

Assembling Data

```

1 #importing packages
2 import nltk
3 nltk.download('gutenberg')
4 nltk.download('punkt')
5 import re
6 from nltk.stem import WordNetLemmatizer
7 from nltk.tokenize import word_tokenize
8 nltk.download('omw-1.4')
9 nltk.download('wordnet')
10 import pandas as pd
11 import string

[nltk_data] Downloading package gutenberg to /root/nltk_data...
[nltk_data] Unzipping corpora/gutenberg.zip.
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data] Downloading package wordnet to /root/nltk_data...

1 nltk.corpus.gutenberg.fileids()

['austen-emma.txt',
 'austen-persuasion.txt',
 'austen-sense.txt',
 'bible-kjv.txt',
 'blake-poems.txt',
 'bryant-stories.txt',
 'burgess-busterbrown.txt',
 'carroll-alice.txt',
 'chesterton-ball.txt',
 'chesterton-brown.txt',
 'chesterton-thursday.txt',
 'edgeworth-parents.txt',
 'melville-moby_dick.txt',
 'milton-paradise.txt',
 'shakespeare-caesar.txt',
 'shakespeare-hamlet.txt',
 'shakespeare-macbeth.txt',
 'whitman-leaves.txt']

1 #preprocessing
2 #grabbing the Shakespeare works
3
4 caesar = nltk.corpus.gutenberg.raw('shakespeare-caesar.txt')
5 hamlet = nltk.corpus.gutenberg.raw('shakespeare-hamlet.txt')
6 mcbeth = nltk.corpus.gutenberg.raw('shakespeare-macbeth.txt')
7
8 #lowercasing everything
9
10 c_lower = caesar.lower()
11 h_lower = hamlet.lower()
12 m_lower = mcbeth.lower()
13
14 #tokenizing the text blobs
15
16 from nltk.tokenize import sent_tokenize
17 caesar_tok = sent_tokenize(c_lower)
18 hamlet_tok = sent_tokenize(h_lower)
19 mcbeth_tok = sent_tokenize(m_lower)
20
21 #https://www.guru99.com/tokenize-words-sentences-nltk.html
22 #https://www.nltk.org/howto/corpus.html

1 #removing punctuation
2 #caesar_tok = re.sub('[%s]' % re.escape(string.punctuation), '', caesar_tok)
3
4 def remove_punct(str_list):
5     no_punct = []
6
7     for sent in str_list:
8         sent = re.sub('[%s]' % re.escape(string.punctuation), '', sent)
9

```

```

10 no_punct.append(sent)
11 return no_punct

1 caesar_tok = remove_punct(caesar_tok)
2 hamlet_tok = remove_punct(hamlet_tok)
3 mcbeth_tok = remove_punct(mcbeth_tok)

1 print(caesar_tok)

['the tragedie of julius caesar by william shakespeare 1599\n\nnactus primus', 'scoena prima', 'enter flavius murellus and certaine com

1 #checking length
2
3 len(caesar)
4 len(hamlet)
5 len(mcbeth)
6
7 #average length of a play is 125,000 lines (aka sentences)

100351

1 #Make Shakespeare plays into dataframes
2 import pandas as pd
3 ceasar_df = pd.DataFrame(caesar_tok)
4 author = 'Shakespeare'
5 ceasar_df['Author'] = author
6 ceasar_df = ceasar_df.rename(columns={0: 'Text'})
7
8 hamlet_df = pd.DataFrame(hamlet_tok)
9 hamlet_df['Author'] = author
10 hamlet_df = hamlet_df.rename(columns={0: 'Text'})
11
12 mcbeth_df = pd.DataFrame(mcbeth_tok)
13 mcbeth_df['Author'] = author
14 mcbeth_df = mcbeth_df.rename(columns={0: 'Text'})
15
16 print(ceasar_df)

      Text      Author
0  the tragedie of julius caesar by william shake...  Shakespeare
1              scoena prima  Shakespeare
2  enter flavius murellus and certaine commoners ...  Shakespeare
3              flavius  Shakespeare
4  hence home you idle creatures get you home\nis...  Shakespeare
...      ...      ...
1547 within my tent his bones to night shall ly\nmo...  Shakespeare
1548              exeunt  Shakespeare
1549              omnes  Shakespeare
1550              finis  Shakespeare
1551          the tragedie of ivliivs caesar  Shakespeare

[1552 rows x 2 columns]

```

Non-Shakespeare Corpus Creation

```

1 #creating a non-Shakespeare corpus
2
3 bible = nltk.corpus.gutenberg.raw('bible-kjv.txt')
4 milton = nltk.corpus.gutenberg.raw('milton-paradise.txt')
5 blake = nltk.corpus.gutenberg.raw('blake-poems.txt')
6
7 b_lower = bible.lower()
8 mil_lower = milton.lower()
9 bl_lower = blake.lower()
10
11 #tokenizing the text-blob
12
13 bible_tok = sent_tokenize(b_lower)
14 milton_tok = sent_tokenize(mil_lower)
15 blake_tok = sent_tokenize(bl_lower)
16
17 bible_tok = remove_punct(bible_tok)

```

```

18 milton_tok = remove_punct(milton_tok)
19 blake_tok = remove_punct(blake_tok)

1 #making the non-Shakespeare all one dataframe.
2
3 author_2 = 'Other'
4
5 bible_df = pd.DataFrame(bible_tok)
6 bible_df['Author'] = author_2
7 bible_df = bible_df.rename(columns={0: 'Text'})
8
9 milton_df = pd.DataFrame(milton_tok)
10 milton_df['Author'] = author_2
11 milton_df = milton_df.rename(columns={0: 'Text'})
12
13 blake_df = pd.DataFrame(blake_tok)
14 blake_df['Author'] = author_2
15 blake_df = blake_df.rename(columns={0: 'Text'})
16
17 noshakespear = bible_df.append([milton_df, blake_df], ignore_index=True)
18 print(noshakespear)
19 len(noshakespear)

```

	Text	Author
0	the king james bible\n\nthe old testament of t...	Other
1	12 and the earth was without form and void and...	Other
2	and the spirit of god moved upon the face of t...	Other
3	13 and god said let there be light and there w...	Other
4	14 and god saw the light that it was good and ...	Other
...
31994	why a tongue impressd with honey from every wind	Other
31995	why an ear a whirlpool fierce to draw creation...	Other
31996	why a nostril wide inhaling terror trembling ...	Other
31997	why a little curtain of flesh on the bed of ou...	Other
31998	the virgin started from her seat with a shrie...	Other

[31999 rows x 2 columns]

<ipython-input-8-828362e671fb>:17: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future vers
noshakespear = bible_df.append([milton_df, blake_df], ignore_index=True)
31999

```

1 #subsampling No Shakespeare
2
3 import random
4 random.seed(1)
5
6 noshakes_sub = noshakespear.sample(n = 3500, random_state=1)
7 len(noshakes_sub)
8 print(noshakes_sub)
9
10 #https://www.geeksforgeeks.org/how-to-randomly-select-rows-from-pandas-dataframe/
11 #https://www.w3schools.com/python/ref_random_seed.asp
12 #https://stackoverflow.com/questions/46638641/how-to-fix-valueerror-expected-2d-array-got-1d-array-instead-in-sklearn-pyth

```

	Text	Author
13667	3111 i was a reproach among all mine enemies b...	Other
1978	1922 and let the priests also which come near ...	Other
1428	4914 issachar is a strong ass couching down be...	Other
17621	3818 for the grave cannot praise thee death ca...	Other
3203	2637 and they shall fall one upon another as i...	Other
...
26792	219 and the same man had four daughters virgin...	Other
22408	verily i say unto you they\nhave their reward	Other
7435	524 and let it be when thou hearest the sound ...	Other
13147	3918 what time she lifteth up herself on high ...	Other
26311	733 then said the lord to him put off thy shoe...	Other

[3500 rows x 2 columns]

```

1 #composite data frames
2
3 no_ceasar = hamlet_df.append([noshakes_sub, mcbeth_df], ignore_index=True)
4 no_hamlet = ceasar_df.append([noshakes_sub, mcbeth_df], ignore_index=True)
5 no_mcbeth = hamlet_df.append([noshakes_sub, ceasar_df], ignore_index=True)

```

<ipython-input-10-47ee7286166d>:3: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future vers
no_ceasar = hamlet_df.append([noshakes_sub, mcbeth_df], ignore_index=True)

```
<ipython-input-10-47ee7286166d>:4: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future vers
no_hamlet = ceasar_df.append([noshakes_sub, mcbeth_df], ignore_index=True)
<ipython-input-10-47ee7286166d>:5: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future vers
no_mcbeth = hamlet_df.append([noshakes_sub, ceasar_df], ignore_index=True)
```

```
1 display(no_ceasar.shape)
2 no_ceasar["Author"].value_counts()/no_ceasar.shape[0]
3
4 #move these down!
```

```
(7273, 2)
Shakespeare    0.518768
Other           0.481232
Name: Author, dtype: float64
```

```
1 display(no_hamlet.shape)
2 no_hamlet["Author"].value_counts()/no_hamlet.shape[0]
```

```
(6472, 2)
Other           0.540791
Shakespeare    0.459209
Name: Author, dtype: float64
```

```
1 display(no_mcbeth.shape)
2 no_mcbeth["Author"].value_counts()/no_mcbeth.shape[0]
```

```
(7405, 2)
Shakespeare    0.527346
Other           0.472654
Name: Author, dtype: float64
```

N-Grams

Preprocessing

```
1 from nltk.corpus import stopwords
2 nltk.download('stopwords')
3 import string

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.

1 #cleaning text: remove stopwords, remove non-text characters
2 def clean_text(str_list, lemmatize = False):
3     clean_list = []
4
5     for text in str_list:
6         text = re.sub('[^a-zA-Z]', ' ', text)
7         text = re.sub(r'\s+', ' ', text)
8         text = re.sub(r'[\w\s]', '', text)
9         text = re.sub('[%s]' % re.escape(string.punctuation), '', text)
10        words = word_tokenize(text)
11        clean_words = []
12
13        for word in words:
14            if word not in stopwords.words('english'):
15                if (len(word) > 1) and (re.match(r'^\w+$', word)):
16
17                    clean_words.append(word)
18        clean_text = ' '.join(clean_words)
19        clean_list.append(clean_text)
20
21    return clean_list
22
23 #https://datagy.io/python-remove-punctuation-from-string/
24 #https://stackoverflow.com/questions/265960/best-way-to-strip-punctuation-from-a-string

1 c_ngrams = clean_text(ceasar_df['Text'])
2 h_ngrams = clean_text(hamlet_df['Text'])
3 m_ngrams = clean_text(mcbeth_df['Text'])
4 n_ngrams = clean_text(noshakes_sub['Text'])
```

```

5
6 #there's a bunch of random blanks???

1 print(c_ngrams)

['tragedie julius caesar william shakespeare actus primus', 'scoena prima', 'enter flauius murellus certaine commoners ouer stage', 'fla

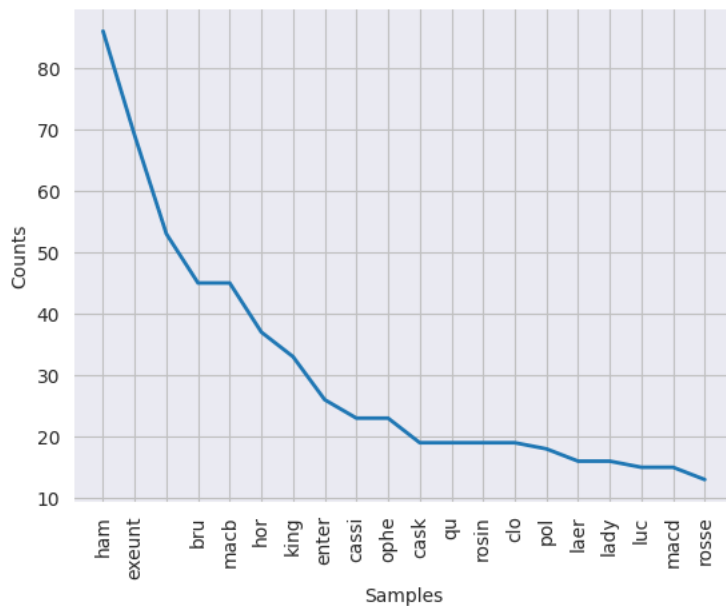
1 #most common words in all Shakespeare
2
3 comb_shakes = ceasar_df.append([hamlet_df, mcbeth_df], ignore_index=True)
4 comb_shakes = clean_text(comb_shakes['Text'])
5
6 from nltk.probability import FreqDist
7
8 fdist_comb_shakes = FreqDist(comb_shakes)
9 fdist_comb_shakes1 = fdist_comb_shakes.most_common(20)
10
11 print(fdist_comb_shakes)
12 print(fdist_comb_shakes1)
13
14 import seaborn as sns
15 sns.set_style('darkgrid')
16 nlp_words=nltk.FreqDist(fdist_comb_shakes)
17 nlp_words.plot(20);

```

```

<ipython-input-17-687fa6e03aba>:1: FutureWarning: The frame.append method is depre
comb_shakes = ceasar_df.append([hamlet_df, mcbeth_df], ignore_index=True)
<FreqDist with 4317 samples and 5325 outcomes>
[('ham', 86), ('exeunt', 69), ('', 53), ('bru', 45), ('macb', 45), ('hor', 37), ('

```



```

1 #counting the most common words in Julius Caesar
2
3 from nltk.probability import FreqDist
4
5 fdist_c = FreqDist(c_ngrams)
6 fdist_c1 = fdist_c.most_common(20)
7 fdist_c
8
9 #https://towardsai.net/p/data-mining/text-mining-in-python-steps-and-examples-78b3f8fd913b

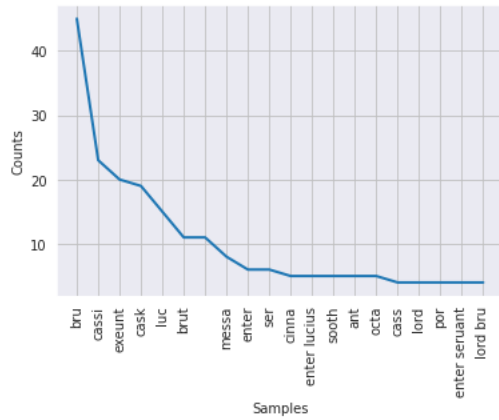
FreqDist({'bru': 45, 'cassi': 23, 'exeunt': 20, 'cask': 19, 'luc': 15, 'brut': 11, '': 11, 'messa': 8, 'enter': 6, 'ser': 6, ...})

1 #graphing the frequency distribution
2
3 import seaborn as sns
4 sns.set_style('darkgrid')
5 nlp_words=nltk.FreqDist(c_ngrams)
6 nlp_words.plot(20);

```

7

8 #https://www.milindsoorya.com/blog/introduction-to-word-frequencies-in-nlp#import-the-needed-libraries



```

1 #counting the most common words in The Tragedy of Hamlet
2
3 fdist_h = FreqDist(h_ngrams)
4 fdist_h1 = fdist_h.most_common(20)
5 fdist_h

```

```

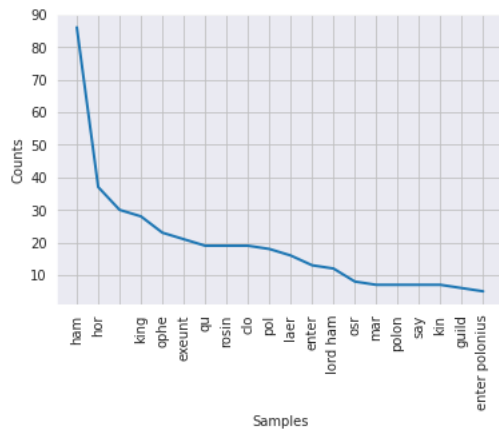
FreqDist({'ham': 86, 'hor': 37, '': 30, 'king': 28, 'ophe': 23, 'exeunt': 21, 'qu': 19, 'rosin': 19, 'clo': 19, 'pol': 18, ...})

```

```

1 sns.set_style('darkgrid')
2 nlp_words=nltk.FreqDist(h_ngrams)
3 nlp_words.plot(20);

```



```

1 #counting the most common words in McBeth
2
3 fdist_m = FreqDist(m_ngrams)
4 fdist_m1 = fdist_m.most_common(20)
5 fdist_m

```

```

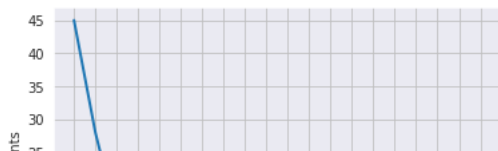
FreqDist({'macb': 45, 'exeunt': 28, 'lady': 15, 'macd': 15, 'rosse': 13, '': 12, 'mal': 10, 'lenox': 9, 'wife': 8, 'enter': 7, ...})

```

```

1 sns.set_style('darkgrid')
2 nlp_words=nltk.FreqDist(m_ngrams)
3 nlp_words.plot(20);

```



```

1 #counting the most common words in Non-Shakespeare
2
3 n_ngrams2 = word_tokenize(str(n_ngrams))
4 fdist_n = FreqDist(n_ngrams2)
5 fdist_n1 = fdist_n.most_common(20)
6 fdist_n

```

```

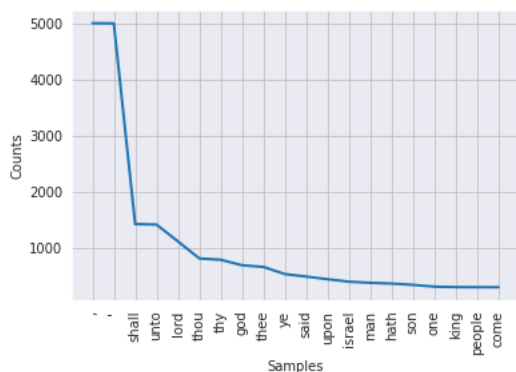
FreqDist({' ': 4999, "'": 4996, 'shall': 1426, 'unto': 1417, 'lord': 1118, 'thou': 815, 'thy': 792, 'god': 694, 'thee': 664, 'ye': 537,
...})

```

```

1 sns.set_style('darkgrid')
2 nlp_words=nltk.FreqDist(n_ngrams2)
3 nlp_words.plot(20);

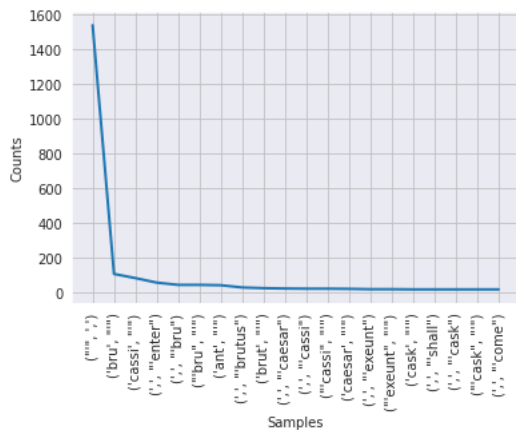
```



```

1 #redoing the frequency distributions with bigrams
2 #Caesar
3
4 c_ngrams2 = word_tokenize(str(c_ngrams))
5 c_bigram = list(nltk.bigrams(c_ngrams2))
6
7 sns.set_style('darkgrid')
8 nlp_words=nltk.FreqDist(c_bigram)
9 nlp_words.plot(20);
10
11 #https://www.tutorialspoint.com/python_text_processing/python_bigrams.htm

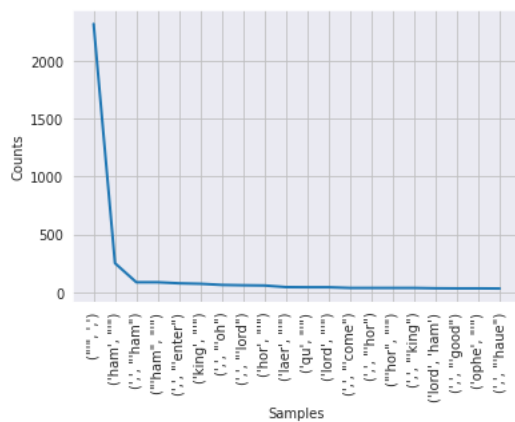
```



```

1 #Hamlet
2
3 h_ngrams2 = word_tokenize(str(h_ngrams))
4 h_bigram = list(nltk.bigrams(h_ngrams2))
5
6 sns.set_style('darkgrid')
7 nlp_words=nltk.FreqDist(h_bigram)
8 nlp_words.plot(20);

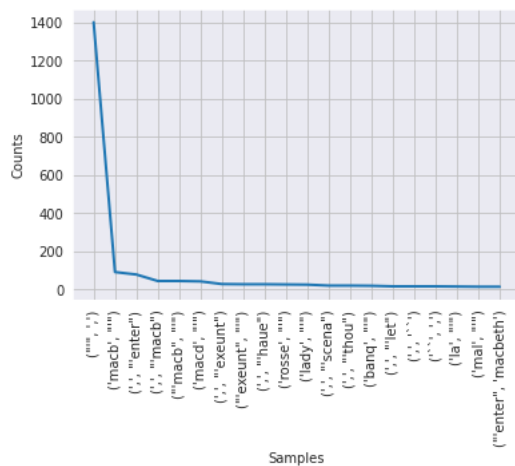
```



```

1 #McBeth
2
3 m_ngrams2 = word_tokenize(str(m_ngrams))
4 m_bigram = list(nltk.bigrams(m_ngrams2))
5
6 sns.set_style('darkgrid')
7 nlp_words=nltk.FreqDist(m_bigram)
8 nlp_words.plot(20);

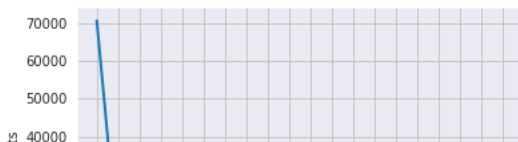
```



```

1 #Non-Shakespeare
2
3 n_ngrams2 = word_tokenize(str(n_ngrams))
4 n_bigram = list(nltk.bigrams(n_ngrams2))
5
6 sns.set_style('darkgrid')
7 nlp_words=nltk.FreqDist(n_bigram)
8 nlp_words.plot(20);

```

Classification Models - Set up

```
1 #we may need to do either one-class training:
2
3 #https://machinelearningmastery.com/one-class-classification-algorithms/
4 #https://towardsdatascience.com/one-class-neural-network-in-keras-249ff56201c0
5 #https://stackoverflow.com/questions/57309958/one-class-classification-using-keras-and-python
6 #https://medium.com/geekculture/sklearn-expects-data-to-be-in-shape-64fbcaf80a8c
7
8 Samples
```

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 from sklearn.feature_extraction import _stop_words
3 from sklearn.model_selection import cross_validate, KFold, GridSearchCV
4
5
6 1 # convert label to a numerical variable
7 no_ceasar['Author'] = no_ceasar['Author'].map({'Shakespeare':1, 'Other':0})
8 no_ceasar.shape
9 print(no_ceasar)
10
11 6 no_hamlet['Author'] = no_hamlet['Author'].map({'Shakespeare':1, 'Other':0})
12 7 no_mcbeth['Author'] = no_mcbeth['Author'].map({'Shakespeare':1, 'Other':0})
```

	Text	Author
0	the tragedie of hamlet by william shakespeare ...	1
1	scoena prima	1
2	enter barnardo and francisco two centinels	1
3	barnardo	1
4	whos there	1
...
7268	this and what need full else\nthat calls vpon ...	1
7269	flourish	1
7270	exeunt omnes	1
7271	finis	1
7272	the tragedie of macbeth	1

[7273 rows x 2 columns]

```
1 #define predictor and response variables
2 c_X = no_ceasar[['Text']]
3 c_y = no_ceasar['Author']
4
5 h_X = no_hamlet[['Text']]
6 h_y = no_hamlet['Author']
7
8 m_X = no_mcbeth[['Text']]
9 m_y = no_mcbeth['Author']
10
11 #define cross-validation method to use
12 cv = KFold(n_splits=10, random_state=1, shuffle=True)
13
14 #https://www.statology.org/k-fold-cross-validation-in-python/
```

```
1 #check proportion of shakespeare and non-shakespeare
```

```
1 #Test Train - The old Fashioned Way
2
3 from sklearn.model_selection import train_test_split
4
5 c_X = no_ceasar.Text
6 c_y = no_ceasar.Author
7 print(c_X.shape, c_y.shape)
8
9 c_X_train, c_X_test, c_y_train, c_y_test = train_test_split(c_X, c_y, random_state=1)
10 print(c_X_train.shape, c_y_train.shape)
11 print(c_X_test.shape, c_y_test.shape)
12
13 #Train/Test and Vectorizing for Hamlet
14
```

```

15 h_X = no_hamlet.Text
16 h_y = no_hamlet.Author
17 print(h_X.shape, h_y.shape)
18
19 h_X_train, h_X_test, h_y_train, h_y_test = train_test_split(h_X, h_y, random_state=1)
20 print(h_X_train.shape, h_y_train.shape)
21 print(h_X_test.shape, h_y_test.shape)
22
23 vect = CountVectorizer()
24 nohamlet_train_vect = vect.fit_transform(h_X_train)
25 nohamlet_test_vect = vect.transform(h_X_test)
26 print(nohamlet_train_vect.shape, nohamlet_test_vect.shape)
27
28 #Train/Test and Vectorizing for McBeth
29
30 m_X = no_mcbeth.Text
31 m_y = no_mcbeth.Author
32 print(m_X.shape, m_y.shape)
33
34 m_X_train, m_X_test, m_y_train, m_y_test = train_test_split(m_X, m_y, random_state=1)
35 print(m_X_train.shape, m_y_train.shape)
36 print(m_X_test.shape, m_y_test.shape)
37
38 vect = CountVectorizer()
39 nomcbeth_train_vect = vect.fit_transform(m_X_train)
40 nomcbeth_test_vect = vect.transform(m_X_test)
41 print(nomcbeth_train_vect.shape, nomcbeth_test_vect.shape)

(7273,) (7273,)
(5454,) (5454,)
(1819,) (1819,)
(6472,) (6472,)
(4854,) (4854,)
(1618,) (1618,)
(4854, 10460) (1618, 10460)
(7405,) (7405,)
(5553,) (5553,)
(1852,) (1852,)
(5553, 11144) (1852, 11144)

1 #Vectorize
2
3 vect = CountVectorizer(min_df=0.29)
4
5 #ceasar
6 noceasar_train_vect = vect.fit_transform(c_X_train)
7 noceasar_test_vect = vect.transform(c_X_test)
8 print(noceasar_train_vect.shape, noceasar_test_vect.shape)
9
10 #hamlet
11 nohamlet_train_vect = vect.fit_transform(h_X_train)
12 nohamlet_test_vect = vect.transform(h_X_test)
13 print(nohamlet_train_vect.shape, nohamlet_test_vect.shape)
14
15 #mcbeth
16 nomcbeth_train_vect = vect.fit_transform(m_X_train)
17 nomcbeth_test_vect = vect.transform(m_X_test)
18 print(nomcbeth_train_vect.shape, nomcbeth_test_vect.shape)
19
20 #consider taking down some of the features pca? but what specifically are we vectorizing here
21 #https://towardsdatascience.com/basics-of-countvectorizer-e26677900f9c

(5454, 3) (1819, 3)
(4854, 3) (1618, 3)
(5553, 3) (1852, 3)

```

Naive Bayes Classification Model

```
1 pip install --upgrade scikit-learn
```

```

Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.9/dist-packages (1.2.2)
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.9/dist-packages (from scikit-learn) (1.22.4)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.9/dist-packages (from scikit-learn) (3.1.0)
Requirement already satisfied: scipy>=1.3.2 in /usr/local/lib/python3.9/dist-packages (from scikit-learn) (1.10.1)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.9/dist-packages (from scikit-learn) (1.2.0)

```

```

1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import confusion_matrix
4 import matplotlib.pyplot as plt
5 import matplotlib as mpl
6 import matplotlib.cm as cm
7 import itertools
8 from sklearn.metrics import roc_auc_score
9 from sklearn import metrics
10 import numpy as np

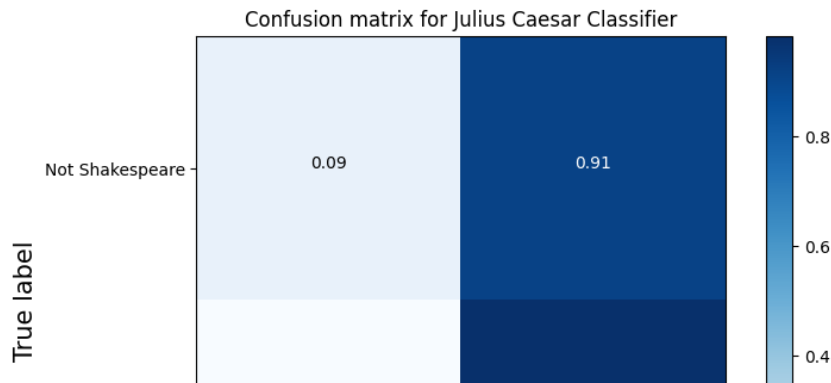
1 def plot_confusion_matrix(cm, classes,
2                           normalize=False,
3                           title='Confusion matrix',
4                           cmap=plt.cm.Blues):
5     if normalize:
6         cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
7
8     plt.imshow(cm, interpolation='nearest', cmap=cmap)
9     plt.title(title)
10    plt.colorbar()
11    tick_marks = np.arange(len(classes))
12    plt.xticks(tick_marks, classes, rotation=45)
13    plt.yticks(tick_marks, classes)
14
15    fmt = '.2f' if normalize else 'd'
16    thresh = cm.max() / 2.
17    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
18        plt.text(j, i, format(cm[i, j], fmt),
19                 horizontalalignment="center",
20                 color="white" if cm[i, j] > thresh else "black")
21
22    plt.tight_layout()
23    plt.ylabel('True label', fontsize=15)
24    plt.xlabel('Predicted label', fontsize=15)

1 #experimental sets
2
3 caesar_experiment = vect.transform(caesar_tok)
4 hamlet_experiment = vect.transform(hamlet_tok)
5 mcbeth_experiment = vect.transform(mcbeth_tok)

1 # Step 3: Train the classifier and predict for test data
2 nb_c = MultinomialNB()
3 nb_c.fit(noceasar_train_vect, c_y_train)
4 noceasar_pred_class = nb_c.predict(noceasar_test_vect)

1 #graph
2
3 ceasar_matrix = confusion_matrix(c_y_test, noceasar_pred_class) #the order probably doesn't matter, but be consistent
4 plt.figure(figsize=(8,6))
5 plot_confusion_matrix(ceasar_matrix, classes=['Not Shakespeare','Shakespeare'],normalize=True,
6                       title='Confusion matrix for Julius Caesar Classifier')
7
8 #check order of first row? should y be first or pred class?

```



```

1 from sklearn.metrics import (accuracy_score, f1_score)
2
3 accuray_c = accuracy_score(noceasar_pred_class, c_y_test)
4 f1_c = f1_score(noceasar_pred_class, c_y_test, average="weighted")
5
6 print("Accuracy:", accuray_c)
7 print("F1 Score:", f1_c)
8
9 #difference in accuracy is because you're calling "weighted", remove average = weighted and see what happens
10 #check why this is really good a finding not shakespeare and not finding anything else
11 #seriously check data balance

```

```

Accuracy: 0.5442550852116548
F1 Score: 0.6606056208769694

```

```

1 #predict
2
3 ceasar_prediction = nb_c.predict(ceasar_experiment)

```

```

1 from numpy import count_nonzero
2
3 count_c_other = np.count_nonzero(ceasar_prediction == 0)
4 count_c_shakes = np.count_nonzero(ceasar_prediction == 1)
5
6 print(count_c_shakes)
7 print(count_c_other)
8
9 #https://www.tutsmake.com/python-count-the-occurrences-in-an-array/
10 #https://www.datacamp.com/tutorial/naive-bayes-scikit-learn

```

```

1541
11

```

```

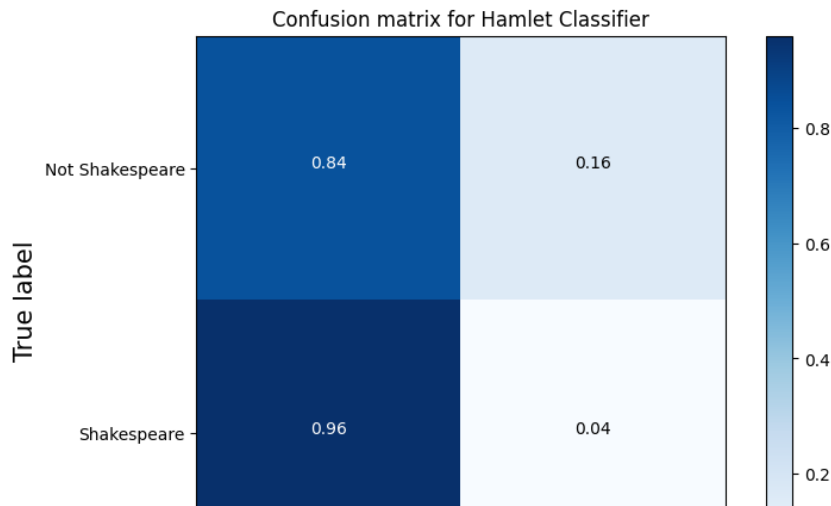
1 #hamlet
2
3 nb_h = MultinomialNB()
4 nb_h.fit(nohamlet_train_vect, h_y_train)
5 nohamlet_pred_class = nb_h.predict(nohamlet_test_vect)

```

```

1 #graph
2
3 hamlet_matrix = confusion_matrix(h_y_test, nohamlet_pred_class)
4 plt.figure(figsize=(8,6))
5 plot_confusion_matrix(hamlet_matrix, classes=['Not Shakespeare','Shakespeare'],normalize=True,
6                       title='Confusion matrix for Hamlet Classifier')

```



```

1 accuray_h = accuracy_score(nohamlet_pred_class, h_y_test)
2 f1_h = f1_score(nohamlet_pred_class, h_y_test, average="weighted")
3
4 print("Accuracy:", accuray_h)
5 print("F1 Score:", f1_h)

```

```

Accuracy: 0.48022249690976515
F1 Score: 0.578537079893798

```

```

1 #predict
2
3 hamlet_prediction = nb_h.predict(hamlet_experiment)
4
5 #at min_def vect = 0.28 or lower, hamlet stops working???

```

```

1 count_h_other = np.count_nonzero(hamlet_prediction == 0)
2 count_h_shakes = np.count_nonzero(hamlet_prediction == 1)
3
4 print(count_h_shakes)
5 print(count_h_other)

```

```

75
2278

```

```

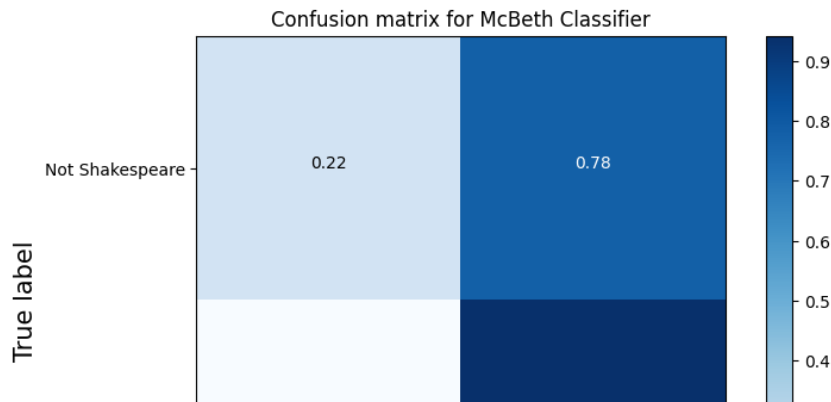
1 #mcbeth
2
3 nb_m = MultinomialNB()
4 nb_m.fit(nomcbeth_train_vect, m_y_train)
5 nomcbeth_pred_class = nb_m.predict(nomcbeth_test_vect)

```

```

1 #graph
2
3 mcbeth_matrix = confusion_matrix(m_y_test, nomcbeth_pred_class)
4 plt.figure(figsize=(8,6))
5 plot_confusion_matrix(mcbeth_matrix, classes=['Not Shakespeare','Shakespeare'],normalize=True,
6                       title='Confusion matrix for McBeth Classifier')

```



```
1 accuracy_m = accuracy_score(nomcbeth_pred_class, m_y_test)
2 f1_m = f1_score(nomcbeth_pred_class, m_y_test, average="weighted")
3
4 print("Accuracy:", accuracy_m)
5 print("F1 Score:", f1_m)
```

```
Accuracy: 0.588012958963283
F1 Score: 0.6497792842010548
```

```
1 #predict
2
3 mcbeth_prediction = nb_m.predict(mcbeth_experiment)
```

```
1 count_m_other = np.count_nonzero(mcbeth_prediction == 0)
2 count_m_shakes = np.count_nonzero(mcbeth_prediction == 1)
3
4 print(count_m_shakes)
5 print(count_m_other)
```

```
1347
73
```

```
1 #metrics?
2
3 auc = nomcbeth_predclass['test_roc_auc']
4 accuracy = nomcbeth_predclass['test_accuracy']
5 f1 = nomcbeth_predclass['test_f1']
6 precision = nomcbeth_predclass['test_precision']
7 recall = nomcbeth_predclass['test_recall']
8 estimators = nomcbeth_predclass['estimator']
9
10 print(accuracy)
11
12 print(precision)
13
14 print(recall)
15
16 print(f1)
```

```
1 #graphs
2
3 #remake the confusion matrixes from below
4 #resource: https://colab.research.google.com/drive/1Hqj4E8SSrq5rMLvKamDPEKB_LikRKPLG
5 #https://stackoverflow.com/questions/49806790/iterable-over-raw-text-documents-expected-string-object-received
6 #https://stackoverflow.com/questions/9884132/what-are-iterator-iterable-and-iteration/9884501#9884501
7 #https://realpython.com/logistic-regression-python/
```

Cross Validation

```
1 m_X2 = no_mcbeth.Text
2 m_y2 = no_mcbeth.Author
3 mx_cv = vect.fit_transform(m_X2)
```

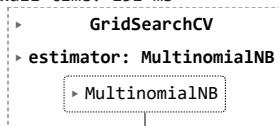
```
1 #mcbet
2 #nomcbeth_train_vect = vect.fit_transform(m_X_train)
```

```

3 #nomcbeth_test_vect = vect.transform(m_X_test)
4 #print(nomcbeth_train_vect.shape, nomcbeth_test_vect.shape)
5
6
7 #print(m_X.shape, m_y.shape)
8
9 from sklearn.model_selection import StratifiedKFold
10
11 params = {}
12
13 #gridsearch searches for the best hyperparameters and keeps the classifier with the highest recall score
14 skf = StratifiedKFold(n_splits=3)
15
16 nb2nb_m = GridSearchCV(MultinomialNB(), cv=skf, param_grid=params)
17 nb2nb_m.fit(mx_cv, m_y2)
18
19 #print(y_scores_nb2)
20
21 #https://stackoverflow.com/questions/51194627/python-naive-bayes-with-cross-validation-using-gaussiannb-classifier

```

CPU times: user 32 ms, sys: 0 ns, total: 32 ms
Wall time: 151 ms



```

1 scores = cross_validate(nb2nb_m,
2                         X = mx_cv,
3                         y = no_mcbeth.Author,
4                         cv = skf,
5                         scoring = ['accuracy', 'f1', 'precision'],
6                         return_estimator = True)

```

```

1 accuracy = scores['test_accuracy']
2 f1 = scores['test_f1']
3 precision = scores['test_precision']
4 #recall = scores['recall']
5
6 print(accuracy)
7 print(f1)
8 print(precision)

```

```

[0.6200891 0.5952188 0.5684765]
[0.72460364 0.69458881 0.70797916]
[0.5865019 0.57694261 0.55029838]

```

```
1 cv_mb_pred = nb2nb_m.predict(mcbeth_experiment)
```

```

1 cv_mb_pred
2
3 count_m_other2 = np.count_nonzero(cv_mb_pred == 0)
4 count_m_shakes2 = np.count_nonzero(cv_mb_pred == 1)
5
6 print(count_m_shakes2)
7 print(count_m_other2)

```

```

1351
69

```

```

1 #Ceasar
2
3 c_X2 = no_ceasar.Text
4 c_y2 = no_ceasar.Author
5 cx_cv = vect.fit_transform(c_X2)
6
7 nb2nb_c = GridSearchCV(MultinomialNB(), cv=skf, param_grid=params)
8 nb2nb_c.fit(cx_cv, c_y2)
9
10 scores = cross_validate(nb2nb_c,
11                         X = cx_cv,
12                         y = no_ceasar.Author,
13                         cv = skf,

```

```

14         scoring = ['accuracy', 'f1', 'precision'],
15         return_estimator = True)
16
17 accuracy = scores['test_accuracy']
18 f1 = scores['test_f1']
19 precision = scores['test_precision']
20 #recall = scores['recall']
21
22 print(accuracy)
23 print(f1)
24 print(precision)

```

```

[0.52618557 0.56683168 0.57343234]
[0.68649386 0.69862227 0.69924375]
[0.52264229 0.54672058 0.55112334]

```

```

1 cv_c_pred = nb2nb_c.predict(caesar_experiment)
2
3 cv_c_pred
4
5 count_c_other2 = np.count_nonzero(cv_c_pred == 0)
6 count_c_shakes2 = np.count_nonzero(cv_c_pred == 1)
7
8 print(count_c_shakes2)
9 print(count_c_other2)

```

```

1543
9

```

```

1 #Hamlet
2
3 h_X2 = no_hamlet.Text
4 h_y2 = no_hamlet.Author
5 hx_cv = vect.fit_transform(h_X2)
6
7 nb2nb_h = GridSearchCV(MultinomialNB(), cv=skf, param_grid=params)
8 nb2nb_h.fit(hx_cv, h_y2)
9
10 scores = cross_validate(nb2nb_h,
11                         X = hx_cv,
12                         y = no_hamlet.Author,
13                         cv = skf,
14                         scoring = ['accuracy', 'f1', 'precision'],
15                         return_estimator = True)
16
17 accuracy = scores['test_accuracy']
18 f1 = scores['test_f1']
19 precision = scores['test_precision']
20 #recall = scores['recall']
21
22 print(accuracy)
23 print(f1)
24 print(precision)
25

```

```

[0.53429101 0.472879    0.45155308]
[0.00593472 0.04694049 0.0497992 ]
[0.15       0.13861386 0.12156863]

```

```

1 cv_h_pred = nb2nb_m.predict(hamlet_experiment)
2
3 cv_h_pred
4
5 count_h_other2 = np.count_nonzero(cv_h_pred == 0)
6 count_h_shakes2 = np.count_nonzero(cv_h_pred == 1)
7
8 print(count_h_shakes2)
9 print(count_h_other2)

```

```

2245
108

```

Resources Used

<https://www.datasciencebytes.com/bytes/2014/12/30/topic-modeling-of-shakespeare-characters/> <https://www.geeksforgeeks.org/adding-new-column-to-existing-dataframe-in-pandas/> <https://www.geeksforgeeks.org/different-ways-to-create-pandas-dataframe/> <https://stackoverflow.com/questions/11346283/renaming-column-names-in-pandas> <https://freelancedatascientist.net/fast-stylometry-tutorial/> <https://towardsdatascience.com/natural-language-processing-count-vectorization-with-scikit-learn-e7804269bb5e> <https://stackoverflow.com/questions/65651544/nameerror-name-plot-confusion-matrix-is-not-defined> <https://stackoverflow.com/questions/67149541/cannot-import-plot-confusion-matrix>

