

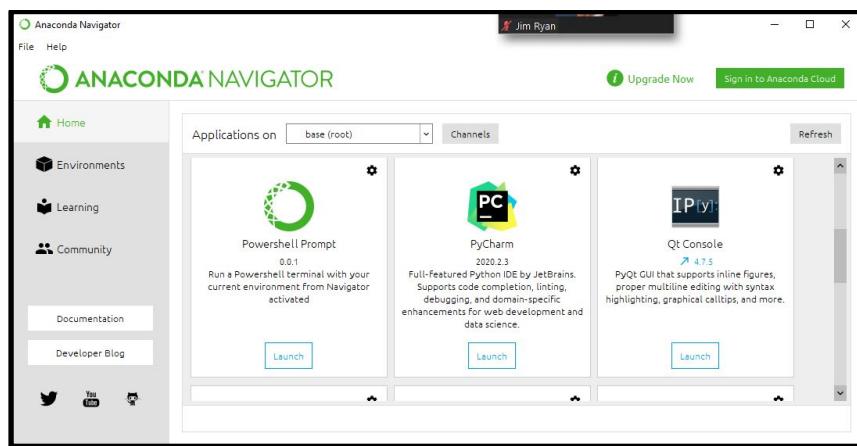
Natural Language Processing (NLP)

Text mining continued via NLP

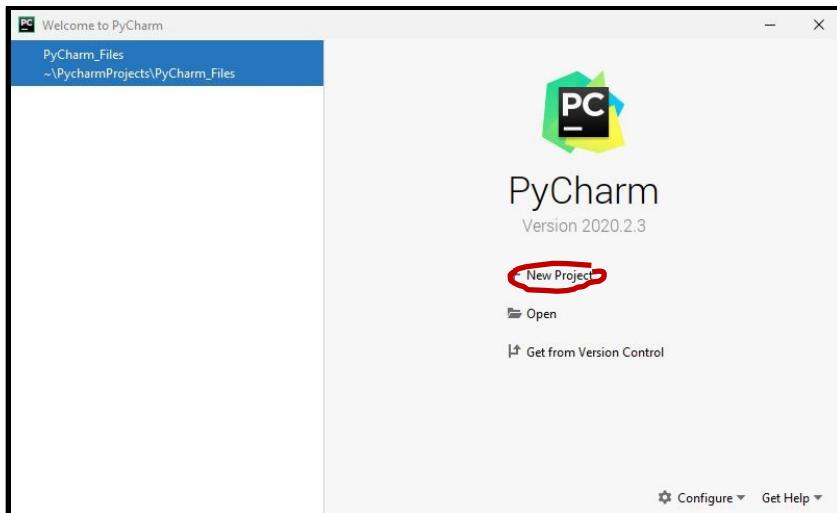
As we saw in Tutorial Eight-B, text mining is a subset of data mining and it implies analyzing data patterns in large batches using one or more software applications. With respect to this tutorial, we will continue our focus on Python data mining processes for text mining via Natural Language Processing (NLP). Text Mining analyzes the text itself, while NLP deals with the underlying/latent metadata to answer questions like - What is the sentiment? - What are the keywords? (using POS tagging & parsers) - What category of content it falls under? - Which are the entities in the sentence? In Tutorial Ten, we will introduce NLP across five sections of Python code in PyCharm: (1) Data Cleaning, (2) Count Vectorizer | TFIDF, (3) Reading a Data Frame using Pandas, and (4) NLP example with Spam data. For more information on NLP outside of this tutorial, please refer to the following URL: <https://medium.com/towards-artificial-intelligence/natural-language-processing-nlp-with-python-tutorial-for-beginners-1f54e610a1a0>

Create the NLP Project in PyCharm

- 0) To ensure your various Python package libraries are available to import, be sure to use Anaconda Navigator to launch PyCharm.

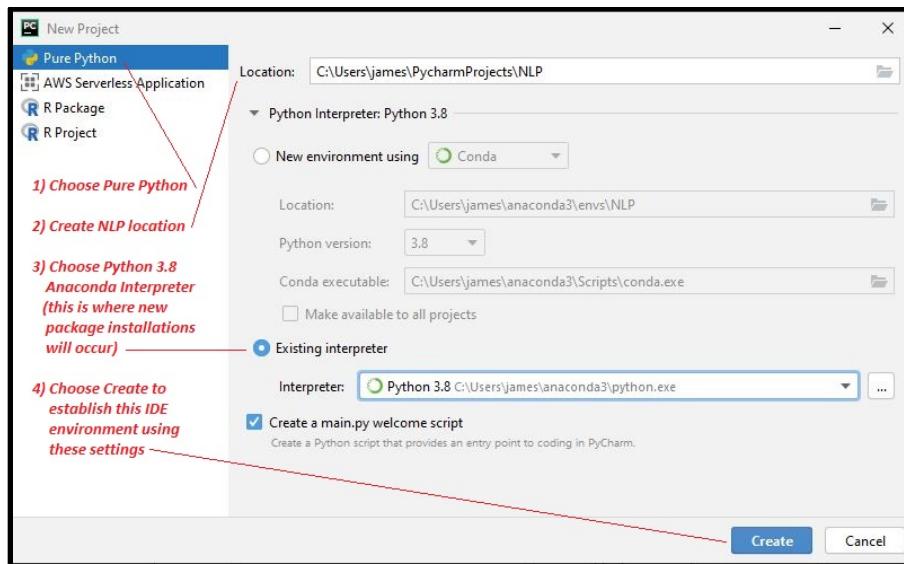


- 1) Launch PyCharm and choose a New Project.

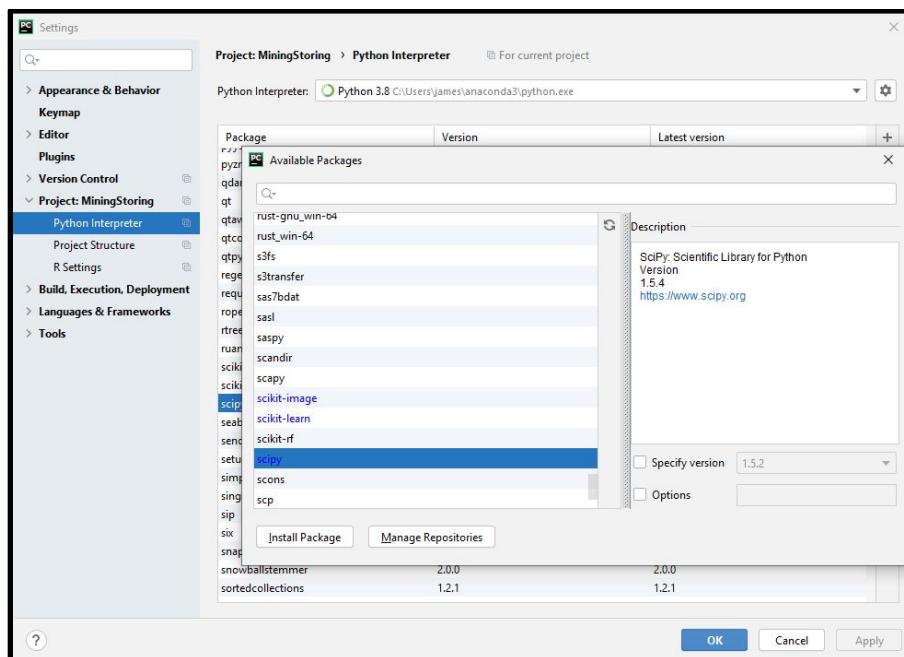


- 1) Choose a Pure Python project,

- 2) Change your project location to create a NLP folder
 - 3) Be sure to choose the Python interpreter as Python 3.8 associated with Anaconda.
 - 4) Choose Create to open PyCharm with these settings.



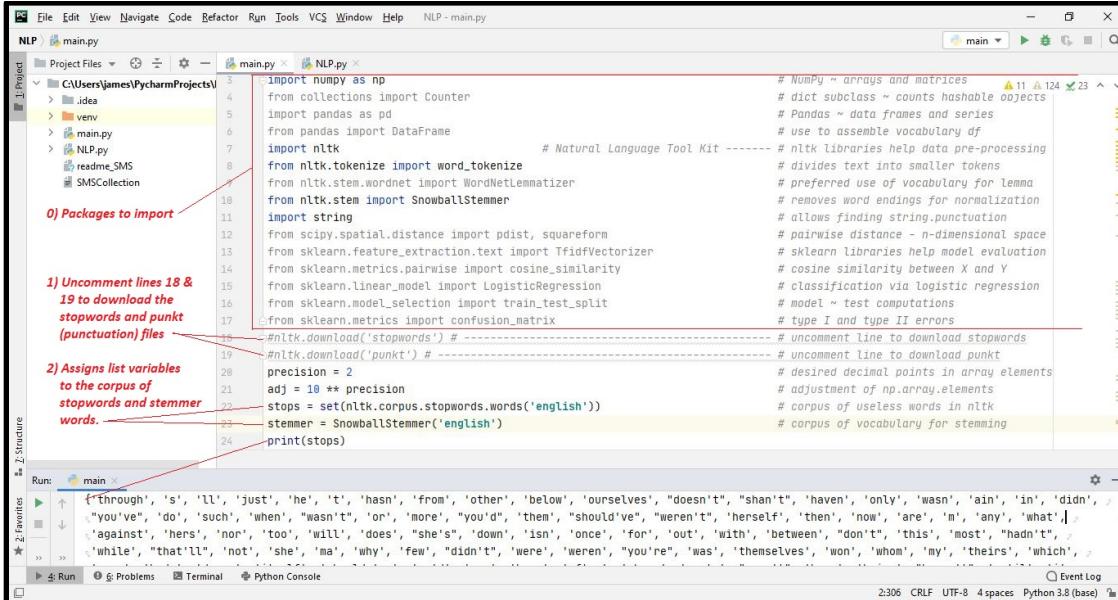
- 5) *Download the SMSCollection and readme_SMS from the Week 12 Canvas Module to the NLP directory created in the last step. The readme_SMS file explains the data in the SMSCollection file.*
 - 6) *When importing a new package library, check the Python Interpreter to ensure the package is installed via Files->Settings->Python Interpreter or left-click Python3.8 in the lower right screen corner choosing Interpreter Settings... Once in the Python Interpreter, pan down to find the particular Python package library and double click on package version. In the resulting window, all installed packages are highlighted in blue.*



Data Cleaning

In this section of Python code, we will use the Natural Language Took Kit (NLTK) as well as the string packages. The PyCharm screenshot below lists Python code for: 0) all the packages used in Tutorial Ten, 1) downloading stopwords and punkt files from NLTK, and 2) assigns the stops list variable to a corpus of English stopwords from NLTK found in the downloaded stopwords file. The last line of code on the screenshot displays all the stops listed in the corpus of stopwords file from NLTK. The Python download code for stopwords and punkt is commented out and the “#” symbol should be removed once to download the files to the Python code directory. Running the download code on lines 17 and 18 more than once will result in unwanted verification notices the files are up-to-date. The punkt file lists punctuation marks to filter out using strings.punctuation in the UDF our_tokenizer(). We introduce stemming and lemmatization in the data cleaning package imports below. For more information on these two topics, please refer to the following URL:

<https://www.datacamp.com/community/tutorials/stemming-lemmatization-python>



```

File Edit View Navigate Code Refactor Run Tools VCS Window Help NLP - main.py
NLP main.py
Project Files C:\Users\James\PycharmProjects\NLP
main.py NLP.py
1) Packages to import
2) Uncomment lines 18 & 19 to download the stopwords and punkt (punctuation) files
3) Assigns list variables to the corpus of stopwords and stemmer words.

0) Packages to import
1) Uncomment lines 18 & 19 to download the stopwords and punkt (punctuation) files
2) Assigns list variables to the corpus of stopwords and stemmer words.

import numpy as np
from collections import Counter
import pandas as pd
from pandas import DataFrame
import nltk
from nltk.tokenize import word_tokenize
from nltk.stem.wordnet import WordNetLemmatizer
from nltk.stem import SnowballStemmer
import string
from scipy.spatial.distance import pdist, squareform
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix

# NumPy ~ arrays and matrices
# dict subclass ~ counts hashable objects
# Pandas ~ data frames and series
# use to assemble vocabulary df
# nltk libraries help data pre-processing
# divides text into smaller tokens
# preferred use of vocabulary for lemma
# removes word endings for normalization
# allows finding string.punctuation
# pairwise distance - n-dimensional space
# sklearn libraries help model evaluation
# cosine similarity between X and Y
# classification via logistic regression
# model ~ test computations
# type I and type II errors

# uncomment line to download stopwords
# nltk.download('stopwords') # -----
# uncomment line to download punkt
# nltk.download('punkt') # -----
precision = 2
adj = 10 ** precision
stops = set(nltk.corpus.stopwords.words('english'))
stemmer = SnowballStemmer('english')
print(stops)

```

All visualizations may be viewed at 400% magnification for clarity.

The next set of Python code in the last data cleaning PyCharm screenshot establishes an user defined function (UDF) named our_tokenizer() that takes each list as a document from the corpus list variable we assign to three sentences and tokenizes each word (doc) in each sentence (corpus[N]) while filtering out punctuation and words listed in stops list of stopwords. These actions clean the corpus data lists and prepares the vocabulary words for vectorization. Specifically, the Python code (3) defines the UDF our_tokenizer() that has parameters list variable, stops, and stemmer; (4) step 1 in the UDF is to convert each doc (word) to lower case; (5) filter out punctuation in each doc; (6) filter out words found in the stops list variable; (7) use stemmer.stem() to reduce a word's conjugation toward normalization; (8) return the tokenized list of words; (9) assigns a list variable corpus to three sentences; (10) call the UDF our_tokenizer for each sentence in corpus; (11) display a corpus sentence; (12) display the

tokenized words; (13) display the tokenized words and check if 'I' is among the stopwords; and (14) display the tokenized words and check if 'drop' is among the stopwords.

```

File Edit View Navigate Code Refactor Run Tools VCS Window Help NLP - main.py
NLP main.py NLP.py
Project Files C:\Users\james\PycharmProject 28 [3] def our_tokenizer(doc, stops=stops, stemmer=stemmer):
    doc = word_tokenize(doc.lower())
    tokens = [''.join([char for char in tok if char not in string.punctuation]) for tok in doc]
    tokens = [tok for tok in tokens if tok]
    if stops:
        tokens = [tok for tok in tokens if (tok not in stops)]
    if stemmer:
        tokens = [stemmer.stem(tok) for tok in tokens]
    return tokens
# user defined function our_tokenizer
# step 1: convert to lower case
# step 2: filter punctuation
# step 3: filter out words from stops
# step 4: stemming cars, car's, cars' >>> car
# a list of documents is a corpus
# one text sentence is a document
# corpus of documents = body of documents
# call our_tokenizer to process corpus
# displays 1st sentence in corpus
# 1st sentence results of tokenized docs
# displays 2nd sentence in corpus
# 2nd sentence results, is I in stops?
# displays 3rd sentence in corpus
# 3rd sentence results, is drop in stops?

[1]  Jim stole my tomato sandwich.
[2]  'Help!', I sobbed, sandwichlessly.
[3]  'Drop the sandwiches!', said the sandwich police.
[4]  print(corpus[0])
[5]  print(tokenized_docs[0])
[6]  print(corpus[1])
[7]  print(tokenized_docs[1], 'Is i a stopword? ', 'i' in stops)
[8]  print(corpus[2])
[9]  print(tokenized_docs[2], 'Is drop a stopword? ', 'drop' in stops)

UDF our_tokenizer() performs the following on each list (doc) in corpus:
step 1: convert each doc to lower case (stopwords are lower case)
step 2: filter out all punctuation found in string.punctuation (punkt)
step 3: filter out all words listed in stops
step 4: reduce word endings (stemmer) for normalization

```

Exercise deliverable: Review the Python code in the prior two PyCharm screenshots. All the Python code listed in these PyCharm screenshots may be downloaded in the [Week 12 file NLP.py](#). Replicate the code exactly as listed, especially the imports and UDF, but create your own sentences to use in the list variable corpus. Capture a PyCharm screenshot showing each sentence of your corpus, followed by the tokenized words.

Count Vectorizer | TFIDF

In this section of Python code, we use the tokenized words (tokenized_docs) to create a sorted vocabulary list of ten unique words. The [ten unique words](#) in the sorted vocabulary list are:

`['drop', 'help', 'jim', 'polic', 'said', 'sandwich', 'sandwich', 'sandwichless', 'sob', 'stole', 'tomato']`

The vocabulary list is the basis for creating NumPy one-dimensional arrays (count vector) that reflect the occurrence of tokenized words in the three documents, where a document is a sentence in the corpus. The Python code adds up the three document arrays to form a count vector. Next, the Python code formulates a term frequency (TF) vector, inverse document frequency (IDF) vector, and their product TFIDF vector. The Python code yields NumPy one-dimensional arrays (vectors) of scores that we can load into many different machine learning algorithms. The different vectors reflect uniqueness of words in the documents and the cosine similarity function from SciKit-Learn Metrics Pairwise package library measures the Euclidean distance between any two of these NumPy array vectors. The following PyCharm screen shot shows Python code for: (0) importing the math libraries for logarithms; (1) creating a vocabulary list by adding the document0, document1, and document2 tokenized_docs lists together; (2) sorting the new vocabulary list; (3) printing the lists and/or arrays; (4) creating a ten element NumPy array for unique words in document0; (5) creating a ten element NumPy array for unique words in document1; (6) creating a ten element NumPy array for unique words in document2; (7) creating a ten element NumPy array that sums the prior three unique word count arrays; and (8) calculating the cosine similarity scores for the three unique word

count arrays (vectors) from the three documents (e.g., 0 vs 1, 0 vs 2, and 1 vs 2). Note document 0 and document 2 have similarity, while the other document combos do not.

Upon calculating the vocabulary word count vector, our next steps are to calculate the term frequency (TF) vectors, the inverse document frequency (IDF) vectors, and use these two vector arrays to calculate the TFIDF vectors. The PyCharm screenshot below contains the Python code to create the TF vectors. The Python code specifically executes: (0) comment statement for the formula to calculate TF vectors for each unique word in each document, (1) assign TF vector array elements for tokenized_docs[0], (2) assign TF vector array elements for tokenized_docs[1], (3) adjust the np.array.elements to two decimal places, (4) assign TF vector array elements for tokenized_docs[2], (5) sum the three document TF arrays to create the vocabulary_TFvector array, (6) calculate the cosine similarity scores comparing the three document TF vector arrays (e.g., 0 vs 1, 0 vs 2, and 1 vs 2), and (7) print out TF results. Note document 0 and document 2 continue the similarity scores while the other document combinations do not.

The screenshot shows the PyCharm IDE interface with the following details:

- Project:** NLP
- File:** main.py
- Code:** Python code for Term Frequency calculation and cosine similarity between TF vectors.
- Output:** The terminal window displays the calculated term frequencies and cosine similarity scores for various document pairs.

```
File Edit View Navigate Code Refactor Run Tool VCS Window Help NLP - main.py
NLP main.py
Project File Structure Run Terminal Python Console
main.py NLP.py
C:\Users\James\PycharmProjects\NLP
.idea
venv
main.py
NLP.py
readme.SMS
SMSCollection
main()
0) # ----- Term Frequency calculation = (# times word appears in document) / (total # of words in document)
1) vocabulary_TFvector0 = np.array([0, 0, 1/4, 0, 0, 1/4, 0, 0, 1/4, 1/4]) # vocabulary_TFvector of unique_tokenized_docs[0]
2) vocabulary_TFvector1 = np.array([0, 1/5, 0, 0, 0, 0, 0, 1/3, 1/3, 0, 0]) # vocabulary_TFvector of unique_tokenized_docs[1]
3) vocabulary_TFvector1 = np.floor(vocabulary_TFvector1*10)/10 # adjust np.array_elements to 2 decimal places
4) vocabulary_TFvector2 = np.array([1/5, 0, 0, 1/5, 1/5, 2/5, 0, 0, 0, 0]) # vocabulary_TFvector of unique_tokenized_docs[2]
5) vocabulary_TFvector0 = vocabulary_TFvector0 + vocabulary_TFvector1 + vocabulary_TFvector2 # addded vocabulary_TF vector
6) vocabulary_TFvector0 = np.floor(vocabulary_TFvector0*10)/10 # adjust np.array_elements to 2 decimal places
7) cs_TS_0vs1 = cosine_similarity([vocabulary_TFvector0,vocabulary_TFvector1]) # cosine similarity array score 0 vs 1
8) cs_TS_0vs2 = cosine_similarity([vocabulary_TFvector0,vocabulary_TFvector2]) # cosine similarity array score 0 vs 2
9) cs_TS_1vs2 = cosine_similarity([vocabulary_TFvector1,vocabulary_TFvector2]) # cosine similarity array score 1 vs 2
10) print(tokenized_docs[0], ' ', vocabulary_TFvector0) # display first sentence words & TF vector
11) print(tokenized_docs[1], ' ', vocabulary_TFvector1) # display second sentence words & TF vector
12) print(tokenized_docs[2], vocabulary_TFvector2) # display third sentence words & TF vector
13) print('vocabulary_TF vector -----> ', vocabulary_TFvector0) # display vocabulary word count vector
14) print(cs_TS_0vs1, '<-----cosine similarity vocabulary_TFvector[0] vs vocabulary_TFvector[1]' ) # cs array scores
15) print(cs_TS_0vs2, '<-----cosine similarity vocabulary_TFvector[0] vs vocabulary_TFvector[2]' ) # cs array scores
16) print(cs_TS_1vs2, '<-----cosine similarity vocabulary_TFvector[1] vs vocabulary_TFvector[2]' ) # cs array scores
main()
[[{"jlm": "stole", "tomato", "Sandwich"}], [0, 0, 0.25, 0, 0.25 0.25], "Calculate term frequency for words in document 0")
[[{"help", "Sob", "sandwichless"}], [0, 0.33, 0, 0, 0, 0.33<0.33>, 0, 0], "Calculate term frequency for words in document 1")
[[{"drop", "sandwich", "said", "sandwich", "polite"}], [0, 0, 0, 0.27<0.27>, 0.4, 0, 0, 0, 0], "Calculate term frequency for words in document 2")
vocabulary_TF vector
[[1, 0, 0.33 0.25 0.2 0.2 0.65 0.33 0.33 0.25 0.25], "Add up three document arrays to generate TF vector
[[0, 1, 0], "cosine similarity ~ document 0 is not similar to document 1")
[[1, 0, 0.37796467], "cosine similarity ~ document 0 similar to document 2")
[[0, 1, 0], "cosine similarity ~ document 1 is not similar to document 2")
Run: main
* D-Structure
* Favorites
* E-Structure
Run: main
[[{"jlm": "stole", "tomato", "Sandwich"}], [0, 0, 0.25, 0, 0.25 0.25], "Calculate term frequency for words in document 0")
[[{"help", "Sob", "sandwichless"}], [0, 0.33, 0, 0, 0, 0.33<0.33>, 0, 0], "Calculate term frequency for words in document 1")
[[{"drop", "sandwich", "said", "sandwich", "polite"}], [0, 0, 0, 0.27<0.27>, 0.4, 0, 0, 0, 0], "Calculate term frequency for words in document 2")
vocabulary_TF vector
[[1, 0, 0.33 0.25 0.2 0.2 0.65 0.33 0.33 0.25 0.25], "Add up three document arrays to generate TF vector
[[0, 1, 0], "cosine similarity ~ document 0 is not similar to document 1")
[[1, 0, 0.37796467], "cosine similarity ~ document 0 similar to document 2")
[[0, 1, 0], "cosine similarity ~ document 1 is not similar to document 2")

```

The following PyCharm screenshot contains the Python code to calculate the document frequency (DF) vector and inverse document frequency (IDF) vectors of the tokenized_docs. Specifically, the Python code executes: (0) a comment statement on the formula to calculate the DF vector array, (1) assigns a NumPy array to vocabulary_DFvector, (2) a comment statement on the formula to calculate the IDF vector arrays, (3) assign IDF vector array elements for tokenized_docs[0], (4) adjust the np.array.elements to two decimal places, (5) assign IDF vector array elements for tokenized_docs[1], (6) assign IDF vector array elements for tokenized_docs[2], (7) sum the three document IDF arrays to create the vocabulary_IDFvector array, (8) calculate the cosine similarity scores comparing the three document IDF vector arrays (e.g., 0 vs 1, 0 vs 2, and 1 vs 2), and (9) print out IDF results. Note document 0 and document 2 continue to have some similarity scores while other document combinations do not.

```

File Edit View Navigate Code Refactor Run Tools VCS Window Help NLP - main.py
NLP > main.py
Project Files > C:\Users\James\PycharmProjects\NLP > main.py > NLP.py
main.py NLP.py
1 # ----- Document Frequency calculation = (# of documents containing word) / (total # of documents)
2 vocabulary_DFvector = np.array([.33,.33,.33,.33,.33,.33,.33,.33]) # vocabulary DF vector of unique tokenized_docs
3 # ----- Inverse Document Frequency calculation = log((total # documents)/(# of documents containing word))
4 IDFV0 = np.array([0,0,math.log(3),0,0,math.log(1.5),0,0,math.log(3)]) # IDF vector ~ tokenized_docs[0]
5 IDFV0 = np.floor(IDFV0 * adj)/adj # adjust np.array.elements to 2 decimal places
6 IDFV1 = np.array([0,math.log(3),0,0,0,0,math.log(3),math.log(3),0,0]) # IDF vector ~ tokenized_docs[1]
7 IDFV1 = np.floor(IDFV1 * adj)/adj # adjust np.array.elements to 2 decimal places
8 IDFV2 = np.array([math.log(3),0,0,math.log(3),math.log(3),math.log(1.5),0,0,0]) # IDF vector ~ tokenized_docs[2]
9 IDFV2 = np.floor(IDFV2 * adj)/adj # adjust np.array.elements to 2 decimal places
10 vocabulary_IDFvector = IDFV0+IDFV1+IDFV2 # vocabulary IDF vector
11 vocabulary_IDFvector = np.floor(vocabulary_IDFvector * adj) # adjust np.array.elements to 2 decimal places
12 cs_IDF_0vs1 = cosine_similarity([IDFV0,IDFV1]) # cosine similarity array score 0 vs 1
13 cs_IDF_0vs2 = cosine_similarity([IDFV0,IDFV2]) # cosine similarity array score 0 vs 2
14 cs_IDF_1vs2 = cosine_similarity([IDFV1,IDFV2]) # cosine similarity array score 1 vs 2
15 print(vocabulary) # display sorted vocabulary list
16 print(tokenized_docs[0],',',IDFV0) # display first sentence words & TF vector
17 print(tokenized_docs[1],',',IDFV1) # display second sentence words & TF vector
18 print(tokenized_docs[2],',',IDFV2) # display third sentence words & TF vector
19 print('vocabulary IDF vector ----->', vocabulary_IDFvector) # display vocabulary word count vector
20 print(cs_IDF_0vs1, '-----cosine similarity IDFV[0] vs IDFV[1]') # cs array scores
21 print(cs_IDF_0vs2, '-----cosine similarity IDFV[0] vs IDFV[2]') # cs array scores
22 print(cs_IDF_1vs2, '-----cosine similarity IDFV[1] vs IDFV[2]') # cs array scores

```

Run: main

Document	Words	IDF Vector	Description
0	drop, help, jim, polic, said, sandwich, sandwich', sandwichless, sob, stole, tomato'	[0. 0. 1.09 0. 0. 0.4 0. 0. 1.09 0.09]	Calculate Inverse Document Frequency for document0
1	sob, sandwich'	[0. 1.09 0. 0. 0. 0. 0. 0. 0. 0.]	Calculate Inverse Document Frequency for document1
2	drop, sandwich, said, sandwich', polic	[0.09 0. 0. 1.09 1.09 0.4 0. 0. 0. 0.]	Calculate Inverse Document Frequency for document2
vocabulary		[1.09 1.09 1.09 1.09 1.09 0.8 1.09 1.09 1.09]	Add up three document arrays for vocabulary IDF vector

[[1. 0.]] -----cosine similarity IDFV[0] vs IDFV[1] cosine similarity score ~ document0 is not similar to document1

[[1. 0.]] -----cosine similarity IDFV[0] vs IDFV[2] cosine similarity score ~ document0 has some similarity to document2

[[1. 0.]] -----cosine similarity IDFV[1] vs IDFV[2] cosine similarity score ~ document1 is not similar to document2

The last PyCharm screenshot for count vectorization identifies Python code to calculate the vocabulary vector array for TFIDF, calculate cosine similarity scores between combinations of the initial vocabulary word count vector and each derived vocabulary vector array (e.g., word count vs TF, word count vs DF, word count vs IDF, and word count vs TFIDF), and then display the results in a summary of count vectorization outcomes and associated cosine similarity. Specifically, the Python code executes: (0) a comment statement on the formula to calculate the TFIDF vector array, (1) calculates the vocabulary TFIDF vector array by multiplying the vocabulary TF vector by the vocabulary IDF vector, (2) adjust the np.array.elements to two decimal places, (3) calculate the cosine similarity scores comparing the vocabulary word count vector array in combination with the derived vocabulary vector arrays (e.g., TF, DF, IDF, and TFIDF), and (4) print out count vectorization and similarity summary results. The annotation on the PyCharm screenshot identifies the uniqueness associated with each tokenized_doc word, based on vocabulary vector arrays for TF and TFIDF. TF and TFIDF vocabulary vector arrays had the most discriminating array elements, where higher is better. Sandwich was the most unique word and said was the least unique word. All four derived vocabulary vector arrays (e.g., TF,

DF, IDF, and TFIDF) show strong similarity scores with TF having the highest and IDF having the lowest.

```

# Project Files
# main.py
# NLP.py

# 0) # TFIDF vector = TF vector * IDF vector
# 1) vocabulary_TFVector = vocabulary_TFvector * vocabulary_IDFvector # TFIDF vector of tokenized_docs
# 2) vocabulary_TFIDFvector = np.floor(vocabulary_TFIDFvector * adj)/adj # adjust np.array.elements to 2 decimal places
# 3) cs_CV_vs_TF = cosine_similarity([vocabulary_vector,vocabulary_TFvector]) # cosine similarity ~ count vs TF
# 4) cs_CV_vs_DF = cosine_similarity([vocabulary_vector,vocabulary_DFvector]) # cosine similarity ~ count vs DF
# 5) cs_CV_vs_IDF = cosine_similarity([vocabulary_vector,vocabulary_IDFvector]) # cosine similarity ~ count vs IDF
# 6) cs_CV_vs_TFIDF = cosine_similarity([vocabulary_vector,vocabulary_TFIDFvector]) # cosine similarity ~ count vs TFIDF
# 7) # display tokenized_docs and vectors
# 8) print(vocabulary)
# 9) print(tokenized_docs[0], vocabulary_vector0) # display first sentence words & vector
# 10) print(tokenized_docs[1], vocabulary_vector1) # display second sentence words & vector
# 11) print(tokenized_docs[2], vocabulary_vector2) # display third sentence words & vector
# 12) print('vocabulary word count vector ----->', vocabulary_word_count_vector) # display vocabulary word count vector
# 13) print('vocabulary TF vector ----->', vocabulary_TFvector) # display vocabulary TF vector array
# 14) print('vocabulary DF vector ----->', vocabulary_DFvector) # display vocabulary DF vector array
# 15) print('vocabulary IDF vector ----->', vocabulary_IDFvector) # display vocabulary IDF vector array
# 16) print('vocabulary TFIDF vector ----->', vocabulary_TFIDFvector) # display vocabulary TFIDF vector
# 17) print(cs_CV_vs_TF, '-----cosine similarity ~ vocabulary count vs vocabulary TF') # cs array scores
# 18) print(cs_CV_vs_DF, '-----cosine similarity ~ vocabulary count vs vocabulary DF') # cs array scores
# 19) print(cs_CV_vs_IDF, '-----cosine similarity ~ vocabulary count vs vocabulary IDF') # cs array scores
# 20) print(cs_CV_vs_TFIDF, '-----cosine similarity ~ vocabulary count vs vocabulary TFIDF') # cs array scores

# Run: main
# Output:
# [drop, 'help', 'jim', 'polic', 'said', 'sandwich', 'sandwich^', 'sandwichless', 'sob', 'stole', 'tomato']
# [^im, 'stole', 'tomato', sandwich] [0 0 1 0 0 1 0 1 1]
# [^im, 'sob', sandwichless] [0 1 0 0 0 0 1 0 0]
# [drop, sandwich, ^said, sandwich, polic] [1 0 0 1 ± 2 0 0 0]
# vocabulary word count vector -----> [1 1 1 1 3 1 1 1]
# vocabulary TF vector -----> [0.2 0.33 0.25 0.2 0.2] [0.65 0.35 0.33 0.25 0.25]
# vocabulary DF vector -----> [0.33 0.33 0.33 0.33 0.33] [0.66 0.33 0.33 0.33 0.33]
# vocabulary IDF vector -----> [1.09 1.09 1.09 1.09 0.8] [1.09 1.09 1.09 1.09 1.09]
# vocabulary TFIDF vector -----> [0.21 0.35 0.27 0.21 0.21] [0.52 0.35 0.35 0.27 0.27]
# [[1. 0.9836604] [1. 0.9836604 1.]] -----cosine similarity ~ vocabulary count vs vocabulary TF
# [[1. 0.9805868] [0.9805868 1.]] -----cosine similarity ~ vocabulary count vs vocabulary DF
# [[1. 0.85488733] [0.85488733 1.]] -----cosine similarity ~ vocabulary count vs vocabulary IDF
# [[1. 0.95993647] [0.95993647 1.]] -----cosine similarity ~ vocabulary count vs vocabulary TFIDF

# Unique words TF TFIDF
# sandwich 0.65 0.52
# help 0.33 0.35
# sandwichless 0.33 0.35
# sob 0.33 0.35
# jim 0.25 0.27
# stole 0.25 0.27
# tomato 0.25 0.27
# drop 0.20 0.21
# polic 0.20 0.21
# said 0.20 0.21

# Most discriminating vector scores TF & TFIDF
# The cosine similarity scores reflect that all vocabulary vector arrays have high similarity with the vocabulary word count vector array.
# Most similarity ~ TF vector | Least similarity ~ IDF vector

```

Exercise deliverable: Review the Python code in the prior four PyCharm screenshots. All the Python code listed in these PyCharm screenshots may be downloaded in the Week 12 file NLP.py. Replicate the code exactly as listed, especially the imports, but use the tokenized_docs from your prior exercise deliverable to build your vocabulary list. Capture a PyCharm screenshot similar to the one above where you summarize your findings.

Reading a text file into a Pandas DataFrame:

We learned about Pandas in prior Tutorials. Again in this tutorial we will look at Data Frames and we will use the SMSCollection.csv file from the text mining example in Tutorial Eight-B. For more specific details about Pandas not covered in this tutorial, please refer to the following URL:

https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.read_table.html

- Anaconda will load the libraries in PyCharm. All the Python code listed in the following PyCharm screenshots may be found in the downloaded Week 12 file NLP_Spam.py. Copy the Python code for the library imports, which are:

```

import numpy as np
import pandas as pd
import nltk
import string
from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
from sklearn.linear_model import LogisticRegression

```

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
#nltk.download('stopwords') # uncomment to download stopwords
#nltk.download('punkt') # uncomment to download punkt
```

- 2) The following Python code formats the run window to handle the size of the csv data file we will be using in Tutorial Ten. You will see this code in the file NLP_Spam.py

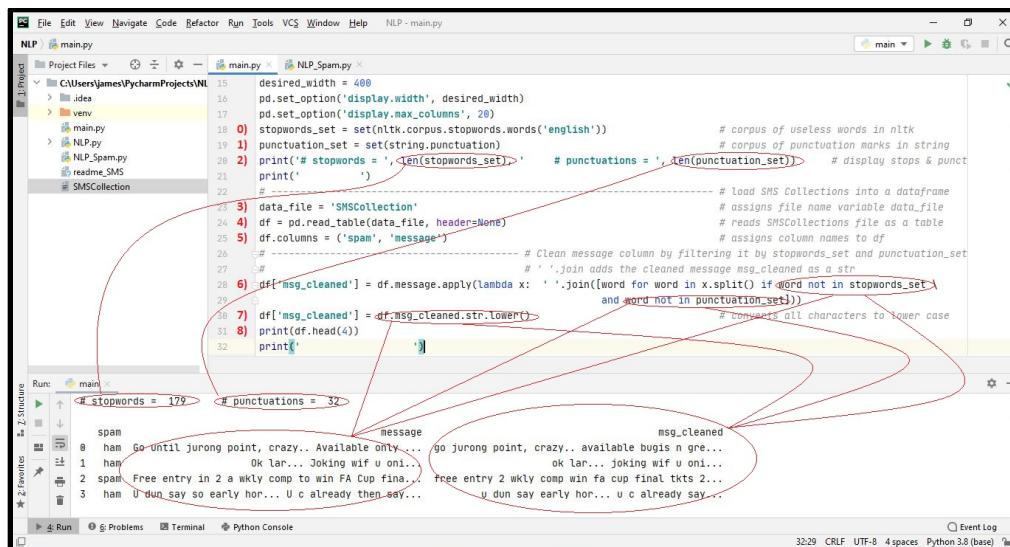
```
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
```

- 3) The following Python code reads the data file SMSCollection into a pandas dataframe as a table. The table does not have a header row, so column labels should be assigned. The Python code is available as file NLP_spam.py with Tutorial Ten.

```
data_file = 'SMSCollection'                      # assigns var data_file
df = pd.read_table(data_file, header=None)        # reads file as a table
df.columns = ('spam', 'message')
```

NLP Example with Spam:

We will continue working with the SMSCollection file we used for the text mining example in Tutorial Eight-B. However, in this Tutorial Ten example we will clean the SMS messages by filtering out the stopwords and punctuation in each SMS message to then use vectorization libraries from sklearn to generate vector arrays (X) to classify the SMS message as spam or ham (y). In this example, we use logistic regression for classification as well as introducing two new classifier functions in RandomForestClassifier() as well as well as GradientBoostingClassifier(). We split the vector arrays (X) using train_test_split from sklearn and use a confusion matrix to review the type I and type II errors. In this NLP example, we tweak NLP classification models by looking at three different vectorization methods combined with three different classification tools. The PyCharm screenshots depict examples of Python code for data cleaning, vectorization, random selection of X,y data points for model training and testing, as well as classification R2 scores and associated type I/type II errors. All Python code in the PyCharm screenshots is available in file NLP_Spam.py.



1. **Data cleaning:** In the preceding PyCharm screenshot, the Python code execution (0) assigns a variable to the set of stopwords from nltk, (1) assigns a variable to the set of punctuation marks from string, (2) displays the number of stopwords in the set and the number of punctuation marks, (3) assigns a file name variable to SMSCollection, (4) reads the file datafile (i.e., SMSCollection) as a table into a Pandas dataframe df, (5) assigns column names of spam and message to the dataframe table df, (6) uses the lambda function to create a new df column msg_cleaned that uses the df column message as the basis and then filters out the stopwords and punctuation, (7) converts each row of the df column msg_cleaned to lower case, and (8) displays the first four rows of the Pandas dataframe df.
2. **Vectorization:** The PyCharm screenshot below shows Python code that executes (0) assigning a variable to the CountVectorizer() function; (1) transforms the df column msg_cleaned into unique word count vectorization and stores it in a NumPy matrix X; (2) displays the type and shape of X; (3) assigns the df column spam as y; (4) randomly splits X and y into model X/y training and test subsets; (5) assigning a variable to the TFIDFVectorizer() function; (6) assigning a variable to the TFIDFVectorizer() function that uses tokenized words, double words, and triple words; (7) transforms the df column msg_cleaned into unique word TFIDF vectorization and stores it in a NumPy matrix X1; (8) transforms the df column msg_cleaned into unique word TFIDF vectorization using bigrams/trigrams and stores it in a NumPy matrix X2; (9) displays the type and shape of X1; (10) displays the type and shape of X2; (11) assign the df column spam rows as y1; (12) randomly splits X1 and y1 into model X1/y1 training and test subsets; and (13) randomly splits X2 and y1 into model X2/y2 training and test subsets For more information on bigrams and trigrams, refer to the following URL: <https://www.geeksforgeeks.org/tf-idf-for-bigrams-trigrams/>

The screenshot shows a PyCharm interface with the following code in main.py:

```

1 df['msg_cleaned'] = df.message.apply(lambda x: ' '.join([word for word in x.split() if word not in stopwords_set and \
2 word not in punctuation_set]))
3 df['msg_cleaned'] = df.msg_cleaned.str.lower() # converts all characters to lower case
4 print(df.head(4))
5
6 # ----- # Use CountVectorizer to create unique word vector array X # -----
7
8 [0] count_vector = CountVectorizer() # assign variable to CountVectorizer()
9 [1] X = count_vector.fit_transform(df.msg_cleaned) # pass df.msg_cleaned for count vectorizer
10 [2] print('count vector X matrix type = ', type(X), 'X matrix shape = ', X.shape) # display type and shape of matrix X
11 [3] y = df.spam # df.spam column is what to predict (spam or ham)
12 [4] X_train, X_test, y_train, y_test = train_test_split(X,y) # split X,y into train/test random subsets
13
14 [5] TFIDF_vector = TfidfVectorizer() # assign var to TfidfVectorizer()
15 [6] TFIDF_vector = TfidfVectorizer(ngram_range=(1,3)) # assign var to Tfidf w/ bigrams & trigrams
16 [7] X1 = TFIDF_vector.fit_transform(df.msg_cleaned) # pass df.msg_cleaned for TFIDF vectorizer
17 [8] X2 = TFIDF_vector1.fit_transform(df.msg_cleaned) # pass df.msg_cleaned for TFIDF w/ bigrams & trigrams
18 [9] print('TFIDF vector X1 matrix type = ', type(X1), 'X1 matrix shape = ', X1.shape) # display type and shape of matrix X1
19 [10] print('TFIDF* vector X2 matrix type = ', type(X2), 'X2 matrix shape = ', X2.shape) # display type and shape of matrix X2
20 [11] y1 = df.spam # df.spam column is what to predict (spam or ham)
21 [12] X1_train, X1_test, y1_train, y1_test = train_test_split(X1,y1) # split X,y into train/test random subsets
22 [13] X2_train, X2_test, y2_train, y2_test = train_test_split(X2,y1)
23 print('')

Number of SMS messages
Number of unique word array vectors

```

The Run tab shows the output of the code execution:

- `count vector X matrix type = <class 'scipy.sparse.csr.csr_matrix'> X matrix shape = (5572, 8703)`
- `TFIDF vector X1 matrix type = <class 'scipy.sparse.csr.csr_matrix'> X1 matrix shape = (5572, 8703)`
- `TFIDF* vector X2 matrix type = <class 'scipy.sparse.csr.csr_matrix'> X2 matrix shape = (5572, 82379)`

3. **Logistic Regression Classification:** In the following PyCharm screenshot, the Python code execution (0) assigns three unique variables the LogisticRegression() function; (1) fits the unique word count vectors (X,y), unique word TFIDF vectors

(X1,y1), and unique word TFIDF vectors using bi/trigrams (X2,y2) using the training data sets; (2) predicts the ham/spam classifications using the three array vectors test subsets; (3) calculates the R2 score for each array vector test subset; (4) displays each models R2 score; (5) calculates the confusion matrix for each array vector predictions; and (6) displays the confusion matrix of type I and type II errors for the array vector predictions. Note: All three of the array vector models yielded high R2 logistic regression scores. However, the highest logistic regression classification R2 score with the minimum type I and type II errors is matrix X that uses the unique word count vectorization.

The screenshot shows the PyCharm IDE interface with the following details:

- Project:** NLP
- File:** main.py
- Code Content:**

```
Project: NLP main.py NLP_Spam.py
C:\Users\James\PycharmProjects\NLP
main.py
venv
NLP.py
NLP_Spam.py
readme_SMS
SMSCollection

[1] 0) logreg = LogisticRegression()
[2] 1) logreg1 = LogisticRegression()
[3] 2) logreg2 = LogisticRegression()
[4] 3) logreg.fit(X_train,y_train)
[5] 4) logreg1.fit(X1_train,y1_train)
[6] 5) logreg2.fit(X2_train,y2_train)
[7] 6) y_pred = logreg.predict(X_test)
[8] 7) y1_pred = logreg1.predict(X1_test)
[9] 8) y2_pred = logreg2.predict(X2_test)
[10] 9) logregX2R = logreg.score(X_test,y_test)
[11] 10) logregX1R2 = logreg1.score(X1_test,y1_test)
[12] 11) logregX2R2 = logreg2.score(X2_test,y2_test)
[13] 12) print('Logistic Regression X R2 -----> ', logregX2R)
[14] 13) print('Logistic Regression X1 R2 -----> ', logregX1R2)
[15] 14) print('Logistic Regression X2 R2 -----> ', logregX2R2)
[16] 15) type_err_matrix = confusion_matrix(y_test,y_pred)
[17] 16) type_err_matrix1 = confusion_matrix(y1_test,y1_pred)
[18] 17) type_err_matrix2 = confusion_matrix(y2_test,y2_pred)
[19] 18) print(type_err_matrix, ' <----- Confusion Matrix [y_test, y_pred]')
[20] 19) print(type_err_matrix1, ' <----- Confusion Matrix [y1_test, y1_pred]')
[21] 20) print(type_err_matrix2, ' <----- Confusion Matrix [y2_test, y2_pred]')
[22] 21) print(' ')
# line separator
```

- Run:** main
- Output:**

 - Highest logistic regression R2
 - Type I errors ~ predicted ham when message was spam
 - Type II errors ~ predicted spam when message as ham

** Best logistic regression model uses count vectorization X,y **

- Event Log:** 21/13 CRLF UTF-8 4 spaces Python 3.8 (base) le

4. ***Random Forest Classification:*** In the following PyCharm screenshot, the Python code execution (0) assigns three unique variables the RandomForestClassifier() function; (1) fits the unique word count vectors (X,y), unique word TFIDF vectors (X1,y1), and unique word TFIDF vectors using bi/trigrams (X2,y2) using the training data sets; (2) predicts the ham/spam classifications using the three array vectors test subsets; (3) calculates the R2 score for each array vector test subset; (4) displays each models R2 score; (5) calculates the confusion matrix for each array vector predictions; and (6) displays the confusion matrix of type I and type II errors for the array vector predictions. Note: All three of the array vector models yielded high R2 random forest classification scores. However, the highest random forest classification R2 score with the minimum type I and type II errors is matrix X2 that uses the unique word TFIDF vectorization. The best model with the highest R2 score and lowest Type I/Type II errors is still the unique word count vectorization using logistic regression classification.

The screenshot shows a PyCharm interface with a project named 'NLP' containing files like main.py, NLP_Spam.py, and SMSCollection. The main.py file contains code for training and testing three different classifiers (RandomForest, Gradient Boosting, and Logistic Regression) on spam and SMS datasets. The execution results in the Run tab show the R2 scores and confusion matrices for each model. The Logistic Regression model (X) consistently shows the highest R2 score and the lowest Type I/Type II errors across all tests.

```

File Edit View Navigate Code Refactor Run Tools VCS Window Help NLP - main.py
NLP main.py
Project Files main.py NLP_Spam.py
CUUsers\james\PycharmProjects\NLP
> idea
> venv
> main.py
> NLP.py
> NLP_Spam.py
> README_SMS
SMSCollection

74 rf = RandomForestClassifier()
75 rf1 = RandomForestClassifier()
76 rf2 = RandomForestClassifier()
77 rf.fit(X_train,y_train)
78 rf1.fit(X1_train,y1_train)
79 rf2.fit(X2_train,y2_train)
80 y3_pred = rf.predict(X_test)
81 y4_pred = rf1.predict(X1_test)
82 y5_pred = rf2.predict(X2_test)
83 print('RandomForest Classifier X R2 -----> ', rf.score(X_test,y_test)) # calculate/display RandomForest X R2
84 print('RandomForest Classifier X1 R2 -----> ', rf1.score(X1_test,y1_test)) # calculate/display RandomForest X1 R2
85 print('RandomForest Classifier X2 R2 -----> ', rf2.score(X2_test,y2_test)) # calculate/display RandomForest X2 R2
86 type_err_matrix3 = confusion_matrix(y_test,y3_pred) # calculate Type I & II errors - X,y
87 type_err_matrix4 = confusion_matrix(y1_test,y4_pred) # calculate Type I & II errors - X1,y1
88 type_err_matrix5 = confusion_matrix(y2_test,y5_pred) # calculate Type I & II errors - X2,y2
89 print(type_err_matrix3, '----- Confusion Matrix [y_test,y3_pred]') # display Type I & II errors - X,y
90 print(type_err_matrix4, '----- Confusion Matrix [y1_test,y4_pred]') # display Type I & II errors - X1,y1
91 print(type_err_matrix5, '----- Confusion Matrix [y2_test,y5_pred]') # display Type I & II errors - X2,y2
92 print('Time: 0.0000000000000001')

Run: main
RandomForestClassifier X R2 -----> 0.9734386216798278
RandomForestClassifier X1 R2 -----> 0.9798994974874372
RandomForestClassifier X2 R2 -----> 0.9662598798724839
Highest RandomForest Classifier R2 uses TFIDF vectorization
Associated Type I/Type II errors
[[11221 1511]
 [11221 1511]
 [11221 1511]] ----- Confusion Matrix [y_test,y3_pred]
[[11221 14611]
 [11221 14611]
 [11221 14611]] ----- Confusion Matrix [y1_test,y4_pred]
[[11221 8]
 [11221 8]
 [11221 8]] ----- Confusion Matrix [y2_test,y5_pred]
** Best model R2 and minimum Type I/Type II errors is still Logistic Classifier X **

75:31 CRLF UTF-8 4 spaces Python 3.8 (base) %

```

5. **Gradient Boosting Classification:** In the following PyCharm screenshot, the Python code execution (0) assigns three unique variables the gradient boosting classification() function; (1) fits the unique word count vectors (X,y), unique word TFIDF vectors (X1,y1), and unique word TFIDF vectors using bi/trigrams (X2,y2) using the training data sets; (2) predicts the ham/spam classifications using the three array vectors test subsets; (3) calculates the R2 score for each array vector test subset; (4) displays each models R2 score; (5) calculates the confusion matrix for each array vector predictions; and (6) displays the confusion matrix of type I and type II errors for the array vector predictions. Note: All three of the array vector models yielded high R2 gradient boosting classification scores. However, the highest gradient boosting classification R2 score with the minimum type I and type II errors is matrix X1 or X2 that uses the unique word TFIDF vectorization or unique word TFIDF using bi/trigrams. The best model with the highest R2 score and lowest Type I/Type II errors is still the unique word count vectorization using logistic regression classification.

The screenshot shows a PyCharm interface with a project named 'NLP' containing files like main.py, NLP_Spam.py, and SMSCollection. The main.py file contains code for training and testing three different classifiers (Gradient Boosting, Random Forest, and Logistic Regression) on spam and SMS datasets. The execution results in the Run tab show the R2 scores and confusion matrices for each model. Similar to the RandomForest example, the Logistic Regression model (X) consistently shows the highest R2 score and the lowest Type I/Type II errors.

```

File Edit View Navigate Code Refactor Run Tools VCS Window Help NLP - main.py
NLP main.py
Project Files main.py NLP_Spam.py
CUUsers\james\PycharmProjects\NLP
> idea
> venv
> main.py
> NLP.py
> NLP_Spam.py
> README_SMS
SMSCollection

93 # ----- run gradient boost classifier with confusion matrix
94 gb = GradientBoostingClassifier()
95 gb1 = GradientBoostingClassifier()
96 gb2 = GradientBoostingClassifier()
97 gb.fit(X_train,y_train)
98 gb1.fit(X1_train,y1_train)
99 gb2.fit(X2_train,y2_train)
100 y6_pred = gb.predict(X_test)
101 y7_pred = gb1.predict(X1_test)
102 y8_pred = gb2.predict(X2_test)
103 print('Gradient Boosting Classifier X R2 -----> ', gb.score(X_test,y_test))
104 print('Gradient Boosting Classifier X1 R2 -----> ', gb1.score(X1_test,y1_test))
105 print('Gradient Boosting Classifier X2 R2 -----> ', gb2.score(X2_test,y2_test))
106 type_err_matrix6 = confusion_matrix(y_test,y6_pred)
107 type_err_matrix7 = confusion_matrix(y1_test,y7_pred)
108 type_err_matrix8 = confusion_matrix(y2_test,y8_pred)
109 print(type_err_matrix6, '----- Confusion Matrix [y_test,y6_pred]')
110 print(type_err_matrix7, '----- Confusion Matrix [y1_test,y7_pred]')
111 print(type_err_matrix8, '----- Confusion Matrix [y2_test,y8_pred]')

Run: main
Gradient Boosting Classifier X R2 -----> 0.966259878724839
Gradient Boosting Classifier X1 R2 -----> 0.968413496951687
Gradient Boosting Classifier X2 R2 -----> 0.968413496951687
[[1175 7]
 [1175 7]
 [1175 7]] ----- Confusion Matrix [y_test,y6_pred]
[[40 171]]
 [[40 171]] ----- Confusion Matrix [y1_test,y7_pred]
[[41 139]]
 [[41 139]] ----- Confusion Matrix [y2_test,y8_pred]
Highest Gradient Boosting Classifier R2 was either X1 or X2 TFIDF vectorization
Lowest Type I/Type II errors were X1 or X2 TFIDF vectorization
** Best model R2 and minimum Type I/Type II errors is still Logistic Classifier X **

94:27 CRLF UTF-8 4 spaces Python 3.8 (base) %

```

6. **Exercise deliverable:** Review the vectorization tools of count, TFIDF, and TFIDF with bi/trigrams, as well as classification tools of logistic regression, random forest classification, and gradient boosting classification. Choose one vectorization and one classification tool on application columns of your project data or on the SMSCollection file. Capture the Python code and results for the data cleaning, vectorization, and classification R2 score/confusion matrix on three separate PyCharm screenshots. Save the PyCharm screenshots for use later in the Tutorial.



Deliverables from Exercises Above:

Take the noted deliverables (i.e., PyCharm screenshots) from each exercise above (i.e., total of 8 screenshots), insert the screenshots into a MS Word document as you label each screenshot deliverable from the corresponding exercise, save the document as one PDF, and submit the PDF to the Tutorial Eight-B Canvas Assignment uplink having the following qualities:

- a. Your screenshots for each of the four exercise deliverables above is included.
- b. Each screenshot is labeled appropriately for each exercise.
 - i. Data Cleaning 1 screenshot
 - ii. Vectorizer | TFIDF 1 screenshot
 - iii. NLP Example with Spam 3 screenshotsTotal PyCharm screenshots to submit 5 screenshots
- c. All screenshots are legible output for each PyCharm exercise deliverable.