

Get the Data

```
In [1]: # load packages
from csv import reader, get_url
import pandas as pd

# store text file
txt = read('resources/data/shakespeare/shakespeare.txt',
           repo='https://github.com/TLP-COL/text-data-course')

# view
print(txt[:250])

First Citizen:
Before we proceed any further, hear me speak.

All:
Speak, speak.

First Citizen:
You are all resolved rather to die than to famish?

All:
Resolved, resolved.

First Citizen:
First, you know Caius Marcius is chief enemy to the people.
```

Part 1

Split the text file into a table, such that

- each row is a single line of dialogue
- there are columns for
  - the speaker
  - the line number
  - the line dialogue (the text)

Hint: you will need to use `Regex` to do this rapidly. See the in-class "markdown" example!

Question(s):

What assumptions have you made about the text that allowed you to do this?

Method 1 - writing a `Regex` pattern

- The speaker name always starts with a capital letter
- Speaker names are always followed by a colon.
- Speaker names are always followed by a line break after the colon and there should be a colon at the end
- The dialogue is always between 'speaker name' and a double line break (`\n\n`)

```
In [2]: # write regex pattern
pat = re.compile(
    # use a non-capturing group to find the start of text file or a double line break
    "(?!(?:[a-z]+[ ]+)*(?:[a-z]+[ ]+)*[a-z]+[ ]+)"
    # match speaker names i.e. first letter of the first word/name should be capital and there should be a colon at the end
    "(?!(?:[a-z]+[ ]+)*[a-z]+[ ]+)"
    # go to the next line, match any character (including line breaks) and stop searching when condition is satisfied
    "(?!(?:[a-z]+[ ]+)*[a-z]+[ ]+)"
    # match till two line breaks or the end of the file
    "(?!(?:[a-z]+[ ]+)*[a-z]+[ ]+)"
    # enable the pattern to handle an input string that consists of multiple lines
    flags=re.M
)

In [3]: # store matches
matches = pat.findall(txt)

# store matches in dataframe
df1 = pd.DataFrame.from_records(matches, columns=['speaker', 'dialogue'])

# add line number column starting from 1
df1['line_number'] = df1.index + 1

# view dataframe
df1.head()
```

```
Out[3]:
```

	speaker	dialogue	line_number
0	First Citizen	Before we proceed any further, hear me speak.	1
1	All	Speak, speak.	2
2	First Citizen	You are all resolved rather to die than to fa...	3
3	All	Resolved, resolved.	4
4	First Citizen	First, you know Caius Marcius is chief enemy...	5

I also tried the following, but realized that one of my assumptions (#2) was incorrect as there were speakers called 'Ghosts of young Princes' and 'Ghost of GREY'

Assumptions:

1. The speaker name always starts with a capital letter.
2. If there is more than one word making up a speaker name e.g. First Citizen, all words start with a capital letter.
3. Speaker names are always followed by a colon.
4. Speaker names are always followed by a line break after the colon and there should be a colon at the end
5. The dialogue is always after 'speaker name' and

`{[A-Z][A-Z-]*}(?:\s+([A-Z][A-Z-]*)(?=[\n\S]))` tried this in <https://regex.com/>

Method 2 - using the `split` function from `re`

The text is structured as follows:

First Speaker\nText\nNext Speaker

Assumptions:

- The speaker's name is always followed by a colon and a line break
- The text (speaker's dialogue) comes after the colon and line break (after the speaker's name) and ends with double line breaks (`\n\n`)

Extract each dialogue spoken and the speaker's name i.e. 'speaker, name:dialogue'

```
In [4]: # split on double line breaks
first_split = re.split("\n\n", txt)

# check length of list
len(first_split)
```

```
Out[4]:
```

7222

Separate speaker name and the dialogue

Assumption:

Speaker name is always followed by a colon.

Since each item in `first_split` contains the speaker's name and the dialogue spoken, splitting on the first colon in the string will separate the speaker name and the dialogue.

```
In [5]: # create list to store split results
second_split = []

# iterate through each item in the list and split on the first colon
for line in first_split:
    second_split.append(line.split(":", 1))

# check length of list
len(second_split)
```

```
Out[5]:
```

7222

Store speakers and dialogues

```
In [6]: # create list to store speakers
speakers = []

# iterate through each item in the list and select the speakers
for line in range(len(second_split)):
    speakers.append(second_split[line][0])

# create list to store speakers
speakers_1 = []

# remove \n from speakers
for i in range(len(speakers)):
    speakers_1.append(speakers[i].replace('\n', ''))

In [7]: # create list to store dialogues
dialogues = []

# iterate through each item in the list and select the dialogues
for line in range(len(second_split)):
    dialogues.append(second_split[line][1])

# create list to store dialogues
dialogues_1 = []

# remove \n from dialogues
for i in range(len(dialogues)):
    dialogues_1.append(dialogues[i].replace('\n', ''))

In [8]: # create empty dataframe
df = pd.DataFrame()
```

```
Out[8]:
```

	speaker	dialogue	line_number
0	First Citizen	Before we proceed any further, hear me speak.	1
1	All	Speak, speak.	2
2	First Citizen	You are all resolved rather to die than to fa...	3
3	All	Resolved, resolved.	4
4	First Citizen	First, you know Caius Marcius is chief enemy...	5

Part 2

You have likely noticed that the lines are not all from the same play! Now, we will add some useful metadata to our table:

- Determine which source title for each line
- add 'title' as a 'play' column in the data table.
- make sure to document your decisions, assumptions, external data sources, etc.

This is fairly open-ended, and you are not being judged completely on accuracy. Instead, think outside the box a bit as to how you might accomplish this, and attempt to justify whatever approximations or assumptions you felt were appropriate.

For the source title, I am going to use a Kaggle dataset (<https://www.kaggle.com/kingburn0660/shakespeare-plays/version4>) on Shakespeare's plays (including character names and dialogues).

```
In [9]: # read in Shakespeare play data
df_shakespeare = pd.read_csv('Shakespeare_data.csv')

# view
df_shakespeare.head()
```

```
Out[9]:
```

	DateLine	Play	PlayerLineNumber	ActSceneLine	Player	PlayerLine
0	1	Henry IV	NAN	NAN	NAN	ACT I
1	2	Henry IV	NAN	NAN	NAN	SCENE I. London. The palace.
2	3	Henry IV	NAN	NAN	NAN	Enter KING HENRY, LORD JOHN OF LANCASTER, the ...
3	4	Henry IV	1.0	1.1	KING HENRY IV	So shaken as we are, so wan with care,
4	5	Henry IV	1.0	1.1.2	KING HENRY IV	Find we a time for frighted peace to pant,

```
In [10]: # make speaker names uppercase (in order to match) and store in df
speaker_df = pd.DataFrame(df_shakespeare['Player']).str.upper()
df_shakespeare['Player'] = df_shakespeare['Player'].str.upper()

# rename column to match shakespeare_df column name for speakers
speaker_df.rename(columns = {'speaker': 'Player'}, inplace = True)
```

I will get the list of plays of a specific speaker e.g. all plays with the speaker 'First Citizen'. Then I will select the most common play from this list as the source title. The assumption here is that the play in which a specific speaker/character had a lot of dialogues is the most likely source title for the dialogues by the speaker.

Limitation of this assumption: some dialogues may be from plays where the speaker has a minor part i.e. it may not be the most common play.

```
In [11]: # create empty list to store source titles
plays = []

# store the most common play for each speaker
for i in range(len(speaker_df)):
    plays_list = list(df_shakespeare.loc[df_shakespeare['Player'] == speaker_df.Player[i], 'Play'])
    if len(plays_list) == 0:
        plays.append('none')
    else:
        plays.append(max(set(plays_list), key = plays_list.count))

# check length of plays
len(plays)
```

```
Out[11]:
```

7222

```
In [12]: # add 'play' column
df['play'] = plays

# view dataframe
df.head()
```

```
Out[12]:
```

	speaker	dialogue	line_number	play
0	First Citizen	Before we proceed any further, hear me speak.	1	Coriolanus
1	All	Speak, speak.	2	macbeth
2	First Citizen	You are all resolved rather to die than to fa...	3	Coriolanus
3	All	Resolved, resolved.	4	macbeth
4	First Citizen	First, you know Caius Marcius is chief enemy ...	5	Coriolanus

```
In [13]: # view row where there was no match
df[df['play'] == 'none']
```

```
Out[13]:
```

	speaker	dialogue	line_number	play
530	Senators, &C	Well surely him.	531	none
538	Senators, &C	Weapons, weapons, weapons! "Thou dost" Pa...	539	none
2142	Ghost of Prince Edward		2143	none
2145	Ghost of King Henry VI		2144	none
2150	Ghosts of young Princes		2151	none
2152	Ghost of BUCKINGHAM		2153	none
6669	ALL SERVING-MEN	Here, here, sir, here, sir.	6670	none

Assumption: for speakers that were not found in an external dataset and for which there is no information, the speaker is assumed to be in the same play as the previous speaker. In other words, for rows where there was not an exact match for the speaker, the preceding row will be used to obtain the play/source title.

Limitation: If the play for the first observation was missing, this would have to be modified.

```
In [14]: # condition: select rows where no play was found
# replace with play in the previous row
s = df['play'].eq("none")
df.loc[s, 'play'] = pd.np.nan
df['play'].ffill(inplace=True)

# check - should not return any matches
df[df['play'] == 'none']

C:\Users\marya\AppData\Local\Temp\ipykernel_16988\1586223887.py:5: FutureWarning: The pandas.np module is deprecated and will be removed from pandas
as in a future version. Import numpy directly instead
df.loc[s, 'play'] = pd.np.nan

Out[14]:
```

	speaker	dialogue	line_number	play
530	Senators, &C	Well surely him.	531	none
538	Senators, &C	Weapons, weapons, weapons! "Thou dost" Pa...	539	none
2142	Ghost of Prince Edward		2143	none
2145	Ghost of King Henry VI		2144	none
2150	Ghosts of young Princes		2151	none
2152	Ghost of BUCKINGHAM		2153	none
6669	ALL SERVING-MEN	Here, here, sir, here, sir.	6670	none

Assumption: for speakers that were not found in an external dataset and for which there is no information, the speaker is assumed to be in the same play as the previous speaker. In other words, for rows where there was not an exact match for the speaker, the preceding row will be used to obtain the play/source title.

Limitation: If the play for the first observation was missing, this would have to be modified.

```
In [14]: # condition: select rows where no play was found
# replace with play in the previous row
s = df['play'].eq("none")
df.loc[s, 'play'] = pd.np.nan
df['play'].ffill(inplace=True)

# check - should not return any matches
df[df['play'] == 'none']

C:\Users\marya\AppData\Local\Temp\ipykernel_16988\1586223887.py:5: FutureWarning: The pandas.np module is deprecated and will be removed from pandas
as in a future version. Import numpy directly instead
df.loc[s, 'play'] = pd.np.nan

Out[14]:
```

	speaker	dialogue	line_number	play
530	Senators, &C	Well surely him.	531	none
538	Senators, &C	Weapons, weapons, weapons! "Thou dost" Pa...	539	none
2142	Ghost of Prince Edward		2143	none
2145	Ghost of King Henry VI		2144	none
2150	Ghosts of young Princes		2151	none
2152	Ghost of BUCKINGHAM		2153	none
6669	ALL SERVING-MEN	Here, here, sir, here, sir.	6670	none

Assumption: for speakers that were not found in an external dataset and for which there is no information, the speaker is assumed to be in the same play as the previous speaker. In other words, for rows where there was not an exact match for the speaker, the preceding row will be used to obtain the play/source title.

Limitation: If the play for the first observation was missing, this would have to be modified.

```
In [14]: # condition: select rows where no play was found
# replace with play in the previous row
s = df['play'].eq("none")
df.loc[s, 'play'] = pd.np.nan
df['play'].ffill(inplace=True)

# check - should not return any matches
df[df['play'] == 'none']

C:\Users\marya\AppData\Local\Temp\ipykernel_16988\1586223887.py:5: FutureWarning: The pandas.np module is deprecated and will be removed from pandas
as in a future version. Import numpy directly instead
df.loc[s, 'play'] = pd.np.nan

Out[14]:
```

	speaker	dialogue	line_number	play
530	Senators, &C	Well surely him.	531	none
538	Senators, &C	Weapons, weapons, weapons! "Thou dost" Pa...	539	none
2142	Ghost of Prince Edward		2143	none
2145	Ghost of King Henry VI		2144	none
2150	Ghosts of young Princes		2151	none
2152	Ghost of BUCKINGHAM		2153	none
6669	ALL SERVING-MEN	Here, here, sir, here, sir.	6670	none

Assumption: for speakers that were not found in an external dataset and for which there is no information, the speaker is assumed to be in the same play as the previous speaker. In other words, for rows where there was not an exact match for the speaker, the preceding row will be used to obtain the play/source title.

Limitation: If the play for the first observation was missing, this would have to be modified.

```
In [14]: # condition: select rows where no play was found
# replace with play in the previous row
s = df['play'].eq("none")
df.loc[s, 'play'] = pd.np.nan
df['play'].ffill(inplace=True)

# check - should not return any matches
df[df['play'] == 'none']

C:\Users\marya\AppData\Local\Temp\ipykernel_16988\1586223887.py:5: FutureWarning: The pandas.np module is deprecated and will be removed from pandas
as in a future version. Import numpy directly instead
df.loc[s, 'play'] = pd.np.nan

Out[14]:
```

	speaker	dialogue	line_number	play
530	Senators, &C	Well surely him.	531	none
538	Senators, &C	Weapons, weapons, weapons! "Thou dost" Pa...	539	none
2142	Ghost of Prince Edward		2143	none
2145	Ghost of King Henry VI		2144	none
2150	Ghosts of young Princes		2151	none
2152	Ghost of BUCKINGHAM		2153	none
6669	ALL SERVING-MEN	Here, here, sir, here, sir.	6670	none

Assumption: for speakers that were not found in an external dataset and for which there is no information, the speaker is assumed to be in the same play as the previous speaker. In other words, for rows where there was not an exact match for the speaker, the preceding row will be used to obtain the play/source title.

Limitation: If the play for the first observation was missing, this would have to be modified.

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```


```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```


```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```

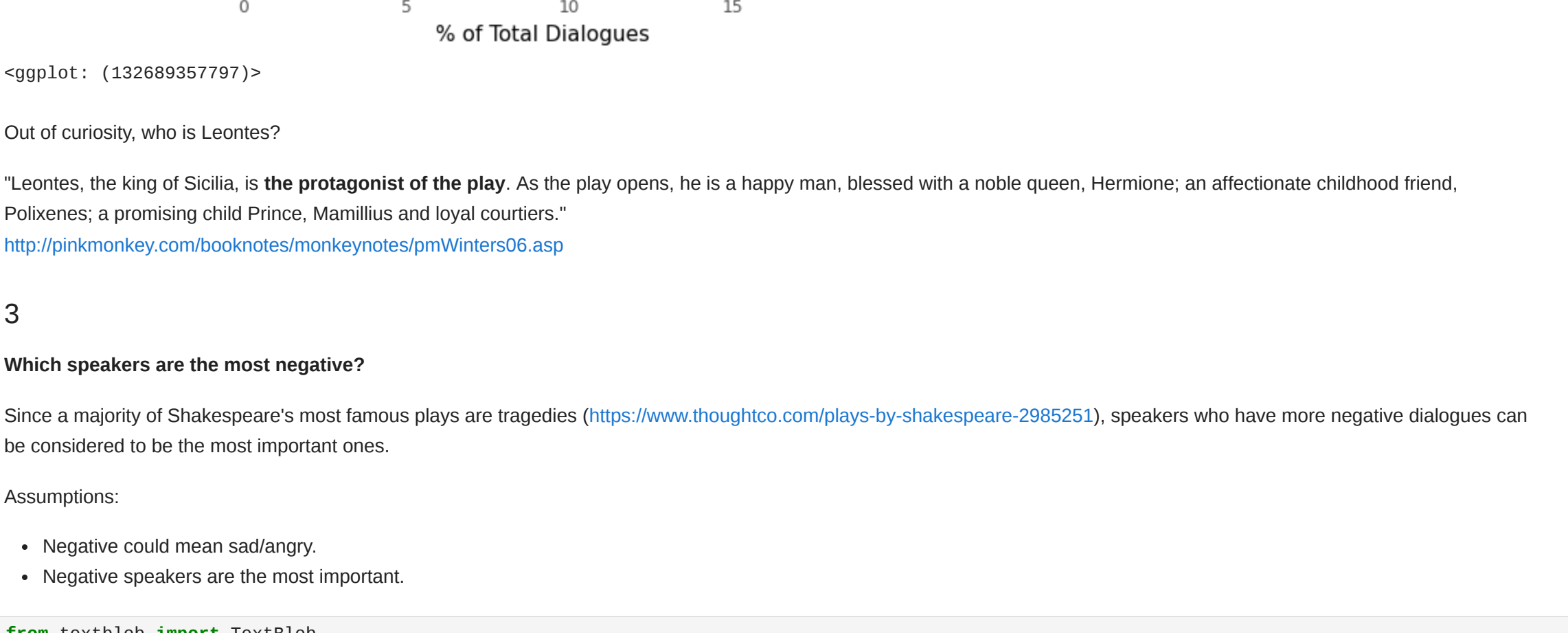
```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```


```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```

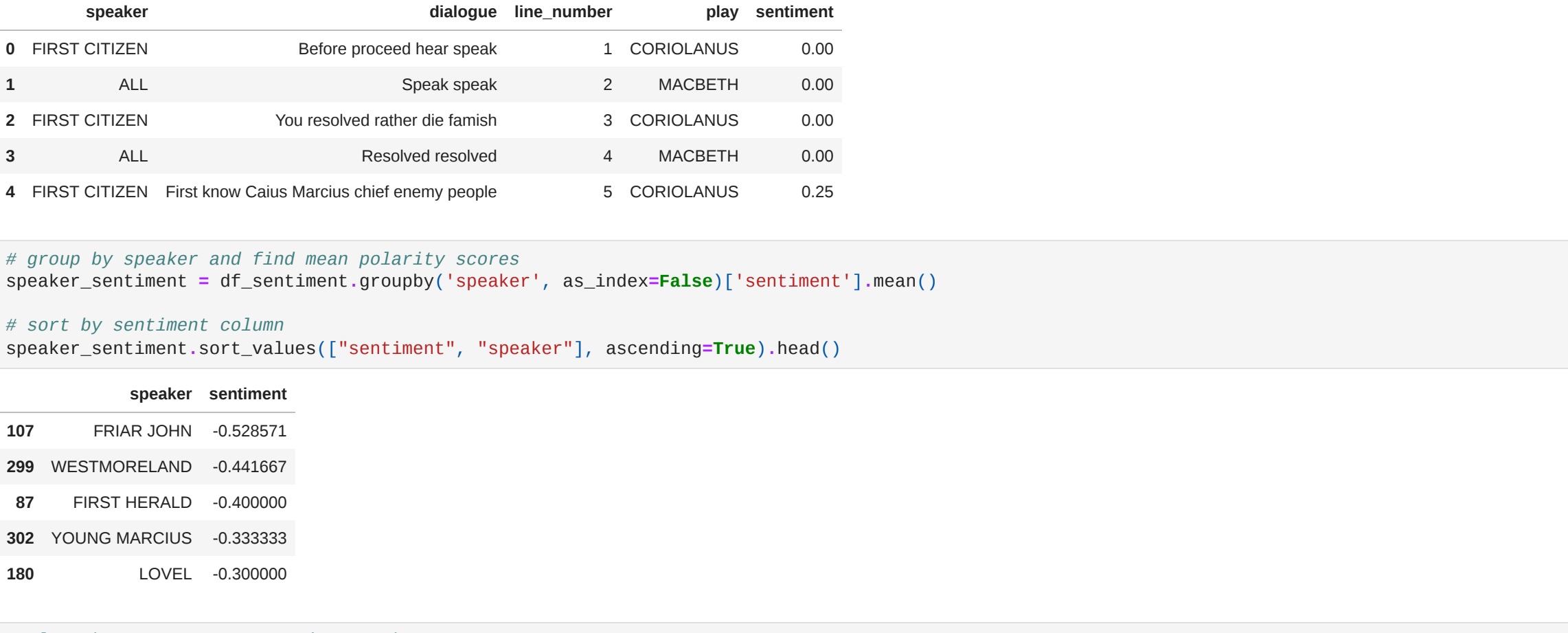
```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```

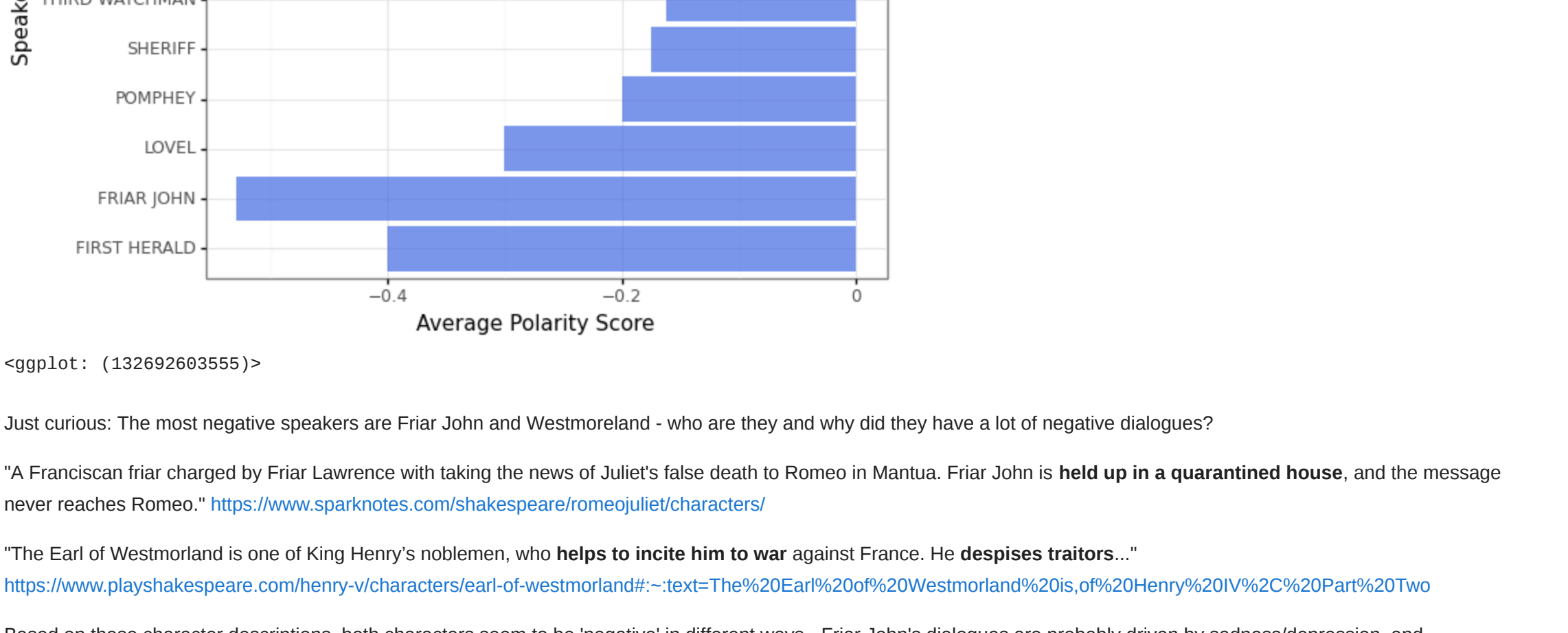
```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```


```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```

```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```

```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```

```
Out[16]:
```

speaker	total_dialogues
0	Coriolanus
1	macbeth
2	Coriolanus
3	macbeth
4	Coriolanus

```
In [17]: # plot the top 10 most important speakers
from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax = df['play'].value_counts().reset_index(name='total_dialogues')[1:10]
ax.set_xlabel('play')
ax.set_ylabel('total_dialogues')
ax.set_title('Speakers With the Most Dialogues')
ax.set_xticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_yticks(['Coriolanus', 'macbeth', 'Coriolanus', 'macbeth', 'Coriolanus'])
ax.set_ylim(0, 200)
ax.set_xlim(0, 200)
```

Speakers With the Most Dialogues



```
Out[17]:
```

<matplotlib.figure.Figure at 132688538647>

```
In [16]: # view speakers and dialogue counts
df.groupby('play').size()
```

```
Out[16]:
```

speaker	total_dialogues
---------	-----------------