

NLP: Yelp Review to Rating

Authors: Tanvee Desai and Tanner Arrizabalaga

Hello! In this project, we will be looking over Yelp reviews (data available here: <https://www.yelp.com/dataset> (<https://www.yelp.com/dataset>)) and utilizing ML/DL to accurately predict what the reviews star rating is based solely on text.

This project is split into the following parts

- Libraries
- EDA
- Data Cleaning
 - Stop word removal, HTML parsing, punctuation removal, etc.
 - Creation of a cleaned *and* stemmed dataset
- Model Implementation
 - Simple BOW Model Neural Network
 - LSTM
 - Bidirectional LSTM
 - One vs. All LSTM Approach
- Exploring Challenges
 - Challenge 5
 - Challenge 6

Importing necessary libraries

```
In [1]: # General Libraries
import json
import sys
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import itertools

# NLP
import nltk
import re
from nltk.corpus import stopwords
from bs4 import BeautifulSoup
from nltk.stem import PorterStemmer

# ML/DL
import tensorflow as tf
import pickle

from sklearn.preprocessing import LabelBinarizer, LabelEncoder
from sklearn.metrics import confusion_matrix, classification_report
from sklearn.model_selection import train_test_split

from tensorflow import keras
from keras import Sequential
from keras.layers import Dense, Activation, Dropout, Embedding, Conv1D, MaxPooling1D, LSTM, BatchNormalization, SpatialDropout1D, Bidirectional
from keras.preprocessing.sequence import pad_sequences
from keras.preprocessing import text, sequence
from keras import utils
from keras import regularizers
from keras.models import load_model
from keras.initializers import Constant
from keras.utils import plot_model
```

Using TensorFlow backend.

```
In [2]: yelp = pd.read_json("./yelp_review_training_dataset.jsonl", lines = True)
yelp.head()
```

Out[2]:

	review_id	text	stars
0	Q1sbwvVQXV2734tPgoKj4Q	Total bill for this horrible service? Over \$8G...	1
1	GJXCdrto3ASJOqKeVWPi6Q	I *adore* Travis at the Hard Rock's new Kelly ...	5
2	2TzJjDVDEuAW6MR5Vuc1ug	I have to say that this office really has it t...	5
3	yi0R0Ugj_xUx_Nek0-_Qig	Went in for a lunch. Steak sandwich was delici...	5
4	11a8sVPMUFtaC7_ABRkmtw	Today was my second out of three sessions I ha...	1

How large is the data?

```
In [3]: yelp.shape
```

```
Out[3]: (533581, 3)
```

EDA - Stars

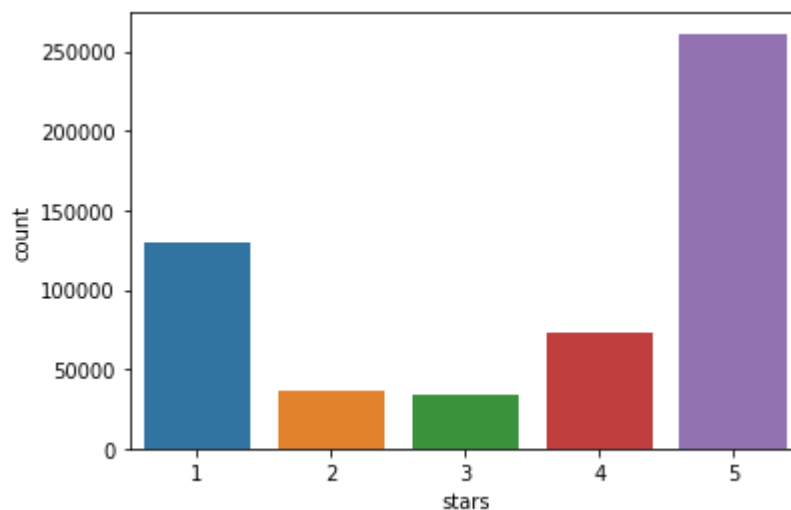
Not too much to go off of, but let's get a general understanding of our data. How many nulls do we have?

```
In [4]: yelp.isna().sum()
```

```
Out[4]: review_id    0  
text              0  
stars             0  
dtype: int64
```

```
In [5]: sns.countplot(yelp['stars'])
```

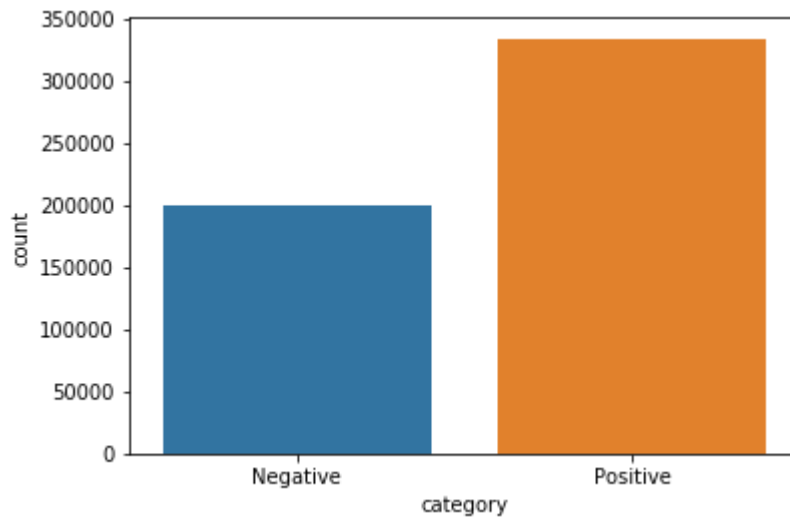
```
Out[5]: <matplotlib.axes._subplots.AxesSubplot at 0x1aeb5c526c8>
```



One thing we can potentially look at is whether or not the reviews are balanced. Let's say ≥ 4 is positive, and < 4 is negative. If we do see a significant difference in positive and negative reviews, we can balance it before training.

```
In [6]: def pos_or_neg(x):  
        if x >= 4:  
            return "Positive"  
        else:  
            return "Negative"  
  
yelp['category'] = yelp['stars'].apply(pos_or_neg)  
  
sns.countplot(yelp['category'])  
num_pos = np.count_nonzero(yelp['category'] == 'Positive')  
num_neg = np.count_nonzero(yelp['category'] == 'Negative')  
print("Positive to negative review ratio: ", num_pos / num_neg)
```

Positive to negative review ratio: 1.6679183395916979



There are roughly 1 and 2/3 times as many positive reviews as negative reviews. We will first try no class balancing when building the model, but may turn to class balancing later on.

Data Cleaning - Text

```

In [7]: REPLACE_BY_SPACE_RE = re.compile('[/(){}\\[\\]\\|@,;]')
BAD_SYMBOLS_RE = re.compile('[^0-9a-z #+_]')
STOPWORDS = set(stopwords.words('english'))
print(STOPWORDS)

def adjust_stopwords(stopwords):
    words_to_keep = set(['nor', 'not', 'very', 'no', 'few', 'too', 'doesn', 'd
    idn', 'wasn', 'ain',
                                "doesn't", "isn't", "hasn't", 'shouldn', "weren't", "d
    on't", "didn't",
                                "shouldn't", "wouldn't", "won't", "above", "below", "h
    aven't", "shan't", "weren"
                                "but", "wouldn", "mightn", "under", "mustn't", "over",
    "won", "aren", "wasn't",
                                "than"])
    return stopwords - words_to_keep

def clean_text(text):
    """
        text: a string

        return: modified initial string
    """
    new_text = BeautifulSoup(text, "lxml").text # HTML decoding
    new_text = new_text.lower() # lowercase text
    new_text = REPLACE_BY_SPACE_RE.sub(' ', new_text) # replace REPLACE_BY_SPACE_RE symbols by space in text
    new_text = BAD_SYMBOLS_RE.sub(' ', new_text) # delete symbols which are in BAD_SYMBOLS_RE from text

    ps = PorterStemmer()

    # new_text = ' '.join(ps.stem(word) for word in new_text.split()) # keeping all words, no stop word removal
    new_text = ' '.join(ps.stem(word) for word in new_text.split() if word not in STOPWORDS) # delete stopwords from text and stem
    return new_text

STOPWORDS = adjust_stopwords(STOPWORDS)
print(STOPWORDS)

```

```
{'this', 'not', 'y', 'ourselves', 'have', 'over', 'with', 'on', 're', 'why',
'if', 'was', 'i', 'up', 'didn', 'does', 'the', 'through', 'll', 'it', 'durin
g', "didn't", 'between', 'were', 'to', 'just', 'under', 'be', 'once', 'off',
'below', 'shouldn', 'such', 'needn', "needn't", 'so', "doesn't", "you'll", 'h
er', 'other', 'some', 'few', 'wasn', 'what', 'by', 'them', "isn't", "must
n't", "hasn't", 'all', 'himself', "haven't", "won't", 'my', 'hers', "don't",
'now', 'above', 'his', 'being', 'their', 'whom', 'o', 'an', 'about', 'more',
"you've", 'from', 'you', 'wouldn', 'of', 'did', 'that', 't', 'mustn', 'as',
"aren't", 've', 'him', 'won', 'most', 'herself', 'but', 'here', 'couldn', 'm
e', 'myself', 'are', 'hadn', 'out', 'no', 'a', 'nor', 'weren', 'd', 'too', 'i
s', 'and', 'there', "wasn't", 'ma', 'down', 'm', "you'd", "weren't", 'these',
'mightn', 'until', 'theirs', 'after', 'don', 'ours', 'hasn', 'doing', 'then',
"shouldn't", 'than', 'we', 'yours', "it's", 'will', "shan't", "you're", 'bee
n', 'am', 'in', "wouldn't", 'both', "she's", 'yourselves', 'he', 'into', 'ow
n', 'those', "mightn't", 'itself', 'do', 'only', 'ain', 'had', 'she', 'whic
h', 'against', 'again', 'themselves', 'aren', 'before', "hadn't", 'at', 'your
self', 'haven', 'or', 'very', 'should', 'having', "couldn't", 'any', 'when',
'isn', 'same', 'they', 's', 'your', 'has', 'for', 'shan', 'further', 'can',
"should've", 'each', 'our', 'its', 'while', "that'll", 'because', 'doesn', 'w
ho', 'where', 'how'}
{'this', 'y', 'ourselves', 'have', 'with', 'on', 're', 'why', 'if', 'was',
'i', 'up', 'does', 'the', 'through', 'll', 'it', 'during', 'between', 'were',
'to', 'just', 'be', 'once', 'off', 'such', 'needn', "needn't", 'so', "you'l
l", 'her', 'other', 'some', 'what', 'by', 'them', 'all', 'himself', 'my', 'he
rs', 'now', 'his', 'being', 'their', 'whom', 'o', 'an', 'about', 'more', "yo
u've", 'from', 'you', 'of', 'did', 'that', 't', 'mustn', 'as', "aren't", 'v
e', 'him', 'most', 'herself', 'but', 'here', 'couldn', 'me', 'myself', 'are',
'hadn', 'out', 'a', 'weren', 'd', 'is', 'and', 'there', 'ma', 'down', 'm', "y
ou'd", 'these', 'until', 'theirs', 'after', 'don', 'ours', 'hasn', 'doing',
'then', 'we', 'yours', "it's", 'will', "you're", 'been', 'am', 'in', 'both',
"she's", 'yourselves', 'he', 'into', 'own', 'those', "mightn't", 'itself', 'd
o', 'only', 'had', 'she', 'which', 'against', 'again', 'themselves', 'befor
e', "hadn't", 'at', 'yourself', 'haven', 'or', 'should', 'having', "could
n't", 'any', 'when', 'isn', 'same', 'they', 's', 'your', 'has', 'for', 'sha
n', 'further', 'can', "should've", 'each', 'our', 'its', 'while', "that'll",
'because', 'who', 'where', 'how'}
```

```
In [ ]: %%time
yelp['text'] = yelp['text'].apply(clean_text)
yelp.to_csv('cleaned_yelp_stemmed.csv')
```

```
In [8]: text_1 = "\"Good morning, cocktails for you?\" \"Wait...what? Oh...it's Vegas!
\n\nDining here, you best not be dieting because this place is literally the d
efinition of excess, but in a good way. I'm a sucker for benedicts so that was
awesome. \"Service was really great too and the staff was so welcoming. It was
our first stop just after landing so really appreciate the service.\n\nBack in
Hawaii this reminds me of Zippys or Anna Millers - that home feeling. Prices a
re a bit high, but for what you get it's totally worth it. Will remember this
place if I ever return to Vegas in the future.\"
text_2 = \"80 bucks, thirty minutes to fix my shattered iPhone screen. Verizon
won't help you so go here\"
text_3 = \"Tr\u00e8s grand caf\u00e9, mais aussi calme et reposant, je m'y suis
arr\u00eat\u00e9 alors que j'\u00e9tais dans le coin.\n\nOn peu y mang\u00e9 l
e midi, prendre une p\u00eatisserie ou un caf\u00e9/th\u00e9. \"J'ai prit un
th\u00e9 qui \u00e9tait vraiment bon, et je me suis pos\u00e9 devant une g
randes baies vitr\u00e9es sur un coussin et j'ai relax\u00e9 compl\u00e8tement
pendant 2 heures. \"Mais c'est aussi une coop\u00e9rative d'artiste, avec un
e estrade etc.\n\nIl y a aussi un magasin Bio \u00e0 l'entr\u00e9e o\u00f9 vou
s retrouverez des savons, huile d'olive et plein d'autres produits.\"
text_4 = \"Sadly, as of July 28, 2016, Silverstein bakery is permanently close
d. I went there today in person and found the bad news posted on their door. :
(\"
text_5 = \"I went here they were about to close but the cashier was especially
helpful ..but I guess they were tired of work...\"

clean_text(text_4)
```

```
Out[8]: 'sadli juli 28 2016 silverstein bakeri perman close went today person found b
ad news post door'
```

Model Implementation

Evaluation

1. Average Star Error (Average Absolute offset between predicted and true number of stars)
2. Accuracy (Exact Match -- Number of exactly predicted star ratings / total samples)

```
In [9]: from keras.losses import mean_absolute_error, binary_crossentropy, categorical_
_crossentropy

def my_custom_loss_ova(y_true, y_pred):
    mse = mean_squared_error(y_true, y_pred)
    crossentropy = binary_crossentropy(y_true, y_pred)
    return mse + crossentropy

def my_custom_loss(y_true, y_pred):
    mse = mean_squared_error(y_true, y_pred)
    crossentropy = categorical_crossentropy(y_true, y_pred)
    return mse + crossentropy

def MAE(y_true, y_pred):
    diffs = np.abs(y_true - y_pred)
    loss = np.mean(diffs)
    return loss

def Accuracy(y_true, y_pred):
    correct = y_true == y_pred
    cor_count = np.count_nonzero(correct)
    return cor_count / len(y_true)

def custom_loss(y_true, y_pred):
    return MAE(y_true, y_pred) + Accuracy(y_true, y_pred)
```

Train/Test Split (Unbalanced and balanced)

```
In [10]: yelp = pd.read_csv('cleaned_yelp_stemmed.csv')
yelp.head()
```

Out[10]:

	Unnamed: 0	review_id	text	stars	category
0	0	Q1sbwvVQXV2734tPgoKj4Q	total bill horribl servic over 8g crook actual...	1	Negative
1	1	GJXCdrto3ASJOqKeVWPi6Q	ador travi hard rock new kelli cardena salon a...	5	Positive
2	2	2TzJjDVDEuAW6MR5Vuc1ug	say offic realli togeth organ friendli dr j ph...	5	Positive
3	3	yi0R0Ugj_xUx_Nek0-_Qig	went lunch steak sandwich delici caesar salad ...	5	Positive
4	4	11a8sVPMUFtaC7_ABRkmtw	today second three session paid although first...	1	Negative

```
In [11]: X = yelp['text'].fillna('').values
y = yelp['stars']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=42)
```



```
In [12]: max_words = 3000
tokenizer = text.Tokenizer(num_words=max_words, char_level=False)

tokenizer.fit_on_texts(X_train)
X_train = tokenizer.texts_to_matrix(X_train)
X_test = tokenizer.texts_to_matrix(X_test)

encoder = LabelEncoder()
encoder.fit(y_train)
y_train = encoder.transform(y_train)
y_test = encoder.transform(y_test)

num_classes = np.max(y_train) + 1
y_train = utils.to_categorical(y_train, num_classes)
y_test = utils.to_categorical(y_test, num_classes)

print('X_train shape:', X_train.shape)
print('X_test shape:', X_test.shape)
print('y_train shape:', y_train.shape)
print('y_test shape:', y_test.shape)

X_train shape: (373506, 3000)
X_test shape: (160075, 3000)
y_train shape: (373506, 5)
y_test shape: (160075, 5)
```

Let's save the tokenizer as well for our test submission file script.

```
In [ ]: # saving
with open('tokenizer.pickle', 'wb') as handle:
    pickle.dump(tokenizer, handle, protocol=pickle.HIGHEST_PROTOCOL)

# Loading
with open('tokenizer.pickle', 'rb') as handle:
    tokenizer = pickle.load(handle)
```

Baseline Sequential Model

Here, we are computing a single model, but in future we will optimize on several parameters, listed below

- Batch size
- Learning rate
- Gradient clipping
- Drop out
- Batch normalization
- Optimizers
- Regularization

After some tests, the main variations I noticed were from the learning rate, regularization, and the choice of the optimizer. With that being said, this baseline model will use **ADAM with a learning rate of .0001 and regularization (kernel, bias, and activity)**

```
In [13]: batch_size = 512
epochs = 10

lr_schedule = keras.optimizers.schedules.ExponentialDecay(
    initial_learning_rate=.0001,
    decay_steps=10000,
    decay_rate=0.9)

optimizer = keras.optimizers.Adam(learning_rate=lr_schedule, beta_1=0.9, beta_2=0.95, amsgrad=False)

baseline = Sequential()
baseline.add(Dense(512, input_shape=(max_words,), kernel_regularizer=regularizers.l1_l2(l1=1e-5, l2=1e-4),
    bias_regularizer=regularizers.l2(1e-4),
    activity_regularizer=regularizers.l2(1e-5)))
baseline.add(BatchNormalization())
baseline.add(Activation('relu'))
baseline.add(Dropout(0.3))
baseline.add(Dense(5))
baseline.add(Activation('softmax'))

baseline.compile(loss='mean_absolute_error',
    optimizer=optimizer,
    metrics=['accuracy'])

history = baseline.fit(X_train, y_train,
    batch_size=batch_size,
    epochs=epochs,
    verbose=1,
    validation_split=0.2)
```

Train on 298804 samples, validate on 74702 samples

Epoch 1/10

298804/298804 [=====] - 32s 106us/step - loss: 0.422
4 - accuracy: 0.6859 - val_loss: 0.2667 - val_accuracy: 0.7296

Epoch 2/10

298804/298804 [=====] - 13s 45us/step - loss: 0.2114
- accuracy: 0.7352 - val_loss: 0.1796 - val_accuracy: 0.7328

Epoch 3/10

298804/298804 [=====] - 12s 41us/step - loss: 0.1613
- accuracy: 0.7449 - val_loss: 0.1530 - val_accuracy: 0.7398

Epoch 4/10

298804/298804 [=====] - 12s 40us/step - loss: 0.1424
- accuracy: 0.7517 - val_loss: 0.1412 - val_accuracy: 0.7409

Epoch 5/10

298804/298804 [=====] - 12s 39us/step - loss: 0.1323
- accuracy: 0.7562 - val_loss: 0.1347 - val_accuracy: 0.7415

Epoch 6/10

298804/298804 [=====] - 12s 40us/step - loss: 0.1263
- accuracy: 0.7595 - val_loss: 0.1310 - val_accuracy: 0.7402

Epoch 7/10

298804/298804 [=====] - 12s 40us/step - loss: 0.1228
- accuracy: 0.7616 - val_loss: 0.1284 - val_accuracy: 0.7424

Epoch 8/10

298804/298804 [=====] - 12s 40us/step - loss: 0.1202
- accuracy: 0.7637 - val_loss: 0.1269 - val_accuracy: 0.7422

Epoch 9/10

298804/298804 [=====] - 12s 41us/step - loss: 0.1181
- accuracy: 0.7663 - val_loss: 0.1258 - val_accuracy: 0.7431

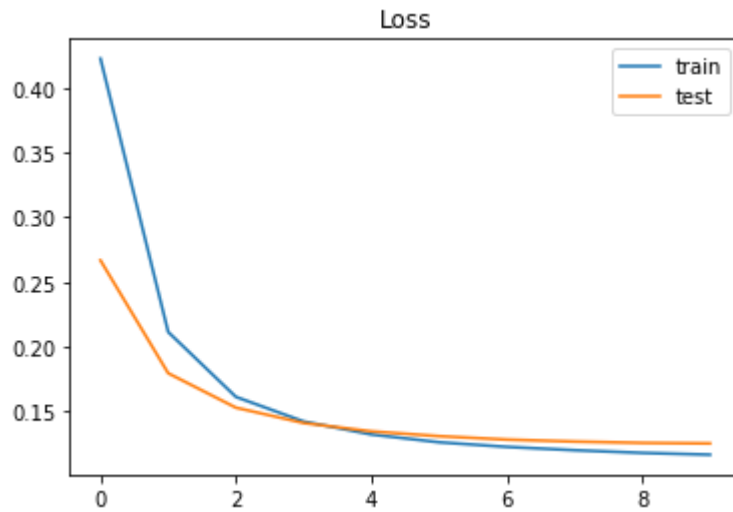
Epoch 10/10

298804/298804 [=====] - 12s 41us/step - loss: 0.1167
- accuracy: 0.7681 - val_loss: 0.1254 - val_accuracy: 0.7428

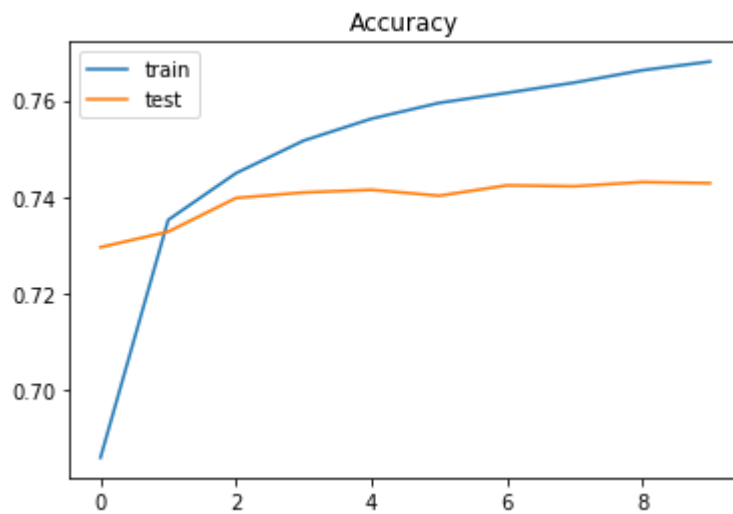
```
In [14]: score = baseline.evaluate(X_test, y_test,
                                   batch_size=batch_size, verbose=1)
print('Test accuracy:', score[1])
```

160075/160075 [=====] - 18s 110us/step
Test accuracy: 0.7449445724487305

```
In [15]: plt.title('Loss')
plt.plot(history.history['loss'], label='train')
plt.plot(history.history['val_loss'], label='test')
plt.legend()
plt.show()
```



```
In [16]: plt.title('Accuracy')
plt.plot(history.history['accuracy'], label='train')
plt.plot(history.history['val_accuracy'], label='test')
plt.legend()
plt.show()
```



```
In [17]: # Get model output
y_pred = baseline.predict(X_test)

cols = [1, 2, 3, 4, 5]

# Creating predictions table
baseline_ps = pd.DataFrame(data=y_pred, columns=cols)
y_pred_true = baseline_ps.idxmax(axis=1)

# Creating truth
baseline_truth = pd.DataFrame(data=y_test, columns=cols)
y_test_true = baseline_truth.idxmax(axis=1)

# Confusion matrix
cm = confusion_matrix(y_pred_true, y_test_true)
pd.DataFrame(cm, index=cols, columns=cols)
```

Out[17]:

	1	2	3	4	5
1	36504	6978	2543	1246	1926
2	0	0	0	0	0
3	536	1820	2723	1108	321
4	389	993	3073	7817	3971
5	1458	952	1924	11590	72203

```
In [18]: print(classification_report(y_pred_true, y_test_true))
```

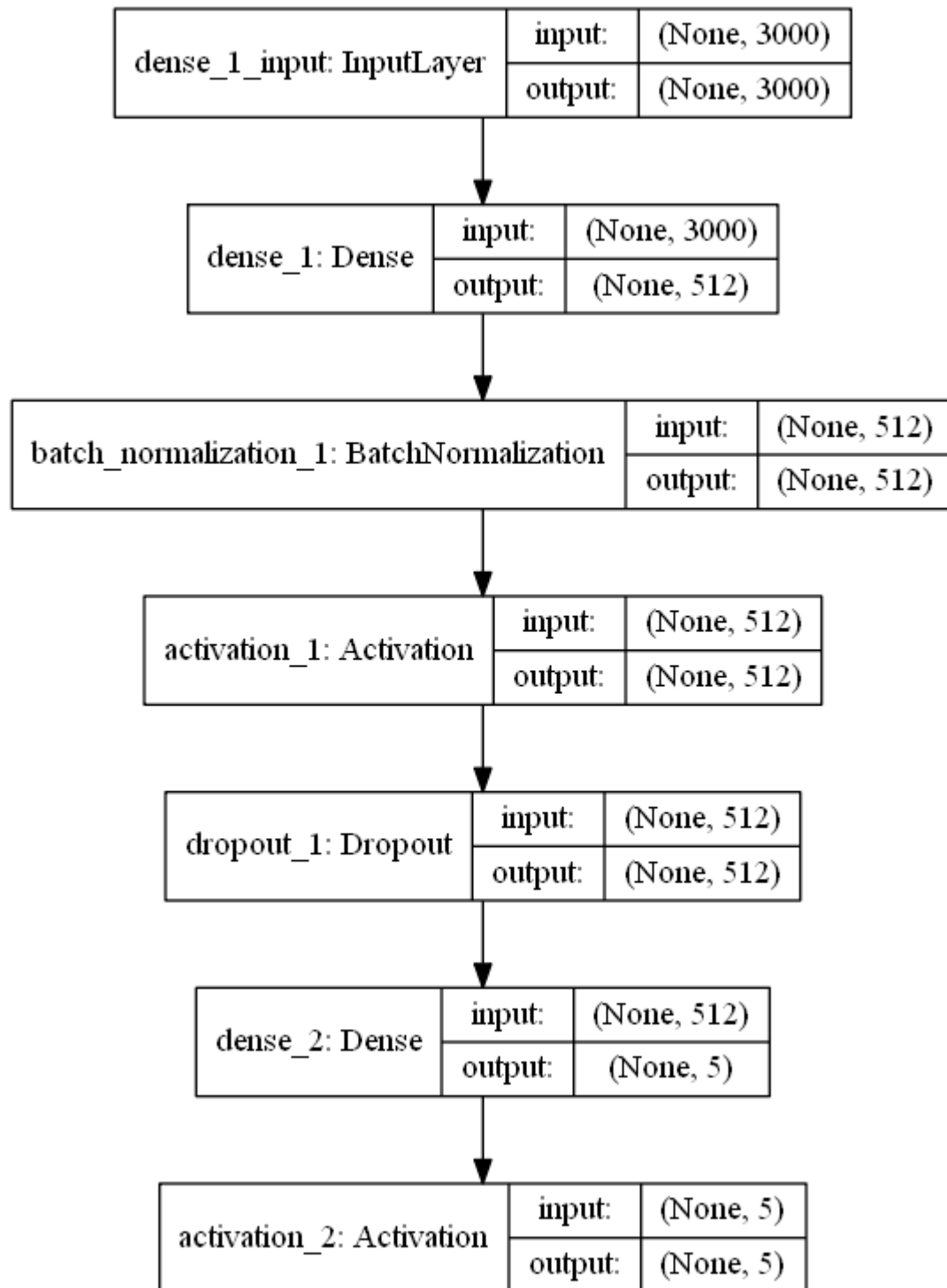
	precision	recall	f1-score	support
1	0.94	0.74	0.83	49197
2	0.00	0.00	0.00	0
3	0.27	0.42	0.32	6508
4	0.36	0.48	0.41	16243
5	0.92	0.82	0.87	88127
accuracy			0.74	160075
macro avg	0.50	0.49	0.49	160075
weighted avg	0.84	0.74	0.79	160075

C:\Users\Tanner\Anaconda3\envs\yelp\lib\site-packages\sklearn\metrics_classification.py:1272: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use `zero_division` parameter to control this behavior.

```
_warn_prf(average, modifier, msg_start, len(result))
```

In [19]: `plot_model(baseline, to_file='baseline.png', show_shapes=True)`

Out[19]:



Let's save this model.

In []: `baseline.save('./models/baseline.h5')`

Now training with several parameter changes

```
In [ ]: batch_sizes = [128, 256, 512]
epochs = [5]
learning_rates = [.01, .001, .0001]
dropout = [False, True]
batch_norm = [False, True]
regularization = [True]
optimizers = ["SGD", "RMSProp", "ADAM"]

all_lists = [batch_sizes, epochs, learning_rates, dropout, batch_norm, regularization, optimizers]

params_to_test = list(itertools.product(*all_lists))
print(len(params_to_test))
```



```
In [ ]: models = {}
        histories = {}
        scores = {}

        for params in params_to_test:
            print(params)
            batch_size, epochs, learning_rate, dropout, batch_norm, regularization, opt = params

            if opt == "SGD":
                optimizer = keras.optimizers.SGD(learning_rate=learning_rate, momentum=0.0, nesterov=False)
            elif opt == "RMSProp":
                optimizer = keras.optimizers.RMSprop(learning_rate=learning_rate, rho=0.9)
            elif opt == "ADAM":
                optimizer = keras.optimizers.Adam(learning_rate=learning_rate, beta_1=0.9, beta_2=0.99, amsgrad=False)
            else:
                optimizer = keras.optimizers.Adadelta(learning_rate=learning_rate, rho=0.95)

            model = Sequential()
            model.add(Dense(512, input_shape=(max_words,), kernel_regularizer=regularizers.l1_l2(l1=1e-5, l2=1e-4)))

            # Check Batch Normalization
            if batch_norm:
                model.add(BatchNormalization())

            model.add(Activation('relu'))

            # Check Dropout
            if dropout:
                model.add(Dropout(0.2))

            model.add(Dense(5))
            model.add(Activation('softmax'))

            model.compile(loss='categorical_crossentropy',
                          optimizer=optimizer,
                          metrics=['accuracy'])

            history = model.fit(X_train, y_train,
                               batch_size=batch_size,
                               epochs=epochs,
                               verbose=0,
                               validation_split=0.1)

            models[params] = model
            histories[params] = history

            score = model.evaluate(X_test, y_test, batch_size=batch_size, verbose=1)
            print(score)

            scores[params] = score
```

LSTM Model

Specific Data Prep

```
In [20]: X = yelp['text'].fillna('').values
y = pd.get_dummies(yelp['stars']).values

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
print(X_train.shape, y_train.shape)
print(X_test.shape, y_test.shape)

max_words = 3000
maxlen = 400

X_train = tokenizer.texts_to_sequences(X_train)
X_test = tokenizer.texts_to_sequences(X_test)

# For the LSTM, we are going to pad our sequences
X_train = pad_sequences(X_train, maxlen=maxlen)
X_test = pad_sequences(X_test, maxlen=maxlen)

(373506,) (373506, 5)
(160075,) (160075, 5)
```

LSTM #1

```

In [22]: batch_size = 512
         epochs = 5

         lr_schedule = keras.optimizers.schedules.ExponentialDecay(
             initial_learning_rate=.001,
             decay_steps=10000,
             decay_rate=0.9)

         optimizer = keras.optimizers.Adam(learning_rate=lr_schedule, beta_1=0.9, beta_2=0.99, amsgrad=False, clipvalue=.3)

         lstm = Sequential()
         lstm.add(Embedding(max_words, 128, input_length=maxlen))
         lstm.add(SpatialDropout1D(0.2))
         lstm.add(Conv1D(64, 5, activation='relu', kernel_regularizer=regularizers.l1_l2(l1=1e-5, l2=1e-4),
             bias_regularizer=regularizers.l2(1e-4)))
         lstm.add(MaxPooling1D(pool_size=4))
         lstm.add(LSTM(128, dropout=0.2, recurrent_dropout=0.2))
         lstm.add(BatchNormalization())
         lstm.add(Dense(5, activation='sigmoid'))

         lstm.compile(loss='mean_absolute_error',
             optimizer=optimizer,
             metrics=['accuracy'])

         history = lstm.fit(X_train, y_train,
             batch_size=batch_size,
             epochs=epochs,
             verbose=1,
             validation_split=0.2)

```

Train on 298804 samples, validate on 74702 samples

Epoch 1/5

298804/298804 [=====] - 96s 320us/step - loss: 0.1543 - accuracy: 0.6931 - val_loss: 0.1097 - val_accuracy: 0.7136

Epoch 2/5

298804/298804 [=====] - 94s 314us/step - loss: 0.1024 - accuracy: 0.7259 - val_loss: 0.0998 - val_accuracy: 0.7304

Epoch 3/5

298804/298804 [=====] - 92s 309us/step - loss: 0.1001 - accuracy: 0.7286 - val_loss: 0.1005 - val_accuracy: 0.7309

Epoch 4/5

298804/298804 [=====] - 92s 308us/step - loss: 0.0990 - accuracy: 0.7308 - val_loss: 0.1001 - val_accuracy: 0.7330

Epoch 5/5

298804/298804 [=====] - 91s 305us/step - loss: 0.0982 - accuracy: 0.7325 - val_loss: 0.0995 - val_accuracy: 0.7356

LSTM #1: Evaluation

```
In [23]: score = lstm.evaluate(X_test, y_test,
                                batch_size=batch_size, verbose=1)
print('Test accuracy:', score[1])
```

160075/160075 [=====] - 12s 72us/step
 Test accuracy: 0.7373855710029602

```
In [24]: lstm.summary()
```

Model: "sequential_3"

Layer (type)	Output Shape	Param #
=====		
embedding_2 (Embedding)	(None, 400, 128)	384000

spatial_dropout1d_2 (Spatial	(None, 400, 128)	0

conv1d_2 (Conv1D)	(None, 396, 64)	41024

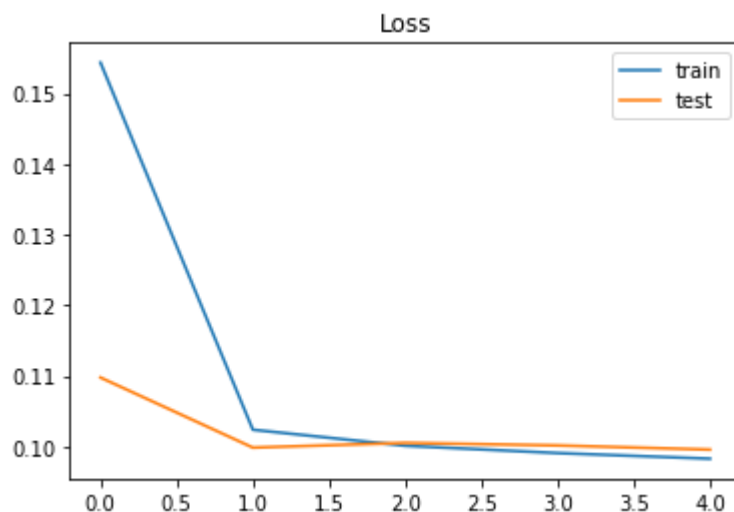
max_pooling1d_2 (MaxPooling1	(None, 99, 64)	0

lstm_2 (LSTM)	(None, 128)	98816

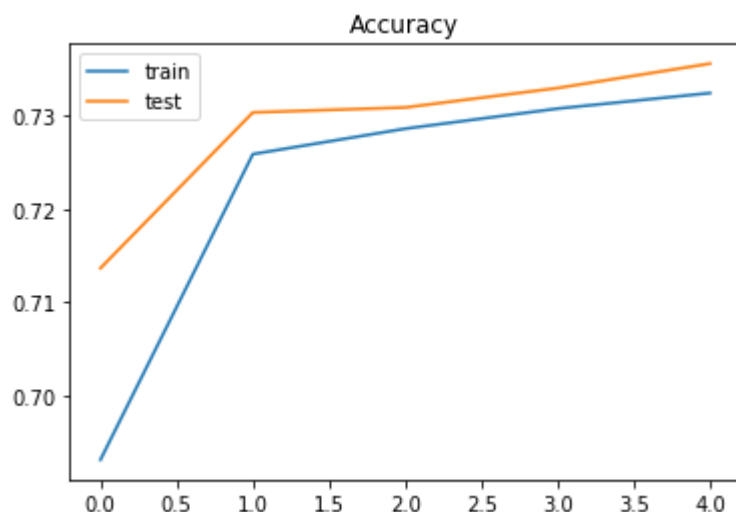
batch_normalization_3 (Batch	(None, 128)	512

dense_4 (Dense)	(None, 5)	645
=====		
Total params: 524,997		
Trainable params: 524,741		
Non-trainable params: 256		

```
In [25]: plt.title('Loss')
plt.plot(history.history['loss'], label='train')
plt.plot(history.history['val_loss'], label='test')
plt.legend()
plt.show()
```



```
In [26]: plt.title('Accuracy')
plt.plot(history.history['accuracy'], label='train')
plt.plot(history.history['val_accuracy'], label='test')
plt.legend()
plt.show()
```



```
In [27]: # Get model output
y_pred = lstm.predict(X_test)
y_pred

cols = [1, 2, 3, 4, 5]

# Creating predictions table
baseline_ps = pd.DataFrame(data=y_pred, columns=cols)
y_pred_true = baseline_ps.idxmax(axis=1)
y_pred_true

# Creating truth
baseline_truth = pd.DataFrame(data=y_test, columns=cols)
y_test_true = baseline_truth.idxmax(axis=1)
y_test_true

# Confusion matrix
cm = confusion_matrix(y_pred_true, y_test_true)
pd.DataFrame(cm, index=cols, columns=cols)
```

Out[27]:

	1	2	3	4	5
1	37048	7473	2950	1617	1952
2	0	0	0	0	0
3	0	0	0	0	0
4	730	2586	5608	8254	3734
5	1109	684	1705	11890	72735

```
In [28]: print(classification_report(y_pred_true, y_test_true))
```

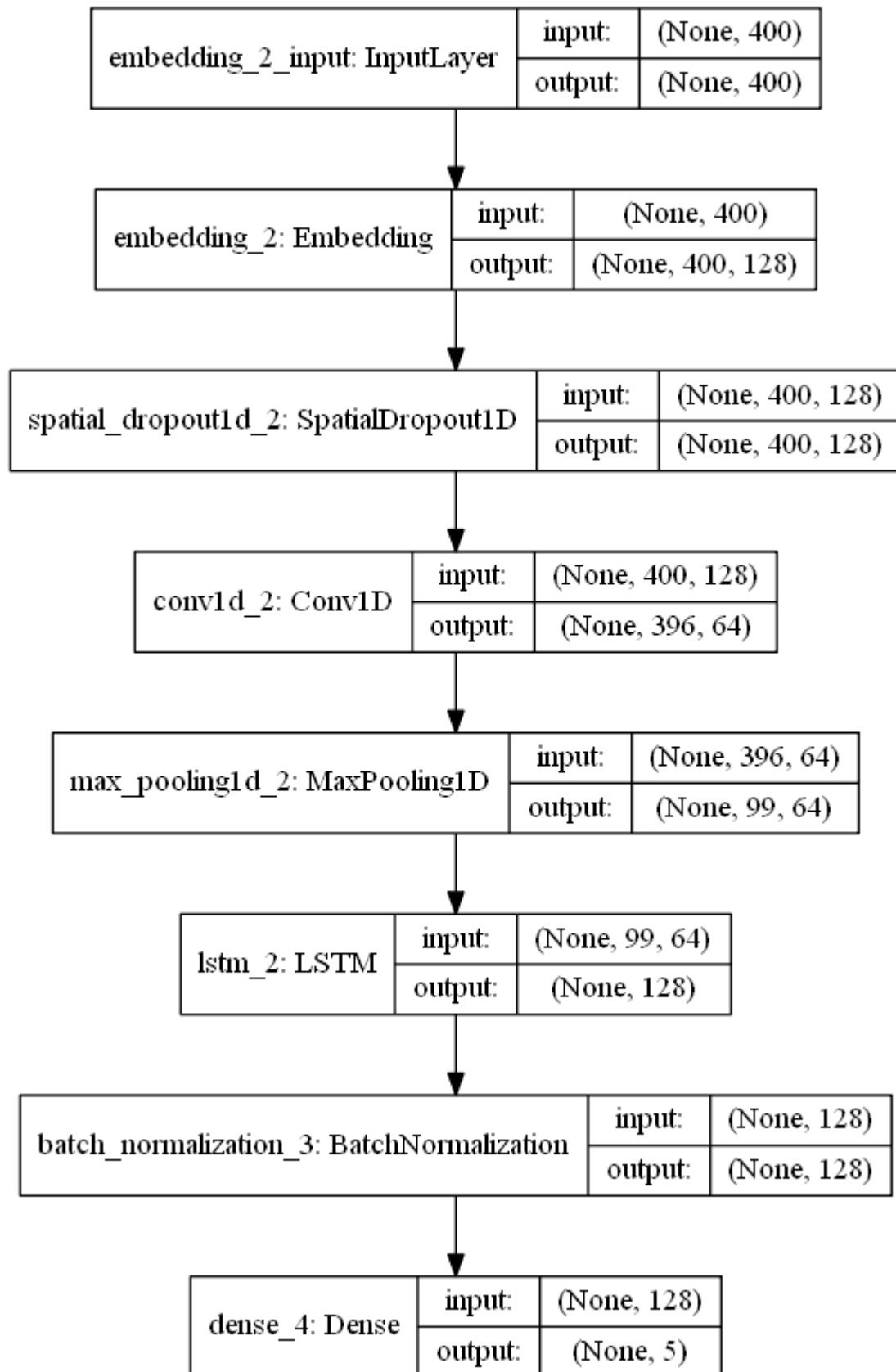
	precision	recall	f1-score	support
1	0.95	0.73	0.82	51040
2	0.00	0.00	0.00	0
3	0.00	0.00	0.00	0
4	0.38	0.39	0.39	20912
5	0.93	0.83	0.87	88123
accuracy			0.74	160075
macro avg	0.45	0.39	0.42	160075
weighted avg	0.86	0.74	0.79	160075

C:\Users\Tanner\Anaconda3\envs\yelp\lib\site-packages\sklearn\metrics_classification.py:1272: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use `zero_division` parameter to control this behavior.

```
_warn_prf(average, modifier, msg_start, len(result))
```

In [29]: `plot_model(lstm, to_file='baseline.png', show_shapes=True)`

Out[29]:



Let's save this model as well.

```
In [ ]: lstm.save('./models/lstm.h5')
```

LSTM #2

```
In [ ]: batch_size = 128
epochs = 5

lr_schedule = keras.optimizers.schedules.ExponentialDecay(
    initial_learning_rate=.001,
    decay_steps=10000,
    decay_rate=0.9)

optimizer = keras.optimizers.Adam(learning_rate=lr_schedule, beta_1=0.9, beta_2=0.99, amsgrad=False, clipvalue=.3)

lstm_v2 = Sequential()
lstm_v2.add(Embedding(max_words, 128, input_length=maxlen))
lstm_v2.add(SpatialDropout1D(0.3))
lstm_v2.add(Bidirectional(LSTM(128, dropout=0.3, recurrent_dropout=0.3)))
lstm_v2.add(Dense(128, activation='relu'))
lstm_v2.add(Dropout(0.2))
lstm_v2.add(Dense(128, activation='relu'))
lstm_v2.add(Dropout(0.2))
lstm_v2.add(Dense(5, activation='sigmoid'))

lstm_v2.compile(loss='categorical_crossentropy',
                optimizer=optimizer,
                metrics=['accuracy'])

history = lstm_v2.fit(X_train, y_train,
                    batch_size=batch_size,
                    epochs=epochs,
                    verbose=1,
                    validation_split=0.2)
```

LSTM #2: Evaluation

```
In [ ]: score = lstm_v2.evaluate(X_test, y_test,
                                batch_size=batch_size, verbose=1)
print('Test accuracy:', score[1])
```

```
In [ ]: lstm_v2.summary()
```

```
In [ ]: plt.title('Loss')
plt.plot(history.history['loss'], label='train')
plt.plot(history.history['val_loss'], label='test')
plt.legend()
plt.show()
```



```
In [ ]: plt.title('Accuracy')
plt.plot(history.history['accuracy'], label='train')
plt.plot(history.history['val_accuracy'], label='test')
plt.legend()
plt.show()
```

Let's save this model as well.

```
In [ ]: lstm.save('./models/lstm_v2.h5')
```

One vs. All Approach

In the one vs. all approach, it goes by the following idea:

- We will have N learners for the multi-class classification problem, where N is the number of classes
- For each learner L , we will train L on our training data X_{Train} and y_{Train} . However, y_{Train} consists of only one label, making it a binary classification problem instead of multinomial
 - For instance, learner L_1 will still use all of X_{Train} , but y_{Train} will now be transformed to be a binary vector v_i where i denotes the star rating we are attempting to predict
- Once we have concluded our training, we will then create an ensemble model (bagging) that does the following
 1. L_1, L_2, \dots, L_5 all assign p_i to each record in X_{Test} , where p_i is the likelihood observation x_n belongs to class i
 2. From there, our prediction is the following: $P_n = \text{argmax}(p_1, p_2, p_3, p_4, p_5)$

After observing the challenge datasets 5 & 6, my partner and I believe this approach is a clever way to tackle the challenges while still having a strong model.

Sources: <https://developers.google.com/machine-learning/crash-course/multi-class-neural-networks/one-vs-all>
(<https://developers.google.com/machine-learning/crash-course/multi-class-neural-networks/one-vs-all>)

```
In [39]: yelp = pd.read_csv('cleaned_yelp_stemmed.csv')

X = yelp['text'].fillna('').values
y = pd.get_dummies(yelp['stars']).values

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=42)

# Loading
# with open('tokenizer.pickle', 'rb') as handle:
#     tokenizer = pickle.load(handle)

max_words = 3000
maxlen = 400

X_train = tokenizer.texts_to_sequences(X_train)
X_test = tokenizer.texts_to_sequences(X_test)
X_train = pad_sequences(X_train, maxlen=maxlen)
X_test = pad_sequences(X_test, maxlen=maxlen)

print('X_train shape:', X_train.shape)
print('X_test shape:', X_test.shape)
print('y_train shape:', y_train.shape)
print('y_test shape:', y_test.shape)

X_train shape: (373506, 400)
X_test shape: (160075, 400)
y_train shape: (373506, 5)
y_test shape: (160075, 5)
```

Buidling all models

```

In [40]: stars = np.arange(1, 6)
models = {}
histories = {}
batch_size = 1024

for star in stars:
    if star in [1]:
        epochs = 2
    elif star in [2, 3, 4]:
        epochs = 3
    else:
        epochs = 4

    print(star)
    y_train_sub = y_train[:, star - 1]

    lr_schedule = keras.optimizers.schedules.ExponentialDecay(
        initial_learning_rate=.001,
        decay_steps=10000,
        decay_rate=0.9)

    optimizer = keras.optimizers.Adam(learning_rate=lr_schedule, beta_1=0.9, b
eta_2=0.99, amsgrad=False, clipvalue=.3)

    sub_lstm = Sequential()
    sub_lstm.add(Embedding(max_words, 128, input_length=maxlen))
    sub_lstm.add(SpatialDropout1D(0.2))
    sub_lstm.add(Conv1D(64, 5, activation='relu', kernel_regularizer=regulariz
ers.l1_l2(l1=1e-5, l2=1e-4),
                bias_regularizer=regularizers.l2(1e-4)))
    sub_lstm.add(MaxPooling1D(pool_size=4))
    sub_lstm.add(LSTM(128))
    sub_lstm.add(BatchNormalization())
    sub_lstm.add(Dense(8))
    sub_lstm.add(Dense(1, activation='sigmoid'))

    sub_lstm.compile(loss='mean_absolute_error',
                    optimizer=optimizer,
                    metrics=['accuracy'])

    history = sub_lstm.fit(X_train, y_train_sub,
                          batch_size=batch_size,
                          epochs=epochs,
                          verbose=1,
                          validation_split=0.2)

    models[star] = sub_lstm
    histories[star] = sub_lstm

```

```
1
Train on 298804 samples, validate on 74702 samples
Epoch 1/2
298804/298804 [=====] - 66s 222us/step - loss: 0.123
0 - accuracy: 0.8994 - val_loss: 0.0978 - val_accuracy: 0.9172
Epoch 2/2
298804/298804 [=====] - 65s 219us/step - loss: 0.089
5 - accuracy: 0.9209 - val_loss: 0.0925 - val_accuracy: 0.9162
2
Train on 298804 samples, validate on 74702 samples
Epoch 1/3
298804/298804 [=====] - 66s 222us/step - loss: 0.113
8 - accuracy: 0.9140 - val_loss: 0.0738 - val_accuracy: 0.9323
Epoch 2/3
298804/298804 [=====] - 65s 216us/step - loss: 0.070
1 - accuracy: 0.9329 - val_loss: 0.0688 - val_accuracy: 0.9323
Epoch 3/3
298804/298804 [=====] - 64s 215us/step - loss: 0.067
5 - accuracy: 0.9329 - val_loss: 0.0678 - val_accuracy: 0.9323
3
Train on 298804 samples, validate on 74702 samples
Epoch 1/3
298804/298804 [=====] - 63s 211us/step - loss: 0.119
2 - accuracy: 0.9051 - val_loss: 0.0775 - val_accuracy: 0.9357
Epoch 2/3
298804/298804 [=====] - 62s 209us/step - loss: 0.072
7 - accuracy: 0.9354 - val_loss: 0.0681 - val_accuracy: 0.9363
Epoch 3/3
298804/298804 [=====] - 63s 209us/step - loss: 0.066
7 - accuracy: 0.9356 - val_loss: 0.0648 - val_accuracy: 0.9363
4
Train on 298804 samples, validate on 74702 samples
Epoch 1/3
298804/298804 [=====] - 64s 213us/step - loss: 0.209
4 - accuracy: 0.8149 - val_loss: 0.1472 - val_accuracy: 0.8618
Epoch 2/3
298804/298804 [=====] - 63s 211us/step - loss: 0.140
7 - accuracy: 0.8655 - val_loss: 0.1395 - val_accuracy: 0.8643
Epoch 3/3
298804/298804 [=====] - 63s 211us/step - loss: 0.136
6 - accuracy: 0.8674 - val_loss: 0.1385 - val_accuracy: 0.8654
5
Train on 298804 samples, validate on 74702 samples
Epoch 1/4
298804/298804 [=====] - 63s 211us/step - loss: 0.175
9 - accuracy: 0.8462 - val_loss: 0.1600 - val_accuracy: 0.8585
Epoch 2/4
298804/298804 [=====] - 63s 209us/step - loss: 0.143
7 - accuracy: 0.8688 - val_loss: 0.3789 - val_accuracy: 0.6311
Epoch 3/4
298804/298804 [=====] - 63s 210us/step - loss: 0.137
3 - accuracy: 0.8722 - val_loss: 0.1468 - val_accuracy: 0.8619
Epoch 4/4
298804/298804 [=====] - 63s 210us/step - loss: 0.133
1 - accuracy: 0.8754 - val_loss: 0.1594 - val_accuracy: 0.8487
```

Building an ensemble model (maximization between learners) for all trained models

Testing

```
In [41]: %%time
# Evaluating the models above (TEST)
y_test_und = pd.DataFrame(y_test)
y_test_true = pd.DataFrame(y_test_und.columns[np.where(y_test_und!=0)[1]]) + 1

# Unload models
lstm_1, lstm_2, lstm_3, lstm_4, lstm_5 = models[1], models[2], models[3], models[4], models[5]

## Predicting the probability for each observation each model
print("Predicting 1 star")
one_star_ps = lstm_1.predict(X_test)
print("Predicting 2 star")
two_star_ps = lstm_2.predict(X_test)
print("Predicting 3 star")
three_star_ps = lstm_3.predict(X_test)
print("Predicting 4 star")
four_star_ps = lstm_4.predict(X_test)
print("Predicting 5 star")
five_star_ps = lstm_5.predict(X_test)

data = [one_star_ps.flatten(), two_star_ps.flatten(), three_star_ps.flatten(),
four_star_ps.flatten(), five_star_ps.flatten()]
cols = [1, 2, 3, 4, 5]
ps = pd.DataFrame(data=data, index=cols).T

ps["pred"] = ps.idxmax(axis=1)
ps.head()

print(MAE(ps["pred"], y_test_true[0]))
print(Accuracy(ps["pred"], y_test_true[0]))
```

```
Predicting 1 star
Predicting 2 star
Predicting 3 star
Predicting 4 star
Predicting 5 star
0.42414493206309545
0.7273215680149929
Wall time: 5min 30s
```

```
In [42]: # Confusion matrix
cm = confusion_matrix(ps["pred"], y_test_true[0])
pd.DataFrame(cm, index=cols, columns=cols)
```

```
Out[42]:
```

	1	2	3	4	5
1	35924	7165	2909	1332	1396
2	0	0	0	0	0
3	691	1994	3478	3566	1559
4	87	333	1205	2115	557
5	2185	1251	2671	14748	74909

```
In [43]: print(classification_report(ps["pred"], y_test_true[0]))
```

	precision	recall	f1-score	support
1	0.92	0.74	0.82	48726
2	0.00	0.00	0.00	0
3	0.34	0.31	0.32	11288
4	0.10	0.49	0.16	4297
5	0.96	0.78	0.86	95764
accuracy			0.73	160075
macro avg	0.46	0.46	0.43	160075
weighted avg	0.88	0.73	0.79	160075

C:\Users\Tanner\Anaconda3\envs\yelp\lib\site-packages\sklearn\metrics_classification.py:1272: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use `zero_division` parameter to control this behavior.

```
_warn_prf(average, modifier, msg_start, len(result))
```

Saving the models

```
In [ ]: # lstm_1.save("./models/one_star.h5")
# lstm_2.save("./models/two_star.h5")
# lstm_3.save("./models/three_star.h5")
# lstm_4.save("./models/four_star.h5")
# lstm_5.save("./models/five_star.h5")
```

Ensemble on Test Set

```

In [44]: yelp = pd.read_csv('cleaned_yelp_stemmed.csv')

X = yelp['text'].fillna('').values
y = pd.get_dummies(yelp['stars'])

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=42)
print(X_train.shape, y_train.shape)
print(X_test.shape, y_test.shape)

max_words = 3000
maxlen = 400

# with open('tokenizer.pickle', 'rb') as handle:
#     tokenizer = pickle.load(handle)

print(y_test)

necc_cols = [1, 2, 3, 4, 5]
for col in necc_cols:
    if col not in y_test.columns:
        y_test[col] = 0

y_test = y_test[necc_cols]
y_test = y_test.values

X_baseline = tokenizer.texts_to_matrix(X_test)
X_lstm = tokenizer.texts_to_sequences(X_test)
X_lstm = pad_sequences(X_lstm, maxlen=maxlen)

(373506,) (373506, 5)
(160075,) (160075, 5)
      1  2  3  4  5
255947 0  0  0  0  1
261035 0  0  0  0  1
355633 0  0  0  0  1
205506 0  0  0  0  1
97222  0  0  0  1  0
...    .. .. .. .. ..
491832 0  0  0  0  1
311959 0  0  0  0  1
140524 1  0  0  0  0
125037 0  0  1  0  0
200135 0  0  0  1  0

[160075 rows x 5 columns]

```

```
In [45]: ## Trying our pretrained models
## Optimizer
# lr_schedule = keras.optimizers.schedules.ExponentialDecay(initial_learning_r
ate=.001, decay_steps=10000, decay_rate=0.9)
# optimizer = keras.optimizers.Adam(learning_rate=lr_schedule, beta_1=0.9, bet
a_2=0.99, amsgrad=False, clipvalue=.3)

## Baseline
# baseline = load_model('./models/baseline.h5')

# baseline.compile(loss='categorical_crossentropy',
#                   optimizer=optimizer,
#                   metrics=['accuracy'])

## LSTM
# lstm = load_model('./models/lstm.h5')

# lstm.compile(loss='categorical_crossentropy',
#              optimizer=optimizer,
#              metrics=['accuracy'])

## One vs. all
# lstm_1 = load_model('./models/one_star.h5')

# lstm_1.compile(loss='binary_crossentropy',
#               optimizer=optimizer,
#               metrics=['accuracy'])

# lstm_2 = load_model('./models/two_star.h5')

# lstm_2.compile(loss='binary_crossentropy',
#               optimizer=optimizer,
#               metrics=['accuracy'])

# lstm_3 = load_model('./models/three_star.h5')

# lstm_3.compile(loss='binary_crossentropy',
#               optimizer=optimizer,
#               metrics=['accuracy'])

# lstm_4 = load_model('./models/four_star.h5')

# lstm_4.compile(loss='binary_crossentropy',
#               optimizer=optimizer,
#               metrics=['accuracy'])

# lstm_5 = load_model('./models/five_star.h5')

# lstm_5.compile(loss='binary_crossentropy',
#               optimizer=optimizer,
#               metrics=['accuracy'])
```



```
In [46]: cols = [1, 2, 3, 4, 5]
# Baseline
print("Baseline")
baseline_preds = pd.DataFrame(baseline.predict(X_baseline), columns=cols)
baseline_preds['baseline_pred'] = baseline_preds.idxmax(axis=1)

# LSTM
print("LSTM")
lstm_preds = pd.DataFrame(lstm.predict(X_lstm), columns=cols)
lstm_preds['lstm_pred'] = lstm_preds.idxmax(axis=1)

# One vs. all
print("OVA")
one_star_ps = lstm_1.predict(X_lstm)
two_star_ps = lstm_2.predict(X_lstm)
three_star_ps = lstm_3.predict(X_lstm)
four_star_ps = lstm_4.predict(X_lstm)
five_star_ps = lstm_5.predict(X_lstm)

data = [one_star_ps.flatten(), two_star_ps.flatten(), three_star_ps.flatten(),
four_star_ps.flatten(), five_star_ps.flatten()]
ova_preds = pd.DataFrame(data=data, index=cols).T

ova_preds["ova_pred"] = ova_preds.idxmax(axis=1)

all_preds = pd.DataFrame([baseline_preds['baseline_pred'], lstm_preds['lstm_pred'],
ova_preds['ova_pred']]).T
all_preds["final_pred"] = all_preds.mode(axis=1)[0]

Baseline
LSTM
OVA
```

```
In [47]: print([MAE(all_preds["final_pred"], pd.DataFrame(data=y_test, columns=cols).idxmax(axis=1)),
Accuracy(all_preds["final_pred"], pd.DataFrame(data=y_test, columns=cols).idxmax(axis=1))])

[0.39061065125722316, 0.7443760737154459]
```

```
In [48]: # Confusion matrix
cm = confusion_matrix(all_preds["final_pred"], pd.DataFrame(data=y_test, columns=cols).idxmax(axis=1))
pd.DataFrame(cm, index=cols, columns=cols)
```

Out[48]:

	1	2	3	4	5
1	37132	7953	3391	1730	1766
2	0	0	0	0	0
3	272	1149	1967	1087	520
4	198	819	2952	6333	2411
5	1285	822	1953	12611	73724

```
In [49]: print(classification_report(y_pred_true, y_test_true))
```

	precision	recall	f1-score	support
1	0.95	0.73	0.82	51040
2	0.00	0.00	0.00	0
3	0.00	0.00	0.00	0
4	0.38	0.39	0.39	20912
5	0.93	0.83	0.87	88123
accuracy			0.74	160075
macro avg	0.45	0.39	0.42	160075
weighted avg	0.86	0.74	0.79	160075

C:\Users\Tanner\Anaconda3\envs\yelp\lib\site-packages\sklearn\metrics_classification.py:1272: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use `zero_division` parameter to control this behavior.

```
_warn_prf(average, modifier, msg_start, len(result))
```

Challenges

Challenge 5

```
In [50]: c5 = pd.read_json("./yelp_challenge_5_with_answers.jsonl", lines = True)
print(c5.shape)
c5.head()
```

```
(500, 3)
```

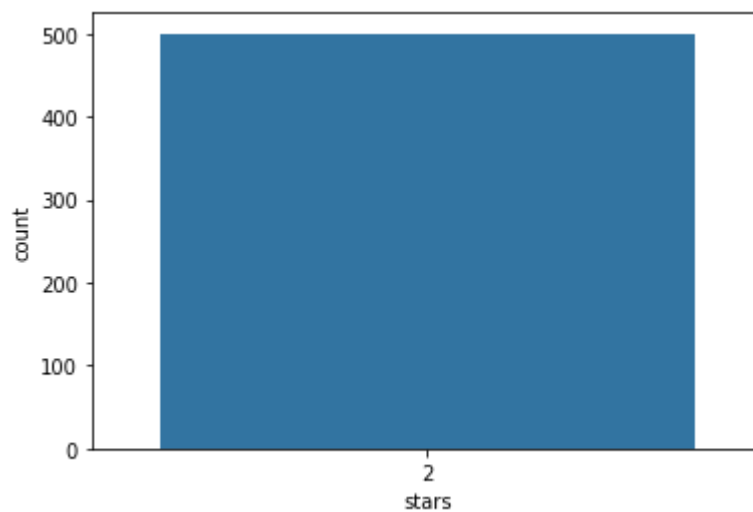
```
Out[50]:
```

	review_id	text	stars
0	50	I went to this campus for 1 semester. I was in...	2
1	51	I have rated it a two star based on its compar...	2
2	52	Just like most of the reviews, we ordered and ...	2
3	53	I only go here if it is an emergency. I HATE i...	2
4	54	Rude staff. I got 60 feeder fish and about 15 ...	2

Quick EDA

```
In [51]: sns.countplot(c5['stars'])
```

```
Out[51]: <matplotlib.axes._subplots.AxesSubplot at 0x1b17a688408>
```



Pre-processing

```
In [52]: c5['text'] = c5['text'].apply(clean_text)
c5.head()
```

```
Out[52]:
```

	review_id	text	stars
0	50	went campu 1 semest busi inform system campu o...	2
1	51	rate two star base comparison shop find staff ...	2
2	52	like review order paid half front door advanc ...	2
3	53	go emerg hate one door enter exit loss prevent...	2
4	54	rude staff got 60 feeder fish 15 dead cashier ...	2

Load previous tokenizer

```
In [53]: X = c5['text'].fillna('').values
y = pd.get_dummies(c5['stars'])

# with open('tokenizer.pickle', 'rb') as handle:
#     tokenizer = pickle.load(handle)

max_words

necc_cols = [1, 2, 3, 4, 5]
for col in necc_cols:
    if col not in y.columns:
        y[col] = 0

y = y[necc_cols]
y = y.values

X_baseline = tokenizer.texts_to_matrix(X)
X_lstm = tokenizer.texts_to_sequences(X)
X_lstm = pad_sequences(X_lstm, maxlen=400)
```

Load and compile models

```
In [ ]: # Baseline
baseline = load_model('./models/baseline.h5')

baseline.compile(loss='categorical_crossentropy',
                 optimizer=optimizer,
                 metrics=['accuracy'])

# LSTM
lstm = load_model('./models/lstm.h5')

lstm.compile(loss='categorical_crossentropy',
             optimizer=optimizer,
             metrics=['accuracy'])

# One vs. all
lstm_1 = load_model('./models/one_star.h5')

lstm_1.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_2 = load_model('./models/two_star.h5')

lstm_2.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_3 = load_model('./models/three_star.h5')

lstm_3.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_4 = load_model('./models/four_star.h5')

lstm_4.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_5 = load_model('./models/five_star.h5')

lstm_5.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])
```

Evaluate Models

```

In [54]: # Baseline
print(baseline.evaluate(X_baseline, y))

# LSTM
print(lstm.evaluate(X_lstm, y))

# One vs. All
one_star_ps = lstm_1.predict(X_lstm)
two_star_ps = lstm_2.predict(X_lstm)
three_star_ps = lstm_3.predict(X_lstm)
four_star_ps = lstm_4.predict(X_lstm)
five_star_ps = lstm_5.predict(X_lstm)

data = [one_star_ps.flatten(), two_star_ps.flatten(), three_star_ps.flatten(),
four_star_ps.flatten(), five_star_ps.flatten()]
cols = [1, 2, 3, 4, 5]
ps = pd.DataFrame(data=data, index=cols).T

ps["ova_pred"] = ps.idxmax(axis=1)

print([MAE(ps["ova_pred"], pd.DataFrame(data=y, columns=cols).idxmax(axis=1)),
Accuracy(ps["ova_pred"], pd.DataFrame(data=y, columns=cols).idxmax(axis=1))])

500/500 [=====] - 0s 76us/step
[0.41860255336761476, 0.0]
500/500 [=====] - 0s 534us/step
[0.29035107016563416, 0.0]
[1.248, 0.0]

```

Attempt Ensemble

```

In [55]: # Baseline
baseline_preds = pd.DataFrame(baseline.predict(X_baseline), columns=cols)
baseline_preds['baseline_pred'] = baseline_preds.idxmax(axis=1)

# LSTM
lstm_preds = pd.DataFrame(lstm.predict(X_lstm), columns=cols)
lstm_preds['lstm_pred'] = lstm_preds.idxmax(axis=1)

# One vs. all
ova_preds = ps

all_preds = pd.DataFrame([baseline_preds['baseline_pred'], lstm_preds['lstm_pr
ed'], ova_preds['ova_pred']]).T
all_preds["final_pred"] = all_preds.mode(axis=1)[0]

print([MAE(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols).idxmax(
axis=1)), Accuracy(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols)
.idxmax(axis=1))])

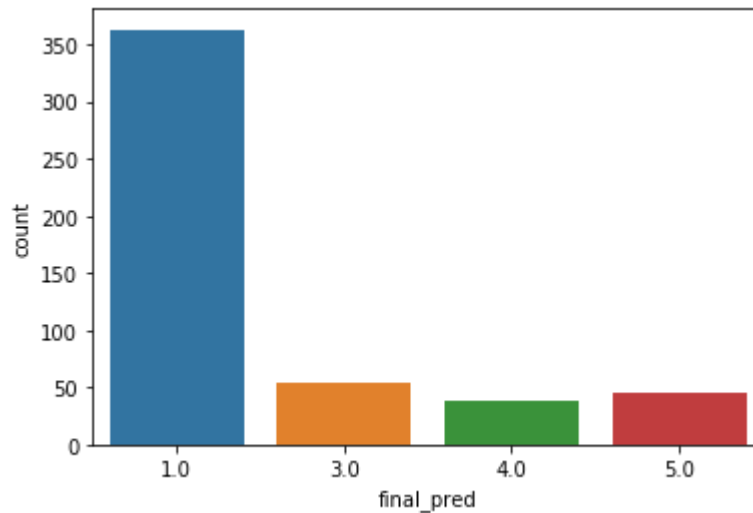
[1.256, 0.0]

```

Misc.

```
In [56]: sns.countplot(all_preds["final_pred"])
```

```
Out[56]: <matplotlib.axes._subplots.AxesSubplot at 0x1aee3aaebc8>
```

**Challenge 6**

```
In [57]: c6 = pd.read_json("./yelp_challenge_6_with_answers.jsonl", lines = True)
print(c6.shape)
c6.head()
```

```
(500, 3)
```

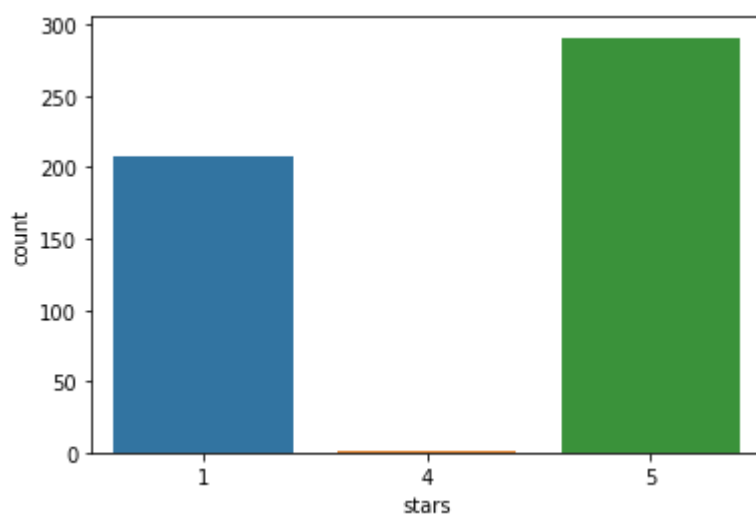
```
Out[57]:
```

	review_id	text	stars
0	60	Amazing for Trees\n\n\$20 for a 5 gallon . I wi...	5
1	61	How the hell can Taco Bell be closed before mi...	5
2	62	I actually had no intention of visiting this p...	5
3	63	Yesterday around 3:30 pm I was driving west on...	5
4	64	DR FITZMAURICE did surgery on both hands on th...	5

Quick EDA

```
In [58]: sns.countplot(c6['stars'])
```

```
Out[58]: <matplotlib.axes._subplots.AxesSubplot at 0x1b1a7a4f948>
```



Pre-processing

```
In [59]: c6['text'] = c6['text'].apply(clean_text)
c6.head()
```

```
Out[59]:
```

	review_id	text	stars
0	60	amaz tree 20 5 gallon never go low home depot ...	5
1	61	hell taco bell close midnight illeg mean pract...	5
2	62	actual no intent visit place disgust next door...	5
3	63	yesterday around 3 30 pm drive west pinnacel re...	5
4	64	dr fitzmauric surgeri hand day 8 plu year ago ...	5

Load previous tokenizer


```
In [60]: X = c6['text'].fillna('').values
y = pd.get_dummies(c6['stars'])

# with open('tokenizer.pickle', 'rb') as handle:
#     tokenizer = pickle.load(handle)

max_words

necc_cols = [1, 2, 3, 4, 5]
for col in necc_cols:
    if col not in y.columns:
        y[col] = 0

y = y[necc_cols]
y = y.values

X_baseline = tokenizer.texts_to_matrix(X)
X_lstm = tokenizer.texts_to_sequences(X)
X_lstm = pad_sequences(X_lstm, maxlen=400)
```

Load and compile models

```
In [ ]: # Baseline
baseline = load_model('./models/baseline.h5')

baseline.compile(loss='categorical_crossentropy',
                 optimizer=optimizer,
                 metrics=['accuracy'])

# LSTM
lstm = load_model('./models/lstm.h5')

lstm.compile(loss='categorical_crossentropy',
             optimizer=optimizer,
             metrics=['accuracy'])

# One vs. all
lstm_1 = load_model('./models/one_star.h5')

lstm_1.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_2 = load_model('./models/two_star.h5')

lstm_2.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_3 = load_model('./models/three_star.h5')

lstm_3.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_4 = load_model('./models/four_star.h5')

lstm_4.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_5 = load_model('./models/five_star.h5')

lstm_5.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])
```

Evaluate Models

```

In [61]: # Baseline
print(baseline.evaluate(X_baseline, y))

# LSTM
print(lstm.evaluate(X_lstm, y))

# One vs. All
one_star_ps = lstm_1.predict(X_lstm)
two_star_ps = lstm_2.predict(X_lstm)
three_star_ps = lstm_3.predict(X_lstm)
four_star_ps = lstm_4.predict(X_lstm)
five_star_ps = lstm_5.predict(X_lstm)

data = [one_star_ps.flatten(), two_star_ps.flatten(), three_star_ps.flatten(),
four_star_ps.flatten(), five_star_ps.flatten()]
cols = [1, 2, 3, 4, 5]
ps = pd.DataFrame(data=data, index=cols).T

ps["ova_pred"] = ps.idxmax(axis=1)

print([MAE(ps["ova_pred"], pd.DataFrame(data=y, columns=cols).idxmax(axis=1)),
Accuracy(ps["ova_pred"], pd.DataFrame(data=y, columns=cols).idxmax(axis=1))])

500/500 [=====] - 0s 74us/step
[0.24863504767417907, 0.4359999895095825]
500/500 [=====] - 0s 482us/step
[0.21767447805404663, 0.4259999990463257]
[2.126, 0.45]

```

Attempt Ensemble

```

In [62]: # Baseline
baseline_preds = pd.DataFrame(baseline.predict(X_baseline), columns=cols)
baseline_preds['baseline_pred'] = baseline_preds.idxmax(axis=1)

# LSTM
lstm_preds = pd.DataFrame(lstm.predict(X_lstm), columns=cols)
lstm_preds['lstm_pred'] = lstm_preds.idxmax(axis=1)

# One vs. all
ova_preds = ps

all_preds = pd.DataFrame([baseline_preds['baseline_pred'], lstm_preds['lstm_pr
ed'], ova_preds['ova_pred']]).T
all_preds["final_pred"] = all_preds.mode(axis=1)[0]

print([MAE(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols).idxmax(
axis=1)), Accuracy(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols)
.idxmax(axis=1))])

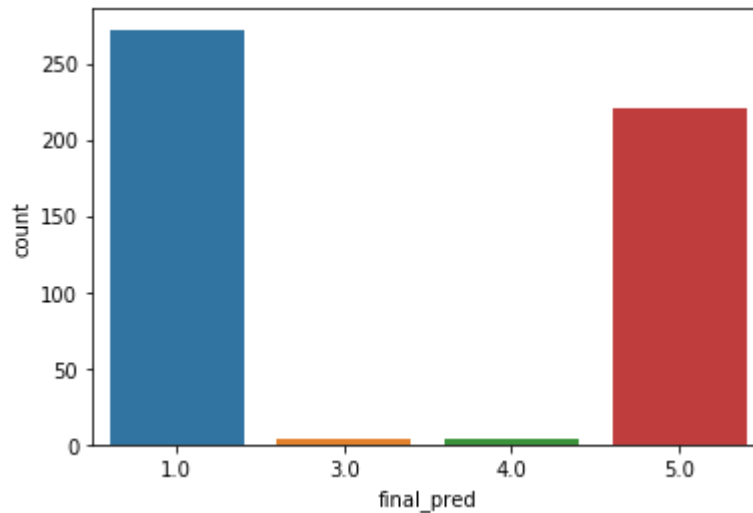
[2.174, 0.448]

```

Misc.

```
In [63]: sns.countplot(all_preds["final_pred"])
```

```
Out[63]: <matplotlib.axes._subplots.AxesSubplot at 0x1b29eefb308>
```

**Challenge 3**

```
In [64]: c3 = pd.read_json("./yelp_challenge_3_with_answers.jsonl", lines = True)
print(c3.shape)
c3.head()
```

```
(534, 3)
```

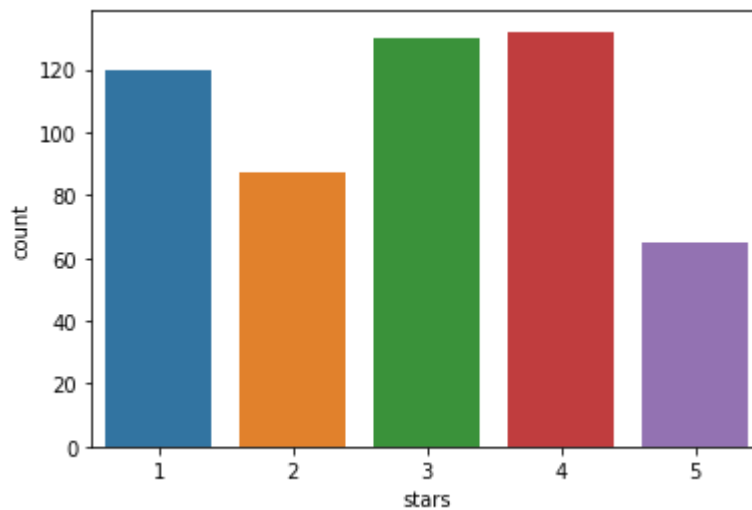
```
Out[64]:
```

	review_id	text	stars
0	30	We stopped here for lunch today and were pleas...	4
1	31	We went for a quick lunch here - it's all reas...	3
2	32	Very bad food, avoid it. We were a group of 4 ...	2
3	33	Bring a friend or two to help open the door. I...	3
4	34	Ukai serves some of the best sushi and sashimi...	4

Quick EDA

```
In [65]: sns.countplot(c3['stars'])
```

```
Out[65]: <matplotlib.axes._subplots.AxesSubplot at 0x1b29edef088>
```



Pre-processing

```
In [66]: c3['text'] = c3['text'].apply(clean_text)
c3.head()
```

```
Out[66]:
```

	review_id	text	stars
0	30	stop lunch today pleasantli surpris great ambi...	4
1	31	went quick lunch reason well price good food n...	3
2	32	veri bad food avoid group 4 veri hungri came o...	2
3	33	bring friend two help open door think weigh 40...	3
4	34	ukai serv best sushi sashimi london bar nobu i...	4

Load previous tokenizer

```
In [67]: X = c3['text'].fillna('').values
y = pd.get_dummies(c3['stars'])

# with open('tokenizer.pickle', 'rb') as handle:
#     tokenizer = pickle.load(handle)

max_words

necc_cols = [1, 2, 3, 4, 5]
for col in necc_cols:
    if col not in y.columns:
        y[col] = 0

y = y[necc_cols]
y = y.values

X_baseline = tokenizer.texts_to_matrix(X)
X_lstm = tokenizer.texts_to_sequences(X)
X_lstm = pad_sequences(X_lstm, maxlen=400)
```

Load and compile models

```
In [ ]: # Baseline
baseline = load_model('./models/baseline.h5')

baseline.compile(loss='categorical_crossentropy',
                 optimizer=optimizer,
                 metrics=['accuracy'])

# LSTM
lstm = load_model('./models/lstm.h5')

lstm.compile(loss='categorical_crossentropy',
             optimizer=optimizer,
             metrics=['accuracy'])

# One vs. all
lstm_1 = load_model('./models/one_star.h5')

lstm_1.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_2 = load_model('./models/two_star.h5')

lstm_2.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_3 = load_model('./models/three_star.h5')

lstm_3.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_4 = load_model('./models/four_star.h5')

lstm_4.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_5 = load_model('./models/five_star.h5')

lstm_5.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])
```

Evaluate Models

```

In [68]: # Baseline
print(baseline.evaluate(X_baseline, y))

# LSTM
print(lstm.evaluate(X_lstm, y))

# One vs. All
one_star_ps = lstm_1.predict(X_lstm)
two_star_ps = lstm_2.predict(X_lstm)
three_star_ps = lstm_3.predict(X_lstm)
four_star_ps = lstm_4.predict(X_lstm)
five_star_ps = lstm_5.predict(X_lstm)

data = [one_star_ps.flatten(), two_star_ps.flatten(), three_star_ps.flatten(),
four_star_ps.flatten(), five_star_ps.flatten()]
cols = [1, 2, 3, 4, 5]
ps = pd.DataFrame(data=data, index=cols).T

ps["ova_pred"] = ps.idxmax(axis=1)

print([MAE(ps["ova_pred"], pd.DataFrame(data=y, columns=cols).idxmax(axis=1)),
Accuracy(ps["ova_pred"], pd.DataFrame(data=y, columns=cols).idxmax(axis=1))])

```

```

534/534 [=====] - 0s 75us/step
[0.2116432555494237, 0.5262172222137451]
534/534 [=====] - 0s 519us/step
[0.17886928136875566, 0.4325842559337616]
[0.7191011235955056, 0.4438202247191011]

```

Attempt Ensemble

```

In [69]: # Baseline
baseline_preds = pd.DataFrame(baseline.predict(X_baseline), columns=cols)
baseline_preds['baseline_pred'] = baseline_preds.idxmax(axis=1)

# LSTM
lstm_preds = pd.DataFrame(lstm.predict(X_lstm), columns=cols)
lstm_preds['lstm_pred'] = lstm_preds.idxmax(axis=1)

# One vs. all
ova_preds = ps

all_preds = pd.DataFrame([baseline_preds['baseline_pred'], lstm_preds['lstm_pr
ed'], ova_preds['ova_pred']]).T
all_preds["final_pred"] = all_preds.mode(axis=1)[0]

print([MAE(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols).idxmax(
axis=1)), Accuracy(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols)
.idxmax(axis=1))])

```

```

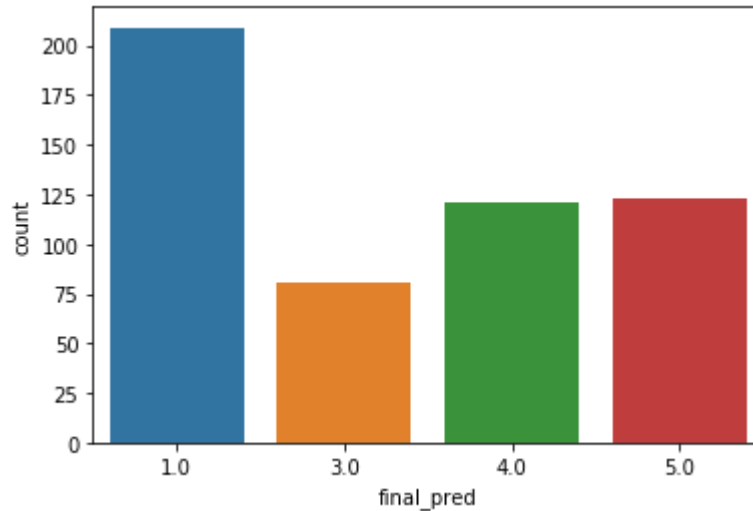
[0.6741573033707865, 0.4887640449438202]

```


Misc.

```
In [70]: sns.countplot(all_preds["final_pred"])
```

```
Out[70]: <matplotlib.axes._subplots.AxesSubplot at 0x1b2a1be0b88>
```

**Challenge 8**

```
In [71]: c8 = pd.read_json("./yelp_challenge_8_with_answers.jsonl", lines = True)
print(c8.shape)
c8.head()
```

```
(500, 3)
```

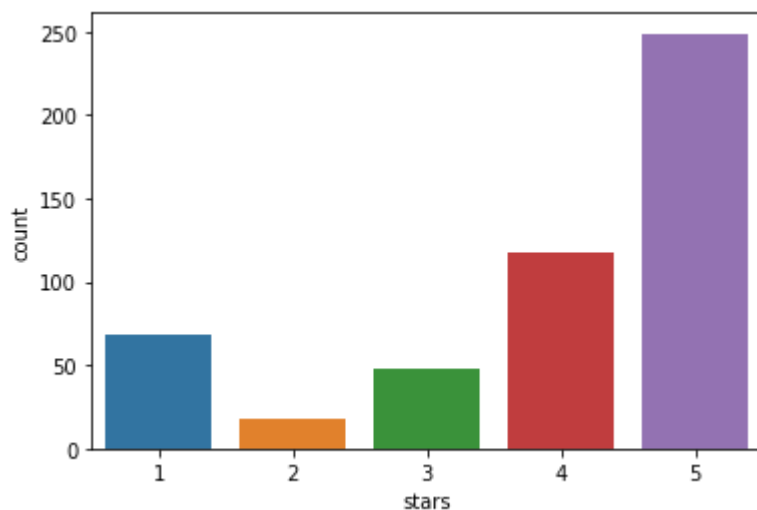
```
Out[71]:
```

	review_id	text	stars
0	qOOv-A-vo3kMT0yi4jlllg	Not bad for fast food.	4
1	uqxoO6B6w_sIDSAGr0k_0A	Une institution du café	4
2	0o_gGSU0m_4QyNLWEHKgug	J ai vraiment aimé !!!!	4
3	BKAj-fKWW5G3yt3xAkbUCQ	They have good poutine.	4
4	fAhp8lwuGNT0ywKmsCs6VQ	Very old and dirty vans.	1

Quick EDA

```
In [72]: sns.countplot(c8['stars'])
```

```
Out[72]: <matplotlib.axes._subplots.AxesSubplot at 0x1b29efc7248>
```



Pre-processing

```
In [73]: c8['text'] = c8['text'].apply(clean_text)
c8.head()
```

C:\Users\Tanner\Anaconda3\envs\yelp\lib\site-packages\bs4__init__.py:398: UserWarning: "https://casetext.com/case/united-states-v-butterbaugh-2" looks like a URL. BeautifulSoup is not an HTTP client. You should probably use an HTTP client like requests to get the document behind the URL, and feed that document to BeautifulSoup.

markup

```
Out[73]:
```

	review_id	text	stars
0	qOOv-A-vo3kMT0yi4jlllg	not bad fast food	4
1	uqxkO6B6w_sIDSAGr0k_0A	une institut du caf	4
2	0o_gGSU0m_4QyNLWEHKgug	j ai vraiment aim	4
3	BKAj-fKWW5G3yt3xAkbUCQ	good poutine	4
4	fAhp8lwuGNT0ywKmsCs6VQ	veri old dirti van	1

Load previous tokenizer

```
In [74]: X = c8['text'].fillna('').values
y = pd.get_dummies(c8['stars'])

# with open('tokenizer.pickle', 'rb') as handle:
#     tokenizer = pickle.load(handle)

max_words

necc_cols = [1, 2, 3, 4, 5]
for col in necc_cols:
    if col not in y.columns:
        y[col] = 0

y = y[necc_cols]
y = y.values

X_baseline = tokenizer.texts_to_matrix(X)
X_lstm = tokenizer.texts_to_sequences(X)
X_lstm = pad_sequences(X_lstm, maxlen=400)
```

Load and compile models

```
In [ ]: # Baseline
baseline = load_model('./models/baseline.h5')

baseline.compile(loss='categorical_crossentropy',
                 optimizer=optimizer,
                 metrics=['accuracy'])

# LSTM
lstm = load_model('./models/lstm.h5')

lstm.compile(loss='categorical_crossentropy',
             optimizer=optimizer,
             metrics=['accuracy'])

# One vs. all
lstm_1 = load_model('./models/one_star.h5')

lstm_1.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_2 = load_model('./models/two_star.h5')

lstm_2.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_3 = load_model('./models/three_star.h5')

lstm_3.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_4 = load_model('./models/four_star.h5')

lstm_4.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])

lstm_5 = load_model('./models/five_star.h5')

lstm_5.compile(loss='binary_crossentropy',
               optimizer=optimizer,
               metrics=['accuracy'])
```

Evaluate Models

```

In [75]: # Baseline
print(baseline.evaluate(X_baseline, y))

# LSTM
print(lstm.evaluate(X_lstm, y))

# One vs. All
one_star_ps = lstm_1.predict(X_lstm)
two_star_ps = lstm_2.predict(X_lstm)
three_star_ps = lstm_3.predict(X_lstm)
four_star_ps = lstm_4.predict(X_lstm)
five_star_ps = lstm_5.predict(X_lstm)

data = [one_star_ps.flatten(), two_star_ps.flatten(), three_star_ps.flatten(),
four_star_ps.flatten(), five_star_ps.flatten()]
cols = [1, 2, 3, 4, 5]
ps = pd.DataFrame(data=data, index=cols).T

ps["ova_pred"] = ps.idxmax(axis=1)

print([MAE(ps["ova_pred"], pd.DataFrame(data=y, columns=cols).idxmax(axis=1)),
Accuracy(ps["ova_pred"], pd.DataFrame(data=y, columns=cols).idxmax(axis=1))])

500/500 [=====] - 0s 78us/step
[0.17441676127910613, 0.6200000047683716]
500/500 [=====] - 0s 484us/step
[0.1450980542898178, 0.587999995231628]
[0.69, 0.6]

```

Attempt Ensemble

```

In [76]: # Baseline
baseline_preds = pd.DataFrame(baseline.predict(X_baseline), columns=cols)
baseline_preds['baseline_pred'] = baseline_preds.idxmax(axis=1)

# LSTM
lstm_preds = pd.DataFrame(lstm.predict(X_lstm), columns=cols)
lstm_preds['lstm_pred'] = lstm_preds.idxmax(axis=1)

# One vs. all
ova_preds = ps

all_preds = pd.DataFrame([baseline_preds['baseline_pred'], lstm_preds['lstm_pr
ed'], ova_preds['ova_pred']]).T
all_preds["final_pred"] = all_preds.mode(axis=1)[0]

print([MAE(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols).idxmax(
axis=1)), Accuracy(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols)
.idxmax(axis=1))])

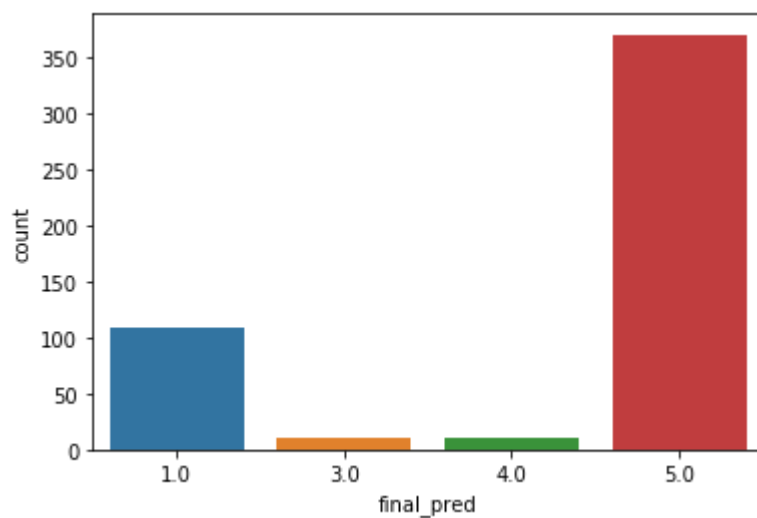
[0.624, 0.616]

```

Misc.

```
In [77]: sns.countplot(all_preds["final_pred"])
```

```
Out[77]: <matplotlib.axes._subplots.AxesSubplot at 0x1b2a1ce8588>
```



```
In [ ]:
```