

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import tensorflow as tf
```

Text Classification

In [2]:

```
1 train = pd.read_csv("C:\\Users\\Acer\\OneDrive\\Desktop\\MKA 9 th Sept. Batch\\Deep L
```



In [3]:

```
1 test = pd.read_csv("C:\\Users\\Acer\\OneDrive\\Desktop\\MKA 9 th Sept. Batch\\Deep L
```



In [4]:

```
1 print(train.shape)
2
3 print(test.shape)
```

(77946, 28)

(42157, 4)

In [5]:

```
1 def _load_and_shuffle_data(data_path,
2                             file_name,
3                             cols,
4                             seed,
5                             separator=',',
6                             header=0):
7
8     np.random.seed(seed)
9     data_path = os.path.join(data_path, file_name)
10    data = pd.read_csv(data_path, usecols=cols, sep=separator, header=header)
11    return data.reindex(np.random.permutation(data.index))
12
13 def _split_training_and_validation_sets(texts, labels, validation_split):
14
15     num_training_samples = int((1 - validation_split) * len(texts))
16     return ((texts[:num_training_samples], labels[:num_training_samples]),
17            (texts[num_training_samples:], labels[num_training_samples:]))
18
19 def load_tweet_weather_topic_classification_dataset(data_path,
20                                                     validation_split=0.2,
21                                                     seed=123):
22
23     columns = [1] + [i for i in range(13, 28)] # 1 - text, 13-28 - topics.
24     data = _load_and_shuffle_data(data_path, 'train.csv', columns, seed)
25
26     # Get tweet text and the max confidence score for the weather types.
27     texts = list(data['tweet'])
28     weather_data = data.iloc[:, 1:]
29
30     labels = []
31     for i in range(len(texts)):
32         # Pick topic with the max confidence score.
33         labels.append(np.argmax(list(weather_data.iloc[i, :].values)))
34
35     return _split_training_and_validation_sets(
36         texts, np.array(labels), validation_split)
```

In [6]:

```
1 from sklearn.model_selection import train_test_split
```

In [9]:

```
1 import os
2 for dirname, _, filenames in os.walk('/content/crowdfower-weather-twitter.zip'):
3     for filename in filenames:
4         print(os.path.join(dirname, filename))
5
```

In [10]:

```
1 (train_data, train_labels),(test_data, test_labels)=load_tweet_weather_topic_classif
```



Tokenize and vectorize

In [11]:

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.feature_selection import SelectKBest
3 from sklearn.feature_selection import f_classif
```

In [12]:

```
1 #i am mouse - i, am, mouse, i am, am mouse
2
3 vectorizer = TfidfVectorizer(dtype='float32',
4                               ngram_range=(1,2),
5                               strip_accents='unicode',
6                               decode_error='replace',
7                               analyzer='word',
8                               min_df=2)
```

In [13]:

```
1 x_train = vectorizer.fit_transform(train_data)
2 x_train.shape
```

C:\Users\Acer\anaconda37\lib\site-packages\sklearn\feature_extraction\tex
t.py:2060: UserWarning: Only (<class 'numpy.float64'>, <class 'numpy.float
32'>, <class 'numpy.float16'>) 'dtype' should be used. float32 'dtype' wil
l be converted to np.float64.
warnings.warn(

Out[13]:

(62356, 85670)

In [14]:

```
1 x_test = vectorizer.transform(test_data)
```

In [15]:

```
1 #feature selection
2 selector = SelectKBest(f_classif, k=20000)
3 selector.fit(x_train, train_labels)
```

Out[15]:

SelectKBest(k=20000)

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [16]:

```
1 x_train = selector.transform(x_train)
2 x_test = selector.transform(x_test)
```

In [17]:

```
1 x_train = x_train.toarray()
2 x_test = x_test.toarray()
```

In [18]:

```
1 x_train.shape
```

Out[18]:

(62356, 20000)

In [19]:

```
1 x_test.shape
```

Out[19]:

(15590, 20000)

In [20]:

```
1 y_train = np.array(train_labels)
2 y_test = np.array(test_labels)
```

Neural network

ANN

In [21]:

```
1 from keras.models import Sequential
```

In [22]:

```
1 from keras.layers import Dense, Flatten, Dropout
```

Model creation - first attempt

Note: each text has 20000 features

In [23]:

```
1 weatherANN = Sequential()
```

In [24]:


```
1 weatherANN.add(Flatten())
2 weatherANN.add(Dense(units=512, activation='relu'))
3 #Add dropout regularization - a percentage of units should be set to 0. Eg: 0.25 - s
4 #be set to 0
5 weatherANN.add(Dropout(rate=0.25))
6 weatherANN.add(Dense(units=1, activation='sigmoid'))
```

In [25]:

```
1 weatherANN.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
```

In [26]:

```
1 from keras.callbacks import ModelCheckpoint, EarlyStopping, ReduceLROnPlateau
2 filename = 'bestmodel_trial1.h5'
3 checkpoint = ModelCheckpoint(filename, monitor='val_accuracy', verbose=1, save_best_
4 #if monitored value doesnt improve for #patience epochs, then stop training
5 es = EarlyStopping(monitor='val_loss', patience=10)
6 #if monitored value doesnt improve for #patience epochs, LR(new) = LR(old)*factor
7 rd = ReduceLROnPlateau(monitor='val_loss', factor=0.1, patience=5)
```



In [27]:

```
1 history = weatherANN.fit(x_train, y_train, epochs=5, validation_split=0.2, callbacks
```

```
Epoch 1/5
1559/1559 [=====] - ETA: 0s - loss: -7202.0425 -
accuracy: 0.0916
Epoch 1: val_accuracy improved from -inf to 0.09469, saving model to bestm
odel_trial1.h5
1559/1559 [=====] - 181s 115ms/step - loss: -720
2.0425 - accuracy: 0.0916 - val_loss: -21491.5469 - val_accuracy: 0.0947 -
lr: 0.0010
Epoch 2/5
1559/1559 [=====] - ETA: 0s - loss: -47311.9570 -
accuracy: 0.0916
Epoch 2: val_accuracy did not improve from 0.09469
1559/1559 [=====] - 180s 115ms/step - loss: -4731
1.9570 - accuracy: 0.0916 - val_loss: -78290.1250 - val_accuracy: 0.0947 -
lr: 0.0010
Epoch 3/5
1559/1559 [=====] - ETA: 0s - loss: -117731.4453
- accuracy: 0.0916
Epoch 3: val_accuracy did not improve from 0.09469
1559/1559 [=====] - 173s 111ms/step - loss: -1177
31.4453 - accuracy: 0.0916 - val_loss: -161819.2031 - val_accuracy: 0.0947
- lr: 0.0010
Epoch 4/5
1559/1559 [=====] - ETA: 0s - loss: -213001.5000
- accuracy: 0.0916
Epoch 4: val_accuracy did not improve from 0.09469
1559/1559 [=====] - 185s 119ms/step - loss: -2130
01.5000 - accuracy: 0.0916 - val_loss: -269040.9688 - val_accuracy: 0.0947
- lr: 0.0010
Epoch 5/5
1559/1559 [=====] - ETA: 0s - loss: -331133.1875
- accuracy: 0.0916
Epoch 5: val_accuracy did not improve from 0.09469
1559/1559 [=====] - 185s 119ms/step - loss: -3311
33.1875 - accuracy: 0.0916 - val_loss: -398364.3750 - val_accuracy: 0.0947
- lr: 0.0010
```

Save the best model (W, b) by comparing accuracy/loss on validation set

In [29]:

```
1 from keras import models
2 #human accuracy - 1.0000
3 #training accuracy: 0.0916
4 #val_accuracy did not improve from 0.9469
5 #Avoidable bias = 1-0.091 = 0.1
6 #Variance = 1 - 0.94 = 0.16
7 finalmodel = models.load_model('bestmodel_trial1.h5')
```

loading best model based on val_accuracy

evaluate on test set

In [30]:

```
1
2 finalmodel.evaluate(x_test, y_test)
```

488/488 [=====] - 5s 10ms/step - loss: -21934.714
8 - accuracy: 0.0921

Out[30]:

[-21934.71484375, 0.09211032837629318]

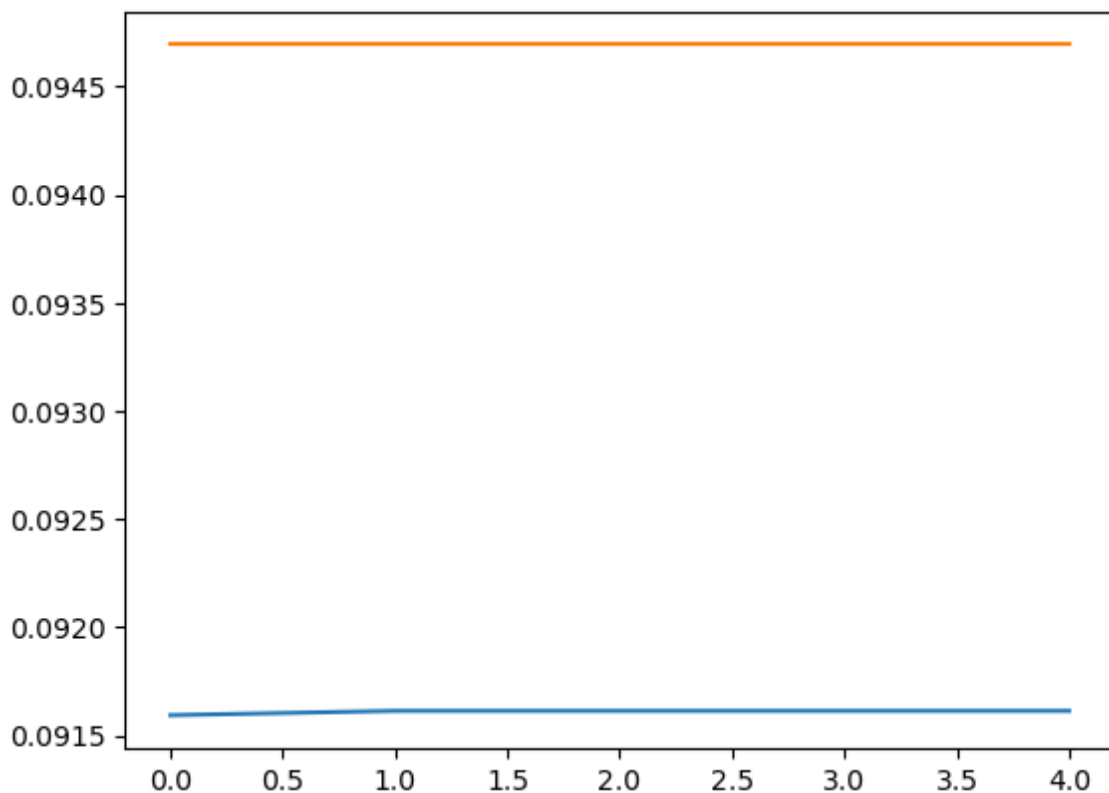
Plots

In [31]:

```
1 import matplotlib.pyplot as plt
2
```

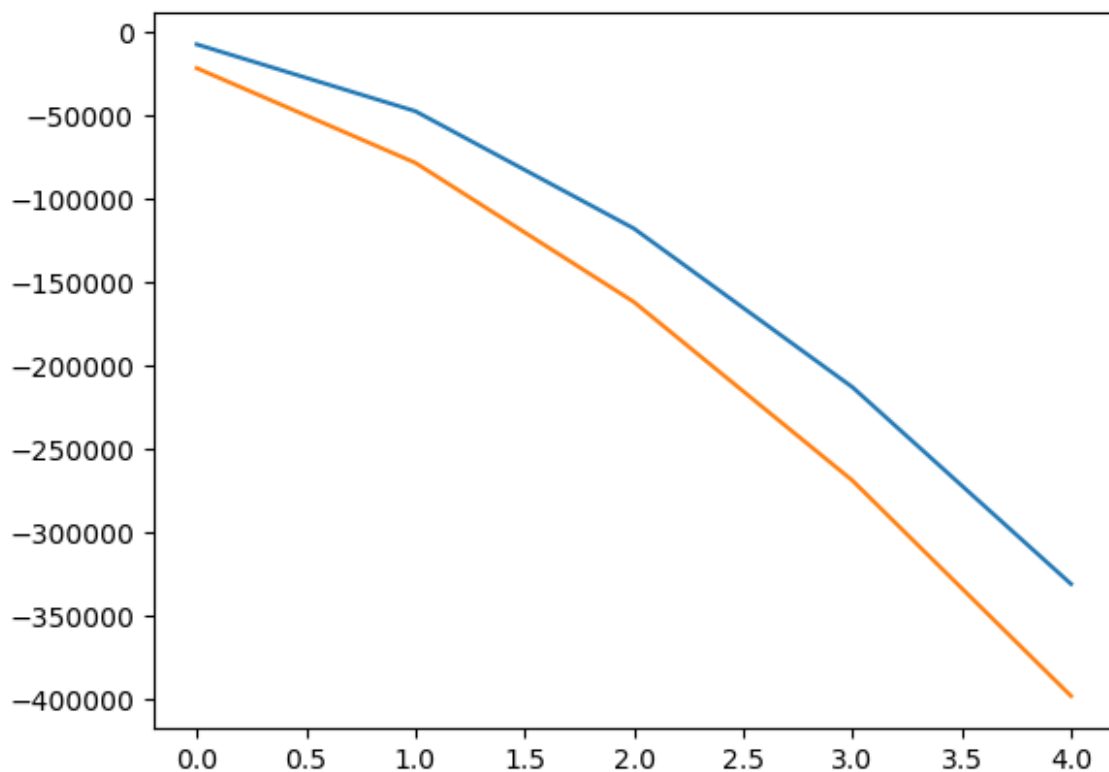
In [38]:

```
1 #plot accuracy with epochs - should increase with epochs
2
3 #training accuracy
4 plt.plot(history.history['accuracy'])
5 #validation accuracy
6 plt.plot(history.history['val_accuracy'])
7 plt.show()
```



In [33]:

```
1 #plot loss
2 import matplotlib.pyplot as plt
3 #training loss
4 plt.plot(history.history['loss'])
5 #validation loss
6 plt.plot(history.history['val_loss'])
7 plt.show()
```



Deployment

In [42]:

```
1 comment = ['At the park with erica, enjoying the weather. :)']
2
3 X_deployment = vectorizer.transform(comment)
4
5 X_deployment_best = selector.transform(X_deployment)
6
7 X_deployment_best_array = X_deployment_best.toarray()
```

In [44]:

```
1 finalmodel.predict(X_deployment_best_array)
```

1/1 [=====] - 0s 128ms/step

Out[44]:

```
array([[1.]], dtype=float32)
```


In []:

1	
---	--