## CE066-Jaydeep Mahajan- ML- LAB 7

```python
1 # import libraries
2 import nltk
3 import re
4 import string
5 import numpy as np
6 import tensorflow as tf
7 from sklearn.preprocessing import StandardScaler
8 from nltk.corpus import twitter_samples
9 from nltk.corpus import stopwords
10 from nltk.stem import PorterStemmer
11 from nltk.tokenize import TweetTokenizer
12 from __future__ import absolute_import, division, print_function
```

```python
1 # download twitter_samples and stopwords dataset
2 nltk.download('twitter_samples')
3 nltk.download('stopwords')
```

```
[nltk_data] Downloading package twitter_samples to /root/nltk_data...
[nltk_data]   Unzipping corpora/twitter_samples.zip.
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
True
```

```python
1 # function for preprocessing task and set train data
2 def process_tweet(tweet):
3   stemmer = PorterStemmer()
4   stopwords_english = stopwords.words('english')
5   tweet = re.sub(r'\$\w*', '', tweet)
6   tweet = re.sub(r'^RT[\s]+', '', tweet)
7   tweet = re.sub(r'https?:\/\/.*[\r\n]*', '', tweet)
8   tweet = re.sub(r'#', '', tweet)
9
10   tokenizer = TweetTokenizer(preserve_case=False, strip_handles=True, reduce_len=True)
```

```
10    tokenizer = TweetTokenizer(preserve_case=False,strip_handles=True,reduce_len=True)
11    tweet_tokens = tokenizer.tokenize(tweet)
12    tweets_clean = []
13    for word in tweet_tokens:
14      if (word not in stopwords_english and word not in string.punctuation):
15        stem_word = stemmer.stem(word)
16        tweets_clean.append(stem_word)
17    return tweets_clean
18
19 def build_freqs(tweets, ys):
20    yslist = np.squeeze(ys).tolist()
21    freqs = {}
22    for y, tweet in zip(yslist, tweets):
23      for word in process_tweet(tweet):
24        pair = (word, y)
25        if pair in freqs:
26          freqs[pair]+=1
27        else:
28          freqs[pair]=1
29    return freqs
30
31 def extract_features(tweet, freqs):
32    word_list = process_tweet(tweet)
33    x = np.zeros((1,2), dtype=np.float32)
34    for word in word_list:
35      if (word,1) in freqs:
36        x[0,0]+=freqs[word,1]
37      if (word,0) in freqs:
38        x[0,1]+=freqs[word,0]
39    assert(x.shape==(1,2))
40    return x
```

```
1 # sample of preprocessed tweet
2 processed_tweet = process_tweet("@Amazon is always #good company")
3 print(processed_tweet)
```

```
['alway', 'good', 'compani']
```

```
1 # Get the positive and negative tweets and create dataset
2 all_positive_tweets = twitter_samples.strings('positive_tweets.json')
3 all_negative_tweets = twitter_samples.strings('negative_tweets.json')
4
5 test_pos = all_positive_tweets[3000:]
6 train_pos = all_positive_tweets[:3000]
7 test_neg = all_negative_tweets[3000:]
8 train_neg = all_negative_tweets[:3000]
9
10 train_x = train_pos + train_neg
11 test_x = test_pos + test_neg
12 train_y = np.append(np.ones((len(train_pos), 1),np.int64), np.zeros((len(train_neg), 1),np.int64), axis=0)
13 test_y = np.append(np.ones((len(test_pos), 1),np.int64), np.zeros((len(test_neg), 1),np.int64), axis=0)
```

```
1 # Get word frequencies for positive and negative sentiment
2 freqs = build_freqs(train_x,train_y)
3 print("type(freqs) = " + str(type(freqs)))
4 print("len(freqs) = " + str(len(freqs.keys())))
```

```
type(freqs) = <class 'dict'>
len(freqs) = 9326
```

```
1 # Define parameters
2 num_classes = 2 # 1 or 0
3 num_features = 2 # positive and negative freqs
4 learning_rate =  0.001
5 training_steps = 1000
6 batch_size = 256
7 display_step = 50
```

```
1 # Get the frequencies of positive and negative word for 2 samples
2 sample_1 = extract_features(train_x[0], freqs)
3 print("sample 1 : ", sample_1)
4 sample_2 = extract_features(train_x[4010], freqs)
5 print("sample 2 : ", sample_2)
```

```
⤷   sample 1 :  [[2276.     47.]]
    sample 2 :  [[   45. 2822.]]
```

```python
1 # Format X_train and X_test
2 X_train = np.zeros((len(train_x),2),dtype=np.float32)
3 X_test = np.zeros((len(test_x),2),dtype=np.float32)
4 for i in range(len(train_x)):
5   X_train[i,:] = extract_features(train_x[i],freqs)
6 for i in range(len(test_x)):
7   X_test[i,:] = extract_features(test_x[i],freqs)
8 sc = StandardScaler()
9 X_train = sc.fit_transform(X_train)
10 X_test = sc.transform(X_test)
11 Y_train = train_y
12 Y_test = test_y
13 print("Train sample : ",X_train[0],Y_train[0])
14 print("Test sample : ",X_test[1500],Y_test[1500])
```

```
⤷   Train sample :  [ 1.0975696  -0.91117305] [1]
    Test sample :  [ 1.2294407  -0.74473023] [1]
```

```python
1 # Intialize weight and bias
2 W = tf.Variable(tf.ones([num_features, num_classes]), name="weight")
3 b = tf.Variable(tf.zeros([num_classes]), name="bias")
4
5 # Use tf.data API to shuffle and batch data.
6 train_data=tf.data.Dataset.from_tensor_slices((X_train,Y_train))
7 train_data=train_data.repeat().shuffle(5000).batch(batch_size).prefetch(1)
```

```python
1 # Main function for perform logistic regression
2 def logistic_regression(x,W,b):
3   return tf.nn.sigmoid(tf.matmul(x,W) + b)
4
5 def cross_entropy(y_pred,y_true):
6   y_true = tf.one_hot(y_true, depth=num_classes)
```

```
 7   y_pred = tf.clip_by_value(y_pred,1e-9,1.)
 8   return tf.reduce_mean(-tf.reduce_sum(y_true*tf.math.log(y_pred)))
 9
10 def accuracy(y_pred, y_true):
11   correct_prediction = tf.equal(tf.argmax(y_pred, 1), tf.cast(y_true, tf.int64))
12   return tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
13
14 def run_optimization(x,y):
15   with tf.GradientTape() as g:
16     pred = logistic_regression(x,W,b)
17     loss = cross_entropy(pred,y)
18   gradients = g.gradient(loss,[W,b])
19   optimizer = tf.optimizers.SGD(learning_rate)
20   optimizer.apply_gradients(zip(gradients, [W,b]))
```

```
1 # Train model for given number of step
2 for step, (batch_x, batch_y) in enumerate(train_data.take(training_steps), 1):
3   run_optimization(batch_x, batch_y)
4   if step % display_step == 0:
5       pred = logistic_regression(batch_x,W,b)
6       loss = cross_entropy(pred, batch_y)
7       acc = accuracy(pred, batch_y)
8       print("step: %i, loss: %f, accuracy: %f" % (step, loss, acc))
```

```
step: 50, loss: 0.028325, accuracy: 0.432709
step: 100, loss: 0.047213, accuracy: 0.452637
step: 150, loss: 0.129363, accuracy: 0.494263
step: 200, loss: 0.015692, accuracy: 0.580933
step: 250, loss: 0.032502, accuracy: 0.591797
step: 300, loss: 0.017747, accuracy: 0.525543
```

```
1 #Final weight
2 print("Weight : ")
3 print(W)
4 #Final bias
5 print("Bias : ")
6 print(b)
```

```
Weight :
<tf.Variable 'weight:0' shape=(2, 2) dtype=float32, numpy=
array([[-0.67513263, -1.3009712 ],
       [-0.8364491 , -1.6109239 ]], dtype=float32)>
Bias :
<tf.Variable 'bias:0' shape=(2,) dtype=float32, numpy=array([14.573064, 21.138458], dtype=float32)>
```

```
1 pred = logistic_regression(X_test,W,b)
2 print("Test accuracy: %f" % accuracy(pred,Y_test))
```

```
Test accuracy: 0.500000
```

```
1
```