# NLP Homework 3

# Sentiment Analysis
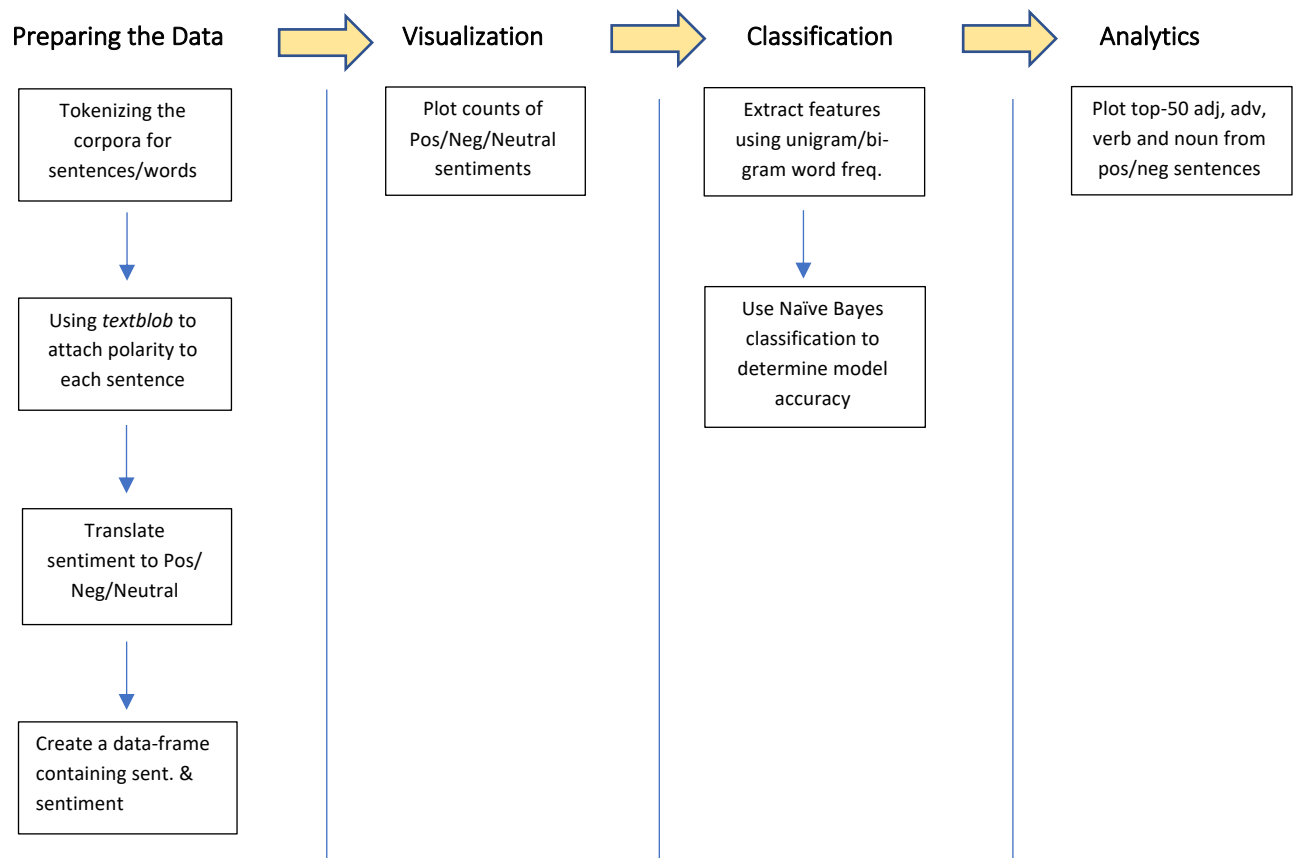
Student: Sharat Sripada (vssripad@syr.edu)

# 1. Introduction

In this week's homework we will perform sentiment analysis on the same corpora of Shakespearean books used in HW-1 and HW-2 namely, Julius Caesar and Hamlet. While the goal was to use knowledge and methodologies related to NLP gathered in IST-664 to bring out any subtle differences in language or style, we have not found any significant deviation thus far (apart from Hamlet just being a longer play and hence a larger corpus).

That is likely to continue with exploration into sentiment analysis as well, since the genre of both plays is tragic. Here are the high-level tasks:
- Sentiment analysis at sentence level
- Build a classifier to classify sentiment polarity of a sentence (as Positive, Negative or Neutral)
- Write to a csv file content of title, author, country, #pos-sent, #neg-sent, #nuetral-sent – This may not be particularly applicable for this case-study (although we will demonstrate the ability to write csv)
- Analyze all the positive sentences to identify top-50 adjective, adverb, noun, or verb phrases and do the same for negative sentences as well.

To better visualize all the tasks or assignments, here is a high-level summary what we are seeking to achieve through this report:

| Preparing the Data | | Visualization | | Classification | | Analytics |
|---|---|---|---|---|---|---|
| Tokenizing the corpora for sentences/words | → | Plot counts of Pos/Neg/Neutral sentiments | → | Extract features using unigram/bi-gram word freq. | → | Plot top-50 adj, adv, verb and noun from pos/neg sentences |
| Using *textblob* to attach polarity to each sentence | | | | Use Naïve Bayes classification to determine model accuracy | | |
| Translate sentiment to Pos/ Neg/Neutral | | | | | | |
| Create a data-frame containing sent. & sentiment | | | | | | |

To begin with, we will continue to extract the corpora as before by indexing books from the Gutenberg corpora:

```
1  # HW-2
2  # Exploratory exercise for sentiment analysis
3  # Finding adverb and adjective phrases, and computing basic statistics
4
5  # importing required nltk libraries
6  import nltk
7  from nltk import sent_tokenize
8
9  # In HW-1, I used books written by Shakespeare - Caesar and Hamlet. We will continue exploratory
10 # sentiment analysis on the same books
11 nltk.corpus.gutenberg.fileids()
12
13 # Get Shakespeare books in the Gutenberh corpus
14 shakespeare_books = [book for book in nltk.corpus.gutenberg.fileids( ) \
15                      if 'shakespeare' in book]
16
17 # Book-1: Caesar (Genre: Tragedy)
18 caesar = nltk.corpus.gutenberg.raw(shakespeare_books[0])
19
20 # Book-2: Hamlet (Genre: Tragedy/Comedy)
21 hamlet = nltk.corpus.gutenberg.raw(shakespeare_books[1])
22
23 print(caesar[:50])
24 print(hamlet[:50])
```

```
[The Tragedie of Julius Caesar by William Shakespe
[The Tragedie of Hamlet by William Shakespeare 159
```

## 2. Preparing the data

Since the corpus mostly comprises contents from the book, the requirement was to first tokenize it into sentences and add a sentiment which could then be used to train and test models. For this purpose, we use a library called *TextBlob* which offers a *sentiment* functionality as below:

```
print (blob)
blob.sentiment
>> Analytics Vidhya is a great platform to learn data science.
Sentiment(polarity=0.8, subjectivity=0.75)
```

The code below shows how TextBlob was adopted for this study:

```
3  # Get polarity/subjectivity for each sent in both corpora
4  caesar_sentiment = [[sent, TextBlob(sent).sentiment] for sent in caesar_sentences[2:]]
5  print(caesar_sentiment[1:10])
6
7  hamlet_sentiment = [[sent, TextBlob(sent).sentiment] for sent in hamlet_sentences[2:]]
8  print(hamlet_sentiment[1:10])
```

Notice that, we walk sentences (tokenized by function *sent_tokenize()*) in both corpora and create a list within a list comprehension to encompass a sentiment polarity or subjectivity tag alongside each sentence. For example:

```
['Hence: home you idle Creatures, get you home:\nIs this a Holiday?', Sentiment(polari
ty=0.0, subjectivity=0.0)]
```

Further, we can also specifically access the polarity or subjectivity of each sentence using an object or instance attribute as shown below. This would later help us translate a polarity range (-1, 1) to Negative (range (-1, <0), or to Neutral (0) or Positive (0, 1):

```
sent = caesar_sentiment[2][0]
polarity = caesar_sentiment[2][1].polarity
subjectivity = caesar_sentiment[2][1].subjectivity
print('Sentence %s has polarity %s and subjectivity %s' %(sent, polarity,
subjectivity))
```

Output showing sentence and importantly its polarity and subjectivity can be extracted:
```
Sentence Hence: home you idle Creatures, get you home:Is this a Holiday? has polarity
0.0 and subjectivity 0.0
```

The translate function, its usage and output:

```
1  # Next we will translate the polarity and subjectivity index to the tags we want to see.
2  import copy
3
4  def translate_sentiment(sent_blob):
5      '''
6      Given a sentence, derive its sentiment
7      '''
8      if sent_blob[1].polarity < 0:
9          sent_blob.append('Negative')
10     elif sent_blob[1].polarity == 0:
11         sent_blob.append('Neutral')
12     elif sent_blob[1].polarity > 0:
13         sent_blob.append('Positive')
14     return sent_blob
15
16 print(translate_sentiment(caesar_sentiment[2]))
17 print(translate_sentiment(hamlet_sentiment[2]))
```

Output:
```
['Hence: home you idle Creatures, get you home:\nIs this a Holiday?', Sentiment(polari
ty=0.0, subjectivity=0.0), 'Neutral'] <- Notice the tag Nuetral which is extracted fro
m the polarity index
```

Finally, we summarize this data into a Pandas data-frame comprising column-names – Sentence, TextBlob-sentiment(raw) and Sentiment. All further study would be based on this data-frame.

```
1  # Let's make a data-frame comprising all the above data
2  import pandas as pd
3
4  def create_df(textblob_raw):
5      df = pd.DataFrame(columns=['Sentence', 'TextBlob-sentiment(raw)', 'Sentiment'])
6      for blob in textblob_raw:
7          new_blob = translate_sentiment(blob)
8          row = {'Sentence': new_blob[0],
9                 'TextBlob-sentiment(raw)': new_blob[1],
10                'Sentiment': new_blob[2]
11         }
12         df = df.append(row, ignore_index=True)
13     return df
14
```

Snapshot of the data-frame:

```
In [7]:   1  df_hamlet = create_df(hamlet_sentiment)
          2  df_hamlet.head()
```
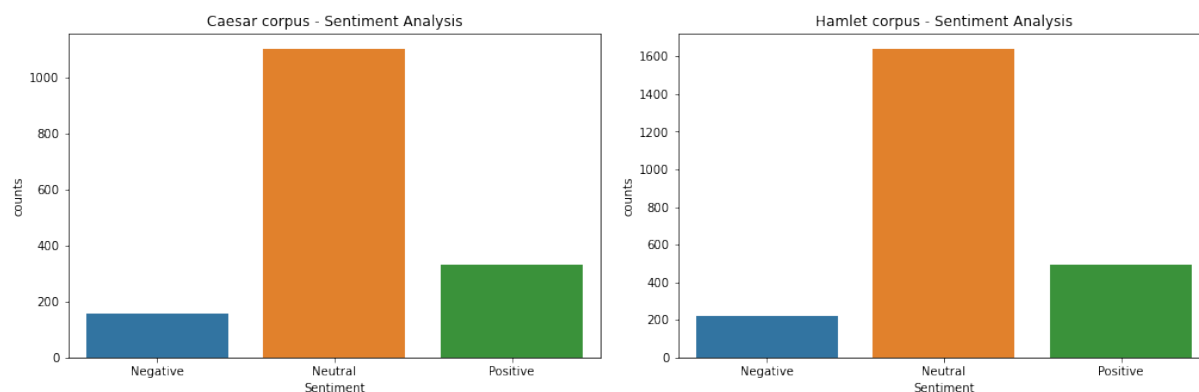
Out[7]:

| | Sentence | TextBlob-sentiment(raw) | Sentiment |
|---|---|---|---|
| 0 | Enter Barnardo and Francisco two Centinels. | (0.0, 0.0) | Neutral |
| 1 | Barnardo. | (0.0, 0.0) | Neutral |
| 2 | Who's there? | (0.0, 0.0) | Neutral |
| 3 | Fran. | (0.0, 0.0) | Neutral |
| 4 | Nay answer me: Stand & vnfold\nyour selfe\n\n ... | (0.0, 0.0) | Neutral |

## 3. Visualization

The translation function based on sentiment polarity, helps us now group and count occurrences of each type viz. Positive, Negative or Neutral. Here is the code that groups the data and plots a bar graph using matplotlib and sns libraries:

```python
1   # Make a bar-plot of sentiment-counts across both the corpora
2
3   import matplotlib.pyplot as plt
4   import seaborn as sns
5
6   def plot_bar_graph(corpus_name, label_groups):
7       plt.figure(figsize=(8,5))
8       ax = sns.barplot(x="Sentiment", y="counts", data=label_groups)
9       ax.set_xticklabels(ax.get_xticklabels(), rotation=0)
10      ax.set_title(label="%s" %corpus_name)
11      plt.show()
12
13
14  caesar_labels = df_caesar.groupby('Sentiment').size().reset_index(name='counts')
15  plot_bar_graph('Caesar corpus - Sentiment Analysis', caesar_labels)
16
17  hamlet_labels = df_hamlet.groupby('Sentiment').size().reset_index(name='counts')
18  plot_bar_graph('Hamlet corpus - Sentiment Analysis', hamlet_labels)
```

Correspondingly, this is the bar-plot visualization:



### Interpreting the charts

We know by now Hamlet is comparatively a longer play or book and would therefore understandably comprise a larger corpus of words.

It is however interesting to see from the plots that although the genre is tragic, there is still largely a positive polarity to sentence sentiment (two times that of negative sentence polarity). This goes to show Shakespeare wrote with a subtle undertone of humor or sarcasm. Perhaps, examining words in greater detail (and their POS references) at a later point would make this evident.

### NOTE

The plots are based on raw text without necessarily filtering for stop words. Neither have we used techniques like regular expressions to clean up for sentence anomalies. Presumably employing such techniques could spill sentiment from Neutral to Positive or Negative; but may not drastically alter the ratio since we have sufficient data when making the inference.

Further, a summary was written to csv (utilizing the *to_csv()* function) with the following code:

```
1  # Also, write a csv with this data
2  df_summary = pd.DataFrame(columns=['Book', 'Negative', 'Neutral', 'Positive'])
3
4  df_summary = df_summary.append({'Book': 'Caesar',
5                                  'Negative': df_caesar_sentiment.loc[0][1],
6                                  'Neutral': df_caesar_sentiment.loc[1][1],
7                                  'Positive': df_caesar_sentiment.loc[2][1]}, ignore_index=True)
8
9  df_summary = df_summary.append({'Book': 'Hamlet',
10                                 'Negative': df_hamlet_sentiment.loc[0][1],
11                                 'Neutral': df_hamlet_sentiment.loc[1][1],
12                                 'Positive': df_hamlet_sentiment.loc[2][1]}, ignore_index=True)
13
14 df_summary.head()
15 df_summary.to_csv('new_data.csv', index=False)
```

Output:
Shows the contents of the csv file with counts of positive, negative and neutral sentences from both corpora:
```
# $ less new_data.csv
# Book,Negative,Neutral,Positive
# Caesar,157,1101,332
# Hamlet,222,1640,491
```

# 4.  Classification Task

With the data collection sufficiently complete that is, every sentence in both corpora have a corresponding sentiment tag, we can now build a model and make predictions.

For this exercise, we will use the Naïve Bayes classification model and use K-fold (k=5) cross-validation method rather than a flat 80-20 split.

## Why is k-fold cross-validation method?
Cross-validation is a powerful tool that helps us better use our data (spread) and therefore provides useful information about the performance of the algorithms we may choose.

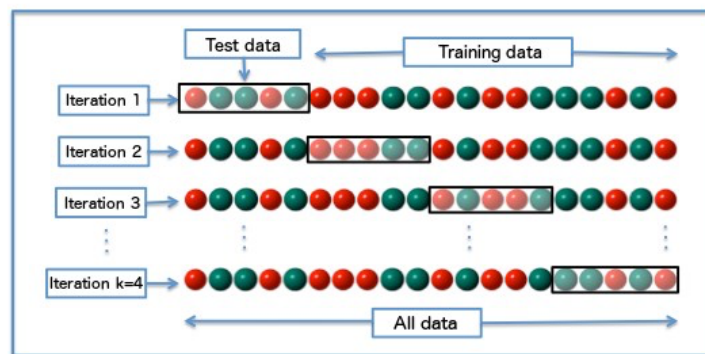Here's a good visualization of how cross-validation works:



Diagram of k-fold cross-validation with k=4.

Data is split more than once (controlled by folds or commonly known as variable k), than in the classic 1-split 80-20, 90-10 or 70-30 representing train and validation/test data.

Each split or fold is a good representation of the whole data and therefore helps train or develop comprehensive and robust models. Accuracy or performance is subsequently derived by averaging the results over several iterations as shown below:

```python
 1  # Using k-folds = 5 we will fairly randomize the train & test data
 2
 3  import numpy as np
 4  from sklearn.model_selection import KFold
 5
 6  def ml_nb(featuresets):
 7      kf = KFold(n_splits = 5)
 8      sum = 0
 9
10      for train, test in kf.split(featuresets):
11          train_data = np.array(featuresets)[train]
12          test_data = np.array(featuresets)[test]
13          classifier = nltk.NaiveBayesClassifier.train(train_data)
14          sum += nltk.classify.accuracy(classifier, test_data)
15
16
17
18      #storing the score in a variable
19      acc1 = sum/5
20
21      return  classifier, acc1
22
23  _, caesar_acc = ml_nb(caesar_featuresets)
24  _, hamlet_acc = ml_nb(hamlet_featuresets)
25
26  print('Accuracy for unigram feature-sets on Caesar corpus %s' %caesar_acc)
27  print('Accuracy for unigram feature-sets on Hamlet corpus %s' %hamlet_acc)
```

See rows 14-19 above on how accuracy is summed over each k-fold iteration and then averaged to provide an overall accuracy.

## Building models

While we chose a classification algorithm, Naïve Bayes to solve our problem here and make a prediction of sentence sentiment we take two fundamental NLP approaches:
1. Uni-gram feature-set

   In this approach, we derive a feature set of single words ordered by frequency and pick the top-2000 to find a match on each sentence.

   Note here, we use a regex to find only words (re.findall(r'\w+')) and eliminate any alpha-numeric words or words that may comprise punctuation or new-line characters.

```python
 1  # Defining set of words that will be used for features
 2
 3  # We'll find the 2000 most common words and used them as an important feature of the whole corpus
 4  def unigram_freq(docs):
 5      all_words = []
 6      # Write a regex to pull only the word portion & leave
 7      # out any punctuation marks etc.
 8      for (sentence, category) in docs:
 9          for word in sentence.split():
10              try:
11                  all_words.append(re.findall(r'\w+', word)[0])
12              except IndexError:
13                  pass
14      top_words = nltk.FreqDist(all_words)
15      most_common_words = top_words.most_common(2000)
16      word_features = [word for (word,count) in most_common_words]
17      return all_words, word_features
```

Once the word features are identified, *document_features()* is called to iterate over the corpus and create a dictionary that would comprise 'contains(<word>): <Boolean>' where Boolean True/False would be dependent on a word match.

```python
1  # now we will use that list of most frequent words in the entire corpus
2  # to iterate over each sentence and check if any of those words are present
3  # in that way, we will see if this unigram corpus feature is present on that particular sentence
4  # using Boolean logic that matches values and returns 'True' or 'False'
5  # we do this by defining a Python "function," i.e.a piece of code writen to be reused
6  def document_features(document, word_features):
7      document_words = set(document)
8      #we open a Pytnon dictionary instead of a list
9      features = {}
10     for word in word_features:
11         #checking if the word from word_features matches a word in the document
12         features['contains({})'.format(word)] = (word in document_words)
13     return features
14
15 caesar_featuresets = [(document_features(d, caesar_word_features), c) for (d, c) in caesar_docs]
16 hamlet_featuresets = [(document_features(d, hamlet_word_features), c) for (d, c) in hamlet_docs]
17
18
19 print("Length of feature set for Caesar = %d" %len(caesar_featuresets))
20 print("Length of feature set for Hamlet = %d" %len(hamlet_featuresets))
```

```
Length of feature set for Caesar = 1590
Length of feature set for Hamlet = 2353
```

The list comprehensions on row 15-16 are organized as a tuple at each index with a dictionary (in the format above) and a sentence sentiment alongside.

Example of the list-comprehension *caesar_featuresets* at index 0 (trimmed for brevity):
```
({'contains(I)': False, 'contains(the)': False, 'contains(and)': False, 'contai
ns(to)': False, 'contains(you)': False, 'contains(of)': False, 'contains(not)':
False, 'contains(a)': True, 'contains(is)': False, 'contains(And)': False, 'con
tains(in)': False, 'contains(that)': False, 'contains(my)': False, 'contains(Ca
esar)': False, 'contains(me)': False, 'contains(it)': False, 'contains(him)': F
alse, 'contains(Brutus)': False, 'contains(Bru)': False, 'contains(his)': False
, 'contains(this)': False, 'contains(your).. False}, 'Neutral')
```

2. Bi-gram feature-set

In this approach, we will group two-words at a time using the NLTK collocations.BigramAssocMeasures() and score frequencies using the score_ngrams() functionality.

```python
1  # Since the accuracy is low, let's also try this with a
2  # bi-gram featureset
3
4  # Re-using the code to clean-up
5  from nltk.collocations import *
6  import re
7
8  def bigram_freq(all_words):
9      #data cleaning and preprocessing
10     stopwords = nltk.corpus.stopwords.words('english')
11
12     #creating bigrams features for the corpus and applying cleaning steps
13     bigram_measures = nltk.collocations.BigramAssocMeasures()
14     finder = BigramCollocationFinder.from_words(all_words)
15     finder.apply_word_filter(lambda w: w in stopwords)
16     scored = finder.score_ngrams(bigram_measures.raw_freq)
17
18     #extracting clean bigrams (no frequency information)
19     bigram_features = [bigram for (bigram, count) in scored[:2000]]
20
21     return bigram_features
```

Correspondingly, the *bi_document_features()* function would walk the word pairs (collocations) in each email, and attempt to find a match with the bi-gram features. Like the unigram document features functionality, the function would return a dictionary comprising 'contains(<word-pairs>): <Boolean>` which would be then placed alongside a category spam or ham at each node or index in the list.

```python
1  def bi_document_features(document, bigram_features):
2      document_words = list(nltk.bigrams(document))
3      features = {}
4      for word in bigram_features:
5          #boolean logic will return 'True' if there is a match, or 'False' if not
6          features['contains({})'.format(word)] = (word in document_words)
7      return features
8
9  # applying the function to our documents
10 caesar_featuresets2 = [(bi_document_features(d, caesar_word_bi_features), c) for (d, c) in caesar_docs]
11 hamlet_featuresets2 = [(bi_document_features(d, hamlet_word_bi_features), c) for (d, c) in hamlet_docs]
```

Example of the list comprehension caeser_featureset2(trimmed for brevity):
({"contains(('I', 'haue'))": False, "contains(('I', 'know'))": False, "contains(('Cass i', 'I'))": False, "contains(('I', 'shall'))": False, "contains(('Bru', 'I'))": False, "contains(('Mark', 'Antony'))": False, "contains(('let', 'vs'))": False, "contains(('M arke', 'Antony'))": False, "contains(('And', 'I'))": False, "contains(('Lord', 'Bru')) ": False, "contains(('I', 'feare'))": False, "contains(('Caesar', 'Caes'))": False, "c ontains(('Enter', 'Brutus'))": False, "contains(('For', 'I'))": False, "contains(('Nob le', 'Brutus'))": False, "contains(('Bru', 'O'))": False, "contains(('I', 'may'))":.. **'Neutral'**)

**NOTE:**
Notice the difference in dictionary keys at index 0 for the unigram and bi-gram feature sets. This will largely influence the models we build hence.

The modeling itself was written in a general manner so it can return a classifier and accuracy score given a feature-set, and is represented by the function *ml_nb():*

```python
1  # Using k-folds = 5 we will fairly randomize the train & test data
2
3  import numpy as np
4  from sklearn.model_selection import KFold
5
6  def ml_nb(featuresets):
7      kf = KFold(n_splits = 5)
8      sum = 0
9
10     for train, test in kf.split(featuresets):
11         train_data = np.array(featuresets)[train]
12         test_data = np.array(featuresets)[test]
13         classifier = nltk.NaiveBayesClassifier.train(train_data)
14         sum += nltk.classify.accuracy(classifier, test_data)
15
16
17
18     #storing the score in a variable
19     acc1 = sum/5
20
21     return  classifier, acc1
22
```

The data is split as train and test before running the NLTK NaiveBayesClassifier() on train-data (seen on row-13). The prediction or classification and accuracy is recorded and averaged over the k-folds.

Uni-gram vs bi-gram results are tabulated here:

| Naïve Bayes classifier k-fold = 5 | Accuracy - Caesar corpus | Accuracy - Hamlet corpus |
|---|---|---|
| Uni-gram frequency | 67.0% | 69.3% |
| Bi-gram frequency | 69.2% | 69.7% |

The results show a slightly improved accuracy when using the bi-gram frequency sets. Specifically, in the case of the Caesar corpus we see a 2.2% improvement while a marginal improvement for the Hamlet corpus.

## 5. Analytics

Finally, we would like to look at the point-of-speech tags and how they vary in sentences conveying positive or negative polarity.

The code is mostly a re-use from HW-2 written generally like this:

```
1  # Utilizing code from HW-2, we will now call the same functions for
2
3  def top_tokens(taggedtext, pos_list):
4      _tokens = []
5      for sentence in taggedtext:
6          for word, pos in sentence:
7              if pos in pos_list:
8                  if len(word)>1:
9                      _tokens.append(word)
10     freq_pos = nltk.FreqDist(_tokens)
11
12     for word, freq in freq_pos.most_common(50):
13         print(word,freq)
```

Essentially, given a POS tagged sentence and pre-defined POS tags the function *top_tokens()* would iterate each sentence and word, sort the words based on frequency and print the top FIFTY.

The Stanford POS tagger *nltk.pos_tag()* will help with tagging and we create list comprehensions comprising tokenized words and the corresponding POS tag:

```
14  # Walk the positive/negative sentences and word tokenize it, before running
15  # it through the pos_tag()
16  hamlet_tokens_pos = [nltk.word_tokenize(sent) for sent in hamlet_sent_pos]
17  hamlet_tags_pos = [nltk.pos_tag(token) for token in hamlet_tokens_pos]
18
19  hamlet_tokens_neg = [nltk.word_tokenize(sent) for sent in hamlet_sent_neg]
20  hamlet_tags_neg = [nltk.pos_tag(token) for token in hamlet_tokens_neg]
```

We then call the *top_tokens()* function for each POS we are interested in viz. Adjectives, Adverbs, Nouns or verbs:

```
 1  # Stats for Hamlet data:
 2  # Top 50 adjective tokens
 3  print('Top 50 Adjective tokens (positive):')
 4  h_adj_pos = top_tokens(hamlet_tags_pos, ['JJ', 'JJR', 'JJS'])
 5
 6  print('Top 50 Adjective tokens (negative):')
 7  h_adj_neg = top_tokens(hamlet_tags_neg, ['JJ', 'JJR', 'JJS'])
 8
 9  # Top 50 adverb tokens
10  print('Top 50 Adverb tokens (positive):')
11  h_adv_pos = top_tokens(hamlet_tags_pos, ['RB', 'RBR', 'RBS'])
12
13  print('Top 50 Adverb tokens (negative):')
14  h_adv_neg = top_tokens(hamlet_tags_neg, ['RB', 'RBR', 'RBS'])
15
16  print('Top 50 Noun tokens: (positive)')
17  h_noun_pos = top_tokens(hamlet_tags_pos, ['NN', 'NNS', 'NNP', 'NNPS']) #Noun, Noun-plural, Noun-Proper, Noun-Pro
18  print('Top 50 Noun tokens: (negative)')
19  h_noun_neg = top_tokens(hamlet_tags_neg, ['NN', 'NNS', 'NNP', 'NNPS'])
20
21
22  print('\nTop 50 Verb tokens: (positive)')
23  h_verb_pos = top_tokens(hamlet_tags_pos, ['VB', 'VBD', 'VBG', 'VBP', 'VBZ']) # Verb, Verb-past-tense, Verb-prese
24                                                  # Verb-past participle, singular present (non-3rd)
25                                                  # singular present (3rd)
26  print('Top 50 Verb tokens: (negative)')
27  h_verb_neg = top_tokens(hamlet_tags_neg, ['VB', 'VBD', 'VBG', 'VBP', 'VBZ'])
```

Following is a summary of the words for different POS tags across the corpora:

1. Caesar corpus (top-50 words)

| Adjectives – Positive sentences | Adjectives – Negative sentences | Adverbs – Positive sentences | Adverbs – Negative sentences | Verbs – Positive sentences | Verbs – Negative sentences | Nouns – Positive sentences | Nouns – Negative sentences |
|---|---|---|---|---|---|---|---|
| good 46 | such 11 | not 85 | not 52 | is 87 | is 45 | Caesar 78 | Caesar 30 |
| Noble 31 | dead 11 | so 46 | then 18 | be 64 | do 24 | Brutus 69 | Cassius 21 |
| great 24 | bad 10 | then 32 | so 16 | haue 47 | haue 24 | Bru 36 | men 20 |
| much 20 | other 10 | more 16 | now 11 | are 38 | be 23 | Antony 29 | Bru 19 |
| Good 18 | thy 7 | now 13 | well 10 | do 38 | are 16 | men 26 | Brutus 16 |
| many 14 | common 6 | well 12 | too 9 | was 28 | did 15 | Cassi 26 | Cassi 13 |
| true 13 | dangerous 6 | yet 11 | Then 9 | am 23 | was 14 | Cassius 24 | Caesars 10 |
| thy 11 | hard 5 | Now 11 | more 8 | know 23 | know 12 | man 24 | man 10 |
| best 10 | strange 5 | too 11 | yet 6 | let 19 | tell 9 | selfe 18 | Rome 9 |
| thou 10 | angry 5 | heere 11 | heere 4 | did 19 | were 8 | day 17 | vs 9 |
| full 10 | true 5 | very 10 | So 3 | come 15 | wrong 8 | Friends 16 | selfe 8 |
| strong 10 | bloody 5 | most 9 | long 3 | say 13 | Let 6 | am 5 | Antony 8 |
| such 9 | wrong 5 | So 9 | once 3 | were 13 | come 6 | go 5 | thou 6 |
| sure 9 | sad 4 | directly 8 | Now 3 | make 13 | see 6 | loues 4 | things 6 |
| worthy 8 | little 4 | there 8 | rather 3 | go 11 | let 5 | hold 4 | day 6 |
| better 8 | Honourable 4 | Then 7 | yong 3 | tell 11 | beare 5 | loue 4 | Romans 6 |
| Most 8 | thou 3 | as 6 | indeed 2 | see 10 | say 5 | had 4 | 'd 6 |
| most 8 | euery 3 | thus 6 | else 2 | had 10 | am 5 | doth 4 | time 5 |
| high 6 | ordinary 2 | first 6 | Well 2 | feare 9 | go 5 | looke 4 | word 5 |
| old 5 | late 2 | still 5 | as 2 | finde 9 | hold 4 | hath 4 | Alas 5 |
| more 5 | gentle 2 | else 5 | alone 2 | take 9 | loue 4 | thee 4 | hands 5 |
| mine 5 | borne 2 | onely 4 | still 2 | stand 9 | had 4 | take 4 | thing 5 |
| gentle 5 | feeble 2 | once 4 | thus 2 | beare 8 | doth 4 | giue 4 | Cas 5 |
| thee 5 | sleepe 2 | much 4 | Indeed 2 | heare 7 | looke 4 | vse 3 | Lord 5 |
| glad 4 | much 2 | wisely 4 | much 2 | fell 7 | hath 4 | perceiue 3 | |
| vs 4 | Such 2 | Therefore 3 | wherefore 1 | thee 7 | thee 4 | | |
| first 4 | third 2 | together 3 | till 1 | Let 7 | take 4 | | |
| welcome 4 | most 2 | Yet 3 | doth 1 | put 6 | giue 4 | | |
| wrong 4 | | yee 3 | therefore 1 | till 6 | vse 3 | | |
| haue 3 | | | very 1 | giue 6 | perceiue 3 | | |
| further 3 | | | | thou 6 | | | |
| other 3 | | | | hath 6 | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| new 3 | terrible 2 | hence 3 | almost 1 | speake 6 | make 3 | Ant 11 | Sir 4 |
| huge 3 | secret 2 | Not 3 | laugh 1 | made 5 | follow 3 | hath 10 | hearts 4 |
| giuen 3 | impossible 2 | along 3 | away 1 | looke 5 | heare 3 | 'd 10 | Will 4 |
| right 3 | better 2 | better 3 | Not 1 | loue 5 | doe 3 | Caska 10 | hand 4 |
| common 3 | good 2 | away 3 | surly 1 | Is 5 | thinke 3 | blood 10 | World 4 |
| seene 3 | constant 2 | Truly 2 | already 1 | keepe 5 | bring 3 | Haue 9 | Caska 4 |
| strange 3 | least 2 | perhaps 2 | Nobly 1 | 'd 5 | speake 3 | Octauius 9 | Cask 4 |
| free 3 | pitty 2 | nod 2 | Quite 1 | thinke 5 | feare 3 | Will 8 | fire 4 |
| enough 3 | cold 2 | downe 2 | hard 1 | comes 5 | keepe 2 | time 8 | Roman 4 |
| dead 3 | more 2 | sayes 2 | further 1 | get 4 | euery 2 | morrow 8 | blood 4 |
| happy 3 | certaine 2 | sicke 2 | impatiently 1 | saw 4 | thou 2 | Lord 8 | Caius 4 |
| mighty 3 | poor 2 | dye 2 | no 1 | doth 4 | shake 2 | heart 8 | death 4 |
| ready 3 | sorry 2 | Thus 2 | litter 1 | run 4 | loose 2 | Octa 8 | night 4 |
| Ambitious 3 | vile 2 | fast 2 | forth 1 | went 4 | write 2 | Did 7 | Which 4 |
| Honourable 3 | worst 2 | smile 2 | deere 1 | lay 4 | get 2 | things 7 | bloody 4 |
| meete 3 | ill 2 | best 2 | marke 1 | vs 4 | thinkes 2 | Friend 7 | thee 4 |
| sir 2 | safe 1 | straight 2 | merry 1 | being 4 | put 2 | Cask 7 | Octauius 4 |
| narrow 2 | senslesse 1 | truly 2 | instantly 1 | appeare 4 | fell 2 | Noble 7 | Trade 3 |
| | many 1 | seene 2 | hence 1 | | went 2 | Thy 7 | Thou 3 |
| | seruile 1 | softly 2 | fast 1 | | | cause 7 | Be 3 |
| | stubborne 1 | partly 2 | | | | thee 7 | eyes 3 |
| | deceiu 1 | | | | | Mark 7 | shew 3 |
| | | | | | | Come 7 | Friend 3 |
| | | | | | | Honor 6 | cold 3 |
| | | | | | | Spirit 6 | hee 3 |
| | | | | | | words 6 | Did 3 |
| | | | | | | Roman 6 | |
| | | | | | | minde 6 | |
| | | | | | | euer 6 | |

2. Hamlet corpus (top-50 words)

| Adjectives – Positive sentences | Adjectives – Negative sentences | Adverbs – Positive sentences | Adverbs – Negative sentences | Verbs – Positive sentences | Verbs – Negative sentences | Nouns – Positive sentences | Nouns – Negative sentences |
|---|---|---|---|---|---|---|---|
| good 74 | dead 20 | not 117 | not 59 | is 129 | is 75 | Ham 60 | Ham 23 |
| more 27 | thy 12 | so 56 | so 26 | be 72 | be 40 | Lord 59 | King 21 |
| most 26 | mad 11 | then 37 | then 15 | haue 55 | haue 16 | King 46 | Hamlet 12 |
| much 23 | more 10 | more 33 | more 13 | are 50 | are 15 | Hamlet 29 | Lord 11 |
| Good 21 | Other 10 | most 31 | very 11 | was 33 | was 13 | Father 23 | selfe 9 |
| true 18 | such 7 | very 30 | too 11 | know 30 | 's 11 | Ile 21 | 'd 8 |
| thy 18 | thou 6 | now 25 | well 11 | am 28 | say 10 | time 21 | Laer 8 |
| great 18 | little 6 | well 22 | thus 10 | do 27 | had 9 | death 20 | Alas 8 |
| such 17 | owne 6 | too 21 | now 8 | did 27 | do 9 | Sir 19 | Qu 8 |
| sweet 14 | other 6 | Then 17 | So 5 | let 25 | let 8 | vs 17 | hath 7 |
| many 14 | long 6 | much 16 | long 5 | make 22 | make 8 | Hor 17 | Hor 7 |
| old 13 | dangerous 6 | heere 13 | Then 5 | had 17 | am 7 | man 17 | death 7 |
| Most 12 | common 5 | So 13 | n't 5 | see 16 | did 7 | Queene 16 | day 7 |
| full 9 | late 5 | thus 13 | away 5 | doe 16 | go 7 | | againe 7 |
| other 9 | wrong 5 | Now 12 | once 4 | come 15 | speake 6 | | time 7 |
| deere 9 | strange 4 | yet 12 | most 4 | tell 15 | know 6 | | Sir 7 |
| excellent 9 | second 4 | as 10 | indeed 4 | 's 15 | come 6 | | |
| | | once 8 | Thus 3 | Let 14 | tell 6 | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Noble 9 | false 4 | there 6 | as 3 | giue 14 | were 5 | Horatio 15 | Heauen 6 |
| first 9 | bad 4 | first 6 | ere 3 | speake 13 | does 5 | life 15 | body 6 |
| welcome 8 | ill 4 | indeed 6 | there 3 | hold 13 | comes 5 | Nature 15 | Father 6 |
| young 8 | guilty 3 | better 6 | Now 3 | put 12 | Let 5 | loue 15 | head 6 |
| best 8 | most 3 | freely 5 | Long 2 | say 11 | liue 4 | thing 14 | Loue 6 |
| free 8 | old 3 | Not 5 | neere 2 | were 11 | hath 4 | Mother 14 | Come 6 |
| fine 8 | hard 3 | else 5 | heartily 2 | hath 10 | thou 4 | Qu 14 | Pol 6 |
| owne 7 | wide 3 | therefore 5 | rather 2 | take 10 | doe 4 | thou 13 | Ile 6 |
| whole 7 | double 3 | neuer 5 | yong 2 | Is 9 | 're 4 | Laer 13 | haue 6 |
| better 7 | true 3 | still 5 | no 2 | set 9 | finde 4 | night 12 | Mother 6 |
| last 6 | horrible 3 | further 4 | heard 2 | said 9 | looke 4 | world 12 | 'T 5 |
| hot 6 | desperate 3 | away 4 | humbly 2 | finde 8 | makes 4 | heart 12 | poore 5 |
| thou 6 | tame 3 | enough 4 | neyther 2 | thinke 8 | thinke 4 | Heauen 12 | winde 5 |
| thee 6 | farre 3 | doth 3 | hence 2 | please 8 | take 4 | Rosin 12 | meanes 5 |
| mine 6 | cold 3 | truly 3 | 'Twere 2 | comes 7 | call 4 | Laertes 11 | man 5 |
| right 6 | absurd 2 | goodly 3 | strangely 2 | pray 7 | thy 3 | Ophelia 11 | Queene 5 |
| particular 5 | flat 2 | Very 3 | sicke 1 | shew 7 | being 3 | Sonne 11 | No 5 |
| strange 5 | pale 2 | poore 3 | twice 1 | selfe 7 | grow 3 | Fathers 11 | thee 5 |
| soft 5 | dull 2 | alone 3 | before 1 | heare 7 | selfe 3 | Giue 10 | vp 5 |
| dead 5 | madnesse 2 | together 3 | vnmanly 1 | go 7 | draw 3 | 'd 10 | nothing 5 |
| selfe 5 | fit 2 | right 3 | still 1 | goes 6 | downe 3 | thee 10 | vs 5 |
| honest 5 | mine 2 | here 3 | meerely 1 | 're 6 | end 3 | Pol 10 | heart 4 |
| vs 5 | Noble 2 | perhaps 3 | Indeed 1 | Take 6 | seeke 3 | eyes 9 | th 4 |
| second 5 | wondrous 2 | quite 3 | coldly 1 | lost 5 | made 3 | day 9 | 't 4 |
| same 4 | violent 2 | twice 2 | together 1 | thy 5 | fell 3 | Which 9 | fault 4 |
| gracious 4 | proper 2 | nightly 2 | exactly 1 | th 5 | beare 3 | head 9 | Horatio 4 |
| farre 4 | hid 2 | n't 2 | stately 1 | stand 5 | put 3 | vp 9 | night 4 |
| late 4 | loose 2 | long 2 | Almost 1 | liue 5 | wrong 2 | purpose 9 | nature 4 |
| glad 4 | tedious 2 | vs 2 | alone 1 | keepe 5 | shewes 2 | youth 9 | Soule 4 |
| oft 4 | capeable 2 | Together 2 | here 1 | vse 5 | Take 2 | blood 9 | Polon 4 |
| sure 4 | heere 2 | Thus 2 | thou 1 | loue 5 | followed 2 | matter 9 | Be 4 |
| fit 4 | good 2 | willingly 2 | shrewdly 1 | does 5 | thine 2 | mine 8 | hast 4 |
| Such 4 | rude 2 | | | | | downe 8 | loue 4 |
| | | | | | | Fortinbras 8 | cause 4 |
| | | | | | | hand 8 | |
| | | | | | | hath 8 | |
| | | | | | | Denmarke 8 | |

Since it is hard to compare the large set of words in the tables via visual inspection, we use Python *sets* to bring out any differences. The syntax is *set1 minus set2 and* would show all elements present in set1 but not set2.

**NOTE:** The prepend *h_<var-name>* is reflective of set comprising words from the Hamlet corpus

There are subtle differences in the words used between the corpora:
- Positive sentiment adjectives:
  print(set(adj_pos) - set(h_adj_pos))
  ```
  {'giuen', 'new', 'enough', 'narrow', 'sir', 'high', 'Ambitious', 'meete', 'gent
  le', 'further', 'wrong', 'Honourable', 'haue', 'ready', 'strong', 'worthy', 'mi
  ghty', 'seene', 'common', 'huge', 'happy'}
  ```

- Negative sentiment adjectives:
  print(set(adj_neg) - set(h_adj_neg))
  ```
  {'better', 'safe', 'least', 'vile', 'bloody', 'sorry', 'feeble', 'constant', 'g
  entle', 'ordinary', 'pitty', 'terrible', 'secret', 'many', 'certaine', 'euery',
  'deceiu', 'Honourable', 'worst', 'stubborne', 'angry', 'senslesse', 'borne', 'm
  uch', 'sad', 'Such', 'impossible', 'seruile', 'sleepe', 'poor', 'third'}
  ```

  *Notice adjectives like bloody, angry, sad etc. a little more prevalent in the Caesar corpus.*

- Nouns would be interesting to compare. We would expect a few different names of protagonists, antagonists, or places as seen below:
  print(set(noun_pos) - set(h_noun_pos))
  ```
  {'Mark', 'Caes', 'Cask', 'Antony', 'morrow', 'Come', 'Caesar', 'Cassi', 'Ant',
  'minde', 'Gods', 'Will', 'Cassius', 'men', 'cause', 'Haue', 'Honor', 'Caska', '
  Octauius', 'Spirit', 'Roman', 'euer', 'Did', 'Brutus', 'Noble', 'Octa', 'things
  ', 'Friends', 'words', 'Rome', 'Thy', 'selfe', 'Bru', 'Caesars', 'Friend'}
  ```

  print(set(h_noun_pos) - set(noun_pos))
  ```
  {'mine', 'Horatio', 'Laertes', 'Qu', 'thing', 'Sir', 'Laer', 'Pol', 'downe', 'F
  athers', 'King', 'Ophelia', 'youth', 'loue', 'matter', 'Hor', 'Hamlet', 'Queene
  ', 'Sonne', 'head', 'vp', 'Heauen', 'purpose', 'life', 'Mother', 'eyes', 'Denma
  rke', 'world', 'Father', 'Nature', 'Giue', 'Fortinbras', 'Rosin', 'Ham', 'Ile'}
  ```

  *Caesar mostly certainly had references to Rome while Hamlet shows references to Denmark, possibly therefore bracing some of that nativity in the plot.*

  *Words like 'Sir' standout in the Hamlet corpus whereas 'Noble' or 'Honor' are prevalent in Caesar.*

- Finally, examining verbs - comparing words with negative sentiment:
  print(set(h_verb_neg) - set(ver_neg))
  ```
  {'thy', 'comes', 'made', 'being', 'end', 'downe', 'thine', 'seeke', 'finde', 'g
  row', 'shewes', 'followed', 'makes', 'Take', 'liue', 'does', 'selfe', 'draw', '
  call'}
  ```

  print(set(ver_neg) - set(h_verb_neg))
  ```
  {'keepe', 'see', 'went', 'shake', 'bring', 'thee', 'giue', 'vse', 'g et', 'hold
  ', 'perceiue', 'follow', 'feare', 'loose', 'heare', 'thinkes', 'loues', 'euery'
  , 'doth', 'loue', 'write'}
  ```

*To corroborate any verbs with adjectives like bloody, angry etc. it seemed worthwhile explore any differences here. Nothing apparent.*

## 6. Conclusion

The case-study at the outset was to analyze two popular Shakespearean plays Julius Caesar and Hamlet. The goal was to be objective and let *Natural Language Processing* analyze corpora and bring out any nuances in style of language, grammar, or sentiment; albeit when the plays were popularly known to be of the same genre.

Turns out through work presented across the past few weeks (HW-1 and HW-2 included), we did not find anything glaringly different and summarize the findings as follows:
- Shakespeare seems to have largely used similar constructs on language, grammar, or sentiment.
- Sentiment was analyzed to be largely neutral but subtly inclining towards a positive polarity and that is particularly interesting given the tragic genre of both plays. It also indicates his style of having an undercurrent of humor or sarcasm to every situation that made his plays popular and entertaining (to both watch or read).
- Finally, we did find Hamlet to be a larger corpus which is in line with any research one would do on the internet.

## Appendix

*Extras - Demonstrate ability to write classification results to a csv file*

The csv will encompass data with the following headings:

**Sentence, pos-sent, neg-sent, neutral-sent**

A function *create_predictions_csv()* was written which given a book-name, associated classifier, sentences, and a feature-set and would correspondingly run the classifier and capture predictions.

The sentence along with a counter to indicate if it was a pos-sent, neg-sent or neutral-sent would be appended to a data-frame df_predict which is eventually written to a .csv called predictions.csv

**NOTE:**
Since results earlier showed better accuracy with bigrams at ~70%, we will use that classifier here.

Here is the code:

```python
 1  # Write classifier predictions to a csv
 2  # Form a test_data with bottom 20% of sentences
 3
 4  def create_predictions_csv(book_name, classifer, book_sentences, bi_features):
 5      test_len = int(0.2 * len(book_sentences))
 6      sentences = book_sentences[-test_len:]
 7
 8      df_predict = pd.DataFrame(columns=['Sentence', 'Book', 'pos-sent', 'neg-sent', 'neutral-sent'])
 9
10      for sents in sentences:
11          senti = caesar_classifier2.classify(bi_document_features(nltk.word_tokenize(sents), bi_features))
12          #adding items to the counter as they are classified
13          if senti.lower() == 'positive':
14              df_predict = df_predict.append({'Sentence': sents,
15                                              'Book': book_name,
16                                              'pos-sent': 1,
17                                              'neg-sent': 0,
18                                              'neutral-sent': 0}, ignore_index=True)
19
20          elif senti.lower() == 'negative':
21              df_predict = df_predict.append({'Sentence': sents,
22                                              'Book': book_name,
23                                              'pos-sent': 0,
24                                              'neg-sent': 1,
25                                              'neutral-sent': 0}, ignore_index=True)
26
27          else:
28              df_predict = df_predict.append({'Sentence': sents,
29                                              'Book': book_name,
30                                              'pos-sent': 0,
31                                              'neg-sent': 0,
32                                              'neutral-sent': 1}, ignore_index=True)
33
34      df_predict.to_csv('predictions.csv', index=False, mode='a')
35      print(df_predict.head())
36
37  print('----Caesar predictions ----')
38  create_predictions_csv('Caesar', caesar_classifier2, caesar_sentences,
39                         caesar_word_bi_features)
40  print('----Hamlet predictions ----')
41  create_predictions_csv('Hamlet', hamlet_classifier2, hamlet_sentences,
42                         hamlet_word_bi_features)
```

See rows 13-32 where we check for a predicted value and compare it with positive, negative, or neutral. Correspondingly we make a row of data and append to the data-frame. As demonstrated earlier we then use the <dataframe>.to_csv() functionality to write the results.

And correspondingly, here are snapshots of the csv comprising data for both corpora:

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | Sentence | Book | pos-sent | neg-sent | neutral-sent |
| 2 | Layest thou thy Leaden Mace vpon my Boy, | Caesar | 0 | 1 | 0 |
| 3 | Gentle knaue good night: | Caesar | 0 | 1 | 0 |
| 4 | Let me see, let me see; is not the Leafe turn'd downe | Caesar | 0 | 1 | 0 |
| 5 | Heere it is I thinke. | Caesar | 0 | 1 | 0 |
| 6 | Enter the Ghost of Caesar. | Caesar | 0 | 0 | 1 |
| 7 | How ill this Taper burnes. | Caesar | 0 | 0 | 1 |
| 8 | Ha! | Caesar | 0 | 0 | 1 |
| 9 | Who comes heere? | Caesar | 0 | 1 | 0 |
| 10 | I thinke it is the weakenesse of mine eyes | Caesar | 0 | 1 | 0 |
| 11 | It comes vpon me: Art thou any thing? | Caesar | 0 | 1 | 0 |
| 12 | Art thou some God, some Angell, or some Diuell, | Caesar | 0 | 1 | 0 |
| 13 | Speake to me, what thou art | Caesar | 0 | 1 | 0 |
| 14 | Thy euill Spirit Brutus? | Caesar | 0 | 0 | 1 |
| 15 | Bru. | Caesar | 0 | 0 | 1 |
| 16 | Why com'st thou? | Caesar | 0 | 0 | 1 |
| 17 | Ghost. | Caesar | 0 | 0 | 1 |
| 18 | To tell thee thou shalt see me at Philippi | Caesar | 0 | 1 | 0 |
| 19 | Well: then I shall see thee againe? | Caesar | 0 | 1 | 0 |
| 20 | Ghost. | Caesar | 0 | 0 | 1 |
| 21 | I, at Philippi | Caesar | 0 | 1 | 0 |

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 321 | So that with ease, | Hamlet | 0 | 0 | 1 |
| 322 | I will doo't. | Hamlet | 0 | 0 | 1 |
| 323 | And for that purpose Ile annoint my Sword: | Hamlet | 0 | 1 | 0 |
| 324 | Let's further thinke of this, | Hamlet | 0 | 0 | 1 |
| 325 | Enter Queene. | Hamlet | 0 | 0 | 1 |
| 326 | Queen. | Hamlet | 0 | 0 | 1 |
| 327 | One woe doth tread vpon anothers heele, | Hamlet | 0 | 0 | 1 |
| 328 | Drown'd! | Hamlet | 0 | 0 | 1 |
| 329 | O where? | Hamlet | 0 | 0 | 1 |
| 330 | Queen. | Hamlet | 0 | 0 | 1 |
| 331 | There is a Willow growes aslant a Brooke, | Hamlet | 0 | 0 | 1 |
| 332 | Alas then, is she drown'd? | Hamlet | 0 | 0 | 1 |
| 333 | Queen. | Hamlet | 0 | 0 | 1 |
| 334 | Drown'd, drown'd | Hamlet | 0 | 0 | 1 |
| 335 | Too much of water hast thou poore Ophelia, | Hamlet | 0 | 1 | 0 |