Homework: State Your Assumptions

```
POLONIUS\ What do you read, my lord?

HAMLET\ Words, words, words

-- Hamlet, Act 2, Scene 2
```

This homework deals with the assumptions made when taking text from its original "raw" form into something more *computable*.

- Assumptions about the shape of text (e.g. how to break a corpus into documents)
- Assumptions about what makes a token, an entity, etc.
- Assumptions about what interesting or important content looks like, and how that informs our analyses.

There are three parts:

- 1. Splitting Lines from Shakespeare
- 2. Tokenizing and Aligning lines into plays
- 3. Assessing and comparing characters from within each play

NB\ This file is merely a *template*, with instructions; do not feel constrained to using it directly if you do not wish to.

Get the Data

Since the class uses dvc , it is possible to get this dataset either using the command line (e.g. dvc import https://github.com/TLP-COI/text-data-course resources/data/shakespeare/shakespeare.txt), or using the python api (if you wish to use python)

```
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.
```

Make sure this works before you continue! Either way, it would likely be beneficial to have the data downloaded locally to keep from needing to re-dowload it every time.

Part 1

Split the text file into a table, such that

- each row is a single line of dialogue
- there are columns for
 - 1. the speaker
 - 2. the line number
 - 3. 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?

```
In [2]: import pandas as pd
   import numpy as np
   import re

In [3]: # Use the first 250 characters as a sample
   sample = txt[:250]
   sample
```

Out[3]: 'First Citizen:\nBefore we proceed any further, hear me speak.\n\nAll:\nSpeak, speak.\n\nFirst Citizen:\nYou are all resolved rather to die than to famish?\n\nAll:\nResolved.resolved.\n\nFirst Citizen:\nFirst, you know Caius Marcius is chief enemy to the people.\n'

Assumptions:

- the speaker begins with a capitalized letter and ends with a colon
- the dialogue begins with \n (new line) and ends with \n\n

```
In [5]: # Find `patt try` such that:
         matches try = patt try.findall(sample)
         pd.DataFrame.from records(matches try, columns=['speaker', 'dialogue'])
Out[5]:
              speaker
                                                     dialogue
         0 First Citizen Before we proceed any further, hear me speak.
                   ΑII
                                                 Speak, speak.
         2 First Citizen
                       You are all resolved rather to die than to fam...
         3
                   ΑII
                                            Resolved, resolved.
In [6]: # Check the last several lines to see whether they meet the assumptions
         print(txt[-330:])
         SEBASTIAN:
         I do; and surely
         It is a sleepy language and thou speak'st
        Out of thy sleep. What is it thou didst say?
         This is a strange repose, to be asleep
         With eyes wide open; standing, speaking, moving,
         And yet so fast asleep.
         ANTONIO:
        Noble Sebastian,
         Thou let'st thy fortune sleep--die, rather; wink'st
         Whiles thou art waking.
In [7]: | txt[-330:]
```

Out[7]: "SEBASTIAN:\nI do; and surely\nIt is a sleepy language and thou speak'st\nOut of thy sle ep. What is it thou didst say?\nThis is a strange repose, to be asleep\nWith eyes wide o pen; standing, speaking, moving,\nAnd yet so fast asleep.\n\nANTONIO:\nNoble Sebastia n,\nThou let'st thy fortune sleep--die, rather; wink'st\nWhiles thou art waking.\n"

Here, the last line will be omitted according to the current assumptions because it ends with \n instead of \n\n.

```
In [8]: # Fix the last line by adding an extra \n
new_line = '\n'
txt = txt + new_line
```

Another particular case that doesn't meet the assumptions is empty dialogue.

```
In [9]: # Example of empty dialogue
  print(txt[11225:11536])
```

```
First Senator:
Your company to the Capitol; where, I know,
Our greatest friends attend us.

TITUS:

COMINIUS:
Noble Marcius!

First Senator:

MARCIUS:
Nay, let them follow:
The Volsces have much corn; take these rats thither
To gnaw their garners. Worshipful mutiners,
Your valour puts well forth: pray, follow.
```

My assumption is that the dialogue begins with \n and ends with \n (a total of three \n 's). However, an empty dialoge comes with only two \n 's.

```
In [10]: # Fix empty dialogue by adding an extra \n
          txt = txt.replace(':\n\n', ':\n\n\n')
In [11]: patt = re.compile(
              "(^[A-Z].+?):$" # the speaker
              "\n{1}(.*?)\n{2}", # the dialogue
              flags = re.S | re.M
In [12]: # Find `patt` such that:
          matches = patt.findall(txt)
          df = pd.DataFrame.from records(matches, columns=['speaker', 'dialogue'])
          df.tail()
Out[12]:
                  speaker
                                                        dialogue
          7217
                 ANTONIO
                            Nor I; my spirits are nimble.\nThey fell toget...
```

```
    7217 ANTONIO Nor I; my spirits are nimble.\nThey fell toget...
    7218 SEBASTIAN What, art thou waking?
    7219 ANTONIO Do you not hear me speak?
    7220 SEBASTIAN I do; and surely\nIt is a sleepy language and ...
    7221 ANTONIO Noble Sebastian,\nThou let'st thy fortune slee...
```

```
In [13]: # Create a new column for line number
df['line'] = df.index + 1
df
```

Out[13]:	speaker		dialogue				
	0 First Citizen		Before we proceed any further, hear me speak.				
	1	All	Speak, speak.				
	2 First Citizen		You are all resolved rather to die than to fam				
	3	All	Resolved. resolved.				
	4 First Citizen		First, you know Caius Marcius is chief enemy t	5			
	7217	ANTONIO	Nor I; my spirits are nimble.\nThey fell toget	7218			
	7218	SEBASTIAN	What, art thou waking?	7219			
	7219 ANTONIO		Do you not hear me speak?	7220			
	7220	SEBASTIAN	I do; and surely\nIt is a sleepy language and	7221			
	7221	ANTONIO	Noble Sebastian,\nThou let'st thy fortune slee	7222			
	7222 rows × 3 columns						

df

```
In [14]: # Split dialogue text into multiple rows
          # Ex. split row 7221 into multiple rows based on \n in the dialogue column
         (df['dialogue'].str.split('\n', expand=True).stack()
                         .reset index(level=1, drop=True).rename('dialogue'))
                     Before we proceed any further, hear me speak.
Out[14]:
                                                      Speak, speak.
         2
                 You are all resolved rather to die than to fam...
         3
                                                Resolved. resolved.
         4
                 First, you know Caius Marcius is chief enemy t...
         7220
                  With eyes wide open; standing, speaking, moving,
         7220
                                            And yet so fast asleep.
         7221
                                                   Noble Sebastian,
         7221
                 Thou let'st thy fortune sleep--die, rather; wi...
         7221
                                            Whiles thou art waking.
         Name: dialogue, Length: 25680, dtype: object
In [15]: # Join the splitted dialogue with the original dataframe
         df = (df.drop('dialogue', axis=1)
                  .join(df['dialogue'].str.split('\n', expand=True).stack()
                        .reset index(level=1, drop=True).rename('dialogue'))
                  .reset index(drop=True))
```

speaker		line	dialogue	
	0	First Citizen	1	Before we proceed any further, hear me speak.
	1	All	2	Speak, speak.
	2	First Citizen	3	You are all resolved rather to die than to fam
	3	All	4	Resolved. resolved.
	4	First Citizen	5	First, you know Caius Marcius is chief enemy t
	•••		•••	
256	675	SEBASTIAN	7221	With eyes wide open; standing, speaking, moving,
256	676	SEBASTIAN	7221	And yet so fast asleep.
256	677	ANTONIO	7222	Noble Sebastian,
256	678	ANTONIO	7222	Thou let'st thy fortune sleepdie, rather; wi
256	679	ANTONIO	7222	Whiles thou art waking.

25680 rows × 3 columns

Part 2

Out[15]:

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 a likely source title for each line
- add the 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.

I found a Kaggle dataset about Shakespeare plays

(https://www.kaggle.com/kingburrito666/shakespeare-plays). It has the following columns:

- 1. Dataline
- 2. Play
- 3. PlayerLinenumber
- 4. ActSceneLine
- 5. Player
- 6. PlayerLine

```
In [16]: # Read in the data
    shakespeare = pd.read_csv('Shakespeare_data.csv')
    shakespeare[92333:92338]
```

Out[16]:	Dataline Play		PlayerLinenumber	ActSceneLine	Player	PlayerLine	
	92333	92334	The Tempest	120.0	2.1.227	SEBASTIAN	With eyes wide open, standing, speaking, moving,
	92334	92335	The Tempest	120.0	2.1.228	SEBASTIAN	And yet so fast asleep.
	92335	92336	The Tempest	121.0	2.1.229	ANTONIO	Noble Sebastian,
	92336	92337	The Tempest	121.0	2.1.230	ANTONIO	Thou let'st thy fortune sleep die, rather, wi
	92337	92338	The Tempest	121.0	2.1.231	ANTONIO	Whiles thou art waking.

I can determine the play by matching speaker to Player and dialogue to PlayerLine. For easier matches, I will drop the empty dialogues and transform the remaining to plain lowercase text.

```
In [17]: | # df: drop the empty dialogues
          df = df[df['dialogue'] != ''].reset index(drop=True)
In [18]:
          # Remove punctuation and lowercase (example)
          (shakespeare[92333:92338]['PlayerLine']
           .str.replace(r'[^\w\s]', '', regex=True)
           .str.lower())
          92333
                    with eyes wide open standing speaking moving
Out[18]:
          92334
                                            and yet so fast asleep
          92335
                                                   noble sebastian
          92336
                   thou letst thy fortune sleepdie rather winkst
          92337
                                            whiles thou art waking
          Name: PlayerLine, dtype: object
          # df: create a new column for transformed dialogue
In [19]:
          df['mat'] = (df['dialogue']
                       .str.replace(r'[^\w\s]', '', regex=True)
                        .str.lower())
          df.tail(3)
Out[19]:
                  speaker
                           line
                                                           dialogue
                                                                                                 mat
          25552 ANTONIO 7222
                                                     Noble Sebastian,
                                                                                        noble sebastian
                                  Thou let'st thy fortune sleep--die, rather;
                                                                        thou letst thy fortune sleepdie rather
          25553 ANTONIO 7222
                                                                                                winkst
          25554 ANTONIO 7222
                                                Whiles thou art waking.
                                                                                   whiles thou art waking
In [20]:
          # shakespeare: create a new column for transformed PlayerLine
          shakespeare['che'] = (shakespeare['PlayerLine']
                                 .str.replace(r'[^\w\s]', '', regex=True)
                                 .str.lower())
          shakespeare[92335:92338]
```

```
Out[20]:
                                                                       Player
                                The
           92335
                                                121.0
                    92336
                                                            2.1.229 ANTONIO Noble Sebastian,
                                                                                              noble sebastian
                            Tempest
                                                                                Thou let'st thy
                                                                                                 thou letst thy
                                The
                                                            2.1.230 ANTONIO
                                                                                fortune sleep-- fortune sleepdie
           92336
                    92337
                                                121.0
                           Tempest
                                                                                die, rather, wi...
                                                                                                 rather winkst
                                                                                Whiles thou art
                                                                                               whiles thou art
                               The
                                                             2.1.231 ANTONIO
           92337
                    92338
                                                121.0
                            Tempest
                                                                                      waking.
                                                                                                      waking
In [21]: # Merge df and shakespeare
           # ['speaker','mat'] == ['Player','che']
           df = (df)
                  .merge(shakespeare[['Play', 'Player', 'che']],
                         how = 'left',
                         left on = ['speaker', 'mat'],
                         right on = ['Player','che'])
                  .drop(['Player','mat','che'], axis = 1)
                  .rename(str.lower, axis='columns'))
           df
```

PlayerLine

che

Play PlayerLinenumber ActSceneLine

Out[21]:

	speaker	line	dialogue	play
0	First Citizen	1	Before we proceed any further, hear me speak.	Coriolanus
1	All	2	Speak, speak.	Coriolanus
2	First Citizen	3	You are all resolved rather to die than to fam	Coriolanus
3	All	4	Resolved. resolved.	Coriolanus
4	First Citizen	5	First, you know Caius Marcius is chief enemy t	Coriolanus
•••				
25633	SEBASTIAN	7221	With eyes wide open; standing, speaking, moving,	The Tempest
25634	SEBASTIAN	7221	And yet so fast asleep.	The Tempest
25635	ANTONIO	7222	Noble Sebastian,	The Tempest
25636	ANTONIO	7222	Thou let'st thy fortune sleepdie, rather; wi	The Tempest
25637	ANTONIO	7222	Whiles thou art waking.	The Tempest

25638 rows × 4 columns

Dataline

```
# Check the unmatched dialogues
In [22]:
         df[df['play'].isnull()]
```

	speaker	line	dialogue	play
789	CORIOLANUS	228	But then Aufidius was within my view,	NaN
942	VOLUMNIA	274	Nay,'tis true.	NaN
1569	BRUTUS	455	And	NaN
1570	BRUTUS	455	Twice being	NaN
1605	LARTIUS	466	On safe-guard he came to me; and did curse	NaN
•••				
21617	ANGELO	5878	Happy return be to your royal grace!	NaN
22975	PETRUCHIO	6265	The youngest daughter whom you hearken for	NaN
22995	GRUMIO	6269	O excellent motion! Fellows, let's be gone.	NaN
23326	GREMIO	6383	Amen, say we: we will be witnesses.	NaN
23860	ALL SERVING-MEN	6570	Here, here, sir; here, sir.	NaN

64 rows × 4 columns

Out[22]:

The number of rows increased from 25554 to 25638, indicating that some dialogues got matched more than once. In addition, there are 64 unmatched dialogues.

Assumptions:

- If an unmatch and a match belong to the same line, they should correspond to the same play.
- If the dialogue before an unmatche has a corresponding play and the dialogue after belong to the same play, the unmatch should come from the same play.

```
In [23]: # Example of the first assumption
# line 455 belong to Coriolanus
# row 1569 and 1570 should come from the same play
df[1568:1573]
```

```
Out [23]:speakerlinedialogueplay1568BRUTUS455That our beat water brought by conduits hither;Coriolanus1569BRUTUS455AndNaN1570BRUTUS455Twice beingNaN1571BRUTUS455Was his great ancestor.Coriolanus1572SICINIUS456One thus descended,Coriolanus
```

```
In [25]: # Example of the second assumption
# row 2746 and row 2748 belong to the same play
# => row 2747 should come from the same play
df[2745:2750]
```

```
Out[25]:
                         speaker
                                  line
                                                                       dialogue
                                                                                      play
           2745
                  First Servingman
                                  843
                                                         Ay, and for an assault too. Coriolanus
           2746
                 Third Servingman 844
                                       O slaves, I can tell you news, -- news, you ras... Coriolanus
                                                    What, what, what? let's partake.
           2747
                  First Servingman
                                  845
                                                                                      NaN
                                        I would not be a Roman, of all nations; I had as Coriolanus
           2748
                 Third Servingman
                                  846
           2749 Third Servingman 846
                                                        lieve be a condemned man. Coriolanus
          for i in range(1,len(df)-1):
In [26]:
               if df.loc[i,'play'] is np.nan:
                    if ((df.loc[i-1, 'play'] is not np.nan) and # i-1 belongs to a play
                         (df.loc[i-1,'play'] == df.loc[i+1,'play'])): # i+1 belongs to the same play
                         df.loc[i,'play'] = df.loc[i-1,'play'] # assign i with the same play
In [27]:
           # Check the unmatched dialogues now
           df[df['play'].isnull()]
Out[27]:
                     speaker line
                                                                  dialogue
                                                                           play
           1823 Senators, &C 539
                                                Weapons, weapons!
                                                                           NaN
           1824 Senators, &C 539
                                    'Tribunes!' 'Patricians!' 'Citizens!' 'What, ho!'
                                                                            NaN
                                      'Sicinius!' 'Brutus!' 'Coriolanus!' 'Citizens!'
           1825
                 Senators, &C 539
                                                                            NaN
           1826 Senators, &C 539
                                      'Peace, peace, peace!' 'Stay, hold, peace!'
                                                                           NaN
           # Manually find the play for the 4 unmatched dialogues
In [28]:
           df = df.fillna('Coriolanus')
           df[1823:1827]
Out[28]:
                     speaker line
                                                                  dialogue
                                                                                 play
           1823 Senators, &C
                              539
                                                Weapons, weapons!
                                                                           Coriolanus
           1824 Senators, &C 539
                                   'Tribunes!' 'Patricians!' 'Citizens!' 'What, ho!'
                                                                           Coriolanus
                                      'Sicinius!' 'Brutus!' 'Coriolanus!' 'Citizens!' Coriolanus
           1825 Senators, &C 539
           1826 Senators, &C 539
                                      'Peace, peace, peace!' 'Stay, hold, peace!' Coriolanus
```

In [29]:

df

Resulting dataframe

	speaker	line	dialogue	play
0	First Citizen	1	Before we proceed any further, hear me speak.	Coriolanus
1	All	2	Speak, speak.	Coriolanus
2	First Citizen	3	You are all resolved rather to die than to fam	Coriolanus
3	All	4	Resolved. resolved.	Coriolanus
4	First Citizen	5	First, you know Caius Marcius is chief enemy t	Coriolanus
•••		•••		
25633	SEBASTIAN	7221	With eyes wide open; standing, speaking, moving,	The Tempest
25634	SEBASTIAN	7221	And yet so fast asleep.	The Tempest
25635	ANTONIO	7222	Noble Sebastian,	The Tempest
25636	ANTONIO	7222	Thou let'st thy fortune sleepdie, rather; wi	The Tempest
25637	ANTONIO	7222	Whiles thou art waking.	The Tempest

25638 rows × 4 columns

Part 3

Out[29]:

Pick one or more of the techniques described in this chapter:

- keyword frequency
- entity relationships
- markov language model
- bag-of-words, TF-IDF
- semantic embedding

make a case for a technique to measure how *important* or *interesting* a speaker is. The measure does not have to be both important *and* interesting, and you are welcome to come up with another term that represents "useful content", or tells a story (happiest speaker, worst speaker, etc.)

Whatever you choose, you must

- 1. document how your technique was applied
- 2. describe why you believe the technique is a valid approximation or exploration of how important, interesting, etc., a speaker is.
- 3. list some possible weaknesses of your method, or ways you expect your assumptions could be violated within the text.

This is mostly about learning to transparently document your decisions, and iterate on a method for operationalizing useful analyses on text. Your explanations should be understandable; homeworks will be peer-reviewed by your fellow students.

```
plt.imshow(wc, interpolation="bilinear")
plt.axis("off")
plt.show()
```

```
Richard IIII
Richard II Henry VI Part 3
A Winters Tale
Coriolanus
Romeo and Juliet
```



```
dialogue
Out[33]:
                  speaker line
                                                                                      play
                                                                                                                             tokens
                       First
                                     Before we proceed any further, hear me
                                                                                                 [before, we, proceed, any, further,
             0
                                                                               Coriolanus
                                                                                                                        hear, me, ...
                    Citizen
                                                                      speak.
             1
                        ΑII
                                2
                                                               Speak, speak. Coriolanus
                                                                                                                     [speak, speak]
                      First
                                    You are all resolved rather to die than to
                                                                                              [you, are, all, resolved, rather, to, die,
             2
                                                                               Coriolanus
                    Citizen
                                                                       fam...
                                                                                                                               tha...
```

```
In [34]: # Remove stopwords from tokens
    from nltk.corpus import stopwords
    sw = stopwords.words("english")
    df['tokens'] = df['tokens'].apply(lambda x: [tkn for tkn in x if tkn not in sw])
    df.head(3)
```

```
Out[34]:
                 speaker line
                                                               dialogue
                                                                              play
                                                                                                      tokens
           O First Citizen
                               Before we proceed any further, hear me speak. Coriolanus
                                                                                         [proceed, hear, speak]
           1
                            2
                                                           Speak, speak. Coriolanus
                      ΑII
                                                                                                [speak, speak]
           2 First Citizen
                            3
                               You are all resolved rather to die than to fam... Coriolanus [resolved, rather, die, famish]
           # Count the number of types
In [35]:
           df.tokens.explode().value counts().head(30)
          thou
                     1405
Out[35]:
          thy
                     1059
          shall
                      845
          thee
                      760
          good
                      671
          lord
                      632
          come
                      628
          sir
                      592
          well
                      565
          would
                      532
          ill
                      467
          hath
                      453
          king
                      448
                      445
          say
          let
                      437
                      423
          one
          qo
                      421
          love
                      410
                      406
          may
                      401
          us
                      397
          make
          upon
                      390
                      386
          yet
          like
                      376
          must
                      370
          know
                      350
                      343
          man
          tis
                      328
          see
                      322
                      304
          death
          Name: tokens, dtype: int64
In [36]:
           # Add more stopwords
           extra = ['thou','thy','shall','thee','would','hath',
                      'let','one','may','us','upon','yet','tis']
           sw = sw + extra
           # Remove extra stopwords from tokens
In [37]:
           df['tokens'] = df['tokens'].apply(lambda x: [tkn for tkn in x if tkn not in sw])
           df.tail(3)
Out[37]:
                   speaker
                              line
                                                             dialogue
                                                                                                        tokens
                                                                            play
                                                                            The
           25635 ANTONIO 7222
                                                      Noble Sebastian,
                                                                                               [noble, sebastian]
                                                                        Tempest
                                       Thou let'st thy fortune sleep--die,
                                                                            The
                                                                                    [letst, fortune, sleepdie, rather,
           25636 ANTONIO 7222
                                                           rather; wi...
                                                                        Tempest
                                                                                                        winkst]
                                                                            The
           25637 ANTONIO 7222
                                                Whiles thou art waking.
                                                                                             [whiles, art, waking]
                                                                        Tempest
```

```
In [38]: # Generate word clouds for the 10 most frequent speakers
speakers = df.speaker.explode().value_counts().head(10).index.tolist()
```

```
Out[38]:
         DUKE VINCENTIO
                            824
         KING RICHARD II
                            760
         LEONTES
                            679
         CORIOLANUS
                             676
         ROMEO
                            604
         MENENIUS
                            587
         PETRUCHIO
                            568
         JULIET
                            551
         QUEEN MARGARET
                           498
         Name: speaker, dtype: int64
In [39]:
         speakers[0]
         'GLOUCESTER'
Out[39]:
In [40]:
         # Generate word cloud for a speaker
         def speaker wordcloud(name):
             return (df[df.speaker == name]
                      .tokens.explode().value counts()
                      .pipe(value ct wordcloud))
In [41]:
         # Visualize keywords of a speaker
         for name in speakers:
             print(name)
             speaker wordcloud(name)
         GLOUCESTER
           time edwards
                           edward •clarence
```



df.speaker.explode().value counts().head(10)

904

GLOUCESTER

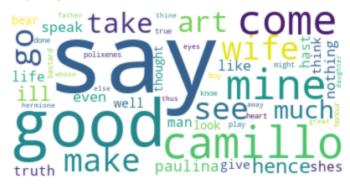
DUKE VINCENTIO



KING RICHARD II



LEONTES



CORIOLANUS



ROMEO



MENENIUS







QUEEN MARGARET



The word clouds above explore how interesting the speakers are. Juliet, for example, is all about love and Romeo, while Queen Margaret mentions a ton about the King and Edward. One weakness of this method is that some keywords do not necessarily represent meanings, but rather the name of the person talking to the speaker.