

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 [285]: # 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
```

```
In [286]: yelp = pd.read_json("./yelp_review_training_dataset.jsonl", lines = True)
yelp.head()
```

Out[286]:

	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...	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 [287]: yelp.shape
```

```
Out[287]: (533581, 3)
```

EDA - Stars

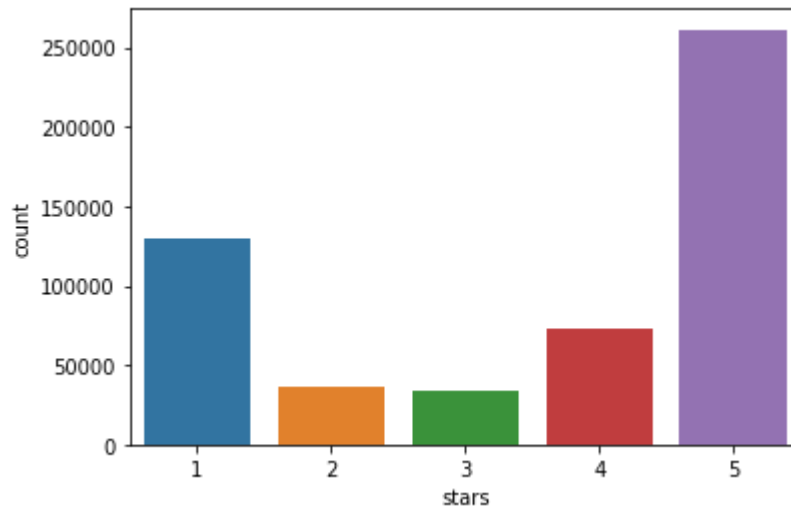
Not too much to go off of, but let's get a general understanding of our data. How many nulls do we have?

```
In [288]: yelp.isna().sum()
```

```
Out[288]: review_id    0  
text              0  
stars             0  
dtype: int64
```

```
In [289]: sns.countplot(yelp['stars'])
```

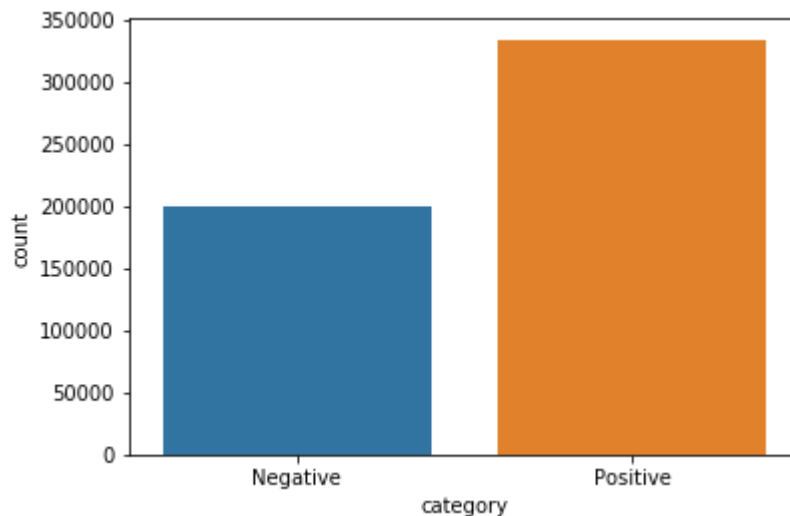
```
Out[289]: <matplotlib.axes._subplots.AxesSubplot at 0x25aee1d0e48>
```



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 [290]: 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 [291]: 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)

```

```
{'their', 'its', 'his', 're', "wouldn't", "you've", 'was', 'we', 'of', "you'r
e", 'a', 'do', 'while', 'been', 'into', 's', 'what', "should've", 'for', 'bef
ore', 'shan', 'o', 'mustn', 'because', 'or', "it's", 'they', 'and', 'off', 'o
ther', 'd', 'your', 'more', "shouldn't", 'during', 'who', "that'll", 'furthe
r', 'didn', 'so', 'from', 'all', 'wouldn', 'about', "mustn't", 'him', 'it',
'am', 'himself', 'doing', 'aren', 'an', 'are', 'being', 'now', 'shouldn', 'ov
er', 'to', 'you', 'y', 'than', 'just', 'with', "mightn't", 'yourselves', 'som
e', 'the', 'be', 'between', 'having', "wasn't", 'same', 'yours', 'down', "nee
dn't", 'were', 'he', 'll', 'how', 'doesn', 'but', 'this', 'ma', 'itself', 'th
emselves', 'once', 'had', 'those', 'is', 'not', 'm', 'ain', "couldn't", "yo
u'll", "aren't", 'mightn', 'by', 'any', 'where', 'own', 'on', 'hasn', 'both',
"doesn't", 'then', "shan't", 'until', 'under', 't', "isn't", 'through', 'was
n', 'did', 'them', 'won', 'up', "don't", 'such', 'after', 'here', 've', 'thes
e', "hasn't", "you'd", "she's", 'most', 'again', 'when', 'ours', 'too', 'abov
e', 'out', 'she', 'myself', 'each', 'below', 'have', 'why', 'will', 'in', 'wh
om', 'herself', 'at', "didn't", 'her', 'which', 'very', "weren't", 'only', 'i
f', 'ourselves', "hadn't", 'me', 'as', 'couldn', 'has', 'few', 'that', 'shoul
d', 'my', 'theirs', 'yourself', 'hers', 'our', 'no', 'can', 'haven', 'nor',
'needn', 'against', 'isn', 'there', 'i', 'does', "won't", "haven't", 'weren',
'hadn', 'don'}
```

```
{'their', 'its', 'his', 're', "wouldn't", "you've", 'was', 'we', 'of', "you'r
e", 'a', 'do', 'while', 'been', 'into', 's', 'what', "should've", 'for', 'bef
ore', 'shan', 'o', 'mustn', 'because', 'or', "it's", 'they', 'and', 'off', 'o
ther', 'd', 'your', 'more', "shouldn't", 'during', 'who', "that'll", 'furthe
r', 'didn', 'so', 'from', 'all', 'wouldn', 'about', "mustn't", 'him', 'it',
'am', 'himself', 'doing', 'aren', 'an', 'are', 'being', 'now', 'shouldn', 'ov
er', 'to', 'you', 'y', 'than', 'just', 'with', "mightn't", 'yourselves', 'som
e', 'the', 'be', 'between', 'having', "wasn't", 'same', 'yours', 'down', "nee
dn't", 'were', 'he', 'll', 'how', 'doesn', 'but', 'this', 'ma', 'itself', 'th
emselves', 'once', 'had', 'those', 'is', 'not', 'm', 'ain', "couldn't", "yo
u'll", "aren't", 'mightn', 'by', 'any', 'where', 'own', 'on', 'hasn', 'both',
"doesn't", 'then', "shan't", 'until', 'under', 't', "isn't", 'through', 'was
n', 'did', 'them', 'won', 'up', "don't", 'such', 'after', 'here', 've', 'thes
e', "hasn't", "you'd", "she's", 'most', 'again', 'when', 'ours', 'too', 'abov
e', 'out', 'she', 'myself', 'each', 'below', 'have', 'why', 'will', 'in', 'wh
om', 'herself', 'at', "didn't", 'her', 'which', 'very', "weren't", 'only', 'i
f', 'ourselves', "hadn't", 'me', 'as', 'couldn', 'has', 'few', 'that', 'shoul
d', 'my', 'theirs', 'yourself', 'hers', 'our', 'no', 'can', 'haven', 'nor',
'needn', 'against', 'isn', 'there', 'i', 'does', "won't", "haven't", 'weren',
'hadn', 'don'}
```

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

```
In [292]: 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\u00e2tisserie ou un caf\u00e9/th\u00e9. \"J'ai prit un
th\u00e9 qui \u00e9tait vraiment bon, et je me suis pos\u00e9 devant une des 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[292]: 'sadli as of juli 28 2016 silverstein bakeri is perman close i went there tod
ay in person and found the bad news post on their 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 [293]: from keras.losses import mean_absolute_error, binary_crossentropy, categorical_
_crossentropy

def my_custom_loss_ova(y_true, y_pred):
    mse = mean_absolute_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_absolute_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 [294]: yelp = pd.read_csv('cleaned_yelp_stemmed.csv')
yelp.head()
```

Out[294]:

	Unnamed: 0	review_id	text	stars	category
0	0	Q1sbwVQXV2734tPgoKj4Q	total bill for thi horribl servic over 8g thes...	1	Negative
1	1	GJXCdrto3ASJOqKeVWPi6Q	i ador travi at the hard rock s new kelli card...	5	Positive
2	2	2TzJjDVDEuAW6MR5Vuc1ug	i have to say that thi offic realli ha it toge...	5	Positive
3	3	yi0R0Ugj_xUx_Nek0-_Qig	went in for a lunch steak sandwich wa delici a...	5	Positive
4	4	11a8sVPMUFtaC7_ABRkmtw	today wa my second out of three session i had ...	1	Negative

```
In [295]: 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, rand
om_state=42)
```



```
In [299]: %%time
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)
Wall time: 1min 18s
```

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

```
In [297]: # # 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 [300]: 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=regulariz
          ers.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', 'mean_absolute_error'])

          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 [=====] - 39s 129us/step - loss: 0.4265 - accuracy: 0.6874 - mean_absolute_error: 0.1358 - val_loss: 0.2666 - val_accuracy: 0.7322 - val_mean_absolute_error: 0.1141

Epoch 2/10

298804/298804 [=====] - 14s 46us/step - loss: 0.2119 - accuracy: 0.7363 - mean_absolute_error: 0.1098 - val_loss: 0.1801 - val_accuracy: 0.7347 - val_mean_absolute_error: 0.1100

Epoch 3/10

298804/298804 [=====] - 12s 39us/step - loss: 0.1619 - accuracy: 0.7414 - mean_absolute_error: 0.1069 - val_loss: 0.1533 - val_accuracy: 0.7343 - val_mean_absolute_error: 0.1093

Epoch 4/10

298804/298804 [=====] - 11s 38us/step - loss: 0.1427 - accuracy: 0.7450 - mean_absolute_error: 0.1052 - val_loss: 0.1413 - val_accuracy: 0.7340 - val_mean_absolute_error: 0.1088

Epoch 5/10

298804/298804 [=====] - 11s 38us/step - loss: 0.1330 - accuracy: 0.7474 - mean_absolute_error: 0.1039 - val_loss: 0.1346 - val_accuracy: 0.7348 - val_mean_absolute_error: 0.1081

Epoch 6/10

298804/298804 [=====] - 11s 38us/step - loss: 0.1274 - accuracy: 0.7500 - mean_absolute_error: 0.1028 - val_loss: 0.1309 - val_accuracy: 0.7365 - val_mean_absolute_error: 0.1077

Epoch 7/10

298804/298804 [=====] - 11s 38us/step - loss: 0.1239 - accuracy: 0.7521 - mean_absolute_error: 0.1019 - val_loss: 0.1286 - val_accuracy: 0.7361 - val_mean_absolute_error: 0.1075

Epoch 8/10

298804/298804 [=====] - 11s 38us/step - loss: 0.1217 - accuracy: 0.7536 - mean_absolute_error: 0.1013 - val_loss: 0.1270 - val_accuracy: 0.7369 - val_mean_absolute_error: 0.1070

Epoch 9/10

298804/298804 [=====] - 11s 38us/step - loss: 0.1199 - accuracy: 0.7554 - mean_absolute_error: 0.1005 - val_loss: 0.1260 - val_accuracy: 0.7370 - val_mean_absolute_error: 0.1070

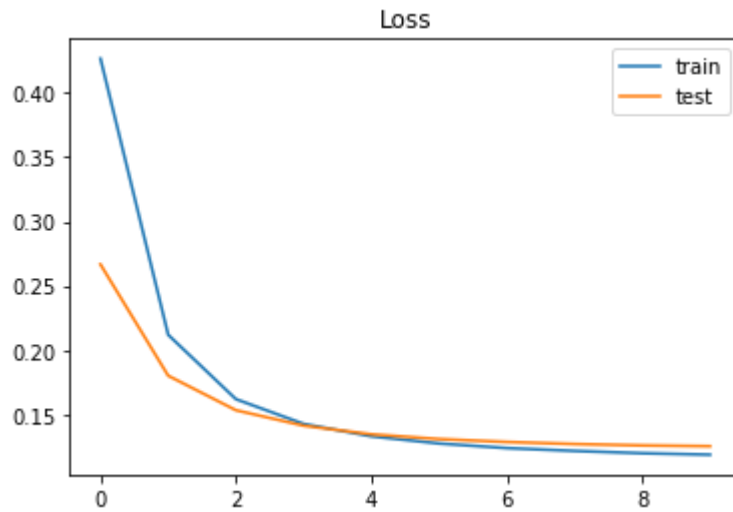
Epoch 10/10

298804/298804 [=====] - 11s 38us/step - loss: 0.1188 - accuracy: 0.7565 - mean_absolute_error: 0.1000 - val_loss: 0.1253 - val_accuracy: 0.7367 - val_mean_absolute_error: 0.1069

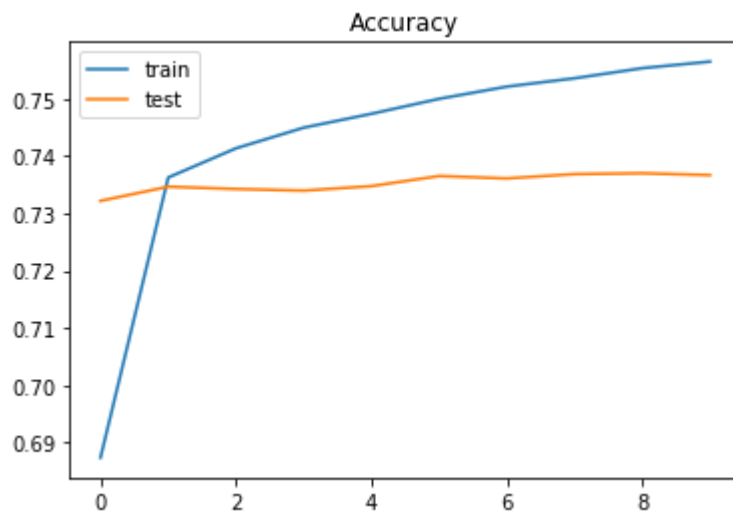
```
In [301]: score = baseline.evaluate(X_test, y_test,
                                     batch_size=batch_size, verbose=1)
print('Test accuracy:', score[1])
```

160075/160075 [=====] - 29s 180us/step
Test accuracy: 0.7380602955818176

```
In [302]: plt.title('Loss')
plt.plot(history.history['loss'], label='train')
plt.plot(history.history['val_loss'], label='test')
plt.legend()
plt.show()
```



```
In [303]: plt.title('Accuracy')
plt.plot(history.history['accuracy'], label='train')
plt.plot(history.history['val_accuracy'], label='test')
plt.legend()
plt.show()
```



```
In [304]: # 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[304]:

	1	2	3	4	5
1	37047	8039	3405	1515	2033
2	0	0	0	0	0
3	0	0	0	0	0
4	474	1738	4872	8133	3423
5	1366	966	1986	12113	72965

```
In [305]: print(classification_report(y_pred_true, y_test_true))
```

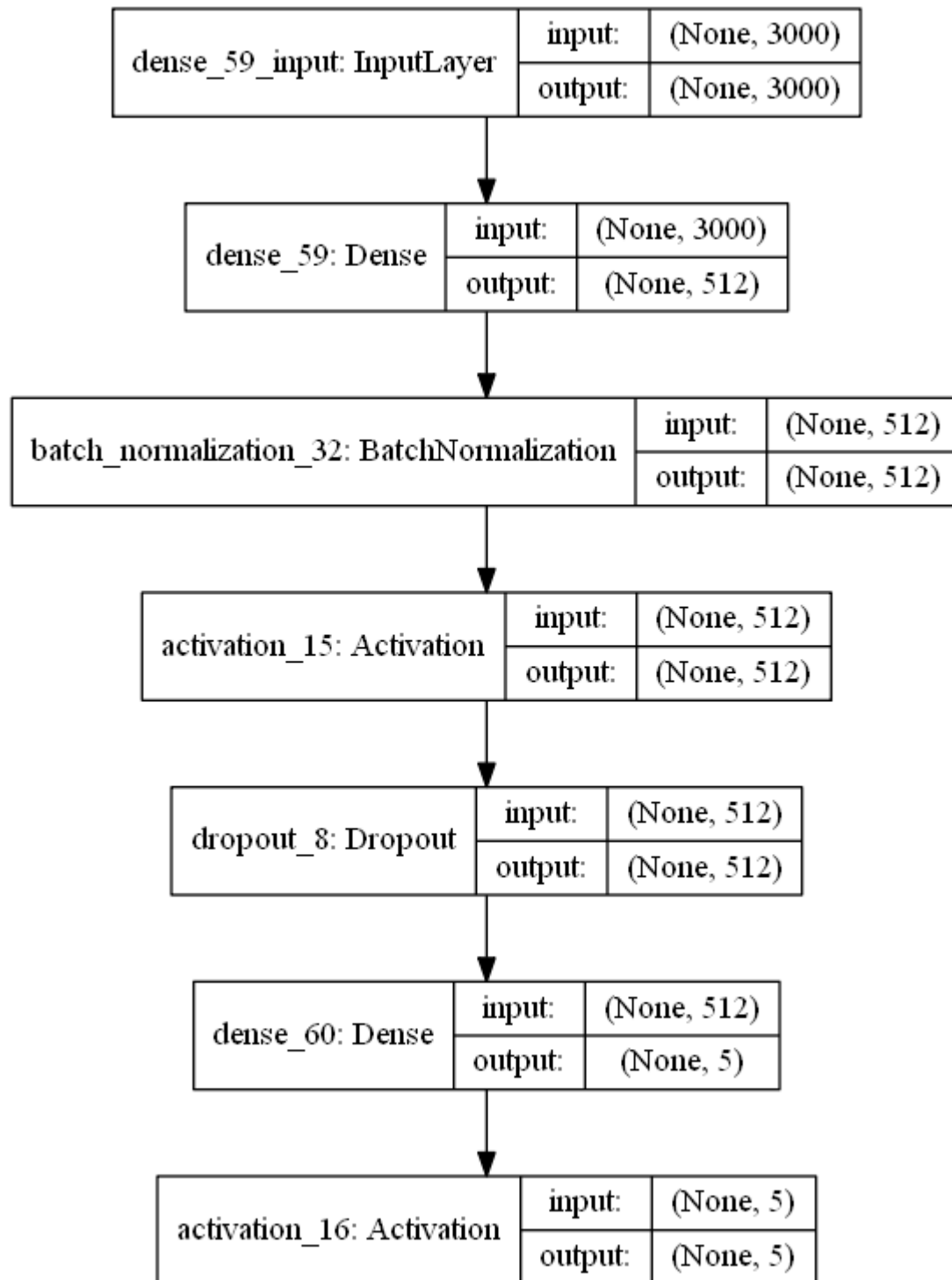
	precision	recall	f1-score	support
1	0.95	0.71	0.81	52039
2	0.00	0.00	0.00	0
3	0.00	0.00	0.00	0
4	0.37	0.44	0.40	18640
5	0.93	0.82	0.87	89396
accuracy			0.74	160075
macro avg	0.45	0.39	0.42	160075
weighted avg	0.87	0.74	0.80	160075

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

Out[306]:



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 [307]: %%time
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)
Wall time: 42.1 s
```

LSTM #1

```

In [308]: 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', 'mean_absolute_error'])

          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 [=====] - 82s 274us/step - loss: 0.1554 - accuracy: 0.6814 - mean_absolute_error: 0.1417 - val_loss: 0.1066 - val_accuracy: 0.7164 - val_mean_absolute_error: 0.1004

Epoch 2/5

298804/298804 [=====] - 81s 270us/step - loss: 0.1025 - accuracy: 0.7247 - mean_absolute_error: 0.0976 - val_loss: 0.1025 - val_accuracy: 0.7280 - val_mean_absolute_error: 0.0982

Epoch 3/5

298804/298804 [=====] - 81s 271us/step - loss: 0.0998 - accuracy: 0.7301 - mean_absolute_error: 0.0954 - val_loss: 0.0981 - val_accuracy: 0.7354 - val_mean_absolute_error: 0.0936

Epoch 4/5

298804/298804 [=====] - 81s 270us/step - loss: 0.0988 - accuracy: 0.7313 - mean_absolute_error: 0.0943 - val_loss: 0.0976 - val_accuracy: 0.7370 - val_mean_absolute_error: 0.0931

Epoch 5/5

298804/298804 [=====] - 81s 270us/step - loss: 0.0979 - accuracy: 0.7332 - mean_absolute_error: 0.0933 - val_loss: 0.1018 - val_accuracy: 0.7275 - val_mean_absolute_error: 0.0970

LSTM #1: Evaluation

```
In [309]: score = lstm.evaluate(X_test, y_test,
                                batch_size=batch_size, verbose=1)
print('Test accuracy:', score[1])
```

160075/160075 [=====] - 10s 66us/step
 Test accuracy: 0.7289021015167236

```
In [310]: lstm.summary()
```

Model: "sequential_33"

Layer (type)	Output Shape	Param #
=====		
embedding_25 (Embedding)	(None, 400, 128)	384000

spatial_dropout1d_25 (SpatialDropout1D)	(None, 400, 128)	0

conv1d_25 (Conv1D)	(None, 396, 64)	41024

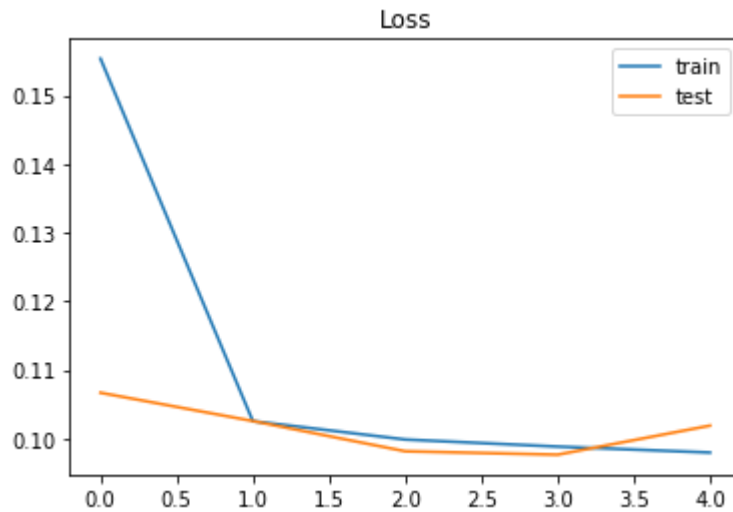
max_pooling1d_25 (MaxPooling1D)	(None, 99, 64)	0

lstm_25 (LSTM)	(None, 128)	98816

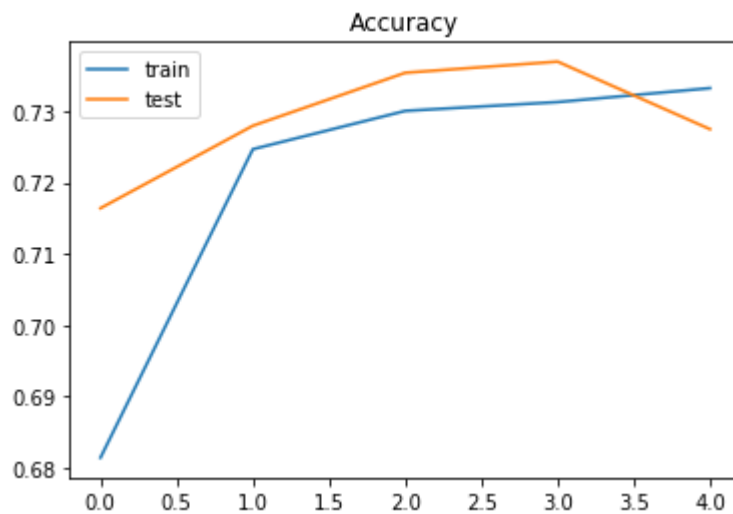
batch_normalization_33 (Batch Normalization)	(None, 128)	512

dense_61 (Dense)	(None, 5)	645
=====		
Total params: 524,997		
Trainable params: 524,741		
Non-trainable params: 256		

```
In [311]: plt.title('Loss')
plt.plot(history.history['loss'], label='train')
plt.plot(history.history['val_loss'], label='test')
plt.legend()
plt.show()
```



```
In [312]: plt.title('Accuracy')
plt.plot(history.history['accuracy'], label='train')
plt.plot(history.history['val_accuracy'], label='test')
plt.legend()
plt.show()
```



```
In [313]: # 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[313]:

	1	2	3	4	5
1	34828	6472	2160	812	609
2	0	0	0	0	0
3	0	0	0	0	0
4	662	2533	5086	5731	1692
5	3397	1738	3017	15218	76120

```
In [314]: print(classification_report(y_pred_true, y_test_true))
```

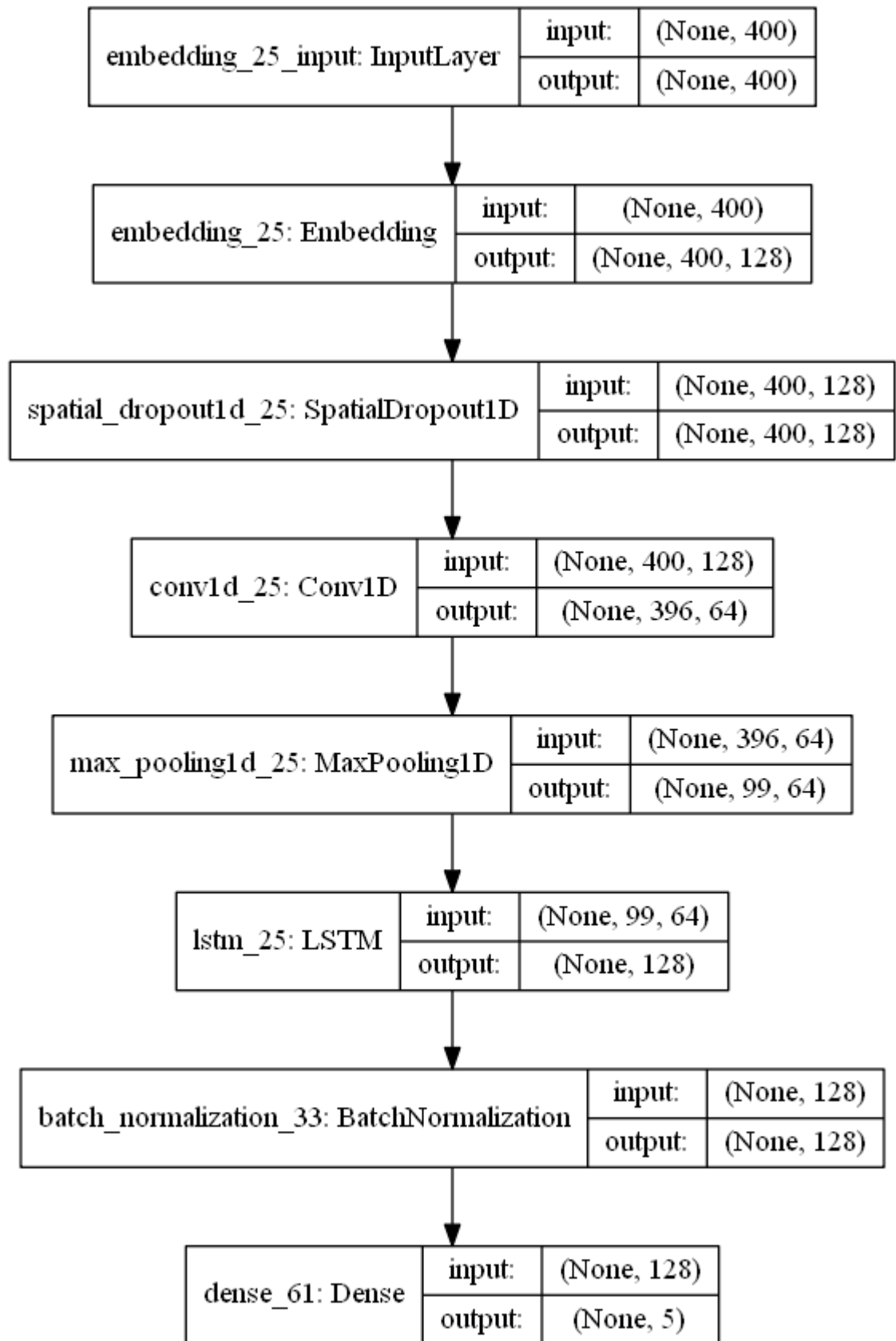
	precision	recall	f1-score	support
1	0.90	0.78	0.83	44881
2	0.00	0.00	0.00	0
3	0.00	0.00	0.00	0
4	0.26	0.36	0.31	15704
5	0.97	0.77	0.86	99490
accuracy			0.73	160075
macro avg	0.43	0.38	0.40	160075
weighted avg	0.88	0.73	0.79	160075

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))
```

In [315]: `plot_model(lstm, to_file='baseline.png', show_shapes=True)`

Out[315]:



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 [316]: 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 [318]: stars = np.arange(1, 6)
models = {}
histories = {}
batch_size = 512

for star in stars:
    if star in [1, 2]:
        epochs = 2
    elif star in [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', 'mean_absolute_error'])

    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 [=====] - 78s 261us/step - loss: 0.119
7 - accuracy: 0.9001 - mean_absolute_error: 0.1019 - val_loss: 0.0945 - val_a
ccuracy: 0.9181 - val_mean_absolute_error: 0.0823
Epoch 2/2
298804/298804 [=====] - 77s 258us/step - loss: 0.095
6 - accuracy: 0.9152 - mean_absolute_error: 0.0851 - val_loss: 0.0937 - val_a
ccuracy: 0.9155 - val_mean_absolute_error: 0.0846
2
Train on 298804 samples, validate on 74702 samples
Epoch 1/2
298804/298804 [=====] - 78s 262us/step - loss: 0.100
3 - accuracy: 0.9198 - mean_absolute_error: 0.0889 - val_loss: 0.0688 - val_a
ccuracy: 0.9323 - val_mean_absolute_error: 0.0678
Epoch 2/2
298804/298804 [=====] - 78s 260us/step - loss: 0.067
3 - accuracy: 0.9329 - mean_absolute_error: 0.0671 - val_loss: 0.0678 - val_a
ccuracy: 0.9323 - val_mean_absolute_error: 0.0677
3
Train on 298804 samples, validate on 74702 samples
Epoch 1/3
298804/298804 [=====] - 78s 262us/step - loss: 0.103
0 - accuracy: 0.9123 - mean_absolute_error: 0.0918 - val_loss: 0.0675 - val_a
ccuracy: 0.9363 - val_mean_absolute_error: 0.0664
Epoch 2/3
298804/298804 [=====] - 78s 260us/step - loss: 0.064
5 - accuracy: 0.9356 - mean_absolute_error: 0.0644 - val_loss: 0.0660 - val_a
ccuracy: 0.9363 - val_mean_absolute_error: 0.0660
Epoch 3/3
298804/298804 [=====] - 77s 258us/step - loss: 0.064
4 - accuracy: 0.9356 - mean_absolute_error: 0.0644 - val_loss: 0.0638 - val_a
ccuracy: 0.9363 - val_mean_absolute_error: 0.0638
4
Train on 298804 samples, validate on 74702 samples
Epoch 1/3
298804/298804 [=====] - 78s 261us/step - loss: 0.170
4 - accuracy: 0.8446 - mean_absolute_error: 0.1588 - val_loss: 0.1640 - val_a
ccuracy: 0.8637 - val_mean_absolute_error: 0.1622
Epoch 2/3
298804/298804 [=====] - 77s 259us/step - loss: 0.135
9 - accuracy: 0.8645 - mean_absolute_error: 0.1355 - val_loss: 0.1362 - val_a
ccuracy: 0.8639 - val_mean_absolute_error: 0.1362
Epoch 3/3
298804/298804 [=====] - 77s 258us/step - loss: 0.135
6 - accuracy: 0.8645 - mean_absolute_error: 0.1355 - val_loss: 0.1362 - val_a
ccuracy: 0.8639 - val_mean_absolute_error: 0.1361
5
Train on 298804 samples, validate on 74702 samples
Epoch 1/4
298804/298804 [=====] - 78s 261us/step - loss: 0.165
5 - accuracy: 0.8532 - mean_absolute_error: 0.1488 - val_loss: 0.1437 - val_a
ccuracy: 0.8686 - val_mean_absolute_error: 0.1319
Epoch 2/4
298804/298804 [=====] - 77s 259us/step - loss: 0.140
5 - accuracy: 0.8701 - mean_absolute_error: 0.1301 - val_loss: 0.1404 - val_a
```

```

ccuracy: 0.8692 - val_mean_absolute_error: 0.1309
Epoch 3/4
298804/298804 [=====] - 77s 259us/step - loss: 0.135
7 - accuracy: 0.8736 - mean_absolute_error: 0.1266 - val_loss: 0.1432 - val_a
ccuracy: 0.8657 - val_mean_absolute_error: 0.1343
Epoch 4/4
298804/298804 [=====] - 77s 259us/step - loss: 0.133
6 - accuracy: 0.8755 - mean_absolute_error: 0.1246 - val_loss: 0.1668 - val_a
ccuracy: 0.8423 - val_mean_absolute_error: 0.1578

```

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

Testing

```

In [319]: %%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], mode
ls[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
0.48825238169608
0.7091113540527878
Wall time: 5min 41s

```

```
In [320]: # Confusion matrix
cm = confusion_matrix(ps["pred"], y_test_true[0])
pd.DataFrame(cm, index=cols, columns=cols)
```

Out[320]:

	1	2	3	4	5
1	31709	3920	980	341	411
2	0	0	0	0	0
3	0	0	0	0	0
4	3980	5312	6473	5945	2153
5	3198	1511	2810	15475	75857

```
In [321]: print(classification_report(ps["pred"], y_test_true[0]))
```

	precision	recall	f1-score	support
1	0.82	0.85	0.83	37361
2	0.00	0.00	0.00	0
3	0.00	0.00	0.00	0
4	0.27	0.25	0.26	23863
5	0.97	0.77	0.86	98851
accuracy			0.71	160075
macro avg	0.41	0.37	0.39	160075
weighted avg	0.83	0.71	0.76	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 [322]: 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 [ ]: ## 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 [323]: 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_pr
ed'], ova_preds['ova_pred']]).T
all_preds["final_pred"] = all_preds.mode(axis=1)[0]

Baseline
LSTM
OVA
```

```
In [324]: print([MAE(all_preds["final_pred"], pd.DataFrame(data=y_test, columns=cols).id
xmax(axis=1)), Accuracy(all_preds["final_pred"], pd.DataFrame(data=y_test, col
umns=cols).idxmax(axis=1))])

[0.4078150866781196, 0.7353927846322037]
```

```
In [325]: # Confusion matrix
cm = confusion_matrix(all_preds["final_pred"], pd.DataFrame(data=y_test, colum
ns=cols).idxmax(axis=1))
pd.DataFrame(cm, index=cols, columns=cols)
```

Out[325]:

	1	2	3	4	5
1	35741	6709	2311	886	766
2	0	0	0	0	0
3	0	0	0	0	0
4	772	2731	5412	6201	1879
5	2374	1303	2540	14674	75776

```
In [326]: print(classification_report(y_pred_true, y_test_true))
```

	precision	recall	f1-score	support
1	0.90	0.78	0.83	44881
2	0.00	0.00	0.00	0
3	0.00	0.00	0.00	0
4	0.26	0.36	0.31	15704
5	0.97	0.77	0.86	99490
accuracy			0.73	160075
macro avg	0.43	0.38	0.40	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))
```

Challenges

Challenge 5

```
In [327]: c5 = pd.read_json("./yelp_challenge_5_with_answers.jsonl", lines = True)
print(c5.shape)
c5.head()
```

```
(500, 3)
```

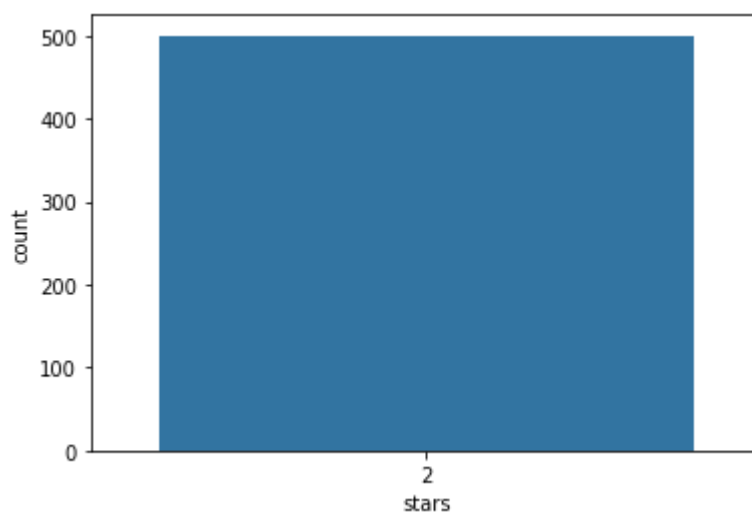
Out[327]:

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

```
Out[328]: <matplotlib.axes._subplots.AxesSubplot at 0x25aedd4de08>
```



Pre-processing

```
In [329]: c5['text'] = c5['text'].apply(clean_text)  
c5.head()
```

```
Out[329]:
```

	review_id	text	stars
0	50	i went to thi campu for 1 semest i wa in busi ...	2
1	51	i have rate it a two star base on it compariso...	2
2	52	just like most of the review we order and paid...	2
3	53	i onli go here if it is an emerg i hate it tha...	2
4	54	rude staff i got 60 feeder fish and about 15 w...	2

Load previous tokenizer

```
In [330]: 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 [331]: ## 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 [332]: # 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 86us/step
[0.4163735044002533, 0.0, 0.3999999165534973]
500/500 [=====] - 0s 526us/step
[0.305676802277565, 0.0, 0.30080291628837585]
[1.802, 0.0]

```

Attempt Ensemble

```

In [333]: # 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))])

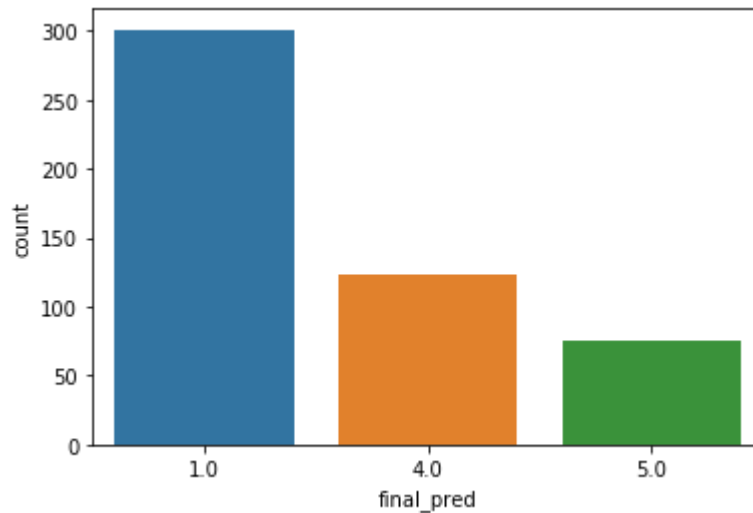
[1.55, 0.0]

```

Misc.

```
In [334]: sns.countplot(all_preds["final_pred"])
```

```
Out[334]: <matplotlib.axes._subplots.AxesSubplot at 0x25a80c5a288>
```

**Challenge 6**

```
In [335]: c6 = pd.read_json("./yelp_challenge_6_with_answers.jsonl", lines = True)
          print(c6.shape)
          c6.head()
```

```
(500, 3)
```

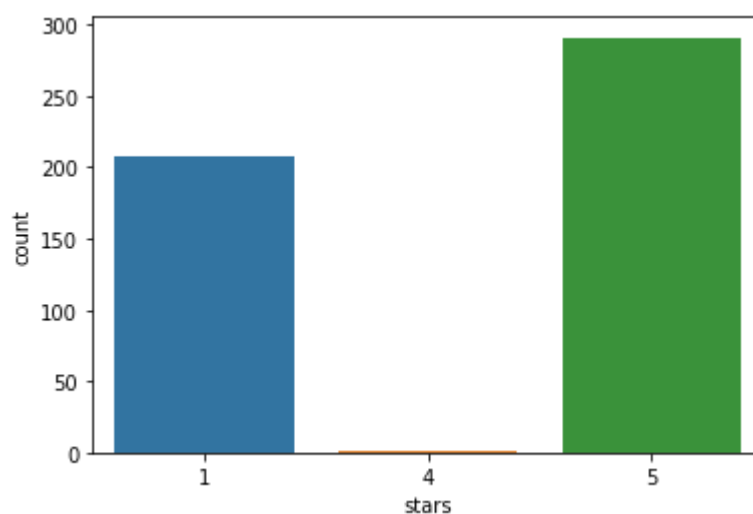
```
Out[335]:
```

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

```
Out[336]: <matplotlib.axes._subplots.AxesSubplot at 0x25a6d15dd88>
```



Pre-processing

```
In [337]: c6['text'] = c6['text'].apply(clean_text)
c6.head()
```

```
Out[337]:
```

	review_id	text	stars
0	60	amaz for tree 20 for a 5 gallon i will never g...	5
1	61	how the hell can taco bell be close befor midn...	5
2	62	i actual had no intent of visit thi place at a...	5
3	63	yesterday around 3 30 pm i wa drive west on pi...	5
4	64	dr fitzmauric did surgeri on both hand on the ...	5

Load previous tokenizer


```
In [338]: 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 [339]: ## 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 [340]: # 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 82us/step
[0.23861910700798034, 0.4399999976158142, 0.2222854048013687]
500/500 [=====] - 0s 496us/step
[0.22340433382987976, 0.43799999356269836, 0.21853041648864746]
[2.218, 0.392]

```

Attempt Ensemble

```

In [341]: # 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))])

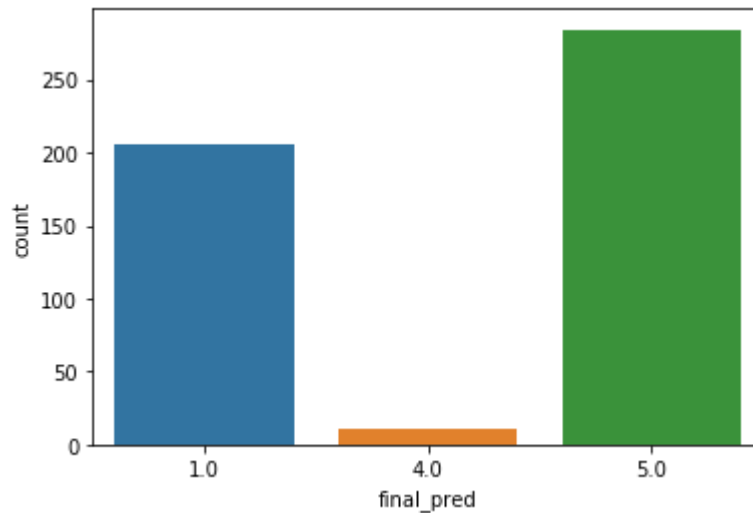
[2.192, 0.442]

```

Misc.

```
In [342]: sns.countplot(all_preds["final_pred"])
```

```
Out[342]: <matplotlib.axes._subplots.AxesSubplot at 0x25af5ff73c8>
```

**Challenge 3**

```
In [343]: c3 = pd.read_json("./yelp_challenge_3_with_answers.jsonl", lines = True)
print(c3.shape)
c3.head()
```

```
(534, 3)
```

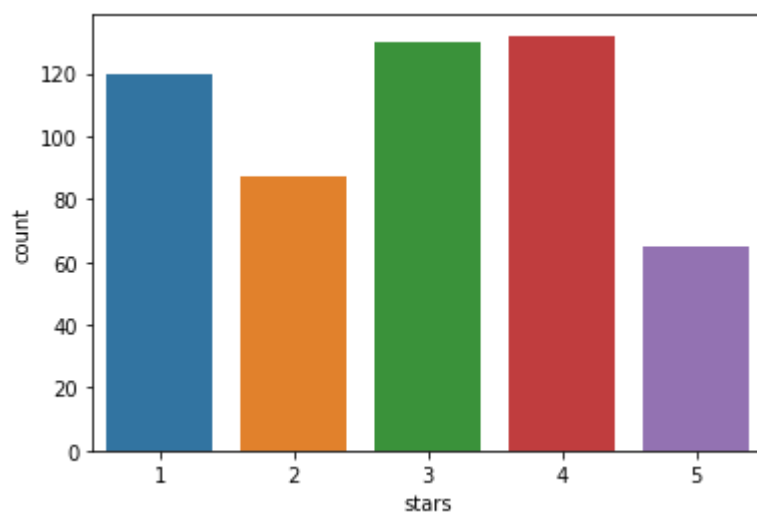
```
Out[343]:
```

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

```
Out[344]: <matplotlib.axes._subplots.AxesSubplot at 0x25a6491cc88>
```



Pre-processing

```
In [345]: c3['text'] = c3['text'].apply(clean_text)
c3.head()
```

```
Out[345]:
```

	review_id	text	stars
0	30	we stop here for lunch today and were pleasant...	4
1	31	we went for a quick lunch here it s all reason...	3
2	32	veri bad food avoid it we were a group of 4 an...	2
3	33	bring a friend or two to help open the door i ...	3
4	34	ukai serv some of the best sushi and sashimi i...	4

Load previous tokenizer

```
In [346]: 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 [347]: ## 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 [348]: # 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 77us/step
[0.24015130800254336, 0.44756555557250977, 0.22375977039337158]
534/534 [=====] - 0s 517us/step
[0.1862835054763694, 0.40449437499046326, 0.18140959739685059]
[1.0, 0.34269662921348315]
```

Attempt Ensemble

```
In [349]: # 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]

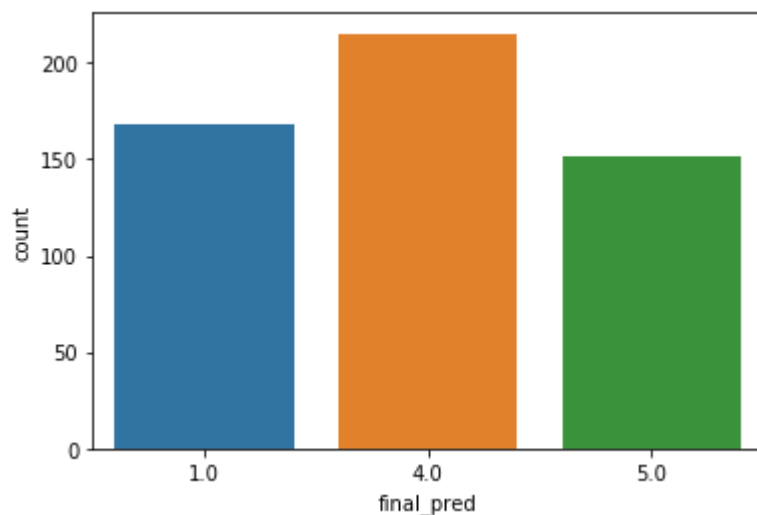
print([MAE(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols).idxmax(
axis=1)), Accuracy(all_preds["final_pred"], pd.DataFrame(data=y, columns=cols)
.idxmax(axis=1))])
```

```
[0.797752808988764, 0.40823970037453183]
```


Misc.

```
In [350]: sns.countplot(all_preds["final_pred"])
```

```
Out[350]: <matplotlib.axes._subplots.AxesSubplot at 0x25b582616c8>
```

**Challenge 8**

```
In [351]: c8 = pd.read_json("./yelp_challenge_8_with_answers.jsonl", lines = True)
print(c8.shape)
c8.head()
```

```
(500, 3)
```

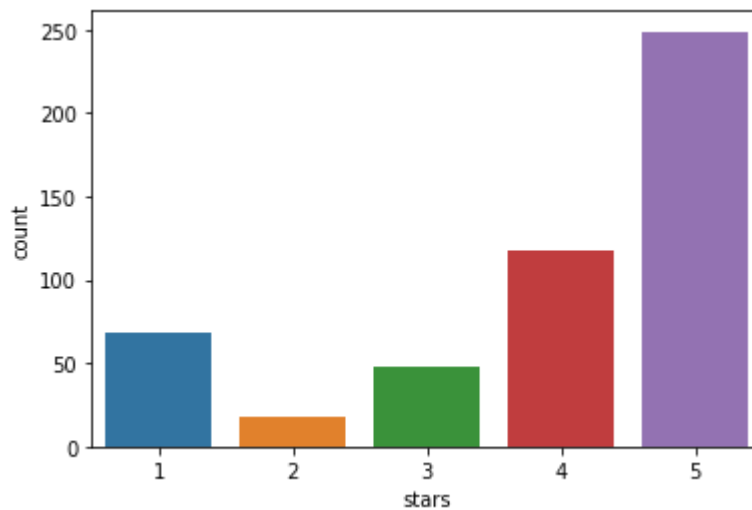
```
Out[351]:
```

	review_id	text	stars
0	qOOv-A-vo3kMT0yi4jlllg	Not bad for fast food.	4
1	uqxo6B6w_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 [352]: sns.countplot(c8['stars'])
```

```
Out[352]: <matplotlib.axes._subplots.AxesSubplot at 0x25b506d5c88>
```



Pre-processing

```
In [353]: 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[353]:
```

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

Load previous tokenizer

```
In [354]: 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 [355]: ## 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 [356]: # 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 82us/step
[0.17698445343971253, 0.6100000143051147, 0.16073599457740784]
500/500 [=====] - 0s 498us/step
[0.14596080994606017, 0.6060000061988831, 0.1410869061946869]
[0.634, 0.598]

```

Attempt Ensemble

```

In [357]: # 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))])

[0.594, 0.626]

```

Misc.

```
In [358]: sns.countplot(all_preds["final_pred"])
```

```
Out[358]: <matplotlib.axes._subplots.AxesSubplot at 0x25acbc85608>
```

