

Web Scraping

with **Jongbin Jung**

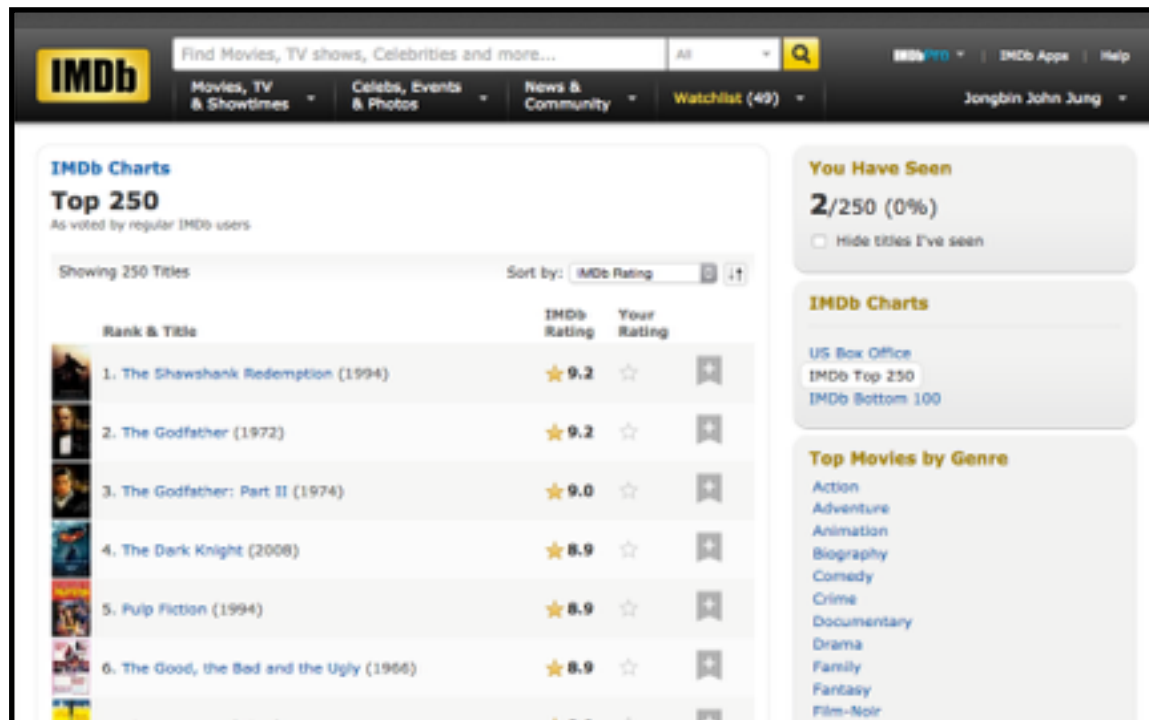
(jongbin@stanford.edu)

PhD student, Decision Analysis, MS&E

Before We Begin

- All material available at:
<https://github.com/jongbinjung/css-python-workshop>
- python packages required:
`requests | BeautifulSoup4 | selenium`
- Firefox browser is preferred
(works better with `selenium`)

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/G/01/2007/1902/24197/Bel13aushRFP120wv020606W02_W1_RK14_C00_0_38_50_R1.jpg" width="134" height="50" />
</td>
</td>
<td class="titleColumn">
<span name="tr" data-value="0.148">114</span>
<a href="/title/tt1049413/?ref=uhdp_t4_114">
<span name="Pete Dinkler (dir.), Edward Zinner, Jordan Kagan" />
<span name="ed" 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></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>
<div class="watchlistColumn"> <div class="wib_ribbon" data-hocount="1104413"></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/G/01/2007/1902/24197/Bel13aushRFP120wv020606W02_W1_RY58_C00_0_38_50_R1.jpg" width="134" height="50" />
</td>

```

Web Scraping?

Different ways to "scrape" the web:

1. **Standard** - load and scrape web page
2. The **Con** - fly like a browser, sting like a robot
3. The **Plumber** - look for the source of data
- * The **Suit** - formal public API (all proper and stuff)

method of choice will likely be determined by how the target website is designed

Word of Caution

- Different people have different ideas regarding what's OK or not - read the fine print!
- It's possible to cause harm if you aren't careful
- Scraping is more art than science

Stanford Data Science Drop-in

- <https://5harad.com/drop-in>
- a **free** educational and consulting service for members of the Stanford community
- many web scraping topics for today's workshop came from questions brought-in by students
- some contents were previously presented at one of the Data Science Drop-in's Mini Tutorials (<https://5harad.com/drop-in/tutorials>)

Standard

load and scrape web page

the **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)

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	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

First, get the Web to python

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

[illegible]

from source (html)
to python

```

✓ jongbinjung@0000a27b2eb:~$ 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/>

Use BeautifulSoup

Make the source (html) into a BeautifulSoup

```
from bs4 import BeautifulSoup

soup = BeautifulSoup(web_page.content)
```

[python]

```
✓ jongbinjung@Mba22b2eb:~$ 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	Cleomenza (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 Python's 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?

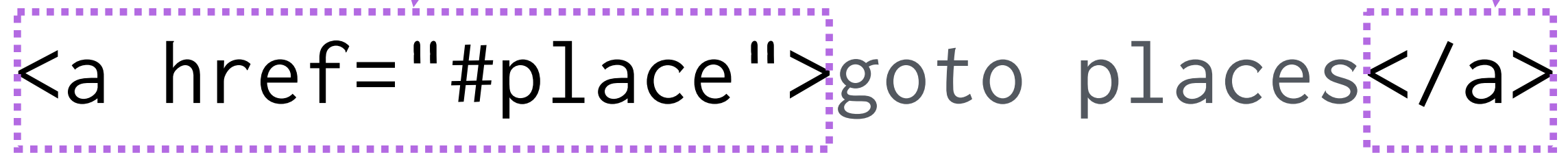
```
21 <div align="center">
22
23 <a href="bs4/download/"><h1>Beautiful Soup</h1></a>
24
25 <p>"A tremendous boon." -- <a
26 href="http://www.awaretek.com/python/index.html">Pyt
27 Podcast</a></p>
28
29 <p>[ <a href="#Download">Download</a> | <a
30 href="bs4/doc/">Documentation</a> | <a href="#Hallof
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>`

Basics of html

html is made of tags that look something like this

opening tags(<a>) and closing tags()

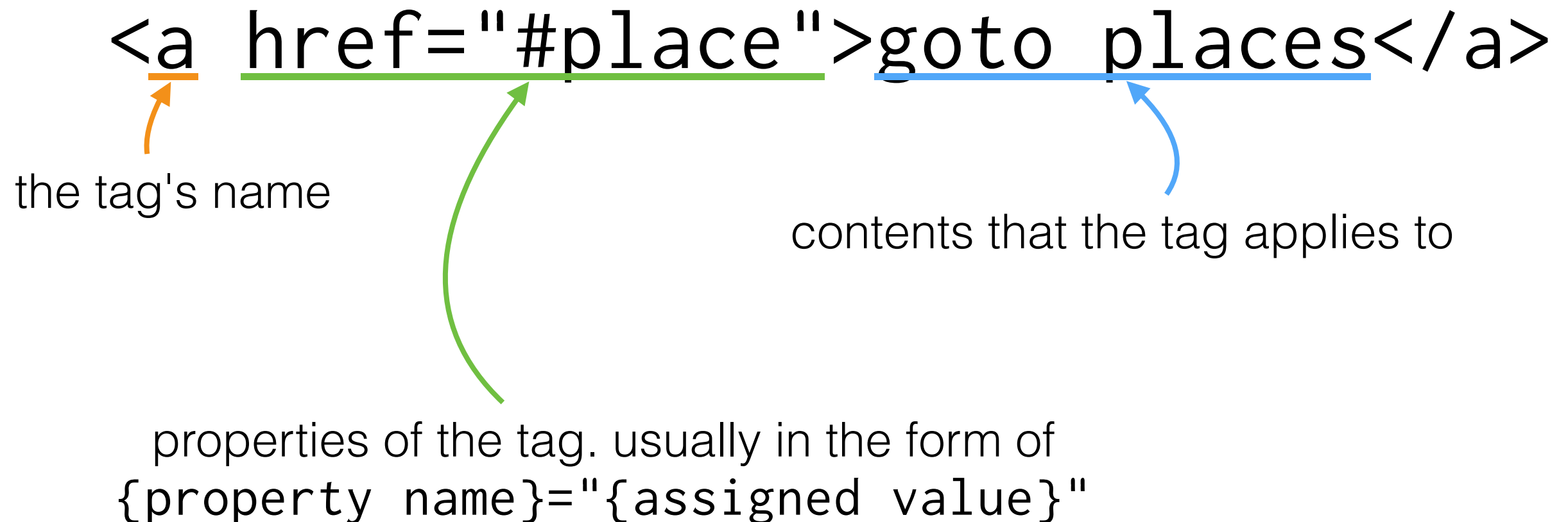


The diagram illustrates the structure of an HTML link. It shows the opening tag `` and the closing tag `` enclosed in dashed purple boxes. A purple arrow points from the text "opening tags(<a>)" to the opening tag box, and another purple arrow points from the text "closing tags()" to the closing tag box. The content "goto places" is placed between the two tags.

```
<a href="#place">goto places</a>
```

Basics of html

html is made of tags that look something like this



Web Scraper's Rule of Thumb

If it looks *different*,
or does something *different*,
it's probably in a *different* <tag>

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 `` 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]

html Structure

html is made of tags that look something like this

`<html>`

`<html>` tag declares the beginning/end of document

`<head></head>`

info/scripts etc. go between `<head></head>`

`<body>`

most of the main content

`</body>`

will be in the `<body></body>`

`</html>`

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 python list
- 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]

...or make a **list** of addresses with list comprehension

```
addresses = [link.get('href') for link in soup.find_all('a')]
```

[python]

Goto IMDb.com

the **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)

Let's start with

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




for just one movie.

We want this

...start by making a soup

Cast [Edit](#)

Cast overview, first billed only:

	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)
```

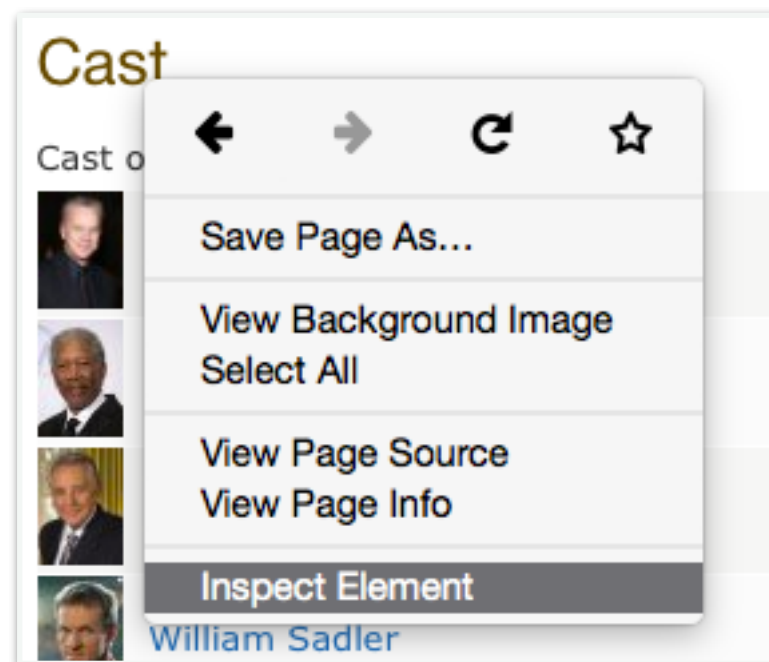
[python]

(http://www.imdb.com/title/tt0111161/?ref_=chttp_tt_1)

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.

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('table', class_='cast_list').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]



Save the values so we can format them appropriately later - in our case, into tab-separated values

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]

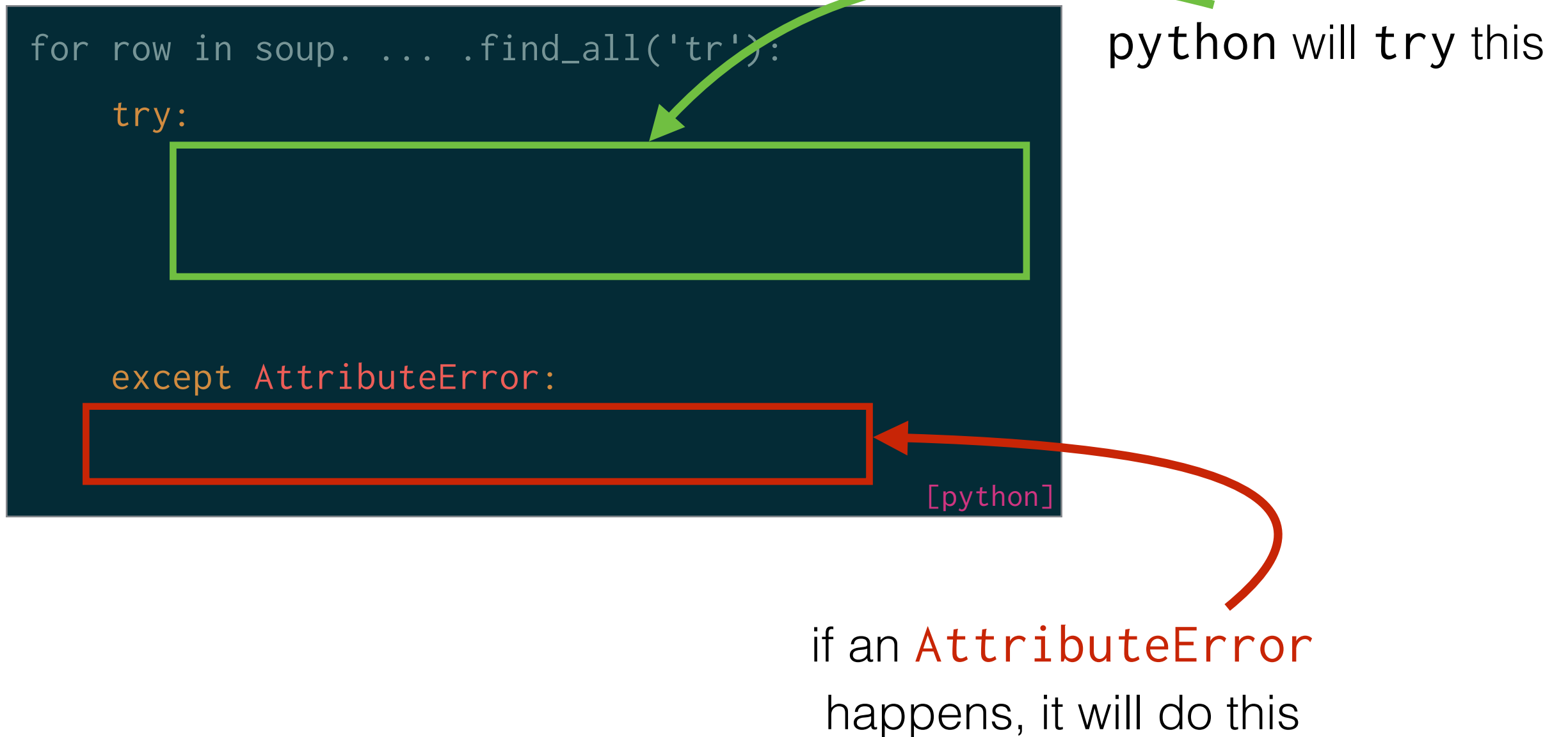
we want python to ignore these errors

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

[output]



Manage with `try-except` Blocks



Manage with `try-except` Blocks

```
for row in soup. ... .find_all('tr'):
```

```
try:
```

```
    actor = row.find(itemprop='name').text  
    role = row.find(class_='character').text
```

```
except AttributeError:
```

```
    pass
```

python will try this

[python]

WARNING!

Practice caution with `try-except`.
Don't pass an error, unless you're certain 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 ... ?

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 = clean_text(row.find(itemprop='name').text)
        role = clean_text(row.find(class_='character').text)

    except AttributeError:
        pass
```

[python]

Let's write a function that does this with any given string. You might find yourself needing this quite often.

```
def clean_text(text):
    return ' '.join(text.split())
```

[python]

Write the Data to a File

Usually, `print` is sufficient in python.

You can make tab-separated values from a `list` of `strings`:

```
print '\t'.join([actor, role])
```

[python]

`print` will send strings to `stdout`, which means you can save the output to a file by redirecting, e.g.:

```
$ python script_name.py > movie_data.tsv
```

[command prompt]

So far ...

```
from bs4 import BeautifulSoup
import requests

def clean_text(text):
    return ' '.join(text.split())

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

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

        print '\t'.join([actor, role])

    except AttributeError:
        pass
```

[python]

Exercise 1

the **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)

Workflow



A screenshot of the IMDb Top 250 movie list. The list contains 10 movies, each with a small thumbnail image on the left and the title followed by the 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

for these links do ***this***

```
from bs4 import BeautifulSoup
import requests

def clean_text(text):
    return ' '.join(text.split())

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

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

        print '\t'.join([actor, role])

    except AttributeError:
        pass
```


Be Nice

Don't harass the servers

Rule-of-thumb(?): act like a human (one request every 3~5 seconds)

Space your requests using `sleep`

```
from time import sleep  
  
...  
for link in movie_list:  
    # request the link and do your thing  
    ...  
    sleep(3)
```

[python]

Unit is seconds.

So, `sleep(3)` = one request every three seconds

DIY

Workflow



A screenshot of the IMDb Top 250 movies list. The list contains 10 items, each with a small movie poster icon on the left and the movie title followed by its release year in parentheses. The items are numbered 1 through 10.

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

for these links do ***this***

```
from bs4 import BeautifulSoup
import requests

def clean_text(text):
    return ' '.join(text.split())

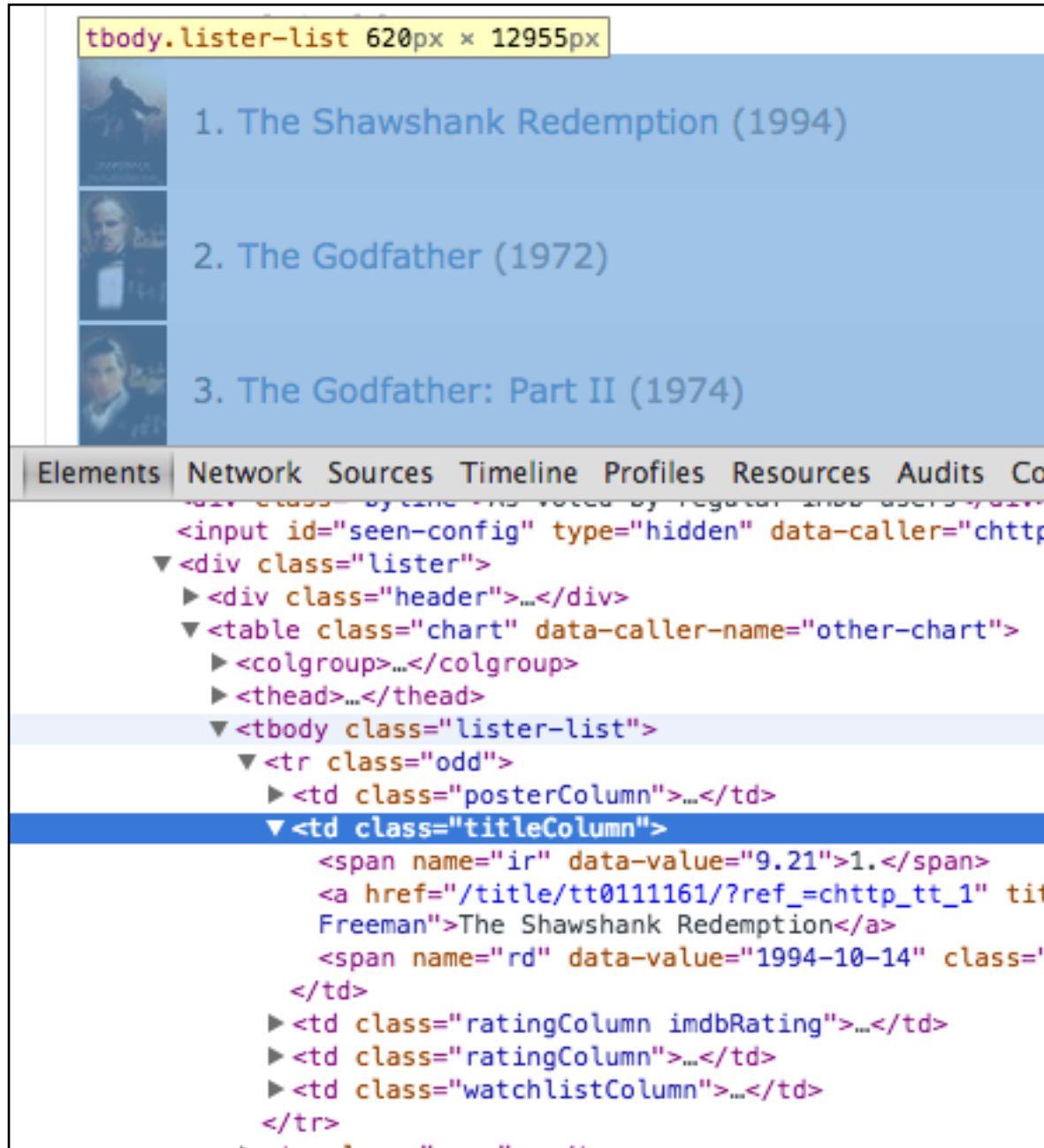
web_page = requests.get('http://...')
soup = BeautifulSoup(web_page.content)

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

        print '\t'.join([actor, role])

    except AttributeError:
        pass
```

The Links



The screenshot shows a web browser's developer tools interface. The top panel displays a list of movies with their posters and titles:

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

The bottom panel shows the HTML structure of the page. The selected element is a table row with the following structure:

```
<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_=chttp_tt_1" tit<br/>Freeman">The Shawshank Redemption</a>  
      <span name="rd" data-value="1994-10-14" class="<br/></td>  
    <td class="ratingColumn imdbRating">...</td>  
    <td class="ratingColumn">...</td>  
    <td class="watchlistColumn">...</td>  
  </tr>
```

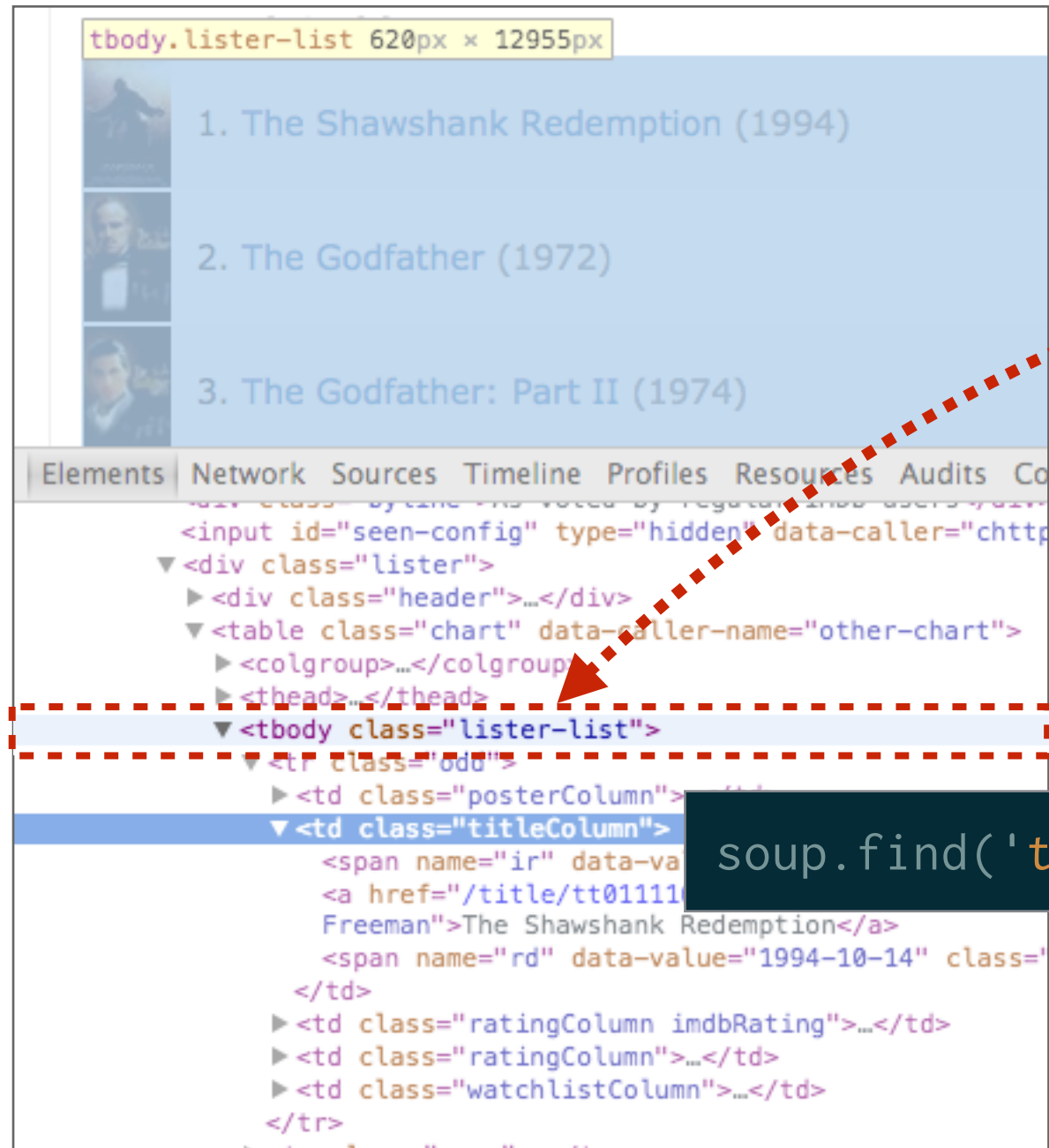
What do we need, and how should we get it?

remember:

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

[python]

The Links



The screenshot shows a web browser displaying a list of movies. The list is titled "tbody.lister-list 620px x 12955px" and contains three items:

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

Below the list, the browser's developer tools are open, showing the HTML structure. The "Elements" panel is selected, and 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 class="titleColumn">
          <span name="ir" data-value="1994-10-14" class="ir">
            <a href="/title/tt0111111" data-caller="http">
              Freeman">The Shawshank Redemption</a>
            <span name="rd" data-value="1994-10-14" class="rd">
              1994-10-14
            </span>
          </td>
          <td class="ratingColumn imdbRating">...</td>
          <td class="ratingColumn">...</td>
          <td class="watchlistColumn">...</td>
        </tr>
      <tr class="even">
        <td class="posterColumn">...
        <td class="titleColumn">
          <span name="ir" data-value="1972-03-09" class="ir">
            <a href="/title/tt0068617" data-caller="http">
              The Godfather</a>
            <span name="rd" data-value="1972-03-09" class="rd">
              1972-03-09
            </span>
          </td>
          <td class="ratingColumn imdbRating">...</td>
          <td class="ratingColumn">...</td>
          <td class="watchlistColumn">...</td>
        </tr>
      <tr class="odd">
        <td class="posterColumn">...
        <td class="titleColumn">
          <span name="ir" data-value="1974-12-13" class="ir">
            <a href="/title/tt0070703" data-caller="http">
              The Godfather: Part II</a>
            <span name="rd" data-value="1974-12-13" class="rd">
              1974-12-13
            </span>
          </td>
          <td class="ratingColumn imdbRating">...</td>
          <td class="ratingColumn">...</td>
          <td class="watchlistColumn">...</td>
        </tr>
    </tbody>
  </table>

```

A red dashed arrow points from the text "<tbody> seems to be a tag indicating the body of a table" to the "<tbody class='lister-list'>" line in the HTML code.

<tbody> seems to be
a tag indicating the
body of a table

```
soup.find('tbody', class_='lister-list') [python]
```

The Links

```
for movie in soup. ....find_all('td', class_='titleColumn')
```

[python]

2. The Godfather (1972)

3. The Godfather: Part II (1974)

Each movie title,
surrounded by an `<a>`
tag, seems to be
contained in `<td>` with
`class='titleColumn'`

`<td class="titleColumn">`

`1.`

`<a href="/title/tt0111161/?ref=chttp_tt_1" tit`
`Freeman">The Shawshank Redemption`

`<span name="rd" data-value="1994-10-14" class="`

`</td>`

`<td class="ratingColumn imdbRating">...</td>`

`<td class="ratingColumn">...</td>`

`<td class="watchlistColumn">...</td>`

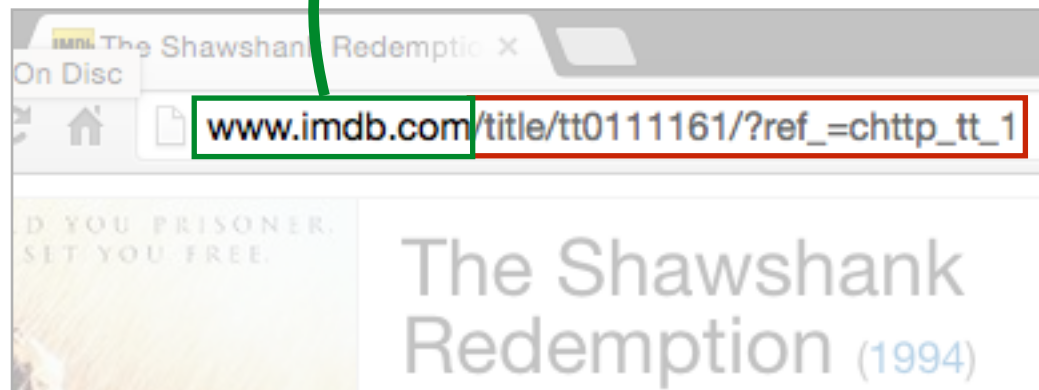
`</tr>`

The Links

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`

Workflow

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

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

soup = BeautifulSoup(web_page.content)

table_body = soup.find('tbody', class_='lister-list')
movie_list = table_body.find_all('td', class_='titleColumn', limit=10)

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

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

...
```

[python]

Tips & Tricks

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]

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

A note on URLs


- Back to IMDb
Most Popular by Genre > Action > 'show more...'

50.  **Fury** (2014) [Add to Watchlist](#)

★★★★★ 7.6/10

April, 1945. As the Allies make their final push in the European Theatre, a battle-hardened Army sergeant named Wardaddy commands a Sherman tank and his five-man crew on a deadly mission behind enemy lines. Outnumbered, out-gunned, and with a rookie soldier thrust into their platoon, Wardaddy and his men face overwhelming odds in their heroic attempts to strike at the heart of Nazi Germany.

Dir: David Ayer With: Brad Pitt, Shia LaBeouf, Logan Lerman
Action | Drama | War

134 mins. 

1-50 of 30,347 titles.

[Next »](#)

Pages? How many? What? Why?

more

Web Scraper's Rule of Thumbs


good web scrapers have many thumbs...

If it was organized by a *robot*,

it can be navigated by a *robot*.

A note on URLs

- Back to IMDb
Most Popular by Genre > Action > 'show more...'


50.  **Fury** (2014) [Add to Watchlist](#)

★★★★★ 7.6/10

April, 1945. As the Allies make their final push in the European Theatre, a battle-hardened Army sergeant named Wardaddy commands a Sherman tank and his five-man crew on a deadly mission behind enemy lines. Outnumbered, out-gunned, and with a rookie soldier thrust into their platoon, Wardaddy and his men face overwhelming odds in their heroic attempts to strike at the heart of Nazi Germany.

Dir: David Ayer With: Brad Pitt, Shia LaBeouf, Logan Lerman

Action | Drama | War

134 mins. 

1-50 of 30,347 titles.

[Next »](#)

Pages? How many? What? Why?

take a look at the link:

http:// ... /title?genres=action&sort=moviemeter,asc&start=51&title_type=feature

A note on URLs

- Look closer ...

```
http:// ... /title?  
    genres=action&  
    sort=moviemeter,asc&  
    start=51&  
    title_type=feature
```

- May require trial-and-error investigation to identify how each parameter's value(s) affect the result

A note on URLs

- Use `requests` to build URLs with parameters

```
import requests
```

```
base_url = 'http://www.imdb.com/search/title'
```

```
params = {'genres': 'action', 'sort': 'moviemeter,asc'}
```

```
page = requests.get(base_url, params=params)
```

```
page.url
```

[python]

```
u'http://www.imdb.com/search/title?genres=action&sort=moviemeter,asc'
```

[output]

Exercise 2

collect the **title**, **release year**, and
grossing US Box Office Income (\$)

for each of the **Top 250 comedy movies** sorted by
US Box Office Gross Income

(http://www.imdb.com/search/title?genres=comedy&sort=boxoffice_gross_us)

The Con

fly like a browser, sting like a robot

the **GOAL**

scrape your own webmail

get comments from today's headline of the guardian

- * Disclaimer: Unless you have some experience developing web pages or at least some kind of experience with javascript, it's unlikely that you'll understand everything in this section at first glance. This section is only included to give you an idea of what's possible, and hopefully you can figure things out more clearly by carefully examining the examples (on github), should you require the methods of this section.

The last resort

- Some web pages are impossible to browse through with python (ok, not impossible, but complicated enough that you don't have the understanding of computer networks required to make it work)
 - proprietary databases that require login
 - secured (encrypted) connections
 - ...
- But you can still see stuff on your browser!
- What if a robot can just emulate the **copy-paste**?

selenium

- `selenium` is a webdriver - it let's you programmatically *drive* your browser
- originally built for automated website testing
- easily repurposed for web scraping
- not many distributions include `selenium` by default, install with pip (in command prompt)

```
$ pip install selenium
```

[command prompt]

selenium

- with `selenium`, you can create a browser object within your code, which will launch a real browser!

```
from selenium import webdriver
```

```
browser = webdriver.Firefox()
```

[python]

- while `selenium` supports a variety of browsers, setting it up to work with anything other than Firefox has proven quite painful (at least to me)
- see the docs here:
<https://selenium-python.readthedocs.org/>

Load a URL

- you can just type a URL into the browser, or
- load a page programmatically with

```
URL = 'http://webmail.stanford.edu'  
browser.get(URL)
```

[python]

- we'll scrape our own Stanford Webmail, which requires two-factor authentication

find elements

- two approaches once a web page is loaded:

1. load the source with `.page_source`, then bs4

```
html = browser.page_source
soup = BeautifulSoup(html)
for sender in soup.find_all('div', class_='class_name'):
    print sender.text
```

[python]

2. use selenium's `.find_elements_by_*`

```
for element in browser.find_elements_by_class_name('class_name'):
    print element.text
```

[python]

scrolling

- you might need to scroll certain elements
- selenium can send javascript to the browser
- javascript is not within scope of this workshop, but let's use just a snippet ...

```
panes = browser.find_elements_by_class_name('scrollContainer')
for pane in panes:
    if '_xlv_e1' in pane.get_attribute('class'):
        browser.execute_script('arguments[0].scrollTop =
                                arguments[0].scrollHeight', pane)
```

[python]

waiting

- once logged in, the rest can be automated
- even waiting for pages to load
- `selenium` can wait for certain elements to load
- this requires a few extra modules

```
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
from selenium.webdriver.common.by import By
from selenium.common.exceptions import TimeoutException
```

[python]

waiting

- the standard syntax look like:

```
containers = WebDriverWait(browser, MAX_WAIT).until(  
    EC.presence_of_all_elements_located(  
        (By.CLASS_NAME, "fc-container__inner"))
```

[python]

- waits until elements with class 'fc-container__inner' are found, and returns the list of elements
- throws `TimeoutException` if `MAX_WAIT` seconds pass without finding at least one element that matches

The **Plumber**

look for the source of data

the **GOAL**

Collect used car price data

* Some critical contents on this section have been masked from the slide deck to protect the target website(s) from potential harm. Source code and related materials can be found on the [github repo](#)

“Why is the source code that I `get` from `requests` different from what I see in the browser's Developer Mode?”

- Everyone who has tried scraping dynamically loaded pages

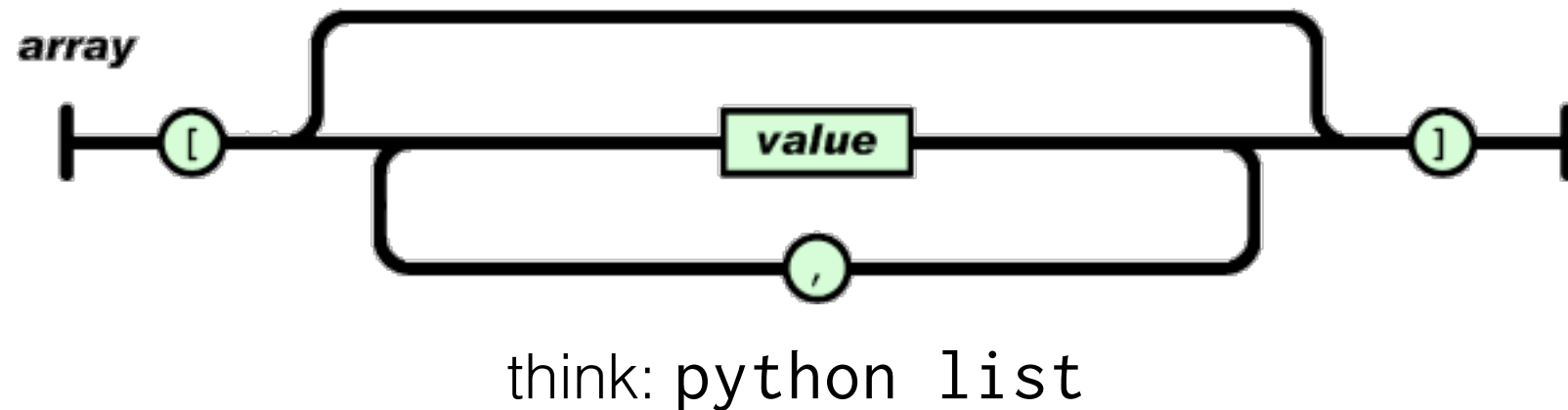
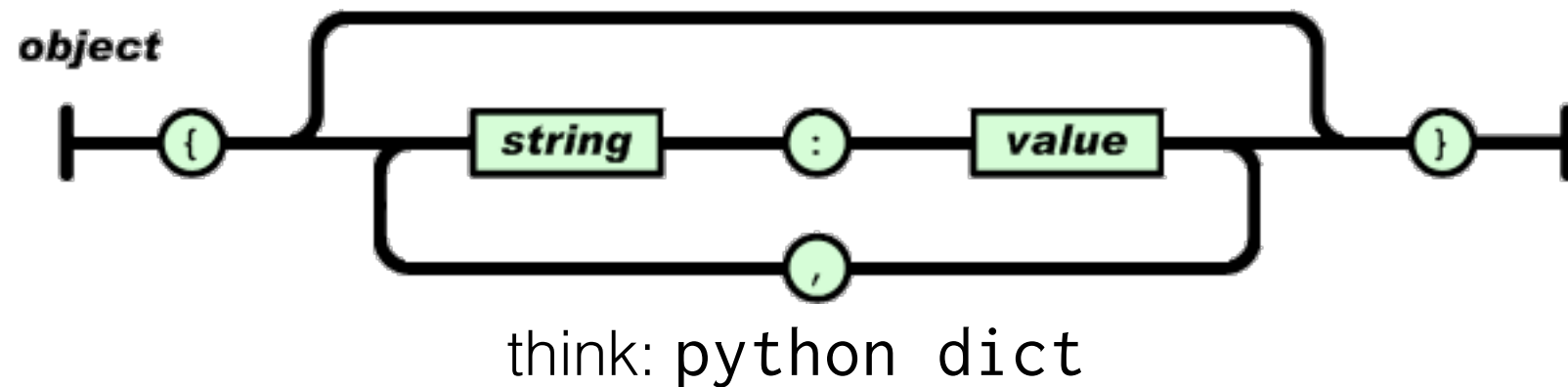
This will only work *sometimes*

- The [Network] tab in your browser's Developer mode will let you investigate from *where* and *how* the web page was loaded
- This may (or may not) give you insight on a better approach to retrieving data that is presented to you
- Sometimes, you might notice that the data on the page is loaded from a `json` file

Introduction to JSON

- a text-base data representation format that can be used to represent complex data structures
- often the data format of choice between server-client communication due to **light-weight** (minimal symbols) and **interchangeable** (text-based) **standardized formatting** (cf. csv)

Introduction to JSON



JSON with requests

- the `requests.Response` object (the object returned by `requests.get`) has a `json()` method, which returns the json-encoded content of a response, if any

```
# assuming the base_url returns json content  
data = requests.get(base_url, params=params).json()
```

[python]

- once the json-encoded content is loaded, you can treat it as a regular python dictionary

Exercise 3

Choose a make/model of interest (e.g., Honda Accord)

Collect the **VIN**, **year**, **trim**, **mileage**, and **price**
for all used cars of selected make/model,
available within 50 miles of your zip-code.

(you should only need to make a single requests)

the **Suit**

will cover if we have time.

no detailed slides

but example codes with comments on git

Using the twitter api with `tweepy`

see: `web-scraping/examples/collect_tweets.py`

additional exercises

if we somehow manage to cover everything,

as *still* have more time

on-site exercise

codes will be synced to git after workshop (if covered)

- Yelp!
notoriously stingy API
- USPTO.gov
organized by humans == difficult for robots

Thank you.