

Web Scraping

Presented by Trixia Belleza

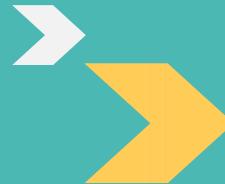




what is web scraping?

**Fetching data from
websites**

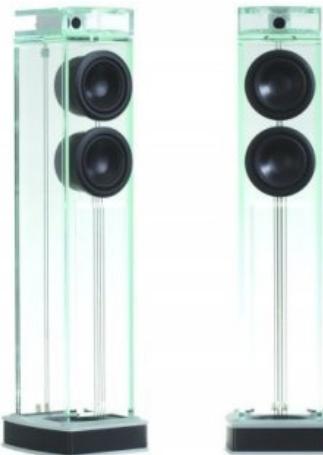
why web scrape?



Big Data
Analytics







★★★★★ FITS PERFECTLY ON ON MY SIDEWALK!

By [Amazon Customer](#) on January 7, 2017

This speaker fits perfectly on my sidewalk, because I sold my house to buy these.

★★★★★ Whoa Baby

By [Patrick J. Hawkins](#) on December 18, 2015

I bought these to put in our newborn's nursery. We've been playing classical music because it is suppose to stimulate brain functionality. It seems to work stupendously. The child is only 8 weeks old and has told me that I'm a freaking moron for buying them.

★★★★☆ Misleading product name...

By [Obi Wan](#) [TOP 100 REVIEWER](#) on February 16, 2016

I was quite despondent to see them shipped in simple wooden and cardboard boxes as if these were being sent to peasants. They should have been packed in hand stained mahogany shipping containers and wrapped tightly in baby panda fur for softness.

The uses of web scraping
for business and personal
requirements are endless.

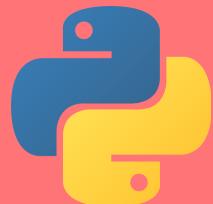
<https://www.webharvy.com/articles/web-scraper-use-cases.html>

How to web scrape?

In this seminar, we will use Python Scrapy

SETTING UP

Python Scrapy



PYTHON

A high-level
programming
language



PIP

PIP is a package
manager for
Python packages.



Scrapy

SCRAPY

Scrapy is a free
and open-source
web-scraping
framework written
in Python

★★★★★ My smart home must have!

By Owner on August 11, 2018

Color: Charcoal | Fabric: Polyester

Amazon Product Reviews

I am a new Echo user, but I love them. I started with a couple Dots because I thought I would just want to use the microphones to turn smart switches on and off. However, when I discovered I could listen to sounds of nature, I wanted the sound quality of the Echo. I now do my own sleep to the sound of waves or rain. That is awesome! I have also listened to music and the Echo sounds surprisingly good, however, I am rarely in

```
graph LR; A[Return a list of text reviews in a page.] --> B[Store the text reviews of a page in a list.]
```

Return a list of text reviews in a page.

Store the text reviews of a page in a list.



In the terminal, write:
scrapy shell "<url of the page>"



```
File Edit View Search Terminal Help
```

```
trix@trix-Satellite-C55-C:~$ scrapy shell "https://www.amazon.com/Certified-Refurbished-All-new-Generation-improved/product-reviews/B071P72R59/ref=cm_cr_dp_d_show_all_btm?ie=UTF8&reviewerType=all_reviews"
```



Find the html element and class that defines a specific text review

span.a-size-base.review-text | 638.5 x 149

Really good Echo plus I saved \$20 by getting a refurbished one. It didn't come in a branded box (just a plain brown box), but it came with everything I needed to get up and running. I got this mostly to play music in my basement. Compared to my original Echo, I think it sounds better. A little more bass. It might be the location though. I love the size. It easily fits on the shelf by my laundry. I do wish it still had the volume ring. It is the one thing I miss. I have this on a high shelf so I can't see the buttons on top.

I have some smart lights in my house. I can control them with the Echo. I also have the Echo Connect so I can answer the house phone with this Echo as well. All in all, I'm very happy with this purchase.

3 people found this helpful

[Helpful](#)

[Not Helpful](#)

[Comment](#) | [Report abuse](#)

★★★★★ so much better than Dot

By James I. Velazquez on August 13, 2018

Color: Silver Finish | [Verified Purchase](#)

The Echo is amazing. It picks up commands and responds very quickly compared to my Dot. I also like the built-in speaker much better. I know the difference from the Dot is only supposed to be the speaker, but the Echo actually seems to function much better than my Dot. The Dot seems to disconnect and have issues with the Amazon App. Especially with multi-room music. The Echo seems better all around

4 people found this helpful

[Helpful](#)

[Not Helpful](#)

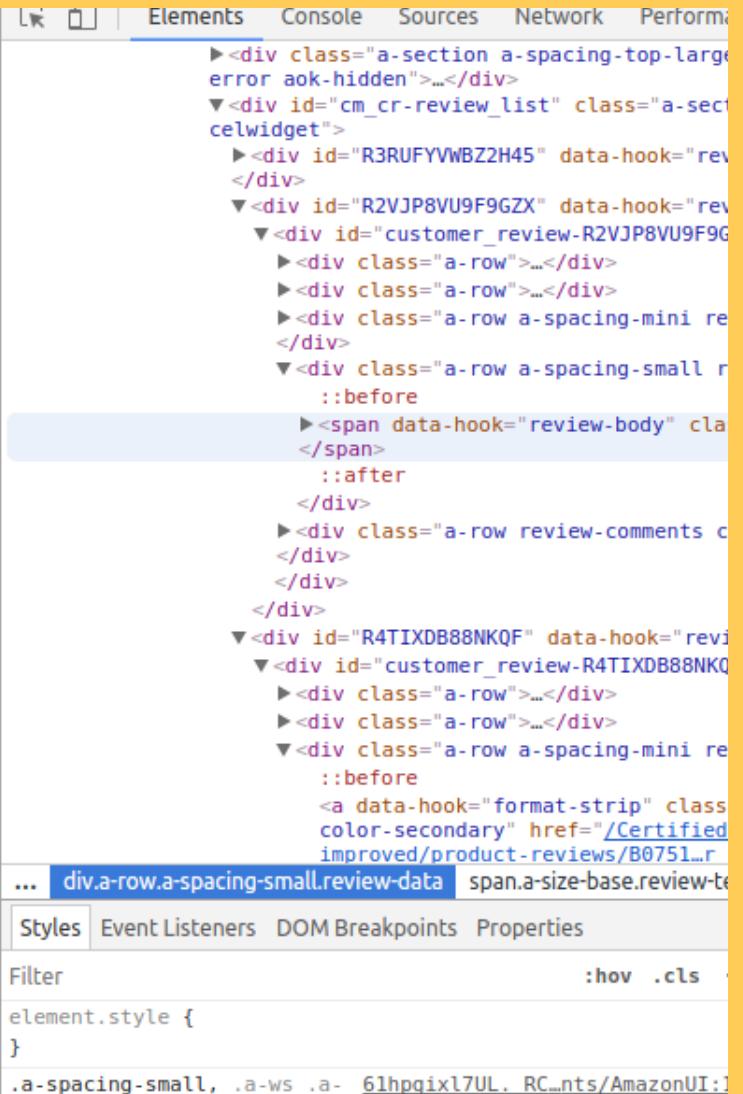
[Comment](#) | [Report abuse](#)

★★★★★ Great for the price!

By Wendy H. on August 17, 2018

Color: Heather Gray Fabric | [Verified Purchase](#)

I already have 2 echo dots and ended up buying speakers for them. I really wish I had skipped all of that and simply bought the 2nd generation echo to begin with. I would never know this was refurbished except for the



```
> <div class="a-section a-spacing-top-large error aok-hidden">...</div>
▼ <div id="cm_cr-review_list" class="a-section celwidget">
  ► <div id="R3RUFYVWBZ2H45" data-hook="review">...
  ▼ <div id="R2VJP8VU9F9GZX" data-hook="review">
    ▼ <div id="customer_review-R2VJP8VU9F9GZX" data-hook="review">
      ► <div class="a-row">...</div>
      ► <div class="a-row">...</div>
      ► <div class="a-row a-spacing-mini review-data">...
      □ <div class="a-row a-spacing-small review-data">...
        :before
        <span data-hook="review-body" class="a-size-base review-text" style="color: #333; font-weight: bold;">...</span>
        :after
      </div>
      ► <div class="a-row review-comments collapse">...
      </div>
    </div>
  </div>
  ▼ <div id="R4TIXDB88NKQF" data-hook="review">...
    ▼ <div id="customer_review-R4TIXDB88NKQF" data-hook="review">
      ► <div class="a-row">...</div>
      ► <div class="a-row">...</div>
      □ <div class="a-row a-spacing-mini review-data">...
        :before
        <a data-hook="format-strip" class="color-secondary" href="/Certified_improved/product-reviews/B0751...>...
      ...</div>
    </div>
  </div>
</div>
...<div class="a-row a-spacing-small review-data" style="background-color: #f0f0f0; border: 1px solid #ccc; padding: 5px; margin-bottom: 10px; border-radius: 5px; font-size: 0.9em; color: #333; font-weight: bold;">...</div>
span.a-size-base.review-text
```

Styles Event Listeners DOM Breakpoints Properties

Filter

element.style {
}

.a-spacing-small, .a-ws .a- 61hpgixl7UL_RComments/AmazonUI:1

Find the html element and class that defines a specific text review

span.a-size-base.review-text | 638.5 x 149

Really good Echo plus I saved \$20 by getting a refund (brown box), but it came with everything I needed to

Element span

Class a-size-base review-text



In the scrapy shell, type:
`response.xpath('//<element>
[contains(@class, "<name of class>")])'`



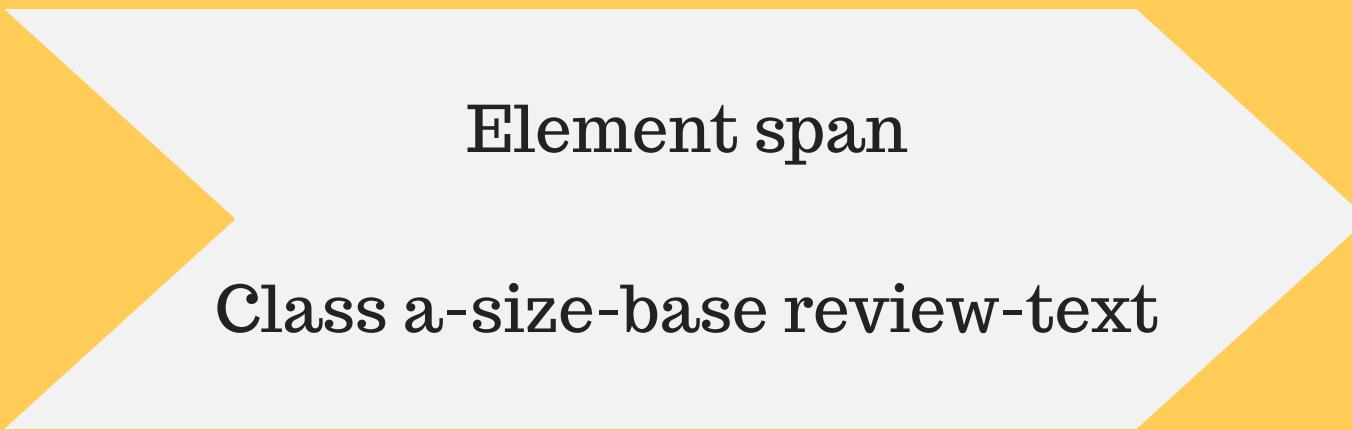
X
P
A
T
H

**XPATH IS A
LANGUAGE FOR
SELECTING NODES IN
XML OR HTML
DOCUMENTS**



```
span.a-size-base.review-text | 638.5 x 149
```

Really good Echo plus I saved \$20 by getting a refurb (in brown box), but it came with everything I needed to



Element span

Class a-size-base review-text





```
>>> response.xpath('//span[contains(@class, "a-size-base review-text")]')
```

**Select nodes from all span elements
that contains the class "a-size-base
review-text"**



Returns a list of selector objects



```
>>>  
>>> response.xpath('//span[contains(@class,"a-size-base review-text")]')  
[<Selector xpath='//span[contains(@class,"a-size-base review-text")]' data=u'<sp  
an data-hook="review-body" class="a-s'>, <Selector xpath='//span[contains(@class  
, "a-size-base review-text")]' data=u'<span data-hook="review-body" class="a-s'>,  
<Selector xpath='//span[contains(@class,"a-size-base review-text")]' data=u'<sp  
an data-hook="review-body" class="a-s'>, <Selector xpath='//span[contains(@class  
, "a-size-base review-text")]' data=u'<span data-hook="review-body" class="a-s'>,  
<Selector xpath='//span[contains(@class,"a-size-base review-text")]' data=u'<sp  
an data-hook="review-body" class="a-s'>, <Selector xpath='//span[contains(@class  
, "a-size-base review-text")]' data=u'<span data-hook="review-body" class="a-s'>,  
<Selector xpath='//span[contains(@class,"a-size-base review-text")]' data=u'<sp  
an data-hook="review-body" class="a-s'>, <Selector xpath='//span[contains(@class  
, "a-size-base review-text")]' data=u'<span data-hook="review-body" class="a-s'>,  
<Selector xpath='//span[contains(@class,"a-size-base review-text")]'
```





Getting the text data from these selectors

In the scrapy shell, type:

```
response.xpath('//<element>  
[contains(@class, "<name of  
class>")]/text()')
```





Getting the text data from these selectors

```
>>>  
>>> response.xpath('//span[contains(@class, "a-size-base")]/text()')  
<Selector xpath='//span[contains(@class, "a-size-base")]/text()' data=u'  
[<Selector xpath='//span[contains(@class, "a-size-base review-text")]/text()' data=u'  
"a-size-base review-text")]/text()' data=u'This arrived  
view-text")]/text()' data=u'It is difficult to get the  
text()' data=u'I thought I was buying an original Vizio'  
'Responds better than the original remote', <Selector x  
ht out of the box, so I hav', <Selector xpath='//span[c  
my 3rd remote', <Selector xpath='//span[contains(@class  
>, <Selector xpath='//span[contains(@class, "a-size-base  
>>>
```



Getting the text data from these selectors

Notice that the list actually
contains data reviews





Finally, extracting the data

In the scrapy shell, type:

```
response.xpath('//<element>  
[contains(@class, "<name of  
class>")]/text()').extract()
```



Finally, extracting the data



Assign it to a list

In the scrapy shell, type:

```
reviews = response.xpath('//<element>  
[contains(@class, "<name of  
class>")]/text()').extract()
```

Then access each review by index:

```
reviews[i]
```





reviews[1]

```
>>> reviews[1]
u"Yay! It works with my TV. Hate, Hate, Hate the remote that came with it.
nd I've had to replace the batteries at least once a month. I've really avo
estick on this TV and have to use it to switch the input. So, thank goodness
```

```
>>>
>>>
>>>
>>
```



”
You can have data without information, but you cannot have information without data.

“

Daniel Keys Moran Quotes. (n.d.). BrainyQuote.com. Retrieved September 14, 2018, from BrainyQuote.com Web site:
https://www.brainyquote.com/quotes/daniel_keys_moran_230935



Presented by:
Trixia Belleza
CMSC 199 - B2