

```
In [1]: from bs4 import BeautifulSoup, NavigableString, Tag
        from datascience import *
        from collections import Counter
```

```
In [2]: data = Table.read_table('scripts_metadata.csv')
        data.show(5)
```

title	Genres	Average user rating	IMSDb rating	IMSDb opinion	Script Date	Movie Release Date	Writers	Submi
10 Things I Hate About You Script	Comedy;Romance;	(8.76 out of 10)	(7 out of 10)	A better- than- most teen film.	: November 1997	nan	Karen McCullah Lutz;Kirsten Smith;William Shakespeare;	
12 Script	Comedy;Read "12" Script;	None available	Not available	None available	nan	nan	Lawrence Bridges;	
12 and Holding Script	Drama;	(7.00 out of 10)	Not available	None available	: April 2004	: May 2006	Anthony Cipriano;	
12 Monkeys Script	Drama;Sci- Fi;Thriller;	(9.25 out of 10)	Not available	None available	: June 1994	nan	David Peoples;Janet Peoples;	
12 Years a Slave Script	Drama;	None available	Not available	None available	nan	: November 2013	John Ridley;	: XXyTi

... (1166 rows omitted)

```
In [3]: data = data.where('script_path', are.not_equal_to('nan'))
data.show(5)
```

title	Genres	Average user rating	IMSDb rating	IMSDb opinion	Script Date	Movie Release Date	Writers	Submi
10 Things I Hate About You Script	Comedy;Romance;	(8.76 out of 10)	(7 out of 10)	A better- than- most teen film.	: November 1997	nan	Karen McCullah Lutz;Kirsten Smith;William Shakespeare;	
12 Script	Comedy;Read "12" Script;	None available	Not available	None available	nan	nan	Lawrence Bridges;	
12 and Holding Script	Drama;	(7.00 out of 10)	Not available	None available	: April 2004	: May 2006	Anthony Cipriano;	
12 Monkeys Script	Drama;Sci- Fi;Thriller;	(9.25 out of 10)	Not available	None available	: June 1994	nan	David Peoples;Janet Peoples;	
12 Years a Slave Script	Drama;	None available	Not available	None available	nan	: November 2013	John Ridley;	: XXyTi

... (1138 rows omitted)

```

In [4]: ## make an empty ditionary then append everthing to it
all_scripts = {}

for fname in data['script_path']:

    print(fname)
    with open(fname, 'r') as f:
        raw = f.read()
        soup = BeautifulSoup(raw, 'html5lib')

    try:
        bolded = soup.find('td', {'class': 'scrtext'}) .find_all('b') #find
        text = soup.find('td', {'class': 'scrtext'}) .text
        b_text = [b.text.strip() for b in bolded]
        bolded_text = [b for b in b_text if len(b) > 0]
        sift_out = ['INT.', "EXT.", "-"] #differenetiaste between scene cues
        characters = []
        scenes = []
        for c in bolded_text:
            character = True
            for s in sift_out:
                if s in c:
                    character = False
            if character == True:
                characters.append(c)
            elif len(c) > 4:
                scenes.append(c)

        characters = [c[0] for c in Counter(characters).most_common() if c[1]
        scenes.extend([c[0] for c in Counter(characters).most_common() if c[1]

        movie_name = fname.split('/')[ -1 ][ :-5 ].replace(' Script', '')

        all_scripts[movie_name] = {}
        all_scripts[movie_name]['cast'] = characters
        all_scripts[movie_name]['scenes'] = scenes
        all_scripts[movie_name]['text'] = text

    except:
        pass

```

```

scripts/10 Things I Hate About You Script.html
scripts/12 Script.html
scripts/12 and Holding Script.html
scripts/12 Monkeys Script.html
scripts/12 Years a Slave Script.html
scripts/127 Hours Script.html
scripts/1492: Conquest of Paradise Script.html
scripts/15 Minutes Script.html
scripts/17 Again Script.html

```

```
-----  
--  
KeyboardInterrupt                                Traceback (most recent call las  
t)  
<ipython-input-4-2a6a504be4e6> in <module>()
```

```
In [5]: all_scripts.keys()
```

```
Out[5]: dict_keys(['10 Things I Hate About You', '12', '12 and Holding', '12 Monk  
eys', '12 Years a Slave', '127 Hours', '1492: Conquest of Paradise', '15  
Minutes'])
```

```
In [6]: import re
```

```
for i in all_scripts.keys():  
    scene_index_list = []  
    for scene in set(all_scripts[i]['scenes']):  
        print(i, scene)  
        indices = [m.start() for m in re.finditer(scene, all_scripts[i]['text'])]  
        scene_index_list.extend(indices)
```

```
10 Things I Hate About You STRATFORD HOUSE - SUNSET  
10 Things I Hate About You EXT. DOWNTOWN STREET - NIGHT  
10 Things I Hate About You INT. STRATFORD HOUSE - DAY  
10 Things I Hate About You INT. DETENTION HALL - DAY  
10 Things I Hate About You INT. BOOK STORE - DAY  
10 Things I Hate About You INT. PROM - NIGHT  
10 Things I Hate About You INT. SOPHOMORE ENGLISH CLASS - DAY  
10 Things I Hate About You EXT. FIELD HOCKEY FIELD - DAY  
10 Things I Hate About You INT. DIVE BAR - NIGHT  
10 Things I Hate About You INT. MATH CLASS - DAY  
10 Things I Hate About You PADUA HIGH SCHOOL - DAY  
10 Things I Hate About You EXT. HOTEL PARKING LOT - NIGHT  
10 Things I Hate About You BOGEY'S KITCHEN - NIGHT  
10 Things I Hate About You INT. STRATFORD HOUSE/DEN - DAY  
10 Things I Hate About You INT. STUDY HALL - DAY  
10 Things I Hate About You EXT. MISS PERKY'S OFFICE - DAY  
10 Things I Hate About You INT. ENGLISH CLASS - DAY  
10 Things I Hate About You INT. KENNY'S THAI FOOD DINER - DAY  
10 Things I Hate About You INT. CAFETERIA - DAY  
10 Things I Hate About You HOTEL - NIGHT
```

```
In [8]: from nltk.util import ngrams
```

```
for i in all_scripts.keys():  
    scene_texts = []  
    for n in ngrams(sorted(scene_index_list), 2):  
        scene_texts.append(all_scripts[i]['text'][n[0]:n[1]])
```

```
In [18]: def make_graph(cast_dict):  
    # setup graph object  
    G = nx.Graph()  
  
    # add nodes with attributes of number of lines and scenes  
    for c in cast_dict.keys():  
        G.add_node(  
            c,  
            scenes = cast_dict[c]  
        )  
  
    # make edges by iterating over all combinations of nodes  
    for (node1, data1), (node2, data2) in itertools.combinations(G.nodes(data=True)):  
        # count scenes together by getting union of their sets  
        scenes_together = len(set(data1['scenes']) & set(data2['scenes']))  
        cast_dict[c]  
  
        if scenes_together:  
            # add more weight for more scenes together  
            G.add_edge(node1, node2, weight=scenes_together)  
  
    return G
```

```

In [24]: for i in all_scripts.keys():
        for i in all_scripts.keys():
            cast_dict = {}
            for c in all_scripts[i]['cast']:
                cast_dict[c] = []
                for i, scene in enumerate(scene_texts):
                    if scene.count(c) > 0:
                        cast_dict[c].append(i)

        G = make_graph(cast_dict)

        node_color = 'blue'

        plt.figure(figsize=(13,8)) # make the figure size a little larger
        plt.axis('off') # remove the axis, which isn't meaningful in this c
        plt.title(i, fontsize=20)

        # The 'k' argument determines how spaced out the nodes will be from
        # one another on the graph.
        pos = nx.spring_layout(G, k=0.5)

        nx.draw_networkx(
            G,
            pos=pos,
            node_color=node_color,
            edge_color='gray', # change edge color
            alpha=0.3, # make nodes more transparent to make labels clearer
            font_size=14,
        )

```

/srv/app/venv/lib/python3.6/site-packages/matplotlib/pyplot.py:524: RuntimeWarning: More than 20 figures have been opened. Figures created through the pyplot interface (`matplotlib.pyplot.figure`) are retained until explicitly closed and may consume too much memory. (To control this warning, see the rcParam `figure.max_open_warning`).

max_open_warning, RuntimeWarning)

```

In [17]: import numpy as np
import networkx as nx
from lxml import etree
import itertools
from datascience import *
import matplotlib.pyplot as plt

for i in all_scripts.keys():
    cast_dict = {}
    for c in all_scripts[i]['cast']:
        cast_dict[c] = []
        for i, scene in enumerate(scene_texts):
            if scene.count(c) > 0:
                cast_dict[c].append(i)

G = make_graph(cast_dict)

network_tab = Table()
network_tab.append_column(label="Characters", values=[c for c in sorted(
dc = [x[1] for x in sorted(nx.degree_centrality(G).items(), key=lambda x:
network_tab.append_column(label="Degree Centrality", values=dc)
network_tab.show()

```

Characters	Degree Centrality
BIANCA	0
BRUCE	0
CAMERON	0
CHASTITY	0
JOEY	0
KAT	0
MANDELLA	0
MICHAEL	0
MISS PERKY	0
MRS. BLAISE	0
PATRICK	0
SHARON	0

In [61]:

```
G = make_graph(cast_dict)

node_size = 0.5
node_color = 'blue'

plt.figure(figsize=(13,8)) # make the figure size a little larger
plt.axis('off') # remove the axis, which isn't meaningful in this case
plt.title("10 Things I Hate About You", fontsize=20)

# The 'k' argument determines how spaced out the nodes will be from
# one another on the graph.
pos = nx.spring_layout(G, k=0.5)

nx.draw_networkx(
    G,
    pos=pos,
    node_size=node_size,
    node_color=node_color,
    edge_color='gray', # change edge color
    alpha=0.3, # make nodes more transparent to make labels clearer
    font_size=14,
)
network_tab = Table()
network_tab.append_column(label="Characters", values=[c for c in sorted(
network_tab.show()
```

File "<ipython-input-61-9aa8a7193e12>", line 2

```
G = make_graph(cast_dict)
```

^

IndentationError: unexpected indent


```
In [40]: cast_dict
```

```
Out[40]: {'ALONSO': [],  
          'AROJAZ': [],  
          'BARTOLOME': [],  
          'BEATRIX': [],  
          'BOBADILLA': [],  
          'BROTHER BUYL': [],  
          'COLUMBUS': [],  
          'CUT TO:': [],  
          'DIEGO': [],  
          'DISSOLVE TO:': [],  
          'FERNANDO': [],  
          'GIACOMO': [],  
          'ISABEL': [],  
          'MARCHENA': [],  
          'MENDEZ': [],  
          'MOXICA': [],  
          'PINZON': [],  
          'SAILOR': [],  
          'SANCHEZ': [],  
          'SANTANGEL': [],  
          'UTAPAN': []}
```

```

In [18]: import numpy as np
import networkx as nx
from lxml import etree
import itertools
from datascience import *
import matplotlib.pyplot as plt

def make_graph(cast_dict):
    '''
    This function accepts a dictionary with number of lines and scenes to create a
    NetworkX graph object
    '''
    # setup graph object
    G = nx.Graph()

    # add nodes with attributes of number of lines and scenes
    for c in cast_dict.keys():
        G.add_node(
            c,
            scenes = cast_dict[c]
        )

    # make edges by iterating over all combinations of nodes
    for (node1, data1), (node2, data2) in itertools.combinations(G.nodes(data=True)):
        # count scenes together by getting union of their sets
        scenes_together = len(set(data1['scenes']) & set(data2['scenes']))
        cast_dict[c]

        if scenes_together:
            # add more weight for more scenes together
            G.add_edge(node1, node2, weight=scenes_together)

    return G

```

```

In [19]: G = make_graph(cast_dict)

```

```
In [20]: import numpy as np
import networkx as nx
from lxml import etree
import itertools
from datascience import *
import matplotlib.pyplot as plt

node_size = 0.5
node_color = 'blue'

plt.figure(figsize=(13,8)) # make the figure size a little larger
plt.axis('off') # remove the axis, which isn't meaningful in this case
plt.title("10 Things I Hate About You", fontsize=20)

# The 'k' argument determines how spaced out the nodes will be from
# one another on the graph.
pos = nx.spring_layout(G, k=0.5)

nx.draw_networkx(
    G,
    pos=pos,
    node_size=node_size,
    node_color=node_color,
    edge_color='gray', # change edge color
    alpha=0.3, # make nodes more transparent to make labels clearer
    font_size=14,
)
```

```
In [21]: network_tab = Table()
network_tab.append_column(label="Characters", values=[c for c in sorted(cast
network_tab.show())
```

Characters
BIANCA
BRUCE
CAMERON
CHASTITY
JOEY
KAT
MANDELLA
MICHAEL
MISS PERKY
MRS. BLAISE
PATRICK
SHARON
WALTER

```
In [23]: dc = [x[1] for x in sorted(nx.degree centrality(G).items(), key=lambda x: x[1])]
network_tab.append_column(label="Degree Centrality", values=dc)
network_tab.show()
```

Characters	Degree Centrality
BIANCA	0.833333
BRUCE	0.25
CAMERON	0.833333
CHASTITY	0.5
JOEY	0.833333
KAT	1
MANDELLA	0.666667
MICHAEL	0.666667
MISS PERKY	0.416667
MRS. BLAISE	0.25
PATRICK	0.833333
SHARON	0.416667
WALTER	0.5

```
In [24]: bc = [x[1] for x in sorted(nx.betweenness centrality(G).items(), key=lambda x: x[1])]
network_tab.append_column(label="Betweenness Centrality", values=bc)
network_tab.show()
```

Characters	Degree Centrality	Betweenness Centrality
BIANCA	0.833333	0.0454545
BRUCE	0.25	0
CAMERON	0.833333	0.0454545
CHASTITY	0.5	0
JOEY	0.833333	0.0671717
KAT	1	0.159091
MANDELLA	0.666667	0.030303
MICHAEL	0.666667	0.0123737
MISS PERKY	0.416667	0
MRS. BLAISE	0.25	0
PATRICK	0.833333	0.0916667
SHARON	0.416667	0
WALTER	0.5	0.0030303

```
In [25]: ec = [x[1] for x in sorted(nx.eigenvector_centrality(G).items(), key=lambda
network_tab.append_column(label="Eigenvector Centrality", values=ec)
network_tab.show()
```

Characters	Degree Centrality	Betweenness Centrality	Eigenvector Centrality
BIANCA	0.833333	0.0454545	0.413741
BRUCE	0.25	0	0.0208809
CAMERON	0.833333	0.0454545	0.385503
CHASTITY	0.5	0	0.115439
JOEY	0.833333	0.0671717	0.304199
KAT	1	0.159091	0.49326
MANDELLA	0.666667	0.030303	0.181087
MICHAEL	0.666667	0.0123737	0.309785
MISS PERKY	0.416667	0	0.0908165
MRS. BLAISE	0.25	0	0.0384913
PATRICK	0.833333	0.0916667	0.417197
SHARON	0.416667	0	0.0626333
WALTER	0.5	0.0030303	0.118897

```
In [26]: def gini(array):
        """Calculate the Gini coefficient of a numpy array."""
        # https://github.com/oliviaguest/gini
        array = np.sort(array) # values must be sorted
        index = np.arange(1, array.shape[0] + 1) # index per array element
        n = array.shape[0] # number of array elements
        return ((np.sum((2 * index - n - 1) * array)) / (n * np.sum(array))) #C
```

```
In [27]: gini(network_tab.column('Eigenvector Centrality'))
```

```
Out[27]: 0.39558396783323707
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [87]: 'hello'.find('e')
```

```
Out[87]: 1
```

```
In [7]: soup = BeautifulSoup(raw, 'html5lib')
```

```
In [8]: bolded = soup.find('td', {'class': 'scrtext'}).find_all('b')
```

```
In [9]: b_text = [b.text.strip() for b in bolded]
```

```
In [10]: bolded_text = [b for b in b_text if len(b) > 0]
```

```
In [11]: sift_out = ['INT.', "EXT.", "-"]
```

```
characters = []
for c in bolded_text:
    character = True
    for s in sift_out:
        if s in c:
            character = False

    if character == True:
        characters.append(c)
```

```
In [12]: from collections import Counter
```

```
In [13]: [c[0] for c in Counter(characters).most_common() if c[1] > 5]
```

```
Out[13]: ['KAT',
          'PATRICK',
          'BIANCA',
          'CAMERON',
          'MICHAEL',
          'JOEY',
          'WALTER',
          'MANDELLA',
          'MISS PERKY',
          'MRS. BLAISE',
          'CHASTITY',
          'SHARON',
          'BRUCE']
```

```
In [ ]:
```

```
In [ ]:
```