] ***Special Control of the product of the produ</td></tr><tr><td>2]:
8]:
8]:
9]:
43
44</td><td>secure from actual Tribe Data only cost in motes tas [1233] ***Secure from a cost find all ("spain", attract" class ("spain" attract state and a spain attract state and a state and a spain attract state and attrac</td></tr><tr><td>2]: 8]: 8]: 9]: 43 44 45</td><td>### ### ### #### #### ################</td></tr><tr><td>2]: 8]: 8]: 9]: 43 45 47 79</td><td>### The action Title Date only exist in what first (1.255) ### Control Contro</td></tr><tr><td>2]:
8]:
8]:
9]:
43
45</td><td>### ### ### ### ### ### ### ### ### ##</td></tr><tr><td>2]:
8]:
9]:
43
45</td><td>### Part Part</td></tr><tr><td>2]:
8]:
9]:
43
45</td><td>### Part of actual Time (Fig. 2017) gas 1. Interest Classes 1 incention and interest 2 incention</td></tr><tr><td>2]:
8]:
8]:
9]:
43
45</td><td>But the actual Title Calco only exist in index set [1.53] ***********************************</td></tr><tr><td>2]: 8]: 8]: 9]: 43 45 47</td><td>The content of the co</td></tr><tr><td>2]: 8]: 8]: 9]: 43 45 47</td><td>Experience Security of Security Security</td></tr><tr><td>2]: 8]: 8]: 9]: 43 45 47 47 47</td><td>The second file of the content of th</td></tr><tr><td>2]: 8]: 8]: 9]: 43 44 45 47 79 79 79</td><td>The second file and an action from the content of t</td></tr><tr><td>2]: 8]: 9]: 43 45 46 47 46</td><td>### Description of the policy flags of the part of</td></tr><tr><td>2]: 8]: 9]: 43 45 46 47 46</td><td>### ### ### ### ### ### ### ### ### ##</td></tr><tr><td>2]: 8]: 8]: 9]: 43 45 46 47 46</td><td>### ### ### ### ### ### ### ### ### ##</td></tr><tr><td>2]: 8]: 8]: 9]: 43 45 45 46 79 79</td><td>### ### ### ### ### ### ### ### ### ##</td></tr><tr><td>2]: 3]: 3]: 3]: 43 45 46 47 46 47 46</td><td>## 18 Part P</td></tr><tr><td>2]: 3]: 3]: 3]: 43 45 46 47 46 47 46</td><td> Company Comp</td></tr><tr><td>2]: 8]: 8]: 9]: 43 45 45 46 93 93 93</td><td> Company Comp</td></tr><tr><td>2]: 8]: 8]: 9]: 43 45 46 47 46</td><td> The content is a content of the co</td></tr><tr><td>2]: 8]: 9]: 93 93 93 93 93 93 93 93</td><td>## 18 Part P</td></tr></tbody></table></script></head></html></pre> |