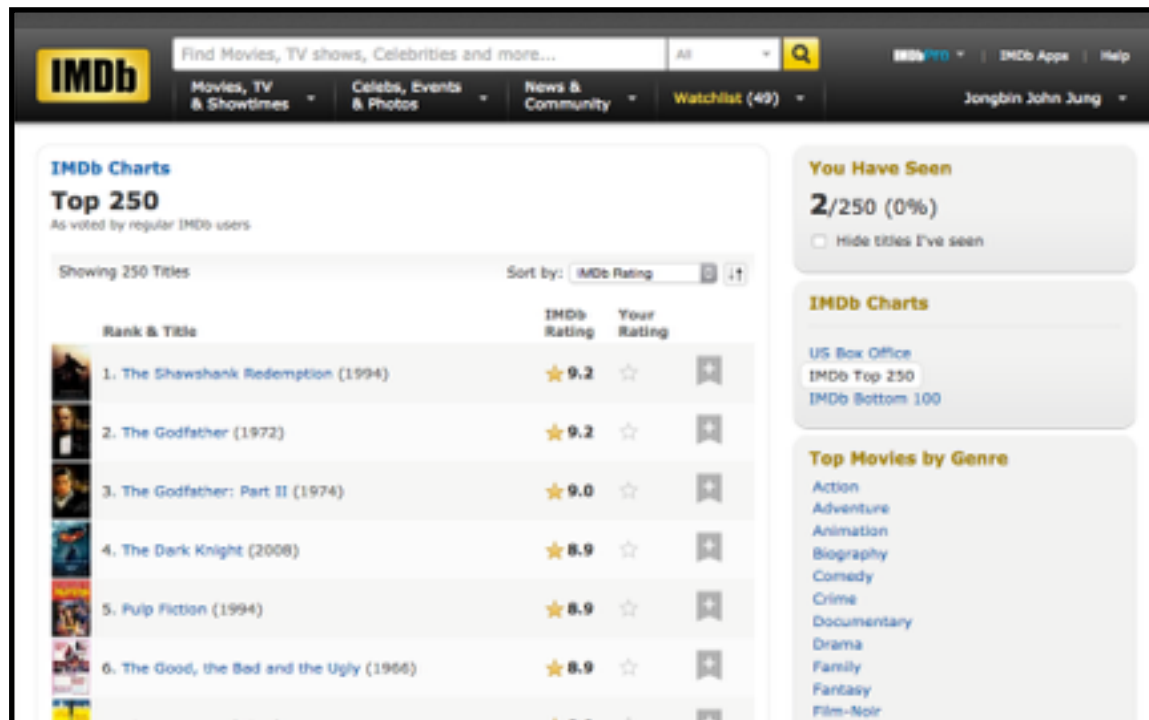


# Web Scraping

a **Data Science Drop-in** Tutorial  
by **Jongbin Jung**  
([jongbin@stanford.edu](mailto:jongbin@stanford.edu))

# Web Scrapping?



	movie_title	actor_name	character
1	The Shawshank Redemption	Tim Robbins	Andy Dufresne
2	The Shawshank Redemption	Morgan Freeman	Ellis Boyd 'Red' Redding
3	The Shawshank Redemption	Bob Gunton	Warden Norton
4	The Shawshank Redemption	William Sadler	Heywood
5	The Shawshank Redemption	Clancy Brown	Captain Hadley
6	The Shawshank Redemption	Gil Bellows	Tommy
7	The Shawshank Redemption	Mark Rolston	Bugs Diamond
8	The Shawshank Redemption	James Whitmore	Brooks Hatlen
9	The Shawshank Redemption	Jeffrey DeMunn	1946 D.A.
10	The Shawshank Redemption	Larry Brandenburg	Skeet
11	The Shawshank Redemption	Neil Giuntoli	Jigger
12	The Shawshank Redemption	Brian Libby	Floyd
13	The Shawshank Redemption	David Proval	Snooze
14	The Shawshank Redemption	Joseph Ragno	Ernie
15	The Shawshank Redemption	Jude Ciccolella	Guard Mert
16	The Godfather	Marlon Brando	Don Vito Corleone
17	The Godfather	Al Pacino	Michael Corleone
18	The Godfather	James Caan	Sonny Corleone
19	The Godfather	Richard S. Castellano	Clemenza (as Richard Castellano)
20	The Godfather	Robert Duvall	Tom Hagen
21	The Godfather	Sterling Hayden	Capt. McCluskey
22	The Godfather	John Marley	Jack Woltz
23	The Godfather	Richard Conte	Barzini

# Get from here...

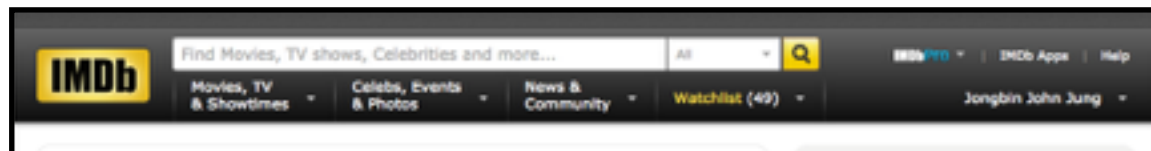
... to **here**

```

<div class="pending"></div>
<div class="unreadable">GET GET RELEASED</div>
<div class="unread"></div>
<div class="rating"></div>
<div class="seen">Seen</div>
</div>
</div>
</div>
</div>
<div class="watchlistColumn"> <div class="wib_ribbon" data-hocount="11044138"></div>
</div>
</div>
<div class="even">
<td class="posterColumn"><a href="/title/tt1049413/?ref=uhdp_t4_114">
<img alt="https://ia.media-imdb.com/images/C6/9Y/3867/1900264197/Bel13aushXFP1Z0wvC0H6dW0E8_VI_8X14_C00_0_38_50_81.jpg" width="360" height="50"/>
</td>
</div>
<div class="titleColumn">
<span name="tr" data-value="8_248">214</span>
<a href="/title/tt1049413/?ref=uhdp_t4_114">
<span name="td" data-value="2009-01-19" class="secondaryInfo">[2009]</span>
</td>
<td class="ratingColumn indRating">
<strong name="sr" data-value="50640" title="8.2 based on 506,460 votes">8.2</strong>
</td>
<td class="ratingColumn"><span name="sr" data-value="0"> <div class="seen-widget seen-widget-tt1049413 pending" data-tikid="tt1049413">
<div class="boundary">
<div class="popover">
<span class="delete"><span></span></span></div>
</div>
<div class="inLine">
<div class="pending"></div>
<div class="unreadable">GET GET RELEASED</div>
<div class="unread"></div>
<div class="rating"></div>
<div class="seen">Seen</div>
</div>
</div>
</div>
</div>
</div>
<div class="watchlistColumn"> <div class="wib_ribbon" data-hocount="11044138"></div>
</div>
</div>
<div class="odd">
<td class="posterColumn"><a href="/title/tt1187043/?ref=uhdp_t4_115">
<img alt="https://ia.media-imdb.com/images/C6/9Y/3867/1900264197/Bel13aushXFP1Z0wvC0H6dW0E8_VI_8Y58_C00_0_38_50_81.jpg" width="360" height="50"/>
</td>

```

# Web Scrapping?



	movie_title	actor_name	character
1	The Shawshank Redemption	Tim Robbins	Andy Dufresne
2	The Shawshank Redemption	Morgan Freeman	Ellis Boyd 'Red' Redding
3	The Shawshank Redemption	Bob Gunton	Warden Norton
4	The Shawshank Redemption	William Sadler	Heywood

# Today's **GOAL**

Collect the **cast overview** (actor and character played) for each of the **Top 10 movies** of **IMDb Charts' Top 250**

([http://www.imdb.com/chart/top?ref\\_nv\\_ch\\_250\\_4](http://www.imdb.com/chart/top?ref_nv_ch_250_4))

[illegible]

# First Things First

- Get on the Stanford corn servers  
(see, <https://farmshare.stanford.edu/>)

- Mac/Linux: Open a terminal and type:

```
> ssh [SUNet ID]@corn.stanford.edu
```

[Terminal]

replacing [SUNet ID] with your ID

- Windows: Use PowerShell / PuTTY / Cygwin ...

# First Things First

- Let's get some python packages!

```
> pip install --user requests beautifulsoup4 selenium
```

[Terminal]

- Clone (i.e., download) example scripts

```
> git clone https://github.com/5harad/datascience.git
```

```
> cd datascience/webscraping
```

[Terminal]

- Start python

```
> python
```

[Terminal]

# Meet BeautifulSoup

<http://www.crummy.com/software/BeautifulSoup/>

[illegible]

from source (html)  
to python

```

✓ jongbinjung@KNSba22b7eb:(master)$ python
Python 2.7.8 [Anaconda 2.0.1 (x86_64)] (default, Aug 21 2014, 15:21:48)
[GCC 4.2.1 (Apple Inc. build 5577)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
Anaconda is brought to you by Continuum Analytics.
Please check out: http://continuum.io/thanks and https://binstar.org
>>>

```

BeautifulSoup  
goes **here**

	movie_title	actor_name	character
1	The Shawshank Redemption	Tim Robbins	Andy Dufresne
2	The Shawshank Redemption	Morgan Freeman	Ellis Boyd 'Red' Redding
3	The Shawshank Redemption	Bob Gunton	Warden Norton
4	The Shawshank Redemption	William Sadler	Heywood
5	The Shawshank Redemption	Clancy Brown	Captain Hadley
6	The Shawshank Redemption	Gil Bellows	Tommy
7	The Shawshank Redemption	Mark Rolston	Bugs Diamond
8	The Shawshank Redemption	James Whitmore	Brooks Hatlen
9	The Shawshank Redemption	Jeffrey DeMunn	1946 D.A.
10	The Shawshank Redemption	Larry Brandenburg	Skeet
11	The Shawshank Redemption	Neil Giuntoli	Jigger
12	The Shawshank Redemption	Brian Libby	Floyd
13	The Shawshank Redemption	David Proval	Snooze
14	The Shawshank Redemption	Joseph Ragno	Ernie
15	The Shawshank Redemption	Jude Ciccolella	Guard Mert
16	The Godfather	Marlon Brando	Don Vito Corleone
17	The Godfather	Al Pacino	Michael Corleone
18	The Godfather	James Caan	Sonny Corleone
19	The Godfather	Richard S. Castellano	Cleenna (as Richard Castellano)
20	The Godfather	Robert Duvall	Tom Hagen
21	The Godfather	Sterling Hayden	Capt. McCluskey
22	The Godfather	John Marley	Jack Woltz
23	The Godfather	Richard Conte	Barzini

This can be done many ways.  
But we use **requests** (for now)

[illegible]

from source (html)  
to python

```

✓ jongbinjung@MBA272eb:~$ python
Python 2.7.8 [Anaconda 2.0.1 (x86_64)] (default, Aug 21 2014, 15:21:46)
[GCC 4.2.1 (Apple Inc. build 5577)] on darwin
Type "help", "copyright", "credits" or "license()" for more information.
Anaconda is brought to you by Continuum Analytics.
Please check out: http://continuum.io/thanks and https://binstar.org
>>>

```

```
import requests
```

```
web_page =  
requests.get('address')  
web_page.content
```

```
[python]
```

Let's try this with the  
BeautifulSoup web page

<http://www.crummy.com/software/BeautifulSoup/>

# Meet BeautifulSoup

Make the source (html) into a BeautifulSoup

```
from bs4 import BeautifulSoup

soup = BeautifulSoup(web_page.content)
```

[python]

```
✓ jongbinjung@Mba22b2eb:(master)$ python
Python 2.7.8 [Anaconda 2.0.1 (x86_64)] (default, Aug 21 2014, 15:21:46)
[GCC 4.2.1 (Apple Inc. build 5577)] on darwin
Type "help", "copyright", "credits" or "license()" for more information.
Anaconda is brought to you by Continuum Analytics.
Please check out: http://continuum.io/thanks and https://binstar.org
>>>
```

	movie_title	actor_name	character
1	The Shawshank Redemption	Tim Robbins	Andy Dufresne
2	The Shawshank Redemption	Morgan Freeman	Ellis Boyd 'Red' Redding
3	The Shawshank Redemption	Bob Gunton	Warden Norton
4	The Shawshank Redemption	William Sadler	Heywood
5	The Shawshank Redemption	Clancy Brown	Captain Hadley
6	The Shawshank Redemption	Gill Bellows	Tommy
7	The Shawshank Redemption	Mark Rolston	Bugs Diamond
8	The Shawshank Redemption	James Whitmore	Brooks Hatlen
9	The Shawshank Redemption	Jeffrey DeMunn	1946 D.A.
10	The Shawshank Redemption	Larry Brandenburg	Skeet
11	The Shawshank Redemption	Neil Giuntoli	Jigger
12	The Shawshank Redemption	Brian Libby	Floyd
13	The Shawshank Redemption	David Proval	Snooze
14	The Shawshank Redemption	Joseph Ragno	Ernie
15	The Shawshank Redemption	Jude Ciccolella	Guard Mert
16	The Godfather	Marlon Brando	Don Vito Corleone
17	The Godfather	Al Pacino	Michael Corleone
18	The Godfather	James Caan	Sonny Corleone
19	The Godfather	Richard S. Castellano	Clemenza (as Richard Castellano)
20	The Godfather	Robert Duvall	Tom Hagen
21	The Godfather	Sterling Hayden	Capt. McCluskey
22	The Godfather	John Marley	Jack Woltz
23	The Godfather	Richard Conte	Barzini



# Use BeautifulSoup

```
soup.h1
```

[python]

```
<h1>Beautiful Soup</h1>
```

[output]

You didn't write that awful page. You're just trying to get some data out of it. BeautifulSoup is here to help. Since 2004, it's been saving programmers hours or days of work on quick-turnaround screen scraping projects.

## Beautiful Soup

"A tremendous boost to Python411 Podcast"

[ [Download](#) | [Documentation](#) | [Hall of Fame](#) | [Source](#) | [Discussion group](#) ]

If BeautifulSoup has saved you a lot of time and money, the best way to pay me back is to check out [Constellation Gamer](#), my sci-fi novel about alien video games.

You can [read the first two chapters for free](#), and the full novel starts at 5 USD. Thanks!

If you have questions, send them to [the discussion group](#). If you find a bug, [file it](#).

Beautiful Soup is a Python library designed for quick turnaround projects like screen-scraping. Three features make it powerful:

1. BeautifulSoup provides a few simple methods and Pythonic idioms for navigating, searching, and modifying a parse tree: a toolkit for dissecting a document and extracting what you need. It doesn't take much code to write an application
2. BeautifulSoup automatically converts incoming documents to Unicode and outgoing documents to UTF-8. You don't have to think about encodings, unless the document doesn't specify an encoding and BeautifulSoup can't detect one. Then you just have to specify the original encoding.
3. BeautifulSoup sits on top of popular Python parsers like [lxml](#) and [html5lib](#), allowing you to try out different parsing strategies or trade speed for flexibility.

Beautiful Soup parses anything you give it, and does the tree traversal stuff for you. You can tell it "Find all the links", or "Find all the links of class `externalLink`", or "Find all the links whose urls match `'foo.com'`", or "Find the table heading that's got bold text, then give me that text."

Valuable data that was once locked up in poorly-designed websites is now within your reach. Projects that would have taken hours take only minutes with BeautifulSoup.

Interested? [Read more](#).



Did we just get this?

# What's Happening?

```
21 <div align="center">
22
23 <a href="bs4/download/"><h1>Beautiful Soup</h1></a>
24
25 <p>"A tremendous boon." -- <a
26 href="http://www.awaretek.com/python/index.html">Pyth
27 Podcast</a></p>
28
29 <p>[ <a href="#Download">Download</a> | <a
30 href="bs4/doc/">Documentation</a> | <a href="#Hallof
href="https://code.launchpad.net/beautifulsoup">Sour
href="https://groups.google.com/forum/?fromgroups#!f
group</a> ]</p>
31
32 <small>If Beautiful Soup has saved you a lot of time
33 best way to pay me back is to check out <a
34 href="http://www.candlemarkandgleam.com/shop/constel
<i>Constellation
35 Games</i>, my sci-fi novel about alien video games</
36 <a
```

- `soup.h1` gives us the content surrounded by `<h1>` and `</h1>`

# What's Happening?

```
21 <div align="center">
22
23 <a href="bs4/download/"><h1>Beautiful Soup</h1></a>
24
25 <p>"A tremendous boon." -- <a
26 href="http://www.awaretek.com/python/index.html">Pyth
27 Podcast</a></p>
28
29 <p>[ <a href="#Download">Download</a> | <a
30 href="bs4/doc/">Documentation</a> | <a href="#Hallof
href="https://code.launchpad.net/beautifulsoup">Sour
href="https://groups.google.com/forum/?fromgroups#!f
group</a> ]</p>
31
32 <small>If Beautiful Soup has saved you a lot of time
33 best way to pay me back is to check out <a
34 href="http://www.candlemarkandgleam.com/shop/constel
<i>Constellation
35 Games</i>, my sci-fi novel about alien video games</
36 <a
```

- The `<a>` and `</a>` make the text between them into a *link*
- The `href="bs4/download/"` indicates where the link should link to

# Get a link and its address

```
soup.a
```

[python]

```
<a href="bs4/download/"><h1>Beautiful Soup</h1></a>
```

[output]

Now let's get the address (the value assigned to href)

```
soup.a.get('href')
```

[python]

```
'bs4/download/ '
```

[output]

You didn't write that awful page. You're just trying to get some data out of it. Beautiful Soup is here to help. Since 2004, it's been saving programmers hours or days of work on quick-turnaround screen scraping projects.

## Beautiful Soup

"A tremendous boon." -- [Python411 Podcast](#)

[ [Download](#) | [Documentation](#) | [Hall of Fame](#) | [Source](#) | [Discussion group](#) ]

If Beautiful Soup has saved you a lot of time and money, the best way to pay me back is to check out [Constellation Games](#), my sci-fi novel about alien video games.

You can [read the first two chapters for free](#), and the full novel starts at 5 USD. Thanks!

*If you have questions, send them to [the discussion group](#). If you find a bug, [file it](#).*

Beautiful Soup is a Python library designed for quick turnaround projects like screen-scraping. Three features make it powerful:

1. Beautiful Soup provides a few simple methods and Pythonic idioms for navigating, searching, and modifying a parse tree: a toolkit for dissecting a document and extracting what you need. It doesn't take much code to write an application
2. Beautiful Soup automatically converts incoming documents to Unicode and outgoing documents to UTF-8. You don't have to think about encodings, unless the document doesn't specify an encoding and Beautiful Soup can't detect one. Then you just have to specify the original encoding.
3. Beautiful Soup sits on top of popular Python parsers like [lxml](#) and [html5lib](#), allowing you to try out different parsing strategies or trade speed for flexibility.

Beautiful Soup parses anything you give it, and does the tree traversal stuff for you. You can tell it "Find all the links", or "Find all the links of class `externalLink`", or "Find all the links whose urls match `'foo.com'`", or "Find the table heading that's got bold text, then give me that text."

Valuable data that was once locked up in poorly-designed websites is now within your reach. Projects that would have taken hours take only minutes with Beautiful Soup.

Interested? [Read more.](#)



What about all these  
other links?



# Get more links and addresses

Lets find\_all the links (a)

```
soup.find_all('a')
```

[python]

```
[<a href="bs4/download/"><h1>Beautiful Soup</h1></a>, <a href="http://  
www.awaretek.com/python/index.html">Python411 Podcast</a>, ...  
...  
... , <a href="http://www.crummy.com/">http://www.crummy.com/</a>, <a  
href="http://www.crummy.com/software/">software/</a>, <a href="http://  
www.crummy.com/software/BeautifulSoup/">BeautifulSoup/</a>]
```

[output]

- Notice [ ..., ..., ... ] is a list in python
- We can `iterate` through a `list` with a `for` loop

# Get more links and addresses

Lets get all the addresses with a for loop

```
for link in soup.find_all('a'):
    link.get('href')
```

[python]

```
'bs4/download/'
'http://www.awaretek.com/python/index.html'
...
'http://www.crummy.com/software/'
'http://www.crummy.com/software/BeautifulSoup/'
```

[output]

Goto IMDb.com

## Today's **GOAL**

Collect the **cast overview** (actor and character played) for each  
of the **Top 10 movies** of **IMDb Charts' Top 250**  
([http://www.imdb.com/chart/top?ref\\_=nv\\_ch\\_250\\_4](http://www.imdb.com/chart/top?ref_=nv_ch_250_4))

Let's start with

“the **cast overview** (actor and character played)”

for just one movie.



We want this

...start by making a soup



The screenshot shows the IMDb 'Cast' page for the movie 'The Shawshank Redemption'. It features a list of actors and their roles, with a small portrait of each actor to the left of their name. The roles are listed to the right of the actor's name, separated by an ellipsis. The page is titled 'Cast' and includes a sub-header 'Cast overview, first billed only:'. An 'Edit' link is visible in the top right corner.

Cast		
	Tim Robbins	... Andy Dufresne
	Morgan Freeman	... Ellis Boyd 'Red' Redding
	Bob Gunton	... Warden Norton
	William Sadler	... Heywood
	Clancy Brown	... Captain Hadley
	Gil Bellows	... Tommy
	Mark Rolston	... Bogs Diamond
	James Whitmore	... Brooks Hatlen
	Jeffrey DeMunn	... 1946 D.A.
	Larry Brandenburg	... Skeet
	Neil Giuntoli	... Jigger
	Brian Libby	... Floyd
	David Proval	... Snooze
	Joseph Ragno	... Ernie
	Jude Ciccolella	... Guard Mert

```
from bs4 import BeautifulSoup
import requests

web_page = requests.get('http://...')
soup = BeautifulSoup(web_page.content)
```

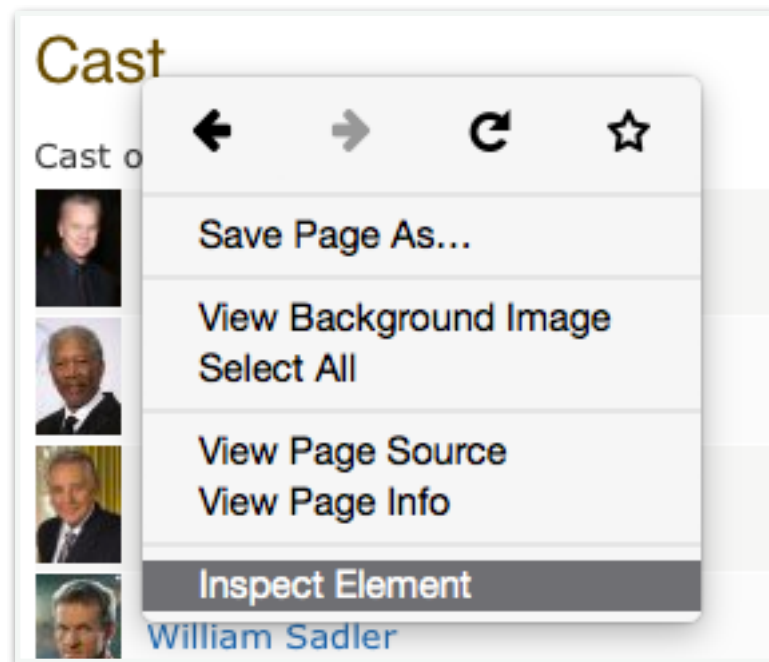
[python]

([http://www.imdb.com/title/tt0111161/?ref\\_=chttp\\_tt\\_1](http://www.imdb.com/title/tt0111161/?ref_=chttp_tt_1))

# Find Stuff in a Soup of html

Using your browser's Developer Mode

(A demo is better than a thousand slides)



But in case you forget, it's  
[right click]  
> [Inspect Element]  
... in most modern browsers

# finding a Specific tag

<table class='cast\_list'>

```
soup.find('table', class_='cast_list')
```

[python]

Note that we use `class_` instead of `class`.

This is because `class` means something else in python.

Anything other than `class`, you should use as is.

# Straining the Soup

When dealing with only a small portion of the entire page (like ourselves), it might speed things up a little to **strain** the soup with `SoupStrainer`, before **finding** things in it.

```
from bs4 import BeautifulSoup, SoupStrainer
import requests

cast_strainer = SoupStrainer('table', class_='cast_list')

web_page = requests.get('http://...')
soup = BeautifulSoup(web_page.content, parse_only=cast_strainer)
```

[python]

# What's in a table?



```
<table>
```

```
  <tr>
```

```
    <td></td>
```

```
    <td></td>
```

```
  </tr>
```

```
  <tr>
```

```
    <td></td>
```

```
  </tr>
```

```
</table>
```

within a `table` tag,

`<tr>` defines rows

while `<td></td>` defines columns within rows.

html tables are written row-by-row

# Iterate a `table` by its rows

We want to `iterate` each row of our strained soup  
(`<table class="cast_list">`)

```
for row in soup.find_all('tr'):
    # do useful things with each row
```

[python]

# Picking out the Cherries

How should we identify the **actor's name** and **character** played, given a single row (<tr>)?

```
▶ <tr>...</tr>
▼ <tr class="odd">
  ▶ <td class="primary_photo">...</td>
  ▼ <td class="itemprop" itemprop="actor" itemscope itemtype="http://
    schema.org/Person">
    ▼ <a href="/name/nm0000209/?ref_=tt_cl_t1" itemprop="url">
      <span class="itemprop" itemprop="name">Tim Robbins</span>
    </a>
    </td>
    <td class="ellipsis">
      ...
    </td>
  ▼ <td class="character">
    ▼ <div>
      <a href="/character/ch0001388/?ref_=tt_cl_t1">Andy Dufresne</a>
    </div>
    </td>
  </tr>
```

An actor's name seems to be uniquely identified by the property `itemprop="name"`

```
▶ <tr>...</tr>
▼ <tr class="odd">
  ▶ <td class="primary_photo">...</td>
  ▼ <td class="itemprop" itemprop="actor" itemsc
    schema.org/Person">
      ▼ <a href="/name/nm0000209/?ref_tt_cl_t1" itemprop="url">
        <span class="itemprop" itemprop="name">Tim Robbins</span>
      </a>
    </td>
    <td class="ellipsis">
      ...
    </td>
  ▼ <td class="character">
    ▼ <div>
      <a href="/character/ch0001388/?ref_tt_cl_t1">Andy Dufresne</a>
    </div>
  </td>
```

The column containing the character name has `class="character"`

*(There's usually more than one way ... )*



An actor's name seems to be uniquely identified by the property `itemprop="name"`

```
for row in soup.find_all('tr'):
    actor = row.find(itemprop='name').text
    role = row.find(class_='character').text
```

[python]

The column containing the character name has `class="character"`

`find_all('tr')` gives an extra row, which doesn't have anything that matches `itemprop='name'`

```
for row in soup.find_all('tr'):
    actor = row.find(itemprop='name').text
    role = row.find(class_='character').text
```

[python]

we want python to ignore these errors

```
...
AttributeError: 'NoneType' object has no
attribute 'text'
```

[output]



python will try this

```
for row in soup.find_all('tr'):
    try:
        actor = row.find(itemprop='name').text
        role = row.find(class_='character').text

    except AttributeError:
        pass
```

[python]

### WARNING!

Practice caution with `try-except`.  
Don't pass an error, unless you're  
certain you it's an error you *want to*  
ignore.

if an `AttributeError`  
happens, it will do this  
(in this case, ignore the error and pass)

```
for row in soup.find_all('tr'):
    try:
        actor = row.find(itemprop='name').text
        role = row.find(class_='character').text

    except AttributeError:
        pass
```

[python]

```
u'Tim Robbins'
u'\n\nAndy Dufresne\n\n'
...
```

[output]

What's with all the  
\n\n\n\n\n\n ... ?

Web designers sometimes use hidden white spaces for layout purposes.

A good way to deal with these in python is to surround your string with ' '.join(string.split())

```
for row in soup.find_all('tr'):
    try:
        actor = ' '.join(row.find(itemprop='name').text.split())
        role = ' '.join(row.find(class_='character').text.split())

    except AttributeError:
        pass
```

[python]

# Write the Data to a File

Usually, `print` is sufficient in `python`.

But `print` doesn't play well with weird characters ...

But the web is full of weird characters!

*One* workaround is to use `codecs`

```
import codecs
...
with codecs.open('file_name.tsv', 'a', encoding='utf-8') as fid:
    print>> fid, '\t'.join([actor, role])
```

[python]

PRO: Works with most languages on the web (Chinese, Korean, ...)

CON: Trickier to pipe/redirect output using command line

# So far ...

```
from bs4 import BeautifulSoup, SoupStrainer
import requests
import codecs

cast_strainer = SoupStrainer('table', class_='cast_list')

web_page = requests.get('http://...')
soup = BeautifulSoup(web_page.content, parse_only=cast_strainer)

for row in soup.find_all('tr'):
    try:
        actor = ' '.join(row.find(itemprop='name').text.split())
        role = ' '.join(row.find(class_='character').text.split())

        with codecs.open('file_name.tsv', 'a', encoding='utf-8') as fid:
            print>> fid, '\t'.join([actor, role])

    except AttributeError:
        pass
```

[python]

# Finally!

## Today's **GOAL**

Collect the **cast overview** (actor and character played) for each  
of the **Top 10 movies** of **IMDb Charts' Top 250**  
([http://www.imdb.com/chart/top?ref\\_=nv\\_ch\\_250\\_4](http://www.imdb.com/chart/top?ref_=nv_ch_250_4))



# Workflow



A screenshot of the IMDb top 250 movies list. It shows a vertical list of 10 movies, each with a small thumbnail image on the left and the movie title followed by its release year in parentheses on the right. The movies are: 1. The Shawshank Redemption (1994), 2. The Godfather (1972), 3. The Godfather: Part II (1974), 4. The Dark Knight (2008), 5. Pulp Fiction (1994), 6. The Good, the Bad and the Ugly (1966), 7. 12 Angry Men (1957), 8. Schindler's List (1993), 9. The Lord of the Rings: The Return of the King (2003), and 10. Fight Club (1999).

1. The Shawshank Redemption (1994)
2. The Godfather (1972)
3. The Godfather: Part II (1974)
4. The Dark Knight (2008)
5. Pulp Fiction (1994)
6. The Good, the Bad and the Ugly (1966)
7. 12 Angry Men (1957)
8. Schindler's List (1993)
9. The Lord of the Rings: The Return of the King (2003)
10. Fight Club (1999)

[http://www.imdb.com/chart/top?ref\\_=nv\\_ch\\_250\\_4](http://www.imdb.com/chart/top?ref_=nv_ch_250_4)

for these links do ***this***

```
from bs4 import BeautifulSoup, SoupStrainer
import requests
import codecs

cast_strainer = SoupStrainer('table', class_='cast_list')

web_page = requests.get('http://...')
soup = BeautifulSoup(web_page.content, parse_only=cast_strainer)

for row in soup.find_all('tr'):
    try:
        actor = ' '.join(row.find(itemprop='name').text.split())
        role = ' '.join(row.find(class_='character').text.split())

        with codecs.open('file_name.tsv', 'a', encoding='utf-8') as fid:
            print>> fid, '\t'.join([actor, role])

    except AttributeError:
        pass
```

[python]

# The Links

What do we need, and how should we get it?

remember:

```
soup.a.get('href')
```

[python]

The screenshot shows a web browser's developer tools interface. At the top, a yellow box highlights the CSS selector `tbody.lister-list` with dimensions 620px x 12955px. Below this, a list of three movies is displayed:

1. The Shawshank Redemption (1994)
2. The Godfather (1972)
3. The Godfather: Part II (1974)

Below the movie list, the 'Elements' tab is active, showing the HTML structure. The following code is visible:

```
<input id="seen-config" type="hidden" data-caller="http">
<div class="lister">
  <div class="header">...</div>
  <table class="chart" data-caller-name="other-chart">
    <colgroup>...</colgroup>
    <thead>...</thead>
    <tbody class="lister-list">
      <tr class="odd">
        <td class="posterColumn">...</td>
        <td class="titleColumn">
          <span name="ir" data-value="9.21">1.</span>
          <a href="/title/tt0111161/?ref_=http_tt_1" tit
            Freeman">The Shawshank Redemption</a>
          <span name="rd" data-value="1994-10-14" class="
        </td>
        <td class="ratingColumn imdbRating">...</td>
        <td class="ratingColumn">...</td>
        <td class="watchlistColumn">...</td>
      </tr>
    </tbody>
  </table>
</div>
```

# Workflow

requests ► BeautifulSoup ► .get('href') ► requests ...

```
web_page = requests.get('http://...')

list_strainer = SoupStrainer('tbody', class_='lister-list')
soup = BeautifulSoup(web_page.content, parse_only=list_strainer)

movie_list = soup.find_all('td', class_='titleColumn')

for movie in movie_list:
    movie_title = movie.a.text

    link = 'http://imdb.com' + movie.a.get('href')
    soup = requests.get(link)
    ...
```

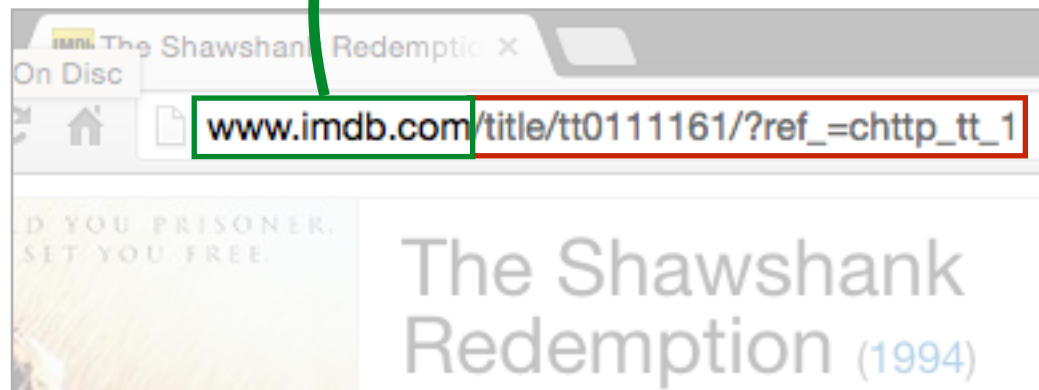
[python]

# Workflow

**Build** a link from the href value

```
link = 'http://imdb.com' + movie.a.get('href')
```

[python]



```
<a href="/title/tt0111161/?ref_=chttp_tt_1" title="Frank Darabont (dir.), Tim Robbins, Morgan Freeman">The Shawshank Redemption</a>
```

**“Top 10 movies of IMDb Charts' Top 250”**

```
soup.find_all('td', class_='titleColumn', limit=10)
```

[python]

The `limit` option will return only the first 10 results of `find_all`

# Be Nice

Don't harass the servers

Rule-of-thumb(?): three requests per second

Space your requests using `sleep`

```
from time import sleep
...
for link in movie_list:
    # request the link and do your thing
    ...
    sleep(0.3)
```

[python]

Unit is seconds.

So, `sleep(0.3)` = (approx.) three requests per second

# Thank you.

**Data Science Drop-in**

<https://5harad.com/drop-in/>

**Mondays @ 4–6 pm in Y2E2 253**