

# Amazon Food Review Rating Prediction

This file experiments with different neural network based models to come up with a great one to submit to the class Kaggle Competition

The pretrained GloVe models used to produce meaningful embeddings were obtained from <https://nlp.stanford.edu/projects/glove/> (<https://nlp.stanford.edu/projects/glove/>). (Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation.)

```
In [2]: 1 import pandas as pd
        2 import numpy as np
        3 from matplotlib import pyplot as plt
        4 plt.style.use("seaborn")
        5 import seaborn as sns
```

## Get Data

```
In [12]: 1 # train data
        2 df_train = pd.read_csv("./data/train.csv")
        3 display(df_train.head(3))
        4 print(df_train.shape)
```

	Score	Review_text
0	5	I received this product early from the seller!...
1	5	***** Numi's Collection Assortment Melang...
2	5	I was very careful not to overcook this pasta,...

(426340, 2)

```
In [167]: 1 # test data
        2 df_test = pd.read_csv("./data/test_new.csv")
        3 display(df_test.head(3))
        4 print(df_test.shape)
```

	Id	Review_text
0	1	I have a very picky German Shephard mix and sh...
1	2	It is hard to believe that this candy is sugar...
2	3	These are delicious cookies but I just cancell...

(142114, 2)

## Clean Data

```
In [168]: 1 import re
        2 from emot.emo_unicode import EMOTIONS_EMO # for emoticon expansion
```

In [169]:

```
1 CONTRACTIONS = {
2     "ain't": "is not",
3     "aren't": "are not",
4     "can't": "cannot",
5     "can't've": "cannot have",
6     "'cause": "because",
7     "could've": "could have",
8     "couldn't": "could not",
9     "couldn't've": "could not have",
10    "didn't": "did not",
11    "doesn't": "does not",
12    "don't": "do not",
13    "hadn't": "had not",
14    "hadn't've": "had not have",
15    "hasn't": "has not",
16    "haven't": "have not",
17    "he'd": "he would",
18    "he'd've": "he would have",
19    "he'll": "he will",
20    "he'll've": "he will have",
21    "he's": "he is",
22    "how'd": "how did",
23    "how'd'y": "how do you",
24    "how'll": "how will",
25    "how's": "how is",
26    "i'd": "i would",
27    "i'd've": "i would have",
28    "i'll": "i will",
29    "i'll've": "i will have",
30    "i'm": "i am",
31    "i've": "i have",
32    "isn't": "is not",
33    "it'd": "it would",
34    "it'd've": "it would have",
35    "it'll": "it will",
36    "it'll've": "it will have",
37    "it's": "it is",
38    "let's": "let us",
39    "ma'am": "madam",
40    "mayn't": "may not",
41    "might've": "might have",
42    "mightn't": "might not",
43    "mightn't've": "might not have",
44    "must've": "must have",
45    "mustn't": "must not",
46    "mustn't've": "must not have",
47    "needn't": "need not",
48    "needn't've": "need not have",
49    "o'clock": "of the clock",
50    "oughtn't": "ought not",
51    "oughtn't've": "ought not have",
52    "shan't": "shall not",
53    "sha'n't": "shall not",
54    "shan't've": "shall not have",
55    "she'd": "she would",
56    "she'd've": "she would have",
57    "she'll": "she will",
58    "she'll've": "she will have",
59    "she's": "she is",
```

```
60     "should've": "should have",
61     "shouldn't": "should not",
62     "shouldn't've": "should not have",
63     "so've": "so have",
64     "so's": "so as",
65     "that'd": "that would",
66     "that'd've": "that would have",
67     "that's": "that is",
68     "there'd": "there would",
69     "there'd've": "there would have",
70     "there's": "there is",
71     "they'd": "they would",
72     "they'd've": "they would have",
73     "they'll": "they will",
74     "they'll've": "they will have",
75     "they're": "they are",
76     "they've": "they have",
77     "to've": "to have",
78     "wasn't": "was not",
79     "we'd": "we would",
80     "we'd've": "we would have",
81     "we'll": "we will",
82     "we'll've": "we will have",
83     "we're": "we are",
84     "we've": "we have",
85     "weren't": "were not",
86     "what'll": "what will",
87     "what'll've": "what will have",
88     "what're": "what are",
89     "what's": "what is",
90     "what've": "what have",
91     "when's": "when is",
92     "when've": "when have",
93     "where'd": "where did",
94     "where's": "where is",
95     "where've": "where have",
96     "who'll": "who will",
97     "who'll've": "who will have",
98     "who's": "who is",
99     "who've": "who have",
100    "why's": "why is",
101    "why've": "why have",
102    "will've": "will have",
103    "won't": "will not",
104    "won't've": "will not have",
105    "would've": "would have",
106    "wouldn't": "would not",
107    "wouldn't've": "would not have",
108    "y'all": "you all",
109    "y'all'd": "you all would",
110    "y'all'd've": "you all would have",
111    "y'all're": "you all are",
112    "y'all've": "you all have",
113    "you'd": "you would",
114    "you'd've": "you would have",
115    "you'll": "you will",
116    "you'll've": "you will have",
117    "you're": "you are",
118    "you've": "you have"
119 }
```

```

120
121 def clean(text):
122     ''' Function that returns a given piece of text after cleaning it.
123     # check for html embeddings and sub them out with a space
124     text = re.sub("<[<>]*>", " ",text)
125
126     # if an emoticon is encountered, replace with word equivalent
127     for emote in EMOTICONS_EMO:
128         text = text.replace(emote, " "+EMOTICONS_EMO[emote])
129
130     # make Lowercase
131     text = text.lower().strip()
132
133     # expand contractions
134     for key, val in CONTRACTIONS.items():
135         text = re.sub(key, val, text)
136
137     # discard 's
138     text = re.sub("'s", "",text)
139
140     # discard if it is not a-z or space
141     text = re.sub('-', " ",text)
142
143     # discard if it is not a-z or space
144     text = re.sub('[^a-z ]*', "",text)
145
146     # discard urls
147     text = re.sub('(https?\w*)|(www\w*)', "",text)
148
149     # discard extra spaces
150     text = re.sub('\s\s+', " ",text)
151
152     # check if final text is empty or not
153     if re.search("[a-z]", text): return text
154     else: return ""

```

**WARNING: Time consuming cell ahead!** (upto 3 mins)

```

In [20]: 1 # clean train data
2 df_train["review"] = df_train["Review_text"].apply(lambda review: clean(review))
3 df_train.drop(["Review_text"], axis=1)
4 df_train = df_train[df_train["review"] != ""]

```

```

In [170]: 1 # clean test data
2 df_test["review"] = df_test["Review_text"].apply(lambda review: clean(review))
3 df_test.drop(["Review_text"], axis=1)
4 df_test = df_test[df_test["review"] != ""]

```

## Feature Engineering

In [171]:

```
1 import nltk
2 from nltk.sentiment.vader import SentimentIntensityAnalyzer
3 nltk.download('vader_lexicon')
4 sid = SentimentIntensityAnalyzer()
5 import math
```

```
[nltk_data] Downloading package vader_lexicon to C:\Users\Gayathri
[nltk_data]   Girish Nair\AppData\Roaming\nltk_data...
[nltk_data]   Package vader_lexicon is already up-to-date!
```

**WARNING: Time consuming cell ahead!** (upto 5 mins)

In [27]:

```
1 # adding features "polarity", "word_count" and "Longest_word_Length" to
2 df_train = pd.DataFrame({
3     "review": df_train.review,
4     "word_count": df_train["review"].apply(lambda review: math.log(len(
5     "longest_word_length": df_train["review"].apply(
6         lambda review: math.log(max([len(word) for word in review.split
7     ),
8     "polarity": df_train["review"].apply(lambda review: sid.polarity_sc
9     "score": df_train.Score
10 })
11 display(df_train.head(3))
```

	review	word_count	longest_word_length	polarity	score
0	i received this product early from the seller ...	3.637586	2.079442	0.9488	5
1	numi collection assortment melange includes h...	5.755742	2.484907	0.9921	5
2	i was very careful not to overcook this pasta ...	5.323010	2.639057	0.9980	5

**WARNING: Time consuming cell ahead!** (upto 3 mins)

In [173]:

```
1 # adding features "polarity", "word_count" and "Longest_word_Length" to
2 df_test = pd.DataFrame({
3     "Id": df_test.Id,
4     "review": df_test.review,
5     "word_count": df_test["review"].apply(lambda review: math.log(len(r
6     "longest_word_length": df_test["review"].apply(
7         lambda review: math.log(max([len(word) for word in review.split
8     ),
9     "polarity": df_test["review"].apply(lambda review: sid.polarity_sco
10 })
11 display(df_test.head(3))
```

	Id	review	word_count	longest_word_length	polarity
0	1	i have a very picky german shephard mix and sh...	5.099866	2.302585	0.9584
1	2	it is hard to believe that this candy is sugar...	4.110874	2.079442	0.8720
2	3	these are delicious cookies but i just cancell...	3.970292	2.484907	0.8663

```
In [174]: 1 # standardize "word_count" and "Longest_word_Length"
2 from sklearn.preprocessing import StandardScaler
3
4 scaler = StandardScaler()
5 scaled_train = scaler.fit_transform(df_train[["word_count", "longest_wo
6 scaled_test = scaler.transform(df_test[["word_count", "longest_word_len
7
8 df_train["word_count"] = scaled_train[:, 0]
9 df_train["longest_word_length"] = scaled_train[:, 1]
10
11 df_test["word_count"] = scaled_test[:, 0]
12 df_test["longest_word_length"] = scaled_test[:, 1]
13
14 del scaled_train
15 del scaled_test
```

## Text Processing And Normalization

```
In [175]: 1 # stop word removal
2 from nltk.corpus import stopwords
3 import nltk
4
5 STOP_WORDS = set(stopwords.words('english'))
6 def remove_stopwords(text): # (M.D. Pietro, 2017)
7     ''' Function that removes stopwords from a given list of word
8         and returns this new possibly smaller list.
9         @param text: The list of words from which to remove stopwords.
10    '''
11    return " ".join([word for word in text.split(" ") if not word in ST
12
13 df_train["review"] = df_train.review.apply(lambda review: remove_stopwo
14 df_test["review"] = df_test.review.apply(lambda review: remove_stopword
```

**WARNING: Time consuming cell ahead!** (upto 30 mins)

```
In [176]: 1 # Lemmatization
2 import spacy
3 nlp = spacy.load('en_core_web_sm', disable=['ner', 'parser']) # !python
4 nlp.add_pipe('sentencizer')
5 nlp.get_pipe('attribute_ruler').add([{"TEXT": "us"}], {"LEMMA": "us"})
6 nlp.get_pipe('attribute_ruler').add([{"TEXT": "them"}], {"LEMMA": "them"}
7 nlp.get_pipe('attribute_ruler').add([{"TEXT": "nt"}], {"LEMMA": "not"})
8
9 def lemmatize(text_list):
10     doc = nlp(text_list)
11     return " ".join([token.lemma_ for token in doc])
12
13 # df_train["review"] = df_train.review.apply(lambda review: Lemmatize(r
14 df_test["review"] = df_test.review.apply(lambda review: lemmatize(revie
```

```
In [177]: 1 # save processed data
2 # df_train.to_csv("data/train_processed_full.csv", index=False)
3 df_test.to_csv("data/test_processed_full.csv", index=False)
```

## Review Vectorization

```
In [178]: 1 from sklearn.feature_extraction.text import CountVectorizer
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 from imblearn.over_sampling import SMOTE
4 from sklearn.preprocessing import label_binarize
5 from imblearn.under_sampling import NearMiss
6 import json
```

```
In [179]: 1 # Load stored processed data
2 df_train = pd.read_csv("data/train_processed_full.csv")
3 df_test = pd.read_csv("data/test_processed_full.csv")
```

```
In [181]: 1 # for resampling
2 def strategy(x, y, threshold, t='majority'):
3     ''' Function that aids in doing over/under sampling. '''
4     targets = ''
5     if t == 'majority': targets = y.value_counts() > threshold
6     elif t == 'minority': targets = y.value_counts() < threshold
7     tc = targets[targets == True].index
8     strategy = {}
9     for target in tc: strategy[target] = threshold
10    return strategy
11
12 def uo_resample(x, y, type):
13     if not type in ["under", "over"]: raise Exception("resample: invalid type")
14
15     if type == "under": # undersampling
16         resample = NearMiss(version=1, n_neighbors=3, sampling_strategy=0.5)
17         x, y = resample.fit_resample(x, y)
18
19     else: # oversampling
20         resample = SMOTE(random_state=30, sampling_strategy=strategy(x, y))
21         x, y = resample.fit_resample(x, y)
22
23     return x, y
```

**WARNING: Time consuming cell ahead!** (upto 5 mins)

```
In [24]: 1 # get bag of words vectors
2 vectorizer = TfidfVectorizer(ngram_range=(2,2), max_features=1000)
3 tfidf_train = vectorizer.fit_transform(df_train['review'])
4 tfidf_test = vectorizer.transform(df_test['review'])
5
6 # get TF-IDF scores
7 tfidf_df_train = pd.DataFrame(data=tfidf_train.toarray(), columns=vectorizer.get_feature_names_out())
8 tfidf_df_test = pd.DataFrame(data=tfidf_test.toarray(), columns=vectorizer.get_feature_names_out())
9
10 # concatenate added features
11 bow_train = pd.concat([tfidf_df_train, df_train[["polarity", "word_count"]]])
12 bow_labels_train = df_train["score"]
13 bow_test = pd.concat([tfidf_df_test, df_test[["polarity", "word_count"]]])
14
15 # resampling
16 bow_train, bow_labels_train = uo_resample(x=bow_train, y=bow_labels_train, type="majority")
17 bow_test, bow_labels_test = uo_resample(x=bow_test, y=bow_labels_test, type="majority")
```

**WARNING: Time consuming cell ahead!** (upto 3 mins)

In [73]:

```
1 # get global vectors
2 from sklearn import preprocessing
3
4 # Load pretrained embeddings
5 GLOVE_MODEL_100 = {}
6 with open("./data/glove.6B.100d.txt", "r", encoding="utf8") as f:
7     for line in f:
8         values = line.split()
9         word = values[0]
10        vector = np.asarray(values[1:], 'float32')
11        GLOVE_MODEL_100[word] = vector
12 f.close()
13
14 # get GloVe embeddings
15 def glove_embed(docs):
16     ''' Function that converts given documents
17         into vectors using a GloVe model.
18         @param docs: Documents to embed as an iterable.
19         @return docs: Embedded documents as np array.
20     '''
21     docs_embedded = []
22     for doc_words in [doc.split(" ") for doc in docs]:
23         doc_vector = np.array([0.0]*100)
24         for word in doc_words:
25             try: word_vec = GLOVE_MODEL_100[word]
26             except: word_vec = np.array([0.0]*100)
27             doc_vector += word_vec
28         docs_embedded.append(doc_vector)
29     docs_embedded = preprocessing.normalize(docs_embedded, axis=1)
30     return np.array(docs_embedded)
31
32 glove_df_train = pd.DataFrame(data=glove_embed(df_train["review"]), columns=["review_embeddings"])
33 glove_df_test = pd.DataFrame(data=glove_embed(df_test["review"]), columns=["review_embeddings"])
34
35 # concatenate added features
36 we_train = pd.concat([glove_df_train, df_train[["polarity", "word_count"]]])
37 we_labels_train = df_train["score"]
38 we_test = pd.concat([glove_df_test, df_test[["polarity", "word_count"]]])
39
40 # resampling
41 we_train, we_labels_train = resample(x=we_train, y=we_labels_train, type="train")
42 we_test, we_labels_test = resample(x=we_test, y=we_labels_test, type="test")
```

**WARNING: Time consuming cell ahead!** (upto 4 mins)

In [250]:

```
1 # save vectorized data
2 bow_train["score"] = bow_labels_train
3 we_train["score"] = we_labels_train
4
5 bow_train.to_csv("./data/bow_train_full.csv", index=False)
6 bow_test.to_csv("./data/bow_test_full.csv", index=False)
7 we_train.to_csv("./data/we_train_full.csv", index=False)
8 we_test.to_csv("./data/we_test_full.csv", index=False)
```



```
In [ ]: 1 sample = pd.read_csv("./SampleSubmission.csv")
```

## Model: LSTM RNN

OBSERVATION: The LSTM model whose predictions are to be entered in the Kaggle competition was build based on learnings obtained from experiments conducted as part of CW2.

- An LSTM variation of RNN models was chosen because LSTM was the model that did best with train and test data among CW2 experiments.
- Additional dropout layers were added to this model since it was observed in CW2 that the LSTM model tends to overfit.
- The embedding layer outputs embeddings as generated by a pre-trained GloVe model since this was seen to produce most meaningful embeddings in CW2 experiments with LSTM and GRU RNN models.
- Since the added feature "polarity" had high correlation with ratings, this feature value shall be concatenated with the output of the LSTM layers and fed to dense layers to produce more fine tuned predictions.
- Activation function chosen for dense layers were "relu" ensuring that this model performs well even when data is not linearly separable as is most likely the case here.

```
In [13]: 1 import tensorflow as tf
2 import keras
3 from sklearn.model_selection import train_test_split
4 from keras.models import Sequential
5 from keras.layers import Embedding, LSTM, Dense, Dropout
6 from keras.initializers import Constant
7 from tensorflow.keras.optimizers import Adam
8 from keras.preprocessing.text import Tokenizer
9 import json
10 from keras.preprocessing.sequence import pad_sequences
11 from keras.models import Sequential
12 from keras.layers import Embedding, LSTM, Dense, Dropout
13 from keras.initializers import Constant
14 from tensorflow.keras.optimizers import Adam
15 import seaborn as sns
```

```
In [17]: 1 import gc
2
3 def clean_memory():
4     ''' Function that cleans unused memory. '''
5     gc.collect()
6     tf.keras.backend.clear_session()
```

```
In [234]: 1 df_train = pd.read_csv("./data/train_processed_full.csv")
2 df_test = pd.read_csv("./data/test_processed_full.csv")
```

```
In [9]: 1 def get_padded_seqmod_input(train_sentences, test_sentences, max_len):
2         # map words in vocabulary to index of corresponding token after tok
3         tokenizer = Tokenizer(oov_token="[UNK]") # replace out of vocabular
4         tokenizer.fit_on_texts(train_sentences)
5         word_index = tokenizer.word_index
6
7         # replace words in each document with its corresponding index in wo
8         train_sequences = tokenizer.texts_to_sequences(train_sentences)
9         test_sequences = tokenizer.texts_to_sequences(test_sentences)
10
11        # pad all document sequences to have max_len words.
12        train_padded = pad_sequences(train_sequences, maxlen=max_len, paddi
13        test_padded = pad_sequences(test_sequences, maxlen=max_len, padding
14
15        return train_padded, test_padded, tokenizer, word_index
```

```
In [10]: 1 def get_embedding_matrix(word_index, glove_dim):
2
3         embedding_dict = {}
4         with open(f"./data/glove.6B.{glove_dim}d.txt", "r", encoding="utf8"
5             for line in f:
6                 values = line.split()
7                 word = values[0]
8                 vector = np.asarray(values[1:], 'float32')
9                 embedding_dict[word] = vector
10        f.close()
11
12        # create GloVe embedding matrix which is a matrix of GloVe embeddin
13        # such that each embedding has index equal to the index of correspo
14        num_words = len(word_index) + 1
15        embedding_matrix = np.zeros((num_words, glove_dim)) # each word vec
16        for word, i in word_index.items():
17            if i < num_words:
18                emb_vec = embedding_dict.get(word)
19                if emb_vec is not None: embedding_matrix[i] = emb_vec
20
21        return embedding_matrix, num_words
```

```
In [237]: 1 train_sentences = df_train["review"]
2 train_labels = pd.get_dummies(df_train["score"]).values
3 test_sentences = df_test["review"]
```

```
In [238]: 1 max_len = 100
2 glove_dim = 100
3 train_padded, test_padded, tokenizer, word_index = get_padded_seqmod_in
4 embedding_matrix, num_words = get_embedding_matrix(word_index, glove_di
```

In [266]:

```
1 # make model
2 input1 = keras.Input(shape=(max_len,)) # documents
3 input2 = keras.Input(shape=(1,)) # polarity
4
5 # Layers
6 embedding = Embedding(
7     name="Embedding_Layer", input_dim = num_words, output_dim = glove_d
8     input_length = max_len, trainable = False,
9     embeddings_initializer = Constant(embedding_matrix)
10 )
11
12 lstm1 = LSTM(name="LSTM_Layer1", units=100, return_sequences=True, drop
13 lstm2 = LSTM(name="LSTM_Layer2", units=100, return_sequences=False, dro
14
15 dropout1 = Dropout(0.1)
16 dropout2 = Dropout(0.1)
17
18 dense1 = Dense(name="Dense_Layer1", units=101, activation='relu')
19 dense2 = Dense(name="Dense_Layer2", units=101, activation='relu')
20 dense3 = Dense(name="Dense_Layer3", units=5, activation='softmax')
21
22 concatenate = keras.layers.Concatenate(axis=1)
23
24 # functional API
25 x = embedding(input1)
26 x = dropout1(x)
27 x = lstm1(x)
28 x = lstm2(x)
29 x = concatenate([x, input2])
30 x = dense1(x)
31 x = dropout2(x)
32 x = dense2(x)
33 outputs = dense3(x)
34
35 # make model
36 model = keras.Model(inputs=[input1, input2], outputs=outputs, name="LST
37
38 # compile model
39 model.compile(loss="categorical_crossentropy", optimizer='adam', metric
```

In [4]:

```
1 import keras
2 model = keras.models.load_model('models/100dlstm80epochs.h5')
3 model.summary()
```

Model: "LSTM\_Functional"

Layer (type)	Output Shape	Param #	Connected to
=====			
input_1 (InputLayer)	[(None, 100)]	0	[]
Embedding_Layer (Embedding)	(None, 100, 100)	13715300	['input_1[0][0]']
dropout (Dropout)	(None, 100, 100)	0	['Embedding_Layer[0][0]']
LSTM_Layer1 (LSTM)	(None, 100, 100)	80400	['dropout[0][0]']
LSTM_Layer2 (LSTM)	(None, 100)	80400	['LSTM_Layer1[0][0]']
input_2 (InputLayer)	[(None, 1)]	0	[]
concatenate (Concatenate)	(None, 101)	0	['LSTM_Layer2[0][0]', 'input_2[0][0]']
Dense_Layer1 (Dense)	(None, 101)	10302	['concatenate[0][0]']
dropout_1 (Dropout)	(None, 101)	0	['Dense_Layer1[0][0]']
Dense_Layer2 (Dense)	(None, 101)	10302	['dropout_1[0][0]']
Dense_Layer3 (Dense)	(None, 5)	510	['Dense_Layer2[0][0]']
=====			
Total params: 13,897,214			
Trainable params: 181,914			
Non-trainable params: 13,715,300			

In [269]:

```
1 # train model (this was ran 8 times)
2 clean_memory()
3 history = model.fit(
4     x=[train_padded, df_train.polarity],
5     y=train_labels,
6     epochs=10,
7     verbose=1,
8     shuffle=True
9 )
```

Epoch 1/10

13324/13324 [=====] - 184s 14ms/step - loss: 0.61  
67 - accuracy: 0.7695

Epoch 2/10

13324/13324 [=====] - 180s 13ms/step - loss: 0.61  
62 - accuracy: 0.7693

Epoch 3/10

13324/13324 [=====] - 179s 13ms/step - loss: 0.61  
56 - accuracy: 0.7692

Epoch 4/10

13324/13324 [=====] - 179s 13ms/step - loss: 0.61  
49 - accuracy: 0.7698

Epoch 5/10

13324/13324 [=====] - 181s 14ms/step - loss: 0.61  
47 - accuracy: 0.7697

Epoch 6/10

13324/13324 [=====] - 189s 14ms/step - loss: 0.61  
47 - accuracy: 0.7701

Epoch 7/10

13324/13324 [=====] - 186s 14ms/step - loss: 0.61  
70 - accuracy: 0.7690

Epoch 8/10

13324/13324 [=====] - 183s 14ms/step - loss: 0.61  
48 - accuracy: 0.7697

Epoch 9/10

13324/13324 [=====] - 183s 14ms/step - loss: 0.61  
49 - accuracy: 0.7701

Epoch 10/10

13324/13324 [=====] - 185s 14ms/step - loss: 0.61  
55 - accuracy: 0.7692

In [270]:

```
1 # save model
2 model.save("models/100dlstm80epochs.h5")
```

In [276]:

1 model.summary()

Model: "LSTM\_Functional"

Layer (type)	Output Shape	Param #	Connected to
=====			
input_1 (InputLayer)	[(None, 100)]	0	[]
Embedding_Layer (Embedding)	(None, 100, 100)	13715300	['input_1[0][0]']
dropout (Dropout)	(None, 100, 100)	0	['Embedding_Layer[0][0]']
LSTM_Layer1 (LSTM)	(None, 100, 100)	80400	['dropout[0][0]']
LSTM_Layer2 (LSTM)	(None, 100)	80400	['LSTM_Layer1[0][0]']
input_2 (InputLayer)	[(None, 1)]	0	[]
concatenate (Concatenate)	(None, 101)	0	['LSTM_Layer2[0][0]', 'input_2[0][0]']
Dense_Layer1 (Dense)	(None, 101)	10302	['concatenate[0][0]']
dropout_1 (Dropout)	(None, 101)	0	['Dense_Layer1[0][0]']
Dense_Layer2 (Dense)	(None, 101)	10302	['dropout_1[0][0]']
Dense_Layer3 (Dense)	(None, 5)	510	['Dense_Layer2[0][0]']
=====			
Total params: 13,897,214			
Trainable params: 181,914			
Non-trainable params: 13,715,300			

```
In [21]: 1 def plot_train_val_curves(history):
2         # plot training and validation results
3         plt.figure(figsize=(18,5))
4
5         plt.subplot(121)
6         plt.title("Loss")
7         plt.plot(history.history["loss"], label="train")
8         plt.plot(history.history["val_loss"], label="validate")
9         plt.legend()
10
11        plt.subplot(122)
12        plt.title("Accuracy")
13        plt.plot(history.history["accuracy"], label="train")
14        plt.plot(history.history["val_accuracy"], label="validate")
15        plt.legend()
16
17        plt.show()
18
19        # plot_train_val_curves(history)
```

```
In [271]: 1 # make prediction
2         pred_prob = model.predict([test_padded, df_test.polarity])
3         pred = np.array([np.argmax(p)+1 for p in pred_prob])
```

```
In [273]: 1 submission_path = "kaggle_submissions/submission16.csv"
```

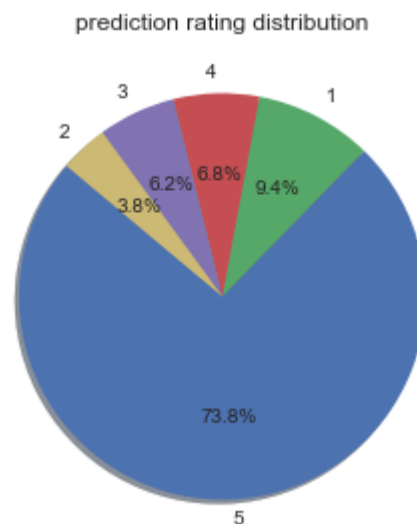
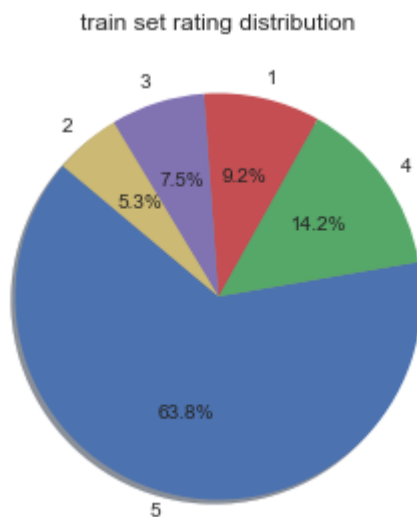
```
In [274]: 1 # save predictions
2         submission = pd.DataFrame({"Id": df_test.Id, "Score": pred})
3         submission.to_csv(submission_path, index=False) # submission to be ente
```

```
In [203]: 1 def plot_distributions(train_df, pred_df):
2         ''' Plot pie charts to view distribution comparison
3             between training set labels and predicted values '''
4
5         train_rating_data = list(train_df.score.value_counts().values)
6         train_rating_labels = list(train_df.score.value_counts().keys())
7         pred_data = list(pred_df.Score.value_counts().values)
8         pred_labels = list(pred_df.Score.value_counts().keys())
9
10        plt.figure(figsize=(10,5))
11        plt.subplot(121)
12        plt.title("train set rating distribution")
13        plt.pie(x=train_rating_data, labels=train_rating_labels, autopct='%1.1f%%')
14        plt.subplot(122)
15        plt.title("prediction rating distribution")
16        plt.pie(x=pred_data, labels=pred_labels, autopct='%1.1f%%', shadow=True)
17        plt.show()
```

In [275]:

```
1 # Load and examine prediction distribution
2 sub = pd.read_csv(submission_path)
3 display(sub.describe())
4
5 # plot pie charts to view distribution comparison between training set
6 plot_distributions(df_train, sub)
```

	Id	Score
count	142114.000000	142114.000000
mean	71057.500000	4.318519
std	41024.922415	1.302047
min	1.000000	1.000000
25%	35529.250000	4.000000
50%	71057.500000	5.000000
75%	106585.750000	5.000000
max	142114.000000	5.000000



**OBSERVATION:** Viewing the distribution of predictions on the test set vs train set labels was found to be a great way to make an educated guess as to whether the model's predictions seem sound.

## Experimentation Zone

In [5]:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 import matplotlib.pyplot as plt
5 from keras.preprocessing.text import Tokenizer
```

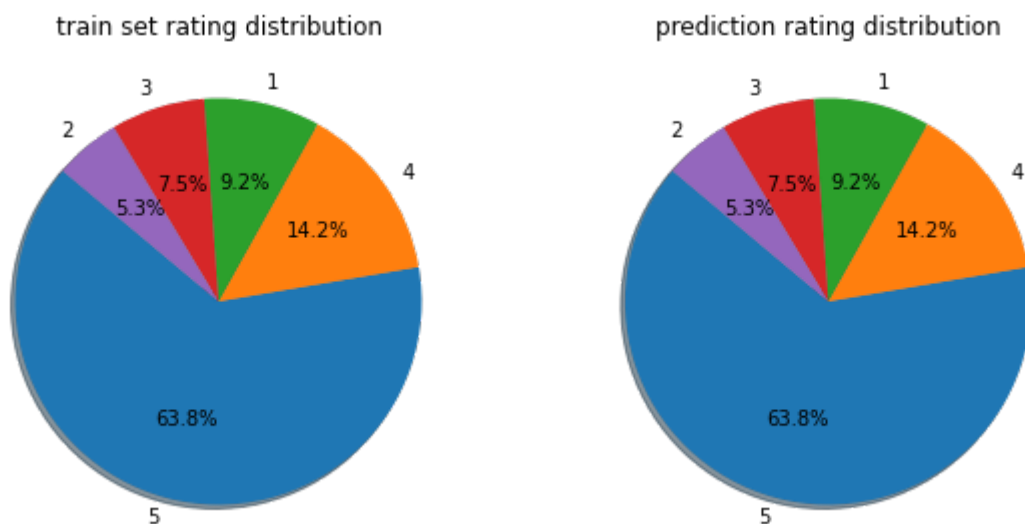
**OBSERVATION:** Since the test set truth values were not available, experiments had to be conducted by splitting the train set into stratified train, test and validation sets such that a large enough portion of the entire dataset was set aside for testing and validation so as to mimic a larger amount of unseen data.



```
In [2]: 1 data_full = pd.read_csv("../data/train_processed_full.csv")
2 x = data_full.drop(["score"], axis=1)
3 y = data_full.score
```

```
In [3]: 1 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.4)
```

```
In [11]: 1 x_score = pd.DataFrame({"score":y_train})
2 y_score = pd.DataFrame({"score":y_test})
3
4 train_rating_data = list(x_score.score.value_counts().values)
5 train_rating_labels = list(x_score.score.value_counts().keys())
6 pred_data = list(y_score.score.value_counts().values)
7 pred_labels = list(y_score.score.value_counts().keys())
8
9 plt.figure(figsize=(10,5))
10 plt.subplot(121)
11 plt.title("train set rating distribution")
12 plt.pie(x=train_rating_data, labels=train_rating_labels, autopct='%1.1f')
13 plt.subplot(122)
14 plt.title("prediction rating distribution")
15 plt.pie(x=pred_data, labels=pred_labels, autopct='%1.1f%%', shadow=True)
16 plt.show()
```



```
In [14]: 1 max_len = 100
2 glove_dim = 100
3 train_padded, test_padded, tokenizer, word_index = get_padded_seqmod_in
4 embedding_matrix, num_words = get_embedding_matrix(word_index, glove_di
5 train_labels = pd.get_dummies(y_train).values
6 test_labels = pd.get_dummies(y_test).values
```

In [19]:

```
1 # make model
2 input1 = keras.Input(shape=(max_len,)) # documents
3 input2 = keras.Input(shape=(1,)) # polarity
4
5 # Layers
6 embedding = Embedding(
7     name="Embedding_Layer", input_dim = num_words, output_dim = glove_d
8     input_length = max_len, trainable = False,
9     embeddings_initializer = Constant(embedding_matrix)
10 )
11
12 lstm1 = LSTM(name="LSTM_Layer1", units=100, return_sequences=True, drop
13 lstm2 = LSTM(name="LSTM_Layer2", units=100, return_sequences=False, dro
14
15 dropout1 = Dropout(0.1)
16 dropout2 = Dropout(0.1)
17
18 dense1 = Dense(name="Dense_Layer1", units=101, activation='relu')
19 dense2 = Dense(name="Dense_Layer2", units=101, activation='relu')
20 dense3 = Dense(name="Dense_Layer3", units=5, activation='softmax')
21
22 concatenate = keras.layers.Concatenate(axis=1)
23
24 # functional API
25 x = embedding(input1)
26 x = dropout1(x)
27 x = lstm1(x)
28 x = lstm2(x)
29 x = concatenate([x, input2])
30 x = dense1(x)
31 x = dropout2(x)
32 x = dense2(x)
33 outputs = dense3(x)
34
35 # make model
36 model = keras.Model(inputs=[input1, input2], outputs=outputs, name="LST
37
38 # compile model
39 model.compile(loss="categorical_crossentropy", optimizer='adam', metric
```

In [20]:

```
1 # train model
2 clean_memory()
3 history = model.fit(
4     x=[train_padded, x_train.polarity],
5     y=train_labels,
6     epochs=10,
7     validation_split=0.3,
8     verbose=1,
9     shuffle=True
10 )
```

Epoch 1/10

5596/5596 [=====] - 98s 17ms/step - loss: 0.8996  
- accuracy: 0.6786 - val\_loss: 0.8115 - val\_accuracy: 0.6983

Epoch 2/10

5596/5596 [=====] - 82s 15ms/step - loss: 0.8049  
- accuracy: 0.7021 - val\_loss: 0.7532 - val\_accuracy: 0.7181

Epoch 3/10

5596/5596 [=====] - 82s 15ms/step - loss: 0.7637  
- accuracy: 0.7153 - val\_loss: 0.7250 - val\_accuracy: 0.7300

Epoch 4/10

5596/5596 [=====] - 81s 14ms/step - loss: 0.7391  
- accuracy: 0.7237 - val\_loss: 0.7094 - val\_accuracy: 0.7356

Epoch 5/10

5596/5596 [=====] - 81s 14ms/step - loss: 0.7242  
- accuracy: 0.7286 - val\_loss: 0.6955 - val\_accuracy: 0.7393

Epoch 6/10

5596/5596 [=====] - 81s 14ms/step - loss: 0.7122  
- accuracy: 0.7337 - val\_loss: 0.7003 - val\_accuracy: 0.7365

Epoch 7/10

5596/5596 [=====] - 81s 15ms/step - loss: 0.7045  
- accuracy: 0.7362 - val\_loss: 0.6841 - val\_accuracy: 0.7433

Epoch 8/10

5596/5596 [=====] - 81s 15ms/step - loss: 0.6977  
- accuracy: 0.7378 - val\_loss: 0.6894 - val\_accuracy: 0.7450

Epoch 9/10

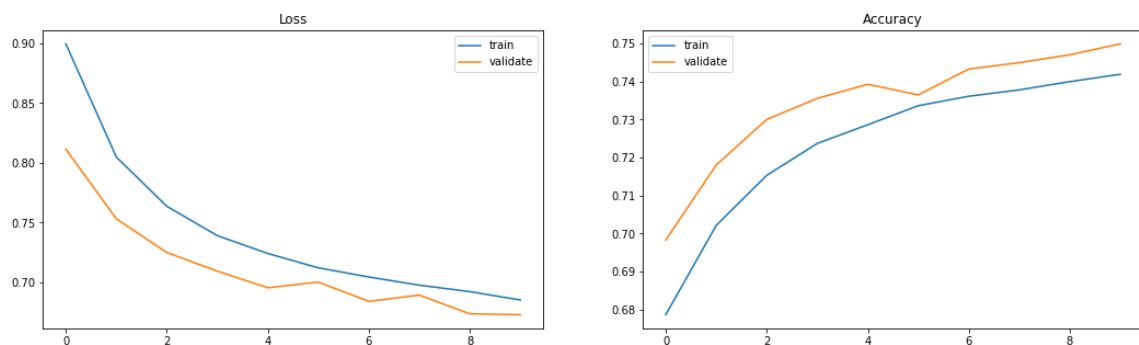
5596/5596 [=====] - 81s 14ms/step - loss: 0.6923  
- accuracy: 0.7400 - val\_loss: 0.6737 - val\_accuracy: 0.7471

Epoch 10/10

5596/5596 [=====] - 81s 14ms/step - loss: 0.6853  
- accuracy: 0.7419 - val\_loss: 0.6730 - val\_accuracy: 0.7500

In [22]:

```
1 plot_train_val_curves(history)
```



In [23]:

```
1 # test model
2 pred_prob = model.predict([test_padded, x_test.polarity])
3 pred = np.array([np.argmax(p)+1 for p in pred_prob])
4 print("\ntrain results")
5 train_loss, train_acc = model.evaluate([train_padded, x_train.polarity])
6 print("test results")
7 test_loss, test_acc = model.evaluate([test_padded, x_test.polarity], te
8 print(f"Train Accuracy = {train_acc}")
9 print(f"Test Accuracy = {test_acc}")
```

train results

7994/7994 [=====] - 46s 6ms/step - loss: 0.6348 - accuracy: 0.7603

test results

5330/5330 [=====] - 31s 6ms/step - loss: 0.6708 - accuracy: 0.7494

Train Accuracy = 0.760260820388794

Test Accuracy = 0.7493549585342407

The above experiment was the last one among various experiments that were conducted. The model obtained above was put to the test on Kaggle because even after 20 epochs, the model continued learn without convergence or overfitting. It is to be noted that while adding more layers or units may have led to better results, this could not be done here since trying it for just a few epochs made it clear that the exponential increase in computing time and resources meant that such a model was not practical with the machines at hand.

## Kaggle Results

**OBSERVATION:** This model performed as expected and produced similar accuracy as obtained here on Kaggle as well. This model was ranked 3rd by kaggle with a provincial score of 0.777 (77.7% accuracy) and a final score of 0.778 (77.8% accuracy). The model submitted on Kaggle had trained for 90 epochs. Even though it continued to learn in the 91st epoch as well, the learning was very slow and thus it was considered to have converged. The predictions of the model after 90 epochs of training was the final Kaggle submission. It is interesting to note that the score when tested on more data (Kaggle final score) was higher than when tested on lesser data (Kaggle provincial score) further indicating this model's good reliability and resistance to overfitting.