

WARNING

Please make sure to "COPY AND EDIT NOTEBOOK" to use compatible library dependencies! DO NOT CREATE A NEW NOTEBOOK AND COPY+PASTE THE CODE - this will use latest Kaggle dependencies at the time you do that, and the code will need to be modified to make it work. Also make sure internet connectivity is enabled on your notebook

✓ Preliminaries

Write requirements to file, anytime you run it, in case you have to go back and recover Kaggle dependencies. **MOST OF THESE REQUIREMENTS WOULD NOT BE NECESSARY FOR LOCAL INSTALLATION**

Latest known such requirements are hosted for each notebook in the companion github repo, and can be pulled down and installed here if needed. Companion github repo is located at <https://github.com/azunre/transfer-learning-for-nlp>

```
1 !pip freeze > kaggle_image_requirements.txt
```

✓ Download IMDB Movie Review Dataset

Download IMDB dataset

```
1 import random
2 import pandas as pd
3
4 ## Read-in the reviews and print some basic descriptions of them
5
6 !wget -q "http://ai.stanford.edu/~amaas/data/sentiment/aclImdb_v1.tar.gz"
7 !tar xzf aclImdb_v1.tar.gz
```

✓ Define Tokenization, Stop-word and Punctuation Removal Functions

Before proceeding, we must decide how many samples to draw from each class. We must also decide the maximum number of tokens per email, and the maximum length of each token. This is done by setting the following overarching hyperparameters

```
1 Nsamp = 1000 # number of samples to generate in each class - 'spam', 'not spam'
2 maxtokens = 200 # the maximum number of tokens per document
3 maxtokenlen = 100 # the maximum length of each token
```

Tokenization

```
1 def tokenize(row):
2     if row is None or row == '':
3         tokens = ""
4     else:
5         tokens = row.split(" ")[ :maxtokens]
6     return tokens
```

Use regular expressions to remove unnecessary characters

Next, we define a function to remove punctuation marks and other nonword characters (using regular expressions) from the emails with the help of the ubiquitous python regex library. In the same step, we truncate all tokens to hyperparameter maxtokenlen defined above.

```
1 import re
2
3 def reg_expressions(row):
4     tokens = []
5     try:
6         for token in row:
7             token = token.lower() # make all characters lower case
8             token = re.sub(r'[\W\d]', "", token)
```

```

9         token = token[:maxtokenlen] # truncate token
10        tokens.append(token)
11    except:
12        token = ""
13        tokens.append(token)
14    return tokens

```

Stop-word removal

Stop-words are also removed. Stop-words are words that are very common in text but offer no useful information that can be used to classify the text. Words such as is, and, the, are are examples of stop-words. The NLTK library contains a list of 127 English stop-words and can be used to filter our tokenized strings.

```

1 import nltk
2
3 nltk.download('stopwords')
4 from nltk.corpus import stopwords
5 stopwords = stopwords.words('english')
6
7 # print(stopwords) # see default stopwords
8 # it may be beneficial to drop negation words from the removal list, as they can change the positive/negative meaning
9 # of a sentence
10 # stopwords.remove("no")
11 # stopwords.remove("nor")
12 # stopwords.remove("not")
13

```

```

[ntlk_data] Downloading package stopwords to /root/nltk_data...
[ntlk_data] Package stopwords is already up-to-date!

```

```

1 def stop_word_removal(row):
2     token = [token for token in row if token not in stopwords]
3     token = filter(None, token)
4     return token

```

✓ Bag-of-words model

For the computer to make inferences of the e-mails, it has to be able to interpret the text by making a numerical representation of it. One way to do this is by using something called a "bag-of-words" model. This model simply counts the frequency of word tokens for each email and thereby represents it as a vector of these counts.

**** Assemble matrices function****

The `assemble_bag()` function assembles a new dataframe containing all the unique words found in the text documents. It counts the word frequency and then returns the new dataframe.

```

1 def assemble_bag(data):
2     used_tokens = []
3     all_tokens = []
4
5     for item in data:
6         for token in item:
7             if token in all_tokens:
8                 if token not in used_tokens:
9                     used_tokens.append(token)
10            else:
11                all_tokens.append(token)
12
13    df = pd.DataFrame(0, index = np.arange(len(data)), columns = used_tokens)
14
15    for i, item in enumerate(data):
16        for token in item:
17            if token in used_tokens:
18                df.iloc[i][token] += 1
19    return df

```

✓ Putting It All Together To Assemble Dataset

Now, putting all the preprocessing steps together we assemble our dataset...

```
1 import os
2 import numpy as np
3
4 # shuffle raw data first
5 def unison_shuffle_data(data, header):
6     p = np.random.permutation(len(header))
7     data = [data[i] for i in p] # Shuffle data as a list
8     header = np.asarray(header)[p]
9     return data, header
10
11 # load data in appropriate form
12 def load_data(path):
13     data, sentiments = [], []
14     for folder, sentiment in (('neg', 0), ('pos', 1)):
15         folder = os.path.join(path, folder)
16         for name in os.listdir(folder):
17             with open(os.path.join(folder, name), 'r') as reader:
18                 text = reader.read()
19                 text = tokenize(text)
20                 text = stop_word_removal(text)
21                 text = reg_expressions(text)
22                 data.append(text)
23                 sentiments.append(sentiment)
24
25     data, sentiments = unison_shuffle_data(data, sentiments)
26     return data, sentiments
27
28 train_path = os.path.join('aclImdb', 'train')
29 test_path = os.path.join('aclImdb', 'test')
30 raw_data, raw_header = load_data(train_path)
31
32 # Aquí ya no intentamos usar np.array con data
33 print(len(raw_data)) # Número total de documentos
34 print(len(raw_header))
35
```

25000
25000

```
1 # Subsample required number of samples
2 random_indices = np.random.choice(range(len(raw_header)), size=(Nsamp*2,), replace=False)
3 data_train = [raw_data[i] for i in random_indices]
4 header = raw_header[random_indices]
5
6 print("DEBUG::data_train::")
7 print(data_train[:5]) # Muestra las primeras 5 filas
```

DEBUG::data_train::
[['first', 'lets', 'agree', 'lorenzo', 'lamas', 'could', 'never', 'considered', 'skilled', 'actor', 'barely', 'even', 'decent', 'sometim

Display sentiments and their frequencies in the dataset, to ensure it is roughly balanced between classes


```
1 unique_elements, counts_elements = np.unique(header, return_counts=True)
2 print("Sentiments and their frequencies:")
3 print(unique_elements)
4 print(counts_elements)
```

Sentiments and their frequencies:
[0 1]
[1002 998]

Featurize and Create Labels


```
1 MixedBagOfReviews = assemble_bag(data_train)
2 # this is the list of words in our bag-of-words model
3 predictors = [column for column in MixedBagOfReviews.columns]
4
5 # expand default pandas display options to make emails more clearly visible when printed
6 pd.set_option('display.max_colwidth', 300)
```

```
7
8 MixedBagOfReviews # you could do print(MixedBagOfReviews), but Jupyter displays this nicer for pandas DataFrames
```




	first	lets	sister	br	actor	acting	lamas	never	the	...	sox	fanbr	uncanny	pupils	hopelessly	brainy	raving	flashy
0	2	2	2	2	3	2	2	2	2	2	...	0	0	0	0	0	0	0
1	0	0	1	0	0	0	0	0	0	1	...	0	0	0	0	0	0	0
2	0	0	1	0	0	0	0	0	1	1	...	0	0	0	0	0	0	0
3	1	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0
4	0	0	0	0	2	0	0	0	1	2	...	0	0	0	0	0	0	0
...
1995	2	0	4	0	0	0	0	0	0	4	...	0	0	0	0	0	0	0
1996	0	0	6	0	6	0	1	0	0	1	...	0	0	0	0	0	0	0
1997	0	0	0	0	0	0	0	0	0	1	...	0	0	0	0	0	0	0
1998	0	0	0	0	1	0	0	0	0	2	...	2	1	1	1	1	1	0
1999	1	0	0	0	0	0	1	0	0	2	...	0	0	0	0	0	0	1

2000 rows x 19598 columns



```
1 # split into independent 70% training and 30% testing sets
2 data = MixedBagOfReviews.values
3
4 idx = int(0.7*data.shape[0])
5
6 # 70% of data for training
7 train_x = data[:idx,:]
8 train_y = header[:idx]
9 # remaining 30% for testing
10 test_x = data[idx:,:]
11 test_y = header[idx:]
12
13 print("train_x/train_y list details, to make sure it is of the right form:")
14 print(len(train_x))
15 print(train_x)
16 print(train_y[:5])
17 print(len(train_y))
```



```
train_x/train_y list details, to make sure it is of the right form:
1400
[[2 2 2 ... 0 0 0]
 [0 0 1 ... 0 0 0]
 [0 0 1 ... 0 0 0]
 ...
 [0 0 1 ... 0 0 0]
 [0 0 1 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
[0 0 0 0 1]
1400
```

✓ How about other vectorization strategies?

We present other vectorization strategies below, for readers who are interested in exploring them...

```
1 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer, HashingVectorizer
2
3 # create the transform - uncomment the one you want to focus on
4 # vectorizer = CountVectorizer() # this is equivalent to the bag of words
5 vectorizer = TfidfVectorizer() # tf-idf vectorizer
6 # vectorizer = HashingVectorizer(n_features=3000) # hashing vectorizer
7
8 # build vocabulary
9 vectorizer.fit([' '.join(sublst) for sublst in data_train])
10 # summarize
11 print(len(vectorizer.vocabulary_))
12 #print(vectorizer.idf_)
13 # encode one document
14 vector = vectorizer.transform([' '.join(data_train[0])])
15 # summarize encoded vector
```

```

8 # Summarize encoded vector
9 print(vector.shape)
10 print(vector.toarray())
11
12 USE = False # set this to 'True' if you want to use the vectorizer featurizers instead of the bag-of-words done before
13 if(USE):
14     data = vectorizer.transform([' '.join(sublst) for sublst in data_train]).toarray()
15     # 70% of data for training
16     train_x = data[idx,: ]
17     # remaining 30% for testing
18     test_x = data[idx[:, :]]
19
20     print("train_x/train_y list details, to make sure it is of the right form:")
21     print(train_x.shape[0])
22     print(train_x)
23     print(train_y[:5])
24     print(len(train_y))
25     predictors = [column for column in vectorizer.vocabulary_]

```

```

➞ 24762
(1, 24762)
[[0. 0. 0. ... 0. 0. 0.]]

```

✓ Logistic Regression Classifier

```

1 from sklearn.linear_model import LogisticRegression
2
3 def fit(train_x,train_y):
4     model = LogisticRegression()
5
6     try:
7         model.fit(train_x, train_y)
8     except:
9         pass
10    return model
11
12 model = fit(train_x,train_y)

```

```

1 predicted_labels = model.predict(test_x)
2
3 # print all labels for full transparency
4 print("DEBUG::The logistic regression predicted labels are::")
5 print(predicted_labels)

```

```

➞ DEBUG::The logistic regression predicted labels are::
[0 1 1 1 0 0 0 1 1 0 1 0 1 1 1 1 0 1 0 1 0 0 1 0 1 0 0 1 1 1 1 1 1 0 0 1 1
 0 1 0 1 1 0 0 1 1 0 0 1 0 1 1 1 1 1 0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 1 0 1 0 0
 0 0 1 1 0 1 0 0 0 1 1 1 1 1 0 1 1 1 1 0 0 1 1 1 0 0 1 0 1 0 1 1 1 1 0 0 1
 0 0 1 1 1 1 0 0 1 1 1 1 0 0 1 1 0 0 1 0 1 1 0 1 0 0 1 0 1 0 0 0 0 1 0 1 0
 0 1 0 0 0 0 1 1 1 1 0 0 1 0 1 0 1 0 0 0 0 0 1 1 1 0 0 0 1 1 0 1 0 0 0 0 1
 1 0 1 1 1 1 0 0 0 0 1 0 1 1 0 1 1 0 0 0 1 1 0 1 1 1 1 0 0 0 0 0 0 1 0 0 0
 1 1 0 0 0 1 1 1 1 0 0 0 1 1 0 0 0 1 1 0 1 0 0 0 1 0 1 1 1 0 0 0 0 1 1 1 1
 1 1 1 1 1 1 1 0 1 0 1 0 1 1 0 1 0 0 1 1 0 0 0 0 1 1 0 0 1 1 1 1 0 0 0 0 0
 1 1 0 0 0 0 0 0 0 0 1 0 0 0 1 1 0 1 1 1 1 1 0 0 1 0 0 1 1 1 1 0 0 1 0 1 0
 1 0 0 0 1 1 0 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 0 1 0 0 0 1 1 1 0 1 0 1 0 0 0
 1 0 1 0 1 1 1 1 0 1 1 0 0 1 1 1 1 0 0 1 1 1 1 0 0 0 1 1 0 0 0 1 0 0 1 0
 0 0 1 1 0 0 1 1 0 1 0 0 1 0 1 1 0 1 1 1 1 1 0 1 0 1 1 0 1 1 0 0 0 0 1 0 0
 1 0 1 0 1 0 0 1 1 1 1 1 0 0 0 1 1 1 1 0 1 0 0 0 1 0 1 0 0 0 0 0 0 0 0 0 1
 0 0 1 0 0 1 1 1 1 1 0 1 0 0 0 1 1 0 0 0 1 1 1 0 1 0 1 1 0 1 1 1 1 0 0
 0 0 1 1 1 1 1 0 0 1 0 0 0 1 0 0 1 0 1 1 1 1 0 0 1 0 0 1 1 0 0 1 1 1 1 1
 0 1 1 1 1 1 0 0 1 0 1 0 0 1 0 1 1 0 0 1 1 1 1 0 0 1 1 1 1 1 0 0 0 0 0 0
 1 0 1 1 1 1 1 0]

```

```

1 from sklearn.metrics import accuracy_score
2
3 acc_score = accuracy_score(test_y, predicted_labels)
4
5 print("The logistic regression accuracy score is::")
6 print(acc_score)

```

```

➞ The logistic regression accuracy score is::
0.7866666666666666

```

✓ Support Vector Machine Classifier

```

1 import time
2 from sklearn.svm import SVC # Support Vector Classification model

3
4 # Create a support vector classifier
5 clf = SVC(C=1, gamma="auto", kernel='linear', probability=False)
6
7 # Fit the classifier using the training data
8 start_time = time.time()
9 clf.fit(train_x, train_y)
10 end_time = time.time()
11 print("Training the SVC Classifier took %3d seconds"%(end_time-start_time))
12
13 # test and evaluate
14 predicted_labels = clf.predict(test_x)
15 print("DEBUG::The SVC Classifier predicted labels are::")
16 print(predicted_labels)
17
18 acc_score = accuracy_score(test_y, predicted_labels)
19 print("The SVC Classifier testing accuracy score is::")
20 print(acc_score)

```

```

↗ Training the SVC Classifier took 22 seconds
DEBUG::The SVC Classifier predicted labels are::
[1 1 1 1 0 0 0 1 1 0 1 0 1 1 1 0 1 1 1 0 0 1 0 0 0 0 1 1 1 1 0 1 0 1 0 0 1 1
 0 0 0 1 1 0 0 1 1 0 0 1 0 0 1 1 1 0 1 0 0 0 1 0 1 1 1 1 0 0 0 1 1 0 1 0 1 0 0
 0 0 1 1 0 0 0 0 0 1 1 1 1 0 0 1 1 1 1 0 0 1 1 1 0 0 1 0 1 0 1 0 1 0 1 1 0 1 1
 0 0 1 1 1 1 0 1 0 1 1 1 0 1 1 1 0 0 1 0 1 1 0 1 0 0 0 0 1 0 0 0 0 1 0 1 0
 0 1 0 0 0 0 1 1 1 0 0 0 1 0 1 0 1 0 0 0 0 0 1 1 1 0 0 0 1 1 0 1 0 0 0 0 1
 1 0 1 1 0 1 1 0 0 0 1 0 1 1 1 1 0 0 0 1 1 0 1 0 1 1 0 0 0 0 1 0 1 0 0 0
 1 1 0 0 0 1 1 1 1 0 0 0 1 1 0 0 0 1 1 0 1 0 0 0 1 0 1 1 0 0 0 0 1 1 1 1
 1 1 1 1 1 1 1 0 1 0 1 0 0 1 0 1 0 0 1 1 0 0 0 0 1 1 0 0 1 0 1 1 0 0 1 0 0
 1 1 0 1 0 0 0 0 1 0 1 0 0 0 1 1 0 1 1 0 1 1 0 0 1 0 0 1 1 1 1 0 0 1 0 1 0
 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1 0 0 0 1 0 1 1 0 1 0 0 0 1 1 1 0 1 1 1 0 0 0
 1 0 1 0 1 1 1 1 0 0 1 0 0 0 1 1 1 0 0 1 1 1 1 0 0 0 1 0 0 0 0 1 1 0 0 0 0
 0 0 1 1 1 0 1 1 0 1 0 0 1 0 1 1 0 1 1 1 1 0 1 0 1 1 0 1 1 0 0 0 0 1 0 0
 1 0 1 0 1 0 0 1 1 0 0 1 0 0 0 1 1 1 1 0 1 0 0 0 1 0 1 0 1 0 0 0 0 0 0 1
 0 0 0 1 0 1 1 1 1 0 1 0 1 0 0 0 1 1 0 0 0 0 1 1 0 1 0 1 1 0 1 0 1 1 1 0 0
 0 0 1 1 1 1 1 0 0 1 0 0 0 0 1 0 1 0 0 1 1 0 1 0 0 1 0 0 1 1 0 0 1 1 1 1
 0 1 1 1 1 1 1 0 0 0 0 0 0 1 0 1 1 0 0 1 1 1 0 0 1 1 1 1 1 0 1 0 0 0 0
 1 0 0 1 1 1 0 0]
The SVC Classifier testing accuracy score is::
0.7533333333333333

```

✓ Random Forests

```

1 # Load scikit's random forest classifier library
2 from sklearn.ensemble import RandomForestClassifier
3
4 # Create a random forest Classifier. By convention, clf means 'Classifier'
5 clf = RandomForestClassifier(n_jobs=1, random_state=0)
6
7 # Train the Classifier to take the training features and learn how they relate
8 # to the training y (spam, not spam?)
9 start_time = time.time()
10 clf.fit(train_x, train_y)
11 end_time = time.time()
12 print("Training the Random Forest Classifier took %3d seconds"%(end_time-start_time))
13
14 predicted_labels = clf.predict(test_x)
15 print("DEBUG::The RF predicted labels are::")
16 print(predicted_labels)
17
18 acc_score = accuracy_score(test_y, predicted_labels)
19
20 print("DEBUG::The RF testing accuracy score is::")
21 print(acc_score)

```

```

↗ Training the Random Forest Classifier took 4 seconds
DEBUG::The RF predicted labels are::
[0 1 1 1 1 0 0 1 1 0 1 0 1 1 1 0 1 0 1 0 0 1 1 1 0 0 1 1 1 1 1 1 0 1 1
 0 0 0 0 1 0 0 1 1 0 0 1 0 0 1 1 1 0 0 0 0 0 0 0 1 0 1 0 0 0 0 1 0 1 1 0 0

```

```

1 0 1 1 1 0 0 0 0 1 1 1 1 0 1 1 1 1 0 0 1 1 1 1 0 1 1 0 1 1 0 1 0 1 1
0 0 1 1 1 1 0 0 0 1 1 1 0 1 1 1 0 0 0 0 0 1 0 1 0 1 1 0 1 0 0 0 0 0 1 1
1 1 0 0 1 1 1 0 1 1 1 0 0 0 1 1 1 0 0 0 0 0 1 1 1 0 0 0 1 1 0 1 0 0 0 1 1
1 0 0 1 0 1 0 0 0 0 1 0 1 1 0 1 1 0 0 0 1 1 0 1 1 1 1 1 0 0 0 0 1 1 0 0 0
1 0 0 0 0 1 1 1 1 0 1 1 1 0 0 1 1 0 0 1 0 0 0 0 1 0 1 1 0 0 0 0 1 1 1 0
1 1 1 1 0 0 0 0 1 1 1 1 1 1 0 1 0 0 1 1 0 0 0 0 0 1 0 0 1 1 1 1 0 0 0 0 0
0 1 0 1 0 0 0 0 0 1 1 0 0 0 1 0 0 1 1 0 1 1 0 0 1 0 0 1 0 1 0 0 0 0 0 1 1
1 0 0 1 1 1 0 0 1 1 0 0 1 1 0 0 1 0 1 0 1 1 0 1 1 1 0 1 1 1 0 1 0 1 0 0 1
1 0 1 0 1 0 1 1 0 1 1 0 1 0 1 1 1 0 0 1 1 0 1 0 1 0 1 0 0 0 0 0 1 0 0 0 0
0 0 1 1 0 0 0 1 0 1 1 0 1 0 1 1 0 1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 0 0
1 0 1 0 0 0 0 1 1 1 0 1 0 0 1 1 0 1 1 0 1 0 0 0 1 0 1 0 0 0 0 0 0 0 0 1
0 0 0 1 0 0 1 0 1 0 1 0 0 0 0 0 1 0 0 0 0 1 1 1 0 1 1 1 0 0 0 1 1 1 1 0
0 0 1 1 1 1 1 0 0 1 0 0 0 1 1 0 0 1 0 1 1 1 1 0 0 1 1 0 1 1 0 0 1 1 1 1
0 0 1 0 0 1 1 0 1 0 0 0 0 0 0 1 0 0 0 0 0 0 1 1 0 1 1 1 0 0 0 0 0 0 0
1 1 1 0 1 1 1 0]
DEBUG::The RF testing accuracy score is::
0.7816666666666666

```

▽ Gradient Boosting Machines

```

1 from sklearn.ensemble import GradientBoostingClassifier # GBM algorithm
2 from sklearn import metrics #Additional sklearn functions
3 from sklearn.model_selection import cross_val_score, GridSearchCV
4
5 def modelfit(alg, train_x, train_y, predictors, test_x, performCV=True, printFeatureImportance=True, cv_folds=5):
6     #Fit the algorithm on the data
7     alg.fit(train_x, train_y)
8
9     #Predict training set:
10    predictions = alg.predict(train_x)
11    predprob = alg.predict_proba(train_x)[:,-1]
12
13    #Perform cross-validation:
14    if performCV:
15        cv_score = cross_val_score(alg, train_x, train_y, cv=cv_folds, scoring='roc_auc')
16
17    #Print model report:
18    print("\nModel Report")
19    print("Accuracy : %.4g" % metrics.accuracy_score(train_y,predictions))
20    print("AUC Score (Train): %f" % metrics.roc_auc_score(train_y, predprob))
21
22    if performCV:
23        print("CV Score : Mean - %.7g | Std - %.7g | Min - %.7g | Max - %.7g" % (np.mean(cv_score),np.std(cv_score),np.min(cv_score),np.max(cv_score)))
24
25    #Print Feature Importance:
26    import matplotlib.pyplot as plt
27
28    if printFeatureImportance:
29        fig,ax = plt.subplots()
30        feat_imp = pd.Series(alg.feature_importances_, predictors).sort_values(ascending=False)
31        feat_imp[:10].plot(kind='bar', title='Feature Importances',ax=ax)
32        plt.ylabel('Feature Importance Score')
33
34        fig.savefig('GBMimportances.eps', format='eps',bbox_inches='tight')
35        fig.savefig('GBMimportances.pdf', format='pdf',bbox_inches='tight')
36        fig.savefig('GBMimportances.png', format='png',bbox_inches='tight')
37        fig.savefig('GBMimportances.svg', format='svg',bbox_inches='tight')
38
39    return alg.predict(test_x)
40
41 gbm = GradientBoostingClassifier(random_state=10)
42
43 start_time = time.time()
44 test_predictions = modelfit(gbm, train_x, train_y, predictors, test_x)
45 end_time = time.time()
46 print("Training the Gradient Boosting Classifier took %3d seconds"%(end_time-start_time))
47
48 predicted_labels = test_predictions
49 print("DEBUG::The Gradient Boosting predicted labels are::")
50 print(predicted_labels)
51
52 acc_score = accuracy_score(test_y, predicted_labels)
53
54 print("DEBUG::The Gradient Boosting testing accuracy score is::")
55 print(acc_score)

```



Model Report

Accuracy : 0.9264

AUC Score (Train): 0.983428

CV Score : Mean - 0.8419613 | Std - 0.01045559 | Min - 0.8242002 | Max - 0.8566508

Training the Gradient Boosting Classifier took 292 seconds

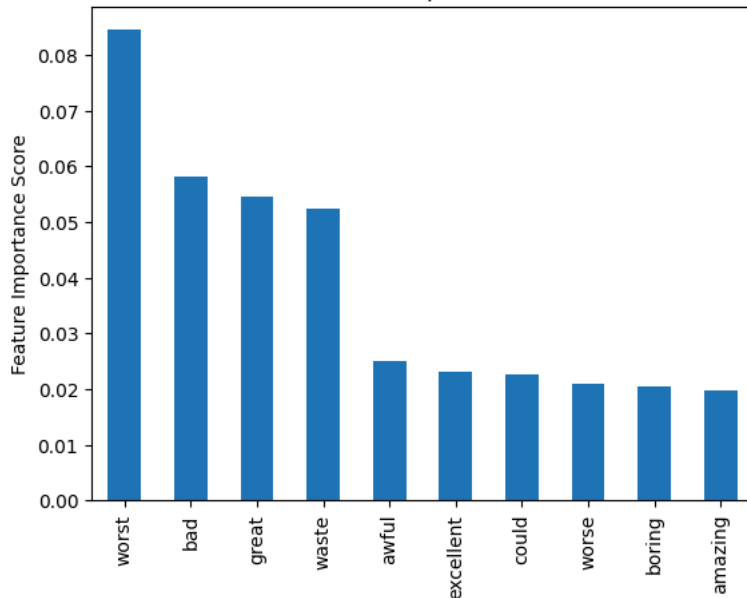
DEBUG::The Gradient Boosting predicted labels are::

```
[0 1 1 1 1 1 0 1 1 0 1 0 1 1 1 1 1 1 1 0 1 1 1 1 1 0 1 1 1 1 0 1 0 1 1
1 0 0 1 1 0 0 1 1 0 0 1 0 1 1 1 1 1 0 1 0 0 0 0 1 1 0 1 0 0 0 0 1 0 0 1 0 1
1 0 1 1 1 1 0 0 0 1 1 1 1 1 1 0 1 1 1 1 0 0 1 1 1 1 0 1 1 1 1 0 1 1 1 0 1 1
0 0 1 1 1 1 0 0 1 0 1 1 0 0 1 0 1 0 0 1 0 1 0 1 1 0 1 1 0 0 0 1 0 1 1
1 1 0 0 0 1 0 1 1 1 1 0 0 0 1 1 1 0 1 0 1 0 0 1 0 0 0 0 1 0 1 1 1 0 0 1 1
1 1 1 1 0 1 0 1 0 0 0 0 1 1 1 1 1 0 0 0 1 1 0 1 1 1 1 0 0 1 0 1 1 0 0 0 0
1 0 0 0 0 1 1 1 1 0 0 1 1 1 0 0 1 1 0 0 1 0 0 0 0 1 1 1 1 0 1 0 0 1 0 1 0
1 1 1 1 0 1 0 0 1 1 1 1 1 0 1 0 0 1 1 0 0 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0
0 1 0 1 1 0 0 0 1 1 1 0 0 0 1 1 0 1 1 0 1 1 0 1 1 0 0 1 0 1 0 0 0 0 0 1 1
1 0 0 1 1 1 0 0 1 1 0 0 1 1 0 0 0 0 1 0 1 1 0 0 0 1 1 1 1 1 0 1 1 1 0 0
1 0 1 1 1 0 1 1 0 1 1 0 1 1 1 1 1 0 0 1 1 1 1 0 1 0 1 0 0 0 0 0 1 0 0 0 0
0 1 0 1 0 0 1 1 0 1 1 0 1 0 1 1 1 1 1 0 1 1 1 0 1 1 0 1 0 0 0 0 1 1 0
1 0 1 0 0 0 0 0 1 1 1 1 1 0 1 1 1 1 1 0 1 0 0 0 0 0 1 1 1 0 0 1 0 1 0 0 1
1 0 0 1 0 0 0 1 1 0 1 0 1 0 0 0 1 1 0 0 0 1 1 1 0 1 1 0 0 0 1 0 1 1 1 1 0
0 0 1 0 1 1 1 1 0 0 1 0 0 0 1 1 0 1 1 0 1 1 1 1 0 1 1 1 0 1 1 1 1 1 1
0 1 1 1 1 1 1 0 1 0 0 0 0 0 1 1 0 0 0 0 1 0 0 1 1 1 1 1 1 1 0 0 0 1 0 1
1 1 1 1 0 1 1 0]
```

DEBUG::The Gradient Boosting testing accuracy score is::

0.76

Feature Importances



Make figures downloadable to local system in interactive mode

```
1 from IPython.display import HTML
2 import base64
3
4 def create_download_link(file_path, title="Download file"):
5     # Open the file in binary mode
6     with open(file_path, "rb") as f:
7         data = f.read()
8
9     # Convert the binary data to a base64 encoded string
10    b64 = base64.b64encode(data).decode()
11
12    # Create a download link
13    html = f'<a download="{file_path}" href="data:application/octet-stream;base64,{b64}">{title}</a>'
14    return HTML(html)
15
16 # Uso con el archivo que deseas descargar
17 create_download_link('GBMimportances.svg')
18
```




```
1 # you must remove all downloaded files - having too many of them on completion will make Kaggle reject your notebook
2 !rm -rf aclImdb
3 !rm aclImdb_v1.tar.gz
```