	Sentiment Analysis Using Naive Bayes Ryan Slipe Grand Canyon University DSC-540-O500 8/23/2023
In [1]:	For this assignment we would like to know if we can use the Naive Bayes theorem to be an accurate classifier on a data set of movie reviews. This dataset was taken from the Kaggle dataset named "Bag of Words Meets Bag of Popcorn". The data consists of 50,000 IMBD movie reviews that do not include any personal information such as the persons' username or adress, or any other identifier that may cause harm. By understanding these are just personal opinions of people that we cannot identify, we deem this data set ethically safe to proceed. Our main goal will be to answer the question, "Can we accurately predict the sentiment of a movie review using Naive Bayes?". We will also try to answer a second question of "Are we able to find words that most commonly belong to each sentiment?". The second question will be answered using the TfidfTransformer that will set weights to each word based on the importance of that word in the corresponding sentiment. We should expect to find positive words in the positive sentiments. import pandas as pd import numpy as np from matplotlib import pyplot as plt
In [4]:	2 7759_3 0 The film starts with a manager (Nicholas Bell) 3 3630_4 0 It must be assumed that those who praised this 4 9495_8 1 Superbly trashy and wondrously unpretentious 8 # check for null df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 25000 entries, 0 to 24999 Data columns (total 3 columns): # Column Non-Null Count Dtype</class>
In [5]: Out[5]:	0 id 25000 non-null object 1 sentiment 25000 non-null int64 2 review 25000 non-null object dtypes: int64(1), object(2) memory usage: 586.1+ KB # view the tail df.tail() id sentiment review 24995 3453_3 0 It seems like more consideration has gone into
In [6]:	24996 5064_1 0 I don't believe they made this film. Completel 24997 10905_3 0 Guy is a loser. Can't get girls, needs to buil 24998 10194_3 0 This 30 minute documentary Buñuel made in the 24999 8478_8 1 I saw this movie as a child and it broke my he We can see that our data is perfectly balanced so we will not have to worry about balancing the data later. This is important because it may make the accuracy of an unbalanced data set seem high when in reality it is not performing very well if one class is much more likely than the other. df['sentiment'].value_counts().sort_index().plot(kind = 'bar', title = 'Distribution of Sentiments')
Out[6]:	CAXES: title={'center': 'Distribution of Sentiments'}> Distribution of Sentiments
Out[7]:	<pre># view value counts for sentiment df['sentiment'].value_counts() 1 12500 0 12500 Name: sentiment, dtype: int64</pre>
In [8]: Out[8]:	df[df['sentiment'] == 1][:10] id sentiment review 0 5814_8
In [14]:	19 9713.10 1 <pre></pre>
Out[14]: In [15]: Out[15]: In [16]:	ragain. Maybe i just want to get a certain insight into this guy who i thought was really cool in the eighties just to maybe make up my mind whether he is guilty or innocent. Moonwalker is part biography, part feature film which i remember going to see at the cinema when it was originally released. Some of it has subtle messages about MJ's fee ling towards the press and also the obvious message of drugs are bad m'kay. Visually impressive but of course this is all about Michael Jackson so unless you remetely like MJ in anyway then you are going to hate this and find it boring. Some may call MJ an egotist for consenting to the making of this movie BUT MJ and most of his fans would say that he made it for the fans which if true is really nice of him. The actual feature film bit when it finally starts is only on for 20 minutes or so excluding the Smooth Cr iminal sequence and Joe Pesci is convincing as a psychopathic all powerful drug lord. Why he wants MJ dead so bad is beyond me. Because MJ overheard his plans? Nah, Joe Pesc i's character ranted that he wanted people to know it is he who is supplying drugs etc so i dunno, maybe he just hates MJ's music. Lots of cool things in this like MJ turn gi nito a car and a robot and the whole Speed Demon sequence. Also, the director must have had the patience of a saint when it came to filming the kiddy Bad sequence as usually directors hate working with one kid let alone a whole bunch of them performing a complex dance scene. Bottom line, this movie is for people who like MJ on one level or which it think is most people). If not, then stay away. It does try and give off a wholesome message and ironically MJ's bestest buddy in this movie is a girl! Micha el Jackson is truly one of the most talented people ever to grace this planet but is he guilty? Well, with all the attention i've gave this subjecthmmm well i don't know because people can be different behind closed doors, i know this for a fact. He is either an extremely nice but stupid guy or one of the most sickest
In [17]:	The first thing we will have to do for our data is 'tokenize' the review. This means we will have to separate the string into separate words. We can do this using the word_tokenize function provided in nltk library. As we see from our example, every word is taken out of the string as its own entity and the punctuation also is included. # example sentence to tokenize example = 'This is an example of word tokenize.' # print the results print(word_tokenize(example)) ['This', 'is', 'an', 'example', 'of', 'word', 'tokenize', '.'] Since we will be using these sentence reviews to classify them as being positive or negative, we can remove words that will appear in both without describing the actual sentence. These are called 'stop words' and we have a few in our example sentence. We can import stop_words from nltk. After importing stopwords we can print the stop words to get an idea of what they are. We can also use the stopwords on our example to see how it is changed. This will clean our data. By looking at some of our reviews we can also see punctuation such as "/" that will have to be cleaned as well. After removing the stop words we should also remove our punctuation since word_tokenize() will separate them like any other word but include them in the tokenized version. This is why we should remove the punctuation before using word_tokenize(). Also, we should do this before using stopwords because they will be working with the tokenized string. In our data we saw there was a \ so after cleaning our data we should make sure to double check those rows to ensure it was removed with the punctuation. If it was not we can set up our own custom stop words.
In [19]:	<pre># set a variable equal to the english stopwords # set a variable equal to the english stopwords # sop_words = set(stopwords.words('english')) # print all english stop words print(stop_words) {"doesn't", 'needn', 'as', 'at', 've', 'some', "isn't", 'mightn', "mustn't", 'nor', 'won', 'of', 'in', 'did', 'above', 'that', 'each', 'had', 'were', 're', 'no', 'itself', 'haven', 'because', 'has', 'our', "it's", 'who', 'after', 'into', 'again', 'same', 'shouldn', 'didn', 'y', 'over', 'both', 'or', "weren't", 'about', "you'll", 'himself', 'it s', 'your', 'during', 'when', 'this', "haven't", 'hadn', 'here', 'having', 'there', 'more', 'them', 'm', 'up', 'she', 'hasn', 'they', 'a', 'l', 'than', 'ourselves', 'doing', 'mustn', 'her', 'those', 'was', 'doin', 'whill, 'aren', 'm', 'doesn', 'gou're", 'too', 'you', 'weren', 'under', "should've", 'are', 'for', 'where', 'hasn', 'they', 'these', 'you'd", 'him', "mightn't", 'down', 'dagainst', 'doesn', "you're", 'too', 'you', 'to', 'the', 'which', 'before', 'my', 'aren', 'worselves', "shouldn't", 'be', 'once', 'should', 'he', 'other', 'all', 'couldn', 'few', 'we', 'll', 'am', "she's", 'now', 'his', 'wasn't", 'bers', 'out', 'me', 'your's', 'thei' r', "won't", 'does', 'such', 'an', 'if', 'while', 'from', "shan't", "that'll", 'it', 'myself', "didn't", 'only', 'isn', 'shan', 'but', 'is', 'been', 'very', 'and', "don't", 'd', 'not', 'why', 'just', 'can', 'have', 'how', 'do', 'until', 'yourself', 'any', "needn't"} # import string to help us remove the punctuation and join letters to form words import string</pre>
	<pre># punctuation removal example punc_example = 'Next, let"s remove the punctuation!' # for loop to remove punctuation punc_example1 = [char for char in punc_example if char not in string.punctuation] # results print(punc_example1) ['N', 'e', 'x', 't', '', 'l', 'e', 't', 's', '', 'r', 'e', 'm', 'o', 'v', 'e', '', 't', 'h', 'e', '', 'p', 'u', 'n', 'c', 't', 'u', 'a', 't', 'i', 'o', 'n'] # join the words together based on the whitespace created from removing punctuation punc_example1 = ''.join(punc_example1) # print results print(punc_example1) Next lets remove the punctuation # use word_tokenize now then use stopwords tokenized = word_tokenize(punc_example1)</pre>
	<pre>tokenized ['Next', 'lets', 'remove', 'the', 'punctuation'] # remove the stop words with a function similar to our punctuation one. ('the') filtered_sentence = [word for word in tokenized if not word.lower() in stop_words] # print the results print(filtered_sentence) ['Next', 'lets', 'remove', 'punctuation']</pre> At this point we tout dots is closed but is not reach for any machine learning elements are also in the result in the result of the sentence of the result o
In [25]:	At this point our text data is cleaned but is not ready for any machine learning analysis. Our machine learning algorithm will have to work with numerical numbers so we will have to create a document term matrix (DTM). This will assign each word to an integer in a matrix throughout our data. We can do this using CountVectorizer from sklearn.feature_extraction.text. We can demonstrate how this works on a single sentence before using it on all of our data. # We can use CountVectorizer so that we are working with numbers from sklearn.feature_extraction.text import CountVectorizer Like many other modules in sklearn we will create an object then fit our data to the object then transform our data using the object. After fitting our test example to the CountVectorizer object we can use .vocabulary_ to see what is happening. This has created a dictionary that stores each word as an integer. Our test sentiment has the word test listed two times so the dictionary only has test listed once. In our example, test is assigned to the value 2. We can also see that the stopwords have not been removed, but the word 'a' is gone because CountVectorizer counts words that are atleast two characters in length.
In [26]:	# sentiment example test = ["This is a test that's called test."] # initiate a CountVectorizer object vectorizer = CountVectorizer() # fit the vectorizer on our data vectorizer.fit(test) # print the fitted data (not transformed yet) print(vectorizer.vocabulary_) {'this': 4, 'is': 1, 'test': 2, 'that': 3, 'called': 0} Once the data is transformed using our vectorizer we can see that a matrix is created. It lists the row as 0 and the dictionary values for each corresponding word. Our vectorizer had 5 words and we see that in the first row, which is the entire sentence, the value (0, 2) corresponds to a 2 while all others correspond to a 1. This is the case because the word 'test' appears in the sentence twice. This will be important information for us when we want to figure out which words are more important than others in our sentiments. # transform our data vector = vectorizer.transform(test)
	print(vector) (0, 0) 1 (0, 1) 1 (0, 2) 2 (0, 3) 1 (0, 4) 1 Preprocessing Our Data Now that we have run through some basic examples on what we have to do to the sentiments in our data we can use these techniques to preprocess our data. We went through numerous steps that included: • removing punctuation • joining the words back together • word_tokenize() • remove stopwords We can create a function that lists these steps and pass our dataframe into it.
In [28]:	<pre>def text_preprocessing(review): # remove punctuation punctuation_removal = [char for char in review if char not in string.punctuation] # join the characters back into words based on whitespace punctuation_removal = ''.join(punctuation_removal) # tokenize our sentiments and return them using .lower so they are all formatted equally return [word for word in word_tokenize(punctuation_removal) if word.lower() not in stopwords.words('english')] # this cell has been commented out because of long run-time</pre>
In [30]:	# create new column of processed sentiment in our data frame by apply and our function # df['preprocesed'] = df.iloc[:, 2].apply(text_preprocessing) After creating the function we can pass our column through using .apply() but since we are working with a large dataset it will take a long time making it computationally expensive. Luckily for us, we can use RegexpTokenizer from nltk.tokenize to speed up this process. When using CountVectorizer there are parameters for both stop words and tokenizer. We can pass through the english stopwords and create a tokenizer from nltk that will tokenize our words based on the pattern. When using the regexptokenizer we will state that we want our pattern to be all lower case and upper case letters and all numbers. This will keep the puncuation out when tokenizing which is what we did when we made the punctuation_removal for loop. from nltk.tokenize import RegexpTokenizer
In [31]:	This can streamline our preprocessing into just a few steps and run much faster than our apply method. Below we have a sample string that we will remove punctuation, stop words, tokenize, and transform into a document term matrix. First, we set up a RegexpTokenizer object with a pattern that we want to look for. This will remove the punctuation from our string. Then we set up our DTM by using a CountVectorizer and remove stopwords by stating 'english' in that parameter. Then, we use the tokenizer that we established as another parameter. The resulting matrix shows there are no stop words and that machine and learning are counted twice as they should be. # create a sample string string = ['I love machine learning, Machine Learning 101 is fun.'] # set up a tokenizer looking for the pattern of letters or numbers toke = RegexpTokenizer('[A-Za-z0-9]+')
	<pre># create a CountVectorizer using the parameters for stop words and our tokenizer vec = CountVectorizer(stop_words='english', tokenizer= toke.tokenize) # fit the CountVectorizer dtm = vec.fit(string) # view the digits dictionary for the matrix print(dtm.vocabulary_) {'love': 3, 'machine': 4, 'learning': 2, '101': 0, 'fun': 1} C:\Users\ryans\anaconda_3\lib\site-packages\sklearn\feature_extraction\text.py:528: UserWarning: The parameter 'token_pattern' will not be used since 'tokenizer' is not Non e' warnings.warn(</pre>
In [32]:	<pre># create a variable for our matrix matrix = dtm.transform(string) # print the matrix to make sure machine and learning are there twice. print(matrix) (0, 0)</pre>
In [33]: In [34]:	Now we will run a test to see how we tokenize data that is in our dataframe. We select the second row of our review column and tokenize that sentiment. This has removed any punctuation and next we would pass this tokenizer into our CountVectorizer. # run a test on the second row (1 index) data_Test = df.iloc[1]['review'] # create a tokenizer with all letters and numbers data_toke = RegexpTokenizer('[A-Za-z0-9]+') tokens = data_toke.tokenize(data_Test) print(tokens) ['The', 'Classic', 'War', 'of', 'the', 'Worlds', 'by', 'Timothy', 'Hines', 'is', 'a', 'very', 'entertaining', 'film', 'that', 'obviously', 'goes', 'to', 'great', 'effort',
	'and', 'lengths', 'to', 'faithfully', 'recreate', 'H', 'G', 'Wells', 'classic', 'book', 'Mr', 'Hines', 'succeeds', 'in', 'doing', 'so', 'I', 'and', 'those', 'who', 'watche d', 'his', 'film', 'with', 'me', 'appreciated', 'the', 'fact', 'that', 'it', 'was', 'not', 'the', 'standard', 'predictable', 'Hollywood', 'fare', 'that', 'comes', 'out', 'every, 'year', 'e', 'g', 'the', 'spielberg', 'version', 'with', 'Tom', 'Cruise', 'that', 'had', 'only', 'the', 'slightest', 'resemblance', 'to', 'the', 'book', 'Obviously', 'everyone', 'looks', 'for', 'different', 'things', 'in', 'a', 'movie', 'Those', 'who', 'envision', 'themselves', 'as', 'amateur', 'critics', 'look', 'only', 'to', 'criticiz', 'e', 'everything', 'they', 'can', 'Others', 'rate', 'a', 'movie', 'on', 'more', 'important', 'bases', 'like', 'being', 'entertained', 'which', 'is', 'why', 'most', 'people', 'never', 'agree', 'with', 'the', 'critics', 'We', 'enjoyed', 'the', 'effort', 'Mr', 'Hines', 'put', 'into', 'being', 'faithful', 'to', 'H', 'G', 'Wells', 'classic', 'novel', 'and', 'we', 'found', 'it', 'to', 'be', 'very', 'entertaining', 'This', 'made', 'it', 'easy', 'to', 'overlook', 'what', 'the', 'critics', 'perceive', 'to', 'be', 'its', 'sho rtcomings'] After running these tests and examples, we know what is happening to the data. After clarifying the results we are ready to implement this preprocessing on our dataframe. Printing the vocabulary dictionary for our dataframe will not give us much useable information because there are too many words for us to gain any insights from, but we can print it to show that the process is the same.
In [91]:	<pre># establish tokenizer we will be using tokenizer = RegexpTokenizer('[A-Za-z0-9]+') # create CountVectorizer using the stopwords and tokenizer count_vect = CountVectorizer(stop_words = 'english', tokenizer = tokenizer.tokenize) # fit our data to the CountVectorizer our_count_vect = count_vect.fit(df['review']) # view (very long) vocabulary dictionary # print(our_count_vect.vocabulary_)</pre>
	C:\Users\ryans\anaconda_3\lib\site-packages\sklearn\feature_extraction\text.py:528: UserWarning: The parameter 'token_pattern' will not be used since 'tokenizer' is not Non e' warnings.warn(# transform our data dtm = our_count_vect.transform(df['review']) We can print the first sentence to see that it is very long and has the row number 0 along with the word that has been transformed into an integer. Many words have only been used one time. # print the first sentence print(dtm[0:1])
	Columbia
	$P(A B) = \frac{P(B A)P(A)}{P(B)}$ The symbol in the equation means "given" so we are solving the Probability (A given B has already occured). In our example we would multiply the probability that the sentiment is a 1 times the probability of each word in the sentence given the sentiment is a 1. The probability of B A is the probability of the word given the sentiment. In a text classifyer the P(B) will be each of our features which are the independent words so we will have $x_1, x_2, \ldots x_n$ for each word in our phrases This means x_0 can be a word such as "good" and we will find the probability of "good" given the probability that the sentiment is a 1. This probability is used on every sentence in our data, not just the sentence that we are looking at so we would look through all the sentences for the probability of "good" appearing the sentiment was positive. As you can tell, this is a lot of arithmetic even though it is simple arithmetic, errors are likely to occur and many words are likely to be in a review, so the sklearn MultinominalNB comes in handy. The first step in model building will be to train our model. We will use a train test split to train our model and then test that model on the test set. The X_train will consist of 80% of our reviews that have been
In [39]:	<pre>from sklearn.model_selection import train_test_split # split our data for training and testing X_train, X_test, y_train, y_test = train_test_split(dtm, df['sentiment'], test_size = 0.2, random_state = 0) # import the Naive Bayes algorithm from sklearn.naive_bayes import MultinomialNB</pre>
In [42]:	# create a multinomial Naive Bayes object MNB = MultinomialNB() # fit to our data model = MNB.fit(X_train, y_train) # Evaluation and Predictions from sklearn.metrics import accuracy_score, confusion_matrix, ConfusionMatrixDisplay, RocCurveDisplay, roc_curve, roc_auc_score from sklearn.metrics import confusion_matrix # make predictions predicted = MNB.predict(X_test) # generate accuracy score accuracy_score = accuracy_score(predicted, y_test) # print our score print(accuracy_score * 100) 86.86
In [45]:	# confusion matrix confusion_matrix = confusion_matrix(y_test, predicted) cm_display = ConfusionMatrixDisplay(confusion_matrix).plot() - 2250 - 2000 - 1750 - 1500
In [46]:	# classification report from sklearn.metrics import classification_report # show results
	print(classification_report(y_test, predicted)) precision recall f1-score support 0 0.86 0.89 0.87 2548 1 0.88 0.85 0.86 2452 accuracy 0.87 5000 macro avg 0.87 0.87 0.87 5000 weighted avg 0.87 0.87 0.87 5000 The ROC curve is a popular metric when working with binary classifiers like we are in this assignment. The ROC curve is a plot of the true positive rate against the false positive rate. This can also be known as plotting sensitivity vs. 1 - specificity. A random classifier would have a diagonal line through the middle of the graph and the firther we are towards the top left corner, the better the results.
In [47]:	plotting sensitivity vs. 1 - specificity. A random classifier would have a diagonal line through the middle of the graph and the firther we are towards the top left corner, the better the results. # ROC Curve and AUC Score fpr, tpr, _ = roc_curve(y_test, predicted) # AUC auc = roc_auc_score(y_test, predicted) # display results roc_display = RocCurveDisplay(fpr = fpr, tpr = tpr, roc_auc = auc).plot()
	O.8 - 0.0 - 0.4 - 0.4 - 0.5 -
	0.0 - AUC = 0.87 0.0 0.2 0.4 0.6 0.8 1.0 False Positive Rate
<pre>In [48]: Out[48]: In [49]: Out[49]:</pre>	We can use our model to make predictions on reviews that have not been made before. To do this, we pass our model followed by .predict(). In the predict function we will need to pass through the review but make sure that we transform it into a dtm first so we will use our_count_vect() and pass the new review in to it. We begin by making a prediction on a sentence that should be received positively and follow it up with a sentence that should be classified as a negative sentence. As a reminder, 1 = positive and 0 = negative. # make predictions of a new sentiment that should be positive model.predict(our_count_vect.transform(['I love this movie'])) array([1], dtype=int64) # make prediction of a new sentiment that should be negative model.predict(our_count_vect.transform(['I hate this movie'])) array([0], dtype=int64) Trying TFIDF TfidfTransformer will transform our words and assign a weight to each word based on importance. The higher the value, the more important the word. For movie reviews we would assume that words like "movie" and
In [50]: In [59]:	"film" are present in many of the reviews whether they are good or bad sentiments. The tfidf assigns weights to words and words like "movie" and "film" would have a very low weight because they will not be adding much information as to what class they belong in when that word is present. In this sentence we see the word that was given the integer value of 14486 has high significance. We can see what this integer is in a word form. This word turns out to be 'cool' which we would think should be used highly in a positive sentiment so we can check if this was indeed a positive sentiment. (1 = positive) # import the TfidfTransformer from sklearn.feature_extraction.text from sklearn.feature_extraction.text import TfidfTransformer Here we create a variable tfidf_transformer and create the transformer. In the same line we also fit and transform our data by passing through the variable dtm. The variable dtm is our CountVectorizer after we removed the stop words, tokenized the data and fit_transformed the data. The data that was transformed was the 'review' column of our data that we stored as a dataframe called df. # fit and transform our data # fit and transform our data
In [60]: Out[60]:	
Out[61]:	# 'cool' sounds like a word that would indicate a positive sentiment. Check it. df.iloc[0]['sentiment'] Using Tfidf in Our Model To use Tfidf in our model we will need to split the data again this time using the tfidf_transformer as the predictor and the sentiment from our dataframe as the target column. To fit the data we can use the MNB model that we created previously and pass through our new X_train, y_train.
	<pre># split the data using the Tfidf X_train, x_test, y_train, y_test = train_test_split(tfidf_transformer, df['sentiment'], test_size = 0.2, random_state = 0) # fit the model to new splits model = MNB.fit(X_train, y_train) # make predictions predicted = model.predict(X_test) from sklearn import metrics # accuracy score score = metrics.accuracy_score(predicted, y_test)</pre>
In [73]: Out[73]:	
Out[73]: In [75]:	# confusion matrix confusion_matrix = metrics.confusion_matrix(y_test, predicted) cm_display = ConfusionMatrixDisplay(confusion_matrix).plot() - 2250 - 2000 - 1750 - 1500 - 1250 - 1000 - 750
In [85]: In [86]:	Cross- Validation Before submitting our model for deployment we should cross validate it using kfolds. We will choose 10 folds to split up our data into training and testing. This ensures that the model is not learning specific information so that it will perform poorly when tested. The train test splits are done ten different times in this case and the results are fairly similar which means we can trust that our model will perform around the 86% accuracy rate that it got. From sklearn.model_selection import KFold, cross_val_score
In [86]: In [89]: In [90]: Out[90]:	
	Grus, J. (2015). Data Science From Scratch. O'Reilly Media, Inc. Harrington, P. (2012). Machine Learning in Action. Manning Publications Co. Patel, D. (2019). Machine Learning Tutorial Python 12 - K Fold Cross Validation[Video]. https://www.youtube.com/watch?v=gJo0uNL-5Qw&t=2s Latrose, D., Latrose, C. (2015). Data Mining and Predictive Analysis. John Wiley & Sons