```python
from sklearn import model_selection, preprocessing, linear_model, naive_bayes, m
from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
from sklearn import decomposition, ensemble

import pandas, xgboost, numpy, textblob, string
from keras.preprocessing import text, sequence
from keras import layers, models, optimizers
```

In [1]:

```
c:\users\mglewis\appdata\local\programs\python\python36\lib\site-packages\sklea
rn\ensemble\weight_boosting.py:29: DeprecationWarning: numpy.core.umath_tests i
s an internal NumPy module and should not be imported. It will be removed in a
future NumPy release.
   from numpy.core.umath_tests import inner1d
Using TensorFlow backend.
```

In [2]:
```python
# load the dataset
data = open("corpus", encoding="utf8")
labels, texts = [], []
for i, line in enumerate(list(data)):
    content = line.split()
    labels.append(content[0])
    texts.append(" ".join(content[1:]))

# create a dataframe using texts and lables
trainDF = pandas.DataFrame()
trainDF['text'] = texts
trainDF['label'] = labels
```

In [3]:
```python
# split the dataset into training and validation datasets
train_x, valid_x, train_y, valid_y = model_selection.train_test_split(trainDF['t

# label encode the target variable
encoder = preprocessing.LabelEncoder()
train_y = encoder.fit_transform(train_y)
valid_y = encoder.fit_transform(valid_y)
```

In [4]:
```python
# create a count vectorizer object
count_vect = CountVectorizer(analyzer='word', token_pattern=r'\w{1,}')
count_vect.fit(trainDF['text'])

# transform the training and validation data using count vectorizer object
xtrain_count =  count_vect.transform(train_x)
xvalid_count =  count_vect.transform(valid_x)
```

```
In [5]:  # word level tf-idf
         tfidf_vect = TfidfVectorizer(analyzer='word', token_pattern=r'\w{1,}', max_featu
         tfidf_vect.fit(trainDF['text'])
         xtrain_tfidf =  tfidf_vect.transform(train_x)
         xvalid_tfidf =  tfidf_vect.transform(valid_x)

         # ngram level tf-idf
         tfidf_vect_ngram = TfidfVectorizer(analyzer='word', token_pattern=r'\w{1,}', ngr
         tfidf_vect_ngram.fit(trainDF['text'])
         xtrain_tfidf_ngram =  tfidf_vect_ngram.transform(train_x)
         xvalid_tfidf_ngram =  tfidf_vect_ngram.transform(valid_x)

         # characters level tf-idf
         tfidf_vect_ngram_chars = TfidfVectorizer(analyzer='char', token_pattern=r'\w{1,}
         tfidf_vect_ngram_chars.fit(trainDF['text'])
         xtrain_tfidf_ngram_chars =  tfidf_vect_ngram_chars.transform(train_x)
         xvalid_tfidf_ngram_chars =  tfidf_vect_ngram_chars.transform(valid_x)
```

```
In [8]:  # load the pre-trained word-embedding vectors
         embeddings_index = {}
         for i, line in enumerate(open('wiki-news-300d-1M.vec', encoding="utf8")):
             values = line.split()
             embeddings_index[values[0]] = numpy.asarray(values[1:], dtype='float32')

         # create a tokenizer
         token = text.Tokenizer()
         token.fit_on_texts(trainDF['text'])
         word_index = token.word_index

         # convert text to sequence of tokens and pad them to ensure equal length vectors
         train_seq_x = sequence.pad_sequences(token.texts_to_sequences(train_x), maxlen=7(
         valid_seq_x = sequence.pad_sequences(token.texts_to_sequences(valid_x), maxlen=7(

         # create token-embedding mapping
         embedding_matrix = numpy.zeros((len(word_index) + 1, 300))
         for word, i in word_index.items():
             embedding_vector = embeddings_index.get(word)
             if embedding_vector is not None:
                 embedding_matrix[i] = embedding_vector
```

```
In [9]:  trainDF['char_count'] = trainDF['text'].apply(len)
         trainDF['word_count'] = trainDF['text'].apply(lambda x: len(x.split()))
         trainDF['word_density'] = trainDF['char_count'] / (trainDF['word_count']+1)
         trainDF['punctuation_count'] = trainDF['text'].apply(lambda x: len("".join(_  for
         trainDF['title_word_count'] = trainDF['text'].apply(lambda x: len([wrd for wrd i
         trainDF['upper_case_word_count'] = trainDF['text'].apply(lambda x: len([wrd for w
```

In [11]:
```python
import nltk
nltk.download('averaged_perceptron_tagger')
pos_family = {
    'noun' : ['NN','NNS','NNP','NNPS'],
    'pron' : ['PRP','PRP$','WP','WP$'],
    'verb' : ['VB','VBD','VBG','VBN','VBP','VBZ'],
    'adj' :  ['JJ','JJR','JJS'],
    'adv' : ['RB','RBR','RBS','WRB']
}

# function to check and get the part of speech tag count of a words in a given se
def check_pos_tag(x, flag):
    cnt = 0
    try:
        wiki = textblob.TextBlob(x)
        for tup in wiki.tags:
            ppo = list(tup)[1]
            if ppo in pos_family[flag]:
                cnt += 1
    except:
        pass
    return cnt

trainDF['noun_count'] = trainDF['text'].apply(lambda x: check_pos_tag(x, 'noun'))
trainDF['verb_count'] = trainDF['text'].apply(lambda x: check_pos_tag(x, 'verb'))
trainDF['adj_count'] = trainDF['text'].apply(lambda x: check_pos_tag(x, 'adj'))
trainDF['adv_count'] = trainDF['text'].apply(lambda x: check_pos_tag(x, 'adv'))
trainDF['pron_count'] = trainDF['text'].apply(lambda x: check_pos_tag(x, 'pron'))
```

```
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]     C:\Users\mglewis\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping taggers\averaged_perceptron_tagger.zip.
```

In [12]:
```python
# train a LDA Model
lda_model = decomposition.LatentDirichletAllocation(n_components=20, learning_met
X_topics = lda_model.fit_transform(xtrain_count)
topic_word = lda_model.components_
vocab = count_vect.get_feature_names()

# view the topic models
n_top_words = 10
topic_summaries = []
for i, topic_dist in enumerate(topic_word):
    topic_words = numpy.array(vocab)[numpy.argsort(topic_dist)][:-(n_top_words+1
    topic_summaries.append(' '.join(topic_words))
```

```
In [13]:  def train_model(classifier, feature_vector_train, label, feature_vector_valid, i:
              # fit the training dataset on the classifier
              classifier.fit(feature_vector_train, label)

              # predict the labels on validation dataset
              predictions = classifier.predict(feature_vector_valid)

              if is_neural_net:
                  predictions = predictions.argmax(axis=-1)

              return metrics.accuracy_score(predictions, valid_y)
```

```
In [15]:  # Naive Bayes on Count Vectors
          accuracy = train_model(naive_bayes.MultinomialNB(), xtrain_count, train_y, xvali
          print("NB, Count Vectors: ", accuracy)
          # Naive Bayes on Word Level TF IDF Vectors
          accuracy = train_model(naive_bayes.MultinomialNB(), xtrain_tfidf, train_y, xvali
          print("NB, WordLevel TF-IDF: ", accuracy)

          # Naive Bayes on Ngram Level TF IDF Vectors
          accuracy = train_model(naive_bayes.MultinomialNB(), xtrain_tfidf_ngram, train_y,
          print("NB, N-Gram Vectors: ", accuracy)

          # Naive Bayes on Character Level TF IDF Vectors
          accuracy = train_model(naive_bayes.MultinomialNB(), xtrain_tfidf_ngram_chars, tr
          print("NB, CharLevel Vectors: ", accuracy)
```

```
NB, Count Vectors:  0.834
NB, WordLevel TF-IDF:  0.8396
NB, N-Gram Vectors:  0.8428
NB, CharLevel Vectors:  0.8064
```

```
In [16]:  # Linear Classifier on Count Vectors
          accuracy = train_model(linear_model.LogisticRegression(), xtrain_count, train_y,
          print("LR, Count Vectors: ", accuracy)

          # Linear Classifier on Word Level TF IDF Vectors
          accuracy = train_model(linear_model.LogisticRegression(), xtrain_tfidf, train_y,
          print("LR, WordLevel TF-IDF: ", accuracy)

          # Linear Classifier on Ngram Level TF IDF Vectors
          accuracy = train_model(linear_model.LogisticRegression(), xtrain_tfidf_ngram, tr
          print("LR, N-Gram Vectors: ", accuracy)

          # Linear Classifier on Character Level TF IDF Vectors
          accuracy = train_model(linear_model.LogisticRegression(), xtrain_tfidf_ngram_cha
          print("LR, CharLevel Vectors: ", accuracy)
```

```
LR, Count Vectors:  0.8564
LR, WordLevel TF-IDF:  0.866
LR, N-Gram Vectors:  0.8364
LR, CharLevel Vectors:  0.8368
```

In [17]:
```python
# SVM on Ngram Level TF IDF Vectors
accuracy = train_model(svm.SVC(), xtrain_tfidf_ngram, train_y, xvalid_tfidf_ngram
print("SVM, N-Gram Vectors: ", accuracy)
```

SVM, N-Gram Vectors:  0.5148

In [ ]: