



CS109A Introduction to Data Science

Homework 1: Data Collection - Web Scraping - Data Parsing ¶

Harvard University

Fall 2018

Instructors: Pavlos Protopapas and Kevin Rader

```
In [58]: ## RUN THIS CELL TO GET THE RIGHT FORMATTING  
import requests  
from IPython.core.display import HTML  
styles = requests.get("https://raw.githubusercontent.com/Harvard-IACS/2018-CS109A/master/content/styles/cs109.css").text  
HTML(styles)
```

Out[58]:

Instructions

- To submit your assignment follow the instructions given in Canvas.
- The deliverables in Canvas are:
 - a) This python notebook with your code and answers, plus a pdf version of it (see Canvas for details),
 - b) the bibtex file you created,
 - c) The CSV file you created,
 - d) The JSON file you created.
- Exercise **responsible scraping**. Web servers can become slow or unresponsive if they receive too many requests from the same source in a short amount of time. Use a delay of 10 seconds between requests in your code. This helps not to get blocked by the target website. Run the webpage fetching part of the homework only once and do not re-run after you have saved the results in the JSON file (details below).
- Web scraping requests can take several minutes. This is another reason why you should not wait until the last minute to do this homework.
- For this assignment, we will use Python 3.5 for grading.

Data Collection - Web Scraping - Data Parsing

In this homework, your goal is to learn how to acquire, parse, clean, and analyze data. Initially you will read the data from a file, and then later scrape them directly from a website. You will look for specific pieces of information by parsing the data, clean the data to prepare them for analysis, and finally, answer some questions.

In doing so you will get more familiar with three of the common file formats for storing and transferring data, which are:

- CSV, a text-based file format used for storing tabular data that are separated by some delimiter, usually comma or space.
- HTML/XML, the stuff the web is made of.
- JavaScript Object Notation (JSON), a text-based open standard designed for transmitting structured data over the web.

```
In [59]: # import the necessary libraries
import matplotlib inline
import numpy as np
import scipy as sp
import matplotlib as mpl
import matplotlib.cm as cm
import matplotlib.pyplot as plt
import pandas as pd
import time
pd.set_option('display.width', 500)
pd.set_option('display.max_columns', 100)
pd.set_option('display.notebook_repr_html', True)
import seaborn as sns
```

Help a professor parse their publications and extract information.

Overview

In this part your goal is to parse the HTML page of a professor containing some of his/her publications, and answer some questions. This page is provided to you in the file `data/publist_super_clean.html`. There are 45 publications in descending order from No. 244 to No. 200.

```
In [60]: # use this file provided  
PUB_FILENAME = 'data/publist_super_clean.html'
```

Question 1 [40 pts]: Parsing and Converting to bibTex and CSV using BeautifulSoup and python string manipulation

A lot of the bibliographic and publication information is displayed in various websites in a not-so-structured HTML files. Some publishers prefer to store and transmit this information in a .bibTex file which looks roughly like this (we've simplified a few things):

```
@article {
    author = "John Doyle"
    title = "Interaction between atoms"
    URL = "Papers/PhysRevB_81_085406_2010.pdf"
    journal = "Phys. Rev. B"
    volume = "81"
}
```

You will notice that this file format is a set of items, each of which is a set of key-value pairs. In the python world, you can think of this as a list of dictionaries. If you think about spreadsheets (as represented by CSV files), they have the same structure. Each line is an item, and has multiple features, or keys, as represented by that line's value for the column corresponding to the key.

You are given an .html file containing a list of papers scraped from the author's website and you are to write the information into .bibTex and .CSV formats. A useful tool for parsing websites is BeautifulSoup (<http://www.crummy.com/software/BeautifulSoup/> (<http://www.crummy.com/software/BeautifulSoup/>)) (BS). In this problem, will parse the file using BS, which makes parsing HTML a lot easier.

1.1 Write a function called `make_soup` that accepts a filename for an HTML file and returns a BS object.

1.2 Write a function that reads in the BS object, parses it, converts it into a list of dictionaries: one dictionary per paper. Each of these dictionaries should have the following format (with different values for each publication):

```
{'author': 'L.A. Agapito, N. Kioussis and E. Kaxiras',
 'title': '"Electric-field control of magnetism in graphene quantum dots:\nAb initio calculations"',
 'URL': 'Papers/PhysRevB_82_201411_2010.pdf',
 'journal': 'Phys. Rev. B',
 'volume': '82'}
```

1.3 Convert the list of dictionaries into standard .bibTex format using python string manipulation, and write the results into a file called `publist.bib`.

1.4 Convert the list of dictionaries into standard tabular .csv format using pandas, and write the results into a file called `publist.csv`. The csv file should have a header and no integer index.

HINT

- Inspect the HTML code for tags that indicate information chunks such as `title` of the paper. The `find_all` method of BeautifulSoup might be useful.
- Question 1.2 is better handled if you break the code into functions, each performing a small task such as finding the author(s) for each paper.

- Question 1.3 is effectively tackled by first using python string formatting on a template string.
- Make sure you catch exceptions when needed.
- Make sure you check for **missing data** and handle these cases as you see fit.

Resources

- [BeautifulSoup Tutorial \(https://www.dataquest.io/blog/web-scraping-tutorial-python/\)](https://www.dataquest.io/blog/web-scraping-tutorial-python/).
- More about the [BibTex format \(http://www.bibtex.org\)](http://www.bibtex.org).

Answers

```
In [61]: # import the necessary libraries  
from bs4 import BeautifulSoup
```

****1.1** Write a function called ``make_soup`` ...

```
In [62]: def make_soup(filename: str) -> BeautifulSoup:  
    '''Open the file and convert into a BS object.  
  
    Args:  
        filename: A string name of the file.  
  
    Returns:  
        A BS object containing the HTML page ready to be parsed.  
    '''  
    with open(filename, 'r', encoding='utf-8') as f:  
        myfile = f.read()  
    return BeautifulSoup(myfile, 'html.parser')  
  
soup = make_soup(PUB_FILENAME)
```

```
In [63]: # check your code - print the BS object, you should get a familiar HTML  
         page as text  
         # clear/remove output before making pdf  
         # print(soup)
```

Your output should look **like** this:

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN"
"http://www.w3.org/TR/html4/loose.dtd">

<title>Kaxiras E journal publications</title>
<head>
<meta content="text/html; charset=utf-8" http-equiv="Content-Type"/>
<link href="../../styles/style_pubs.css" rel="stylesheet" type="text/css"/>
<meta content="" name="description"/>
<meta content="Kaxiras E, Multiscale Methods, Computational Materials" name
="keywords"/>
</head>
<body>
<ol start="244">
<li>
<a href="Papers/2011/PhysRevB_84_125411_2011.pdf" target="paper244">
"Approaching the intrinsic band gap in suspended high-mobility graphene nano
ribbons"</a>
<br/>Ming-Wei Lin, Cheng Ling, Luis A. Agapito, Nicholas Kioussis, Yiyang Zh
ang, Mark Ming-Cheng Cheng,
<i>PHYSICAL REVIEW B </i> <b>84</b>, 125411 (2011)
<br/>
</li>
</ol>
<ol start="243">
<li>
<a href="Papers/2011/PhysRevB_84_035325_2011.pdf" target="paper243">
"Effect of symmetry breaking on the optical absorption of semiconductor nano
particles"</a>
<br/>JAdam Gali, Efthimios Kaxiras, Gergely T. Zimanyi, Sheng Meng,
<i>PHYSICAL REVIEW B </i> <b>84</b>, 035325 (2011)
<br/>
</li>
</ol>

...
```

1.2 Write a function that reads in the BS object, parses it, converts it into a list of dictionaries...

```
In [64]: def make_dict(soup):
    my_list = []
    li_tags = soup.find_all('li')
    for i in li_tags:
        URL = i.find('a')['href'].strip('\n').strip('').rstrip(',')
        title = i.find('a').contents[0].strip('\n').rstrip(',')
        author = i.contents[4].strip('\n').strip('').strip().rstrip(',')
        journal = i.find('i').contents[0].strip('\n').strip('').rstrip(',')
        if i.find('b') != None:
            volume = i.find('b').contents[0].strip('\n').strip('').rstrip(',')
        else:
            volume = ''
        my_dict = {'URL':URL, 'title':title, 'author':author, 'journal':journal, 'volume':volume}
        my_list.append(my_dict)
    return my_list
```

```
In [65]: my_list = make_dict(soup)
    #print(my_list)
```

1.3 Convert the list of dictionaries into the .bibTex format using python string manipulation (python string formatting on a template string is particularly useful)..

```
In [66]: # TO to MK: exception handlings for empty values is done in make_dict()
    - there is not empty value after make_dict()

    from string import Template
    def make_bibtex (my_list):
        template = '@article{\n\tauthor = "$author",\n\ttitle = "$title",\n\tURL = "$url",\n\tjournal = "$journal",\n\tvolume = $volume\n}'
        bibtex = ""
        for i in my_list:
            bibtex = bibtex + Template(template).substitute(author = i['author'], title = i['title'].strip('').replace('\n', ''), url = i['URL'], journal = i['journal'], volume = i['volume']) + '\n\n'
        bibtex = bibtex.replace(',\n\tvolume = \n', '\n')
        return bibtex
```

```
In [67]: # your code here
    with open('publist.bib','w', encoding='utf-8') as my_output:
        my_output.write(make_bibtex(my_list))
```

```
In [68]: # check your answer - print the bibTex file
# clear/remove output before making pdf
f = open('publist.bib','r')
#print (f.read())
```

Your output should look like this

```
@article{
  author = "Ming-Wei Lin, Cheng Ling, Luis A. Agapito, Nicholas Kiuoussis,
  Yiyang Zhang, Mark Ming-Cheng Cheng",
  title = "Approaching the intrinsic band gap in suspended high-mobility
  graphene nanoribbons",
  URL = "Papers/2011/PhysRevB_84_125411_2011.pdf",
  journal = "PHYSICAL REVIEW B",
  volume = 84
}

...

@article{
  author = "E. Kaxiras and S. Succi",
  title = "Multiscale simulations of complex systems: computation meets r
  eality",
  URL = "Papers/SciModSim_15_59_2008.pdf",
  journal = "Sci. Model. Simul.",
  volume = 15
}
```

1.4 Convert the list of dictionaries into the .csv format using pandas, and write the data into publist.csv. The csv file should have a header and no integer index...


```
In [69]: # make sure you use head() when printing the dataframe
# your code here
df = pd.DataFrame(my_list)
df.head()
```

Out[69]:

	URL	author	journal	title
0	Papers/2011/PhysRevB_84_125411_2011.pdf	Ming-Wei Lin, Cheng Ling, Luis A. Agapito, Nic...	PHYSICAL REVIEW B	"Approaching the intrinsic band gap in suspend...
1	Papers/2011/PhysRevB_84_035325_2011.pdf	JAdam Gali, Efthimios Kaxiras, Gergely T. Zima...	PHYSICAL REVIEW B	"Effect of symmetry breaking on the optical ab...
2	Papers/2011/PhysRevB_83_054204_2011.pdf	Jan M. Knaup, Han Li, Joost J. Vlassak, and Ef...	PHYSICAL REVIEW B	"Influence of CH2 content and network defects ...
3	Papers/2011/PhysRevB_83_045303_2011.pdf	Martin Heiss, Sonia Conesa-Boj, Jun Ren, Hsian...	PHYSICAL REVIEW B	"Direct correlation of crystal structure and o...
4	Papers/2011/PhilTransRSocA_369_2354_2011.pdf	Simone Melchionna, Efthimios Kaxiras, Massimo ...	Phil. Trans. R. Soc. A	"Endothelial shear stress from large-scale blo...

```
In [70]: # your code here
df.to_csv('publist.csv', index=False, encoding='utf-8', quoting = 1)
```

```
In [71]: # your code here - testing if the csv works
test_csv = pd.read_csv('publist.csv')
with pd.option_context('display.max_rows', None, 'display.max_columns',
None):
    print(test_csv.head())
```

	author	URL	journal
0	Papers/2011/PhysRevB_84_125411_2011.pdf	Ming-Wei Lin, Cheng Li ng, Luis A. Agapito, Nic...	PHYSICAL REVIEW B "Approaching the i ntrinsic band gap in suspend...
1	Papers/2011/PhysRevB_84_035325_2011.pdf	JAdam Gali, Efthimios Kaxiras, Gergely T. Zima...	PHYSICAL REVIEW B "Effect of symmetr y breaking on the optical ab...
2	Papers/2011/PhysRevB_83_054204_2011.pdf	Jan M. Knaup, Han Li, Joost J. Vlassak, and Ef...	PHYSICAL REVIEW B "Influence of CH2 content and network defects ...
3	Papers/2011/PhysRevB_83_045303_2011.pdf	Martin Heiss, Sonia Co nesa-Boj, Jun Ren, Hsian...	PHYSICAL REVIEW B "Direct correlatio n of crystal structure and o...
4	Papers/2011/PhilTransRSocA_369_2354_2011.pdf	Simone Melchionna, Eft himios Kaxiras, Massimo ...	Phil. Trans. R. Soc. A "Endothelial shear stress from large-scale blo...

Follow the stars in IMDb's list of "The Top 100 Stars for 2017"

Overview

In this part, your goal is to extract information from IMDb's Top 100 Stars for 2017

(<https://www.imdb.com/list/ls025814950/> (<https://www.imdb.com/list/ls025814950/>)) and perform some analysis on each star in the list. In particular we are interested to know: a) how many performers made their first movie at 17? b) how many performers started as child actors? c) who is the most proliferate actress or actor in IMDb's list of the Top 100 Stars for 2017? . These questions are addressed in more details in the Questions below.

When data is not given to us in a file, we need to fetch them using one of the following ways:

- download a file from a source URL
- query a database
- query a web API
- scrape data from the web page

Question 2 [52 pts]: Web Scraping using BeautifulSoup and exploring using Pandas

2.1 Download the webpage of the "Top 100 Stars for 2017" (<https://www.imdb.com/list/ls025814950/>) into a `requests` object and name it `my_page`. Explain what the following attributes are:

- `my_page.text`,
- `my_page.status_code`,
- `my_page.content`.

2.2 Create a BeautifulSoup object named `star_soup` using `my_page` as input.

2.3 Write a function called `parse_stars` that accepts `star_soup` as its input and generates a list of dictionaries named `starlist` (see definition below; order of dictionaries does not matter). One of the fields of this dictionary is the `url` of each star's individual page, which you need to scrape and save the contents in the `page` field. Note that there is a ton of information about each star on these webpages.

```
name: the name of the actor/actress as it appears at the top
gender: 0 or 1: translate the word 'actress' into 1 and 'actor' into '0'
url: the url of the link under their name that leads to a page with details
page: BS object with html text acquired by scraping the above 'url' page'
```

2.4 Write a function called `create_star_table` which takes `starlist` as an input and extracts information about each star (see function definition for the exact information to be extracted and the exact output definition). Only extract information from the first box on each star's page. If the first box is acting, consider only acting credits and the star's acting debut, if the first box is Directing, consider only directing credits and directorial debut.

2.5 Now that you have scraped all the info you need, it's good practice to save the last data structure you created to disk. Save the data structure to a JSON file named `starinfo.json` and submit this JSON file in Canvas. If you do this, if you have to restart, you won't need to redo all the requests and parsings from before.

2.6 We provide a JSON file called `data/staff_starinfo.json` created by CS109 teaching staff for consistency, which you should use for the rest of the homework. Import the contents of this JSON file into a pandas dataframe called `frame`. Check the types of variables in each column and clean these variables if needed. Add a new column to your dataframe with the age of each actor when they made their first appearance, movie or TV, (name this column `age_at_first_movie`). Check some of the values of this new column. Do you find any problems? You don't need to fix them.

2.7 You are now ready to answer the following intriguing questions:

- **2.7.1** How many performers made their first appearance (movie or TV) when he/she was 17 years old?
- **2.7.2** How many performers started as child actors? Define child actor as a person younger than 12 years old.

2.8 Make a plot of the number of credits against the name of actor/actress. Who is the most prolific actress or actor in IMDb's list of the Top 100 Stars for 2017? Define **most prolific** as the performer with the most credits.

Hints

- Create a variable that groups actors/actresses by the age of their first movie. Use pandas' `.groupby` to divide the dataframe into groups of performers that for example started performing as children (age < 12). The grouped variable is a `GroupBy` pandas object and this object has all of the information needed to then apply operations to each of the groups.
- When cleaning the data make sure the variables with which you are performing calculations are in numerical format.
- The column with the year has some values that are double, e.g. '2000-2001' and the column with age has some empty cells. You need to deal with these in a reasonable fashion before performing calculations on the data.
- You should include both movies and TV shows.

Resources

- The `requests` library makes working with HTTP requests powerful and easy. For more on the `requests` library see <http://docs.python-requests.org/> (<http://docs.python-requests.org/>).

Answers

```
In [72]: import requests
```

2.1 Download the webpage of the "Top 100 Stars for 2017 ...

```
In [73]: my_page = requests.get('https://www.imdb.com/list/ls025814950/')  
  
your answer here
```

2.2 Create a BeautifulSoup object named `star_soup` giving `my_page` as input.

```
In [74]: star_soup = BeautifulSoup(my_page.content, 'html.parser')
```

```
In [75]: # check your code - you should see a familiar HTML page  
# clear/remove output before making pdf  
# print (star_soup.prettify()[:])
```

2.3 Write a function called `parse_stars` that accepts `star_soup` as its input ...

Function

parse_stars

Input

star_soup: the soup object with the scraped page

Returns

a list of dictionaries; each dictionary corresponds to a star profile and has the following data:

name: the name of the actor/actress as it appears at the top
 gender: 0 or 1: translate the word 'actress' into 1 and 'actor' into '0'
 url: the url of the link under their name that leads to a page with details
 page: BS object with 'html text acquired by scraping the above 'url' page'

Example:

```
{'name': Tom Hardy,
  'gender': 0,
  'url': https://www.imdb.com/name/nm0362766/?ref_=nm1s_hd,
  'page': BS object with 'html text acquired by scraping the 'url' page'
}
```

```
In [76]: base_url = 'https://www.imdb.com'
def parse_stars(star_soup):
    starlist = []
    all_stars = star_soup.find_all(class_="liister-item mode-detail")
    for star in all_stars:
        name = star.find(class_='liister-item-content').find('a').get_text().strip().title()
        if star.find(class_='text-muted text-small').contents[0].strip().lower() == 'actress':
            gender = 1
        elif star.find(class_='text-muted text-small').contents[0].strip().lower() == 'actor':
            gender = 0
        else:
            gender = ''
        url = base_url + star.find(class_='liister-item-content').find('a')['href']
        #time.sleep(10)
        get_page = requests.get(url)
        page = BeautifulSoup(get_page.content, 'html.parser')

        my_dict = {'name':name, 'gender':gender, 'url':url, 'page':page}
        starlist.append(my_dict)
    return starlist

starlist = parse_stars(star_soup)
```

```
In [77]: len(starlist)
```

```
Out[77]: 100
```

This should give you 100

```
In [78]: # check your code
# this list is large because of the html code into the `page` field
# to get a better picture, print only the first element
# clear/remove output before making pdf
# print(starlist[34])
```

Your output should look like this:

```
{'name': 'Gal Gadot',
 'gender': 1,
 'url': 'https://www.imdb.com/name/nm2933757?ref_=nm1s_hd',
 'page':
<!DOCTYPE html>

<html xmlns:fb="http://www.facebook.com/2008/fbml" xmlns:og="http://ogp.me/
ns#">
<head>
<meta charset="utf-8"/>
<meta content="IE=edge" http-equiv="X-UA-Compatible"/>
<meta content="app-id=342792525, app-argument=imdb:///name/nm2933757?src=md
ot" name="apple-itunes-app"/>
<script type="text/javascript">var IMDbTimer={starttime: new Date().getTime
(),pt:'java'};</script>
<script>
    if (typeof uet == 'function') {
        uet("bb", "LoadTitle", {wb: 1});
    }
</script>
<script>(function(t){ (t.events = t.events || {})[ "csm_head_pre_title" ] = n
ew Date().getTime(); })(IMDbTimer);</script>

...
```

2.4 Write a function called `create_star_table` to extract information about each star ...

Function

create_star_table

Input

the starlist

Returns

a list of dictionaries; each dictionary corresponds to a star profile and has the following data:

star_name: the name of the actor/actress as it appears at the top
gender: 0 or 1 (1 for 'actress' and 0 for 'actor')
year_born : year they were born
first_movie: title of their first movie or TV show
year_first_movie: the year they made their first movie or TV show
credits: number of movies or TV shows they have made in their career.

Example:

```
{ 'star_name': Tom Hardy,  
  'gender': 0,  
  'year_born': 1997,  
  'first_movie' : 'Batman',  
  'year_first_movie' : 2017,  
  'credits' : 24}
```



```
In [79]: def create_star_table(starlist: list) -> list:
    stars_info_list = []
    for i in starlist:
        my_object = i['page']
        name = i['name']
        gender = i['gender']
        try:
            if my_object.find(id='name-born-info') != None:
                year_born = my_object.find(id='name-born-info').find_all(
('a'))[1].get_text().strip(' ')
            else:
                year_born = None
            helper_object = my_object.find(class_='filmo-category-section').find_all('div')[-1]
            if helper_object.find('a') != None:
                first_movie = helper_object.find('a').get_text().strip()
            if helper_object.find('span') != None:
                year_first_movie = helper_object.find('span').get_text().strip()
            else:
                year_first_movie = None
            if my_object.find(id='filmography').find(class_='head') != None:
                credits = my_object.find(id='filmography').find(class_='head').contents[6].split()[0].split('(')[1]
            else:
                credits = None
        except Exception:
            pass
        my_dict={'name':name, 'gender':gender, 'year_born':year_born, 'first_movie':first_movie, 'year_first_movie':year_first_movie, 'credits':credits}
        stars_info_list.append(my_dict)
    return stars_info_list

star_table = create_star_table(starlist)
```

```
In [80]: # check your code
# clear/remove output before making the pdf file
star_table[0]
```

```
Out[80]: {'name': 'Gal Gadot',
'gender': 1,
'year_born': '1985',
'first_movie': 'Bubot',
'year_first_movie': '2007',
'credits': '26'}
```

2.5 Now that you have scraped all the info you need, it's a good practice to save the last data structure you ...

```
In [81]: import json

with open('starinfo.json', 'w', encoding='utf-8') as f:
    json.dump(star_table, f)
```

To check your JSON saving, re-open the JSON file and reload the code

```
In [82]: with open("starinfo.json", "r") as fd:
        star_table = json.load(fd)

# output should be the same
# clear/remove output before making the pdf file
star_table[0]
```

```
Out[82]: {'name': 'Gal Gadot',
          'gender': 1,
          'year_born': '1985',
          'first_movie': 'Bubot',
          'year_first_movie': '2007',
          'credits': '26'}
```

2.6 Import the contents of the staff's JSON file (data/staff_starinfo.json) into a pandas dataframe.

...

```
In [83]: # your code here
with open('data/staff_starinfo.json', 'r') as f:
    my_data = json.load(f)
```

```
In [84]: # your code here
df = pd.DataFrame(my_data)
with pd.option_context('display.max_rows', None, 'display.max_columns',
None):
    print(df.head())
```

	credits	first_movie	gender	name	year_born	year_f
0	25	Bubot	1	Gal Gadot	1985	
	2007					
1	55	Tommaso	0	Tom Hardy	1977	
	2001					
2	17	Doctors	1	Emilia Clarke	1986	
	2009					
3	51	All My Children	1	Alexandra Daddario	1986	
	2002-2003					
4	30	Järngänget	0	Bill Skarsgård	1990	
	2000					

```
In [85]: #clean up data
def split_string (x):
    if '-' in x:
        my_list = x.split('-')
        return (int(my_list[0]))
    elif '/' in x:
        my_list = x.split('/')
        return (int(my_list[0]))
    else:
        return (int(x))

def clean_data (df):
    df['year_first_movie'] = df['year_first_movie'].apply(lambda x: split_string(x))
    df['year_born'] = df['year_born'].astype('int')
    df['year_first_movie'] = df['year_first_movie'].astype('int')
    df['credits'] = df['credits'].astype('int')
    return df

df = clean_data(df)
```

```
In [86]: # your code here
with pd.option_context('display.max_rows', None, 'display.max_columns',
None):
    print(df.head())
```

	credits	first_movie	gender	name	year_born	year_first_movie
0	25	Bubot	1	Gal Gadot	1985	2007
1	55	Tommaso	0	Tom Hardy	1977	2001
2	17	Doctors	1	Emilia Clarke	1986	2009
3	51	All My Children	1	Alexandra Daddario	1986	2002
4	30	Järngänget	0	Bill Skarsgård	1990	2000

your answer here

2.7 You are now ready to answer the following intriguing questions:

2.7.1 How many performers made their first movie at 17?

```
In [87]: # your code here
df['age'] = df['year_first_movie'] - df['year_born']

"{0} performers made their first movie at 17".format(df['age'].value_counts()[17].shape[0])
```

```
Out[87]: '8 performers made their first movie at 17'
```

Your output should look like this:

8 performers made their first movie at 17

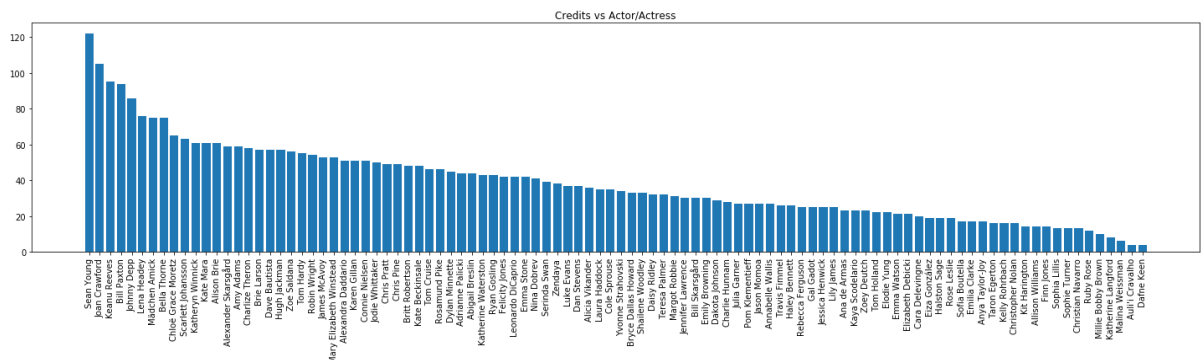
2.7.2 How many performers started as child actors? Define child actor as a person less than 12 years old.

```
In [88]: # your code here
"{0} performers started as child actors".format(df['age'].value_counts()[<12].shape[0])
```

```
Out[88]: '20 performers started as child actors'
```

2.8 Make a plot of the number of credits versus the name of actor/actress.

```
In [89]: # your code here
df = df.sort_values(by='credits', ascending=False)
plt.figure(figsize=(25,5))
plt.bar(df['name'], df['credits'])
plt.xticks(rotation='vertical')
plt.yticks(rotation='horizontal')
plt.title('Credits vs Actor/Actress')
plt.show()
```



```
In [90]: # your code here
#TO: using groupby is a suggestion, i think using fitlers as done below
is much easier and cleaner
prolific = df[df['credits']==max(df['credits'])]
'{0} is the most prolific actor with {0} credits'.format(prolific['name'].values[0], prolific['credits'].values[0])
```

```
Out[90]: 'Sean Young is the most prolific actor with 122 credits'
```

Going the Extra Mile

Be sure to complete problems 1 and 2 before tackling this problem...it is worth only 8 points.

Question 3 [8 pts]: Parsing using Regular Expressions (regex)

Even though scraping HTML with regex is sometimes considered bad practice, you are to use python's **regular expressions** to answer this problem. Regular expressions are useful to parse strings, text, tweets, etc. in general (for example, you may encounter a non-standard format for dates at some point). Do not use BeautifulSoup to answer this problem.

3.1 Write a function called `get_pubs` that takes an `.html` filename as an input and returns a string containing the HTML page in this file (see definition below). Call this function using `data/publist_super_clean.html` as input and name the returned string `prof_pubs`.

3.2 Calculate how many times the author named 'C.M. Friend' appears in the list of publications.

3.3 Find all unique journals and copy them in a variable named `journals`.

3.4 Create a list named `pub_authors` whose elements are strings containing the authors' names for each paper.

Hints

- Look for patterns in the HTML tags that reveal where each piece of information such as the title of the paper, the names of the authors, the journal name, is stored. For example, you might notice that the journal name(s) is contained between the `<l>` HTML tag.
- Learning about your domain is always a good idea: you want to check the names to make sure that they belong to actual journals. Thus, while journal name(s) is contained between the `<l>` HTML tag, please note that *all* strings found between `<l>` tags may not be journal names.
- Each publication has multiple authors.
- C.M. Friend also shows up as Cynthia M. Friend in the file. Count just C. M. Friend.
- There is a comma at the end of the string of authors. You can choose to keep it in the string or remove it and put it back when you write the string as a BibTex entry.
- You want to remove duplicates from the list of journals. Duplicates may also occur due to misspellings or spaces, such as: Nano Lett., and NanoLett. You can assume that any journals with the same initials (e.g., NL for NanoLett.) are the same journal.

Resources

- **Regular expressions:** a) <https://docs.python.org/3.3/library/re.html> (<https://docs.python.org/3.3/library/re.html>), b) <https://regexone.com> (<https://regexone.com>), and c) <https://docs.python.org/3/howto/regex.html> (<https://docs.python.org/3/howto/regex.html>).
- **HTML:** if you are not familiar with HTML see <https://www.w3schools.com/html/> (<https://www.w3schools.com/html/>) or one of the many tutorials on the internet.
- **Document Object Model (DOM):** for more on this programming interface for HTML and XML documents see https://www.w3schools.com/js/js_htmldom.asp (https://www.w3schools.com/js/js_htmldom.asp).

Answers

3.1 Write a function called `get_pubs` that takes an `.html` filename as an input and returns a string ...

```
In [91]: # first import the necessary reg expr library  
import re
```

```
In [92]: # use this file provided  
PUB_FILENAME = 'data/publist_super_clean.html'
```

```
In [93]: # your code here  
# TO to MK: I think they explicitly require us to u regex, that's why we  
import r in 3.1, I figured out a way to make it work  
def get_pubs(filename):  
    with open(filename, 'r') as f:  
        content = f.read()  
        data = re.sub(r'</?w+s+[\^]*>', '', content)  
    return data
```

```
In [94]: # your code here  
prof_pubs = get_pubs(PUB_FILENAME)
```

```
In [95]: # checking your code  
# clear/remove output before creating the pdf file  
#print(prof_pubs)
```

You should see an HTML page that looks like this (colors are not important)

```
<LI>
<A HREF="Papers/2011/PhysRevB_84_125411_2011.pdf" target="paper244">
&quot;Approaching the intrinsic band gap in suspended high-mobility graphene
  nanoribbons&quot;</A>
<BR>Ming-Wei Lin, Cheng Ling, Luis A. Agapito, Nicholas Kioussis, Yiyang Zha
ng, Mark Ming-Cheng Cheng,
<I>PHYSICAL REVIEW B </I> <b>84</b>, 125411 (2011)
<BR>
</LI>
</OL>

<OL START=243>
<LI>
<A HREF="Papers/2011/PhysRevB_84_035325_2011.pdf" target="paper243">
&quot;Effect of symmetry breaking on the optical absorption of semiconductor
  nanoparticles&quot;</A>
<BR>JAdam Gali, Efthimios Kaxiras, Gergely T. Zimanyi, Sheng Meng,
<I>PHYSICAL REVIEW B </I> <b>84</b>, 035325 (2011)
<BR>
</LI>
</OL>

<OL START=242>
<LI>
<A HREF="Papers/2011/PhysRevB_83_054204_2011.pdf" target="paper242">
&quot;Influence of CH2 content and network defects on the elastic properties
  of organosilicate glasses&quot;</A>
<BR>Jan M. Knaup, Han Li, Joost J. Vlassak, and Efthimios Kaxiras,
<I>PHYSICAL REVIEW B </I> <b>83</b>, 054204 (2011)
<BR>
</LI>
</OL>
```


3.2 Calculate how many times the author ...

```
In [96]: # your code here
len(re.findall(r'(C.M. Friend)|(Cynthia M. Friend)', prof_pubs))
```

Out[96]: 8

3.3 Find all unique journals and copy ...

```
In [97]: # your code here
journals = re.findall(r'<i>(.*?) </i>', prof_pubs, re.IGNORECASE)
journals = set(journals)
journals.remove('NanoLett.')
journals.remove('New J. Phys.')
len(journals)
```

Out[97]: 27

```
In [98]: # check your code
#journals
```

Your output should look like this (no duplicates):

```
{'2010 ACM/IEEE International Conference for High Performance',  
 'ACSNano.',  
 'Ab initio',  
 'Acta Mater.',  
 'Catal. Sci. Technol.',  
 'Chem. Eur. J.',  
 'Comp. Phys. Comm.',  
 'Concurrency Computat.: Pract. Exper.',  
 'Energy & Environmental Sci.',  
 'Int. J. Cardiovasc. Imaging',  
 'J. Chem. Phys.',  
 'J. Chem. Theory Comput.',  
 'J. Phys. Chem. B',  
 'J. Phys. Chem. C',  
 'J. Phys. Chem. Lett.',  
 'J. Stat. Mech: Th. and Exper.',  
 'Langmuir',  
 'Molec. Phys.',  
 'Nano Lett.',  
 'New Journal of Physics',  
 'PHYSICAL REVIEW B',  
 'Phil. Trans. R. Soc. A',  
 'Phys. Rev. E - Rap. Comm.',  
 'Phys. Rev. Lett.',  
 'Sci. Model. Simul.',  
 'Sol. St. Comm.',  
 'Top. Catal.'}
```

3.4 Create a list named `pub_authors...`

```
In [99]: # your code here
pub_authors = re.findall(r'<br>\s?(.*)', prof_pubs, re.IGNORECASE)
set(pub_authors)
```

```

Out[99]: {'A. Gali and E. Kaxiras',
'A. Gali, E. Janzen, P. Deak, G. Kresse and E. Kaxiras',
'A. Peters, S. Melchionna, E. Kaxiras, J. Latt, J. Sircar, S. Succi',
'Bingjun Xu, Jan Haubrich, Thomas A. Baker, Efthimios Kaxiras, and Cyn-
thia M. Friend',
'C.E. Lekka, J. Ren, S. Meng and E. Kaxiras',
'C.L. Chang, S.K.R.S. Sankaranarayanan, D. Ruzmetov, M.H. Engelhard,
E. Kaxiras and S. Ramanathan',
'E. Kaxiras and S. Succi',
'E. Manousakis, J. Ren, S. Meng and E. Kaxiras',
'E.M. Kotsalis, J.H. Walther, E. Kaxiras and P. Koumoutsakos',
'F.J. Rybicki, S. Melchionna, D. Mitsouras, A.U. Coskun, A.G. Whitmor-
e, E. Kaxiras, S. Succi, P.H. Stone and C.L. Feldman',
'H. Chen, W.G. Zhu, E. Kaxiras, and Z.Y. Zhang',
'H. Li, J.M. Knaup, E. Kaxiras and J.J. Vlassak',
'H.P. Chen, R.K. Kalia, E. Kaxiras, G. Lu, A. Nakano, K. Nomura',
'J R Maze, A Gali, E Togan, Y Chu, A Trifonov',
'J. Ren, E. Kaxiras and S. Meng',
'JAdam Gali, Efthimios Kaxiras, Gergely T. Zimanyi, Sheng Meng',
'Jan Haubrich, Efthimios Kaxiras, and Cynthia M. Friend',
'Jan M. Knaup, Han Li, Joost J. Vlassak, and Efthimios Kaxiras',
'Jun Ren, Sheng Meng, Yi-Lin Wang, Xu-Cun Ma, Qi-Kun Xue, Efthimios Ka-
xiras',
'Kejie Zhao, Wei L. Wang, John Gregoire, Matt Pharr, Zhigang Suo',
'L.A. Agapito, N. Kioussis and E. Kaxiras',
'M. Bernaschi, M. Fatica, S. Melchionna, S. Succi and E. Kaxiras',
'M. Bernaschi, S. Melchionna, S. Succi, M. Fyta',
'M. Fyta, S. Melchionna, M. Bernaschi, E. Kaxiras and S. Succi',
'Martin Heiss, Sonia Conesa-Boj, Jun Ren, Hsiang-Han Tseng, Adam Gal-
i',
'Masataka Katono, Takeru Bessho, Sheng Meng, Robin Humphry-Baker, Guid-
o Rothenberger',
'Ming-Wei Lin, Cheng Ling, Luis A. Agapito, Nicholas Kioussis, Yiyang
Zhang, Mark Ming-Cheng Cheng',
'S. Melchionna, M. Bernaschi, M. Fyta, E. Kaxiras and S. Succi',
'S. Melchionna, M. Bernaschi, S. Succi, E. Kaxiras, F.J. Rybicki, D. M-
itsouras, A.U. Coskun and C.L. Feldman',
'S. Meng and E. Kaxiras',
'S.K.R.S. Sankaranarayanan, E. Kaxiras and S. Ramanathan',
'S.K.R.S. Sankaranarayanan, E. Kaxiras, S. Ramanathan',
'Sheng Meng, Efthimios Kaxiras, Md. K. Nazeeruddin, and Michael Gratze-
l',
'Simone Melchionna, Efthimios Kaxiras, Massimo Bernaschi and Sauro Suc-
ci',
'T.A. Baker, B.J. Xu, X.Y. Liu, E. Kaxiras and C.M. Friend',
'T.A. Baker, C.M. Friend and E. Kaxiras',
'T.A. Baker, E. Kaxiras and C.M. Friend',
'Thomas A. Baker, Bingjun Xu, Stephen C. Jensen, Cynthia M. Friend and
Efthimios Kaxiras',
'Thomas D. Kuhne, Tod A. Pascal, Efthimios Kaxiras, and Yousung Jung',
'W.L. Wang and E. Kaxiras',
'W.L. Wang, O.V. Yazyev, S. Meng and E. Kaxiras',
'Youdong Mao, Wei L. Wang, Dongguang Wei, Efthimios Kaxiras, and Josep-
h G. Sodroski'}

```

```
In [100]: # check your code: print the list of strings containing the author(s)' names
          # for item in pub_authors:
          #     print (item + ',')
```

Your output should look like this (a line for each paper's authors string of names)

```
Ming-Wei Lin, Cheng Ling, Luis A. Agapito, Nicholas Kioussis, Yiyang Zhang,
Mark Ming-Cheng Cheng,
JAdam Gali, Efthimios Kaxiras, Gergely T. Zimanyi, Sheng Meng,
Jan M. Knaup, Han Li, Joost J. Vlassak, and Efthimios Kaxiras,
Martin Heiss, Sonia Conesa-Boj, Jun Ren, Hsiang-Han Tseng, Adam Gali,
```

```
...
```

```
T.A. Baker, C.M. Friend and E. Kaxiras,
T.A. Baker, C.M. Friend and E. Kaxiras,
E. Kaxiras and S. Succi,
E. Manousakis, J. Ren, S. Meng and E. Kaxiras,
```