

# NLP - Assignment 2

## Sentiment analysis with Deep Learning

Dhivya Chandrasekaran  
*Department of Computer Science)*  
*Student ID : 1111017*  
dchandra@lakeheadu.ca

**Abstract**—This article describes the various steps taken to implement Sentiment Analysis using Rotten Tomatoes movie review dataset using 1-Dimensional Convolutional Neural Network. Preprocessing of text data is carried out by removing punctuations, special characters and stop words. Lemmatization is done with WordNet Lemmatizer and the vectorisation is carried out using pretrained Glove vectors. A 1D Conv model is built using Keras and classification metrics like accuracy, recall, precision and f1score are calculated. This report records the various steps along the process and the achieved results in detail.

### I. INTRODUCTION

Sentiment Analysis of text data is an important Natural Language Processing Task. Various business ventures adopt sentiment analysis as a good marketing tool to identify the acceptance rate of the product in the market and to fine tune the efforts to attract a specific group of people. In this assignment, the given task is to perform Sentiment Analysis on the Rotten Tomatoes Movie Review Dataset using a 1 dimensional Convolutional Neural Network(Conv1D). The given text data is converted to real value vectors called word-embeddings and the vectors are used to predict the class to which the given sentence/phrase is aligned to. The final prediction is performed using softmax activation. The article provides a brief overview of Text data processing in Section II and the experiment analysis of the model is carried out in Section III and section IV records the results obtained.

### II. BACKGROUND

The major challenge in processing Natural Language is the process by which we convert the given text data to machine readable format. Sentiment Analysis of text data maybe considered as a classification problem where the given text is assigned

to a set of class variables. Earlier probabilistic models were proposed to handle the problem of text classification. However due to the versatility of natural language, there was never a dataset big enough to encompass all the possible combinations of the text data. Machine learning approaches like Naive-Bayes, Random-Forrests etc., were proposed for text classification. The given sentences were converted to vectors using techniques like Bag-of-Words(BoW), Term frequency-Inverse Document Frequency(Tf-Idf) that used frequency of a word over a given corpus as the vector values and then cosine similarity between these vectors were calculated and the similar words were grouped to a defined class based on their distribution in the training dataset. However these techniques did not capture the semantic properties of the text data. Hence word-embedding using neural networks like Word2Vec and using word-co occurrence matrix like GloVe vectors. These vectors were formed over a high dimensional space and they captured the semantic properties for example the vectors of the words that are similar are closely aligned.

#### A. *word2vec*

*word2vec*: [1] is developed from Google News Dataset containing approximately 3 million vector representations of words and phrases. *word2vec* is a neural network model used to produce distributed vector representation of words based on an underlying corpus. There are two different models of *word2vec* proposed: The Combined Bag of Words(CBOW) and the Skip-gram model. The architecture of the network is rather simple and contains an input layer, one hidden layer and an output layer. The network is fed with a large text

corpus as the input, and the output of the model is a vector representation of words. The CBOW model predicts the current word using the previous words while the Skip-gram model predicts the neighboring context words given a target word.

### B. GloVe:

*GloVe* [2] was developed by Stanford University, relies on a global word co-occurrence matrix formed based on the underlying corpus. It estimates similarity based on the principle that words similar to each other occur together. The co-occurrence matrix is populated with occurrence values by doing a single pass over the underlying large corpora. *GloVe* model was trained using five different corpora mostly Wikipedia dumps.

## III. EXPERIMENTAL ANALYSIS AND COMPARISON

The entire code for the assignment is provided in Appendix A

### A. Dataset:

For this assignment the *Rotten Tomatoes Movie Review Dataset* is used. This dataset for originally by Pang and Lee [3]. The training dataset contains a total of 8544 sentences that are parsed using Stanford Parser to form 156060 phrases. Each phrase has a phrase id and a sentence id assigned as a feature. The phrases repeated appear only once in the dataset. The dataset contains 5 different classes as in the Table I

Sentiment Id	Sentiment
0	Negative
1	Somewhat Negative
2	Neutral
3	Somewhat Positive
4	Positive

TABLE I  
VARIOUS CLASSES IN THE ROTTEN TOMATOES MOVIE REVIEW DATASET

The dataset has significant class imbalance towards class 2 with nearly 50% of the phrases attributed to class 2.

### B. Model building:

The code snippet for the proposed model is show in A. The proposed model has three convolution blocks. Each convolution block has one convolution layer, one batch normalisation layer and one dropout layer. The first convolution has 512 filters, the second convolution has 256 filters and the third convolution layer has 64 filters with a kernel size of 3. Five fully connected layers are attached to the end with 128,100,50,10 and 5 output dimensions respectively. The following table shows the combination of hyper-parameters the model was tuned with.

Hyper-parameters	Option 1	Option 2	Option 3	Best Option
Vocabulary size	5000	10,000	13,000	<b>10,000</b>
Number of Epochs	10	25	50	<b>50</b>
Filter size	3	5	7	<b>3</b>
Maximum sentence length	9	15	20	<b>15</b>
Embedding dimension	50	200	300	<b>300</b>
Activation function	Relu	Sigmoid	Tanh	<b>Relu</b>
Optimizer	Adam	RMSProp	NAdam	<b>Adam</b>

TABLE II  
HYPER-PARAMETER TUNING.

The dataset as mentioned earlier has class imbalance, however oversampling or under-sampling did not result in good results hence the train\_test\_split is done stratified thus maintaining the same ration of test samples which resulted in improved performance. And early stopping is used to avoid overfitting. The loss function used is categorical\_cross\_entropy

## IV. RESULTS

The data is read from the internet using the code shown in A . The optimal solution for the model is obtained using Adam Optimizer, ReLU activation function, learning rate 0.001, repeated for 50 epochs with early stopping. The training accuracy is 0.7073 . The testing accuracy is 0.6650, precision is 0.5963, Recall is 0.5373, F1 Score is 0.5582 . The loss curve and the accuracy curve are shown in Fig 1 and Fig 2 respectively.

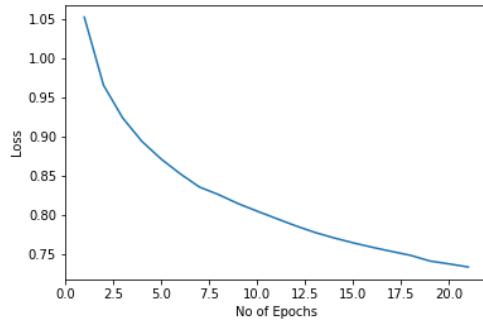


Fig. 1. Loss Curve

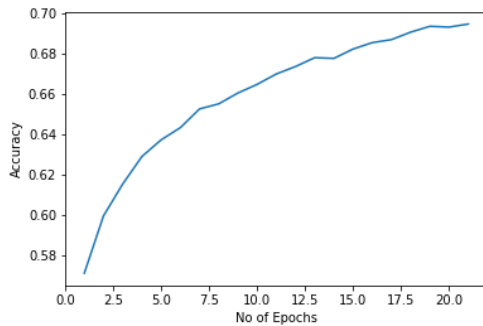


Fig. 2. Accuracy Curve

## REFERENCES

- [1] Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean. "Efficient estimation of word representations in vector space." arXiv preprint arXiv:1301.3781 (2013).
- [2] Pennington, Jeffrey, Richard Socher, and Christopher D. Manning. "Glove: Global vectors for word representation." In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP), pp. 1532-1543. 2014.
- [3] Pang, Bo, and Lillian Lee. "Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales." In Proceedings of the 43rd annual meeting on association for computational linguistics, pp. 115-124. Association for Computational Linguistics, 2005.

## APPENDIX

### A. Reading data from source to a pandas dataframe

```
1 import pandas as pd
2 URL_Tr = 'https://raw.githubusercontent.com/
  cacoderquan/Sentiment-Analysis-on-the-
  Rotten-Tomatoes-movie-review-dataset/
  master/train.tsv'
3 dataset = pd.read_csv(URL_Tr, sep='\t')
```

### B. Importing necessary libraries for data preprocessing

```
1 import re
2 import nltk
3 nltk.download('stopwords')
4 from nltk.corpus import stopwords
5 nltk.download('wordnet')
6 from nltk.stem import WordNetLemmatizer,
  PorterStemmer, LancasterStemmer
7 porter = PorterStemmer()
8 lancaster=LancasterStemmer()
9 wordnet_lemmatizer = WordNetLemmatizer()
10 stop_words = stopwords.words("english")
```

### C. Function which performs data preprocessing

```
1 def preprocessing(review,remove_stopwords =
  True, stem_lem = 'word_net'):
2     temp = review
3     temp.loc[:, "Phrase"] = temp.Phrase.apply(
4         lambda x : str.lower(x))
5     temp.loc[:, "Phrase"] = temp.Phrase.apply(
6         lambda x : " ".join(re.findall('[a-zA-Z]+'
7         ,x)))
8     if remove_stopwords == True:
9         temp.loc[:, "Phrase"] = temp.Phrase.apply(
10             lambda x : ' '.join(word for word in x.
11             split() if word not in stop_words))
12     if stem_lem == 'word_net':
13         temp.loc[:, "Phrase"] = temp.Phrase.apply(
14             lambda x: " ".join([wordnet_lemmatizer.
15             lemmatize(word) for word in x.split()])))
16     elif stem_lem == 'lancaster':
17         temp.loc[:, "Phrase"] = temp.Phrase.apply(
18             lambda x: " ".join([lancaster.stem(word)
19             for word in x.split()])))
20     elif stem_lem == 'porter':
21         temp.loc[:, "Phrase"] = temp.Phrase.apply(
22             lambda x: " ".join([porter.stem(word) for
23             word in x.split()])))
24     return temp
```

### D. Store preprocessed dataset in a pandas DataFrame and remove null values.

```
1 preprocessed_review = preprocessing(review,
2     remove_stopwords=True, stem_lem='word_net')
3 preprocessed_review.head()
4
5 null_index = preprocessed_review[
6     preprocessed_review['Phrase']==''].index
7 preprocessed_review.drop(null_index,inplace=
8     True)
```

### E. Get values of columns as list and perform train\_test\_split at 70:30 ration and random state of 2003.

```
1 phrases = preprocessed_review['Phrase'].values
2 y = preprocessed_review['Sentiment'].values
3
4 from sklearn.model_selection import
5     train_test_split
```

```

5 phrases_train, phrases_test, y_train, y_test =
  train_test_split(phrases, y, test_size
    =0.30, random_state=2003, stratify = y)

```

### F. Code for trying to under-sampling and over-sampling

```

1 phrases_train_df = pd.DataFrame()
2 class_4 = phrases_buffer['Sentiment'] == 4
3 class_4_dataframe = phrases_buffer[class_4]
4 phrases_train_df = phrases_train_df.append([
  class_4_dataframe]*1, ignore_index=True)
5 class_0 = phrases_buffer['Sentiment'] == 0
6 class_0_dataframe = phrases_buffer[class_0]
7 phrases_train_df = phrases_train_df.append([
  class_0_dataframe]*1, ignore_index=True)
8 class_1 = phrases_buffer['Sentiment'] == 1
9 class_1_dataframe = phrases_buffer[class_1]
10 phrases_train_df = phrases_train_df.append([
  class_1_dataframe]*1, ignore_index=True)
11 class_3 = phrases_buffer['Sentiment'] == 3
12 class_3_dataframe = phrases_buffer[class_3]
13 phrases_train_df = phrases_train_df.append([
  class_3_dataframe]*1, ignore_index=True)
14 class_2 = phrases_buffer['Sentiment'] == 2
15 class_2_dataframe = phrases_buffer[class_2]
16 phrases_train_df = phrases_train_df.append([
  class_2_dataframe]*1, ignore_index=True)

```

### G. Import keras tokenizer and tokenize with X\_train data.

```

1 from keras.models import Sequential
2 from keras import layers
3 from keras.preprocessing.text import Tokenizer
4
5 tokenizer = Tokenizer(num_words=13000)
6 tokenizer.fit_on_texts(X_train)
7
8 X_train = tokenizer.texts_to_sequences(X_train)
9
10 X_test = tokenizer.texts_to_sequences(
  phrases_test)

```

### H. Code for padding and converting class variable to one-hot encoded values

```

1 from keras.preprocessing.sequence import
  pad_sequences
2 sentence_length = 15
3 X_train = pad_sequences(X_train, padding='post',
  maxlen=sentence_length)
4 X_test = pad_sequences(X_test, padding='post',
  maxlen=sentence_length)
5
6 X_train.shape
7
8 from keras.utils import to_categorical
9 y_train = to_categorical(y_train, 5)
10 y_test = to_categorical(y_test, 5)
11
12 import numpy as np

```

### I. Function for creating the embedding matrix

```

1 def create_embedding_matrix(filepath,
  word_index, embedding_dim):
2   vocab_size = len(word_index) + 1 # Adding
  again 1 because of reserved 0 index
3   embedding_matrix = np.zeros((vocab_size,
  embedding_dim))
4
5   with open(filepath) as f:
6     for line in f:
7       word, *vector = line.split()
8       if word in word_index:
9         idx = word_index[word]
10        embedding_matrix[idx] = np.
  array(
11          vector, dtype=np.float32)
12
13   return embedding_matrix
14
15 embedding_dim = 300
16 embedding_matrix = create_embedding_matrix('//
  content/glove.6B.300d.txt', tokenizer.
  word_index, embedding_dim)

```

### J. Define vocabulary size

```

1 vocab_size = len(tokenizer.word_index) + 1
2 vocab_size

```

### K. Define the classification model

```

1
2 model = Sequential()
3 model.add(layers.Embedding(vocab_size,
  embedding_dim, weights=[embedding_matrix],
  input_length=sentence_length, trainable=
  False))
4 model.add(layers.Conv1D(512, 3, activation='
  relu'))
5 model.add(layers.BatchNormalization())
6 model.add(layers.Dropout(0.5))
7 model.add(layers.Conv1D(256, 3, activation='
  relu'))
8 model.add(layers.BatchNormalization())
9 model.add(layers.Dropout(0.5))
10 model.add(layers.Conv1D(128, 3, activation='
  relu'))
11 model.add(layers.BatchNormalization())
12 model.add(layers.Dropout(0.5))
13 model.add(layers.MaxPool1D())
14 model.add(layers.Conv1D(64, 3, activation='
  relu'))
15 model.add(layers.GlobalMaxPooling1D())
16 model.add(layers.Dense(128, activation = 'relu'
  ))
17 model.add(layers.Dense(100, activation='relu'
  ))
18 model.add(layers.Dense(50, activation='relu'))
19 model.add(layers.Dense(10, activation='relu'))
20 model.add(layers.Dense(5, activation='softmax'
  ))

```

#### *L. Build the model*

```
1 model.compile(optimizer='adam',loss='
    categorical_crossentropy',metrics=['
    accuracy'])
2 model.summary()
```

```
12 print("The Recall is : %.2f" % recall)
13 fl_score = fl_score(rounded_labels, y_pred ,
    average="macro")
14 print("The F1_score is : %.2f" % fl_score)
```

#### *M. Import and Initialize Early Stopping from keras*

```
1 from keras.callbacks import EarlyStopping
2 es = EarlyStopping(monitor='val_loss', mode='
    min', verbose=1, patience=10)
```

#### *N. Train the model*

```
1 history = model.fit(X_train, y_train,epochs
    =50,verbose=True,validation_data=(X_test,
    y_test),batch_size=25,callbacks=[es])
```

#### *O. Save the model*

```
1 model.save('1111017_1dconv_reg')
```

#### *P. Load the saved model*

```
1 from keras.models import load_model
2 model = load_model('1111017_1dconv_reg')
```

#### *Q. Plot the results*

```
1 training_dict = history.history
2 accuracy_val = training_dict['acc']
3 loss_val = training_dict['loss']
4 epochs = range(1, len(loss_val) + 1)
5 plt.figure()
6 plt.xlabel("No of Epochs")
7 plt.ylabel("Loss")
8 plt.plot(epochs,loss_val)
9 plt.figure()
10 plt.xlabel("No of Epochs")
11 plt.ylabel("Accuracy")
12 plt.plot(epochs,accuracy_val)
```

#### *R. Test the model*

#### *S. Calculate the accuracy, recall, precision and f1\_score for testing sample.*

```
1 from sklearn.metrics import fl_score,
    precision_score, recall_score,
    confusion_matrix
2 import numpy as np
3
4 from collections import Counter
5 rounded_labels=np.argmax(y_test, axis=1)
6 y_pred1 = model.predict(X_test)
7 y_pred = np.argmax(y_pred1, axis=1)
8 print(y_pred)
9 precision = precision_score(rounded_labels,
    y_pred , average="macro")
10 print("The Precision is : %.2f" % precision)
11 recall = recall_score(rounded_labels, y_pred ,
    average="macro")
```