

Multi-class Sentiment Analysis using Deep Learning

Prerak Jayeshkumar Patel
Department of Computer Science
Lakehead University
Thunder Bay, Canada
Student Id: 1117675

Abstract—The dataset of movie reviews from Rotten Tomatoes is a collection of phrases and their corresponding sentiments in five different categories. We have used the Convolutional Neural Network for task of this sentiment classification. But for improving the efficiency of the model we have made use of preprocessing methods.

I. INTRODUCTION

In the problem of sentiment analysis from text data, preprocessing of the text plays an important role in getting most efficient performance. There are multiple methods available for the preprocessing of the text data. Methods include lemmatization, tokenization, removal of stop words and punctuations. All of these methods have a sole motive of making the text more informative. Another method like TF-IDF (Term Frequency – Inverse Document Frequency) is used to decrease the sparsity of the data. It is used to remove the unwanted words which are not important for the particular domain. For example, in this dataset we have movie review's so the word "movie" has more frequency and has less importance. Whereas, the word like "good" or "bad" may have less frequency but more importance. This kind of understanding is helpful for the machine to learn and get more meaning from the data given as input. The model used is Convolutional Neural Network (CNN) for text classification and sentiment analysis.

II. LITERATURE REVIEW

The dataset of movie reviews from Rotten Tomatoes contains totally 156000 phrases in total. There are 5 sentiments labels used for all of them. They are numbered from 0 to 4 with negative, somewhat negative, neutral, somewhat positive and positive as their sentiments. The preprocessing methods used include, TF-IDF (Term Frequency – Inverse Document Frequency). This method is used to differentiate the words which appears frequently and has great importance with the words which appears frequently in all documents and hence has less importance in the particular domain. So, it is a simple mining technique to categorize the document and have a frequency count for each word. For getting an accurate computation of TF-IDF there has been some preprocessing done on the input text data like removal of most frequent words and rare words. Also, bi-grams made more sense in this dataset. The model used for text classification is Convolutional Neural Network (CNN) [1].

The preprocessing part of text plays a significant role in any Natural Language Processing system. Some of the most

well-known text preprocessing methods include tokenization, lowercasing, multiword grouping and Lemmatization. There is a technique named Word2Vec to generate the word embedding for grouping a particular set of words. This technique is nothing but a pre-trained neural network with two layers which generates the linguistic context among words. Tokenization is a simple division of word units with a white space. The motive of tokenization is to separate every unit of the input text. Lowercasing is just converting the input text to lowercase. Lemmatization is a technique of converting the words into their respective lemma. This helps in reducing the sparsity of the input text but with a cost of sometimes ignoring the syntactic nuances. Having a great choice of selection among the preprocessing methods has always been a matter of conflict. Experimental results show that sometimes only a simple tokenization can generate great performance accuracy but, in any domain specific dataset like medical dataset used in the experiment, selecting only tokenization generated very poor results. So, sometime we need to combine multiple methods to get a good accuracy. The experimental results are generated for the best performing model for text classification. There has been, on an average, high variance of $\pm 2.4\%$ -symbol on the results depending on the selection of preprocessing methods. [2].

III. PROPOSED MODEL

The proposed model has a single-layered one-dimensional convolutional neural network. The preprocessing task is shown below.

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3
4 data = pd.read_csv("https://raw.githubusercontent.com/cacoderquan/Sentiment-Analysis-on-the-Rotten-
5                   -Tomatoes-movie-review-dataset/master/train.tsv"
6                   , sep = '\t')
7
8 # Import the numpy library to work with and
9   # manipulate the data
10 import numpy as np
11
12 # Check the head of the dataframe
13 data.head(20)
14
15 # check the shape of df
16 data.shape
17
18 # Get number of unique sentences
19 numSentences = data['SentenceId'].max()
```

```

18 # extract full sentences only from the dataset
19 fullSentences = []
20 curSentence = 0
21 for i in range(data.shape[0]):
22     if data['SentenceId'][i] > curSentence:
23         fullSentences.append((data['Phrase'][i], data['
24             Sentiment'][i]))
25         curSentence = curSentence + 1
26 len(fullSentences)
27
28 # put data into a df
29 fullSentDf = pd.DataFrame(fullSentences,
30                             columns=['Phrase', '
31                                 Sentiment'])
32 # Check class imbalance in tokenized sentences
33 data['Sentiment'].value_counts()
34
35 # Check class imbalance in full sentences
36 fullSentDf['Sentiment'].value_counts()
37
38 import nltk
39 import random
40 nltk.download('punkt')
41 nltk.download('stopwords')
42 nltk.download('wordnet')
43 from nltk.tokenize import word_tokenize
44
45 documents = []
46 # Use only complete sentences
47 for i in range(fullSentDf.shape[0]):
48     tmpWords = word_tokenize(fullSentDf['Phrase'][i])
49     documents.append((tmpWords, fullSentDf['Sentiment'
50         ][i]))
51
52 random.seed(9001)
53 random.shuffle(documents)
54 print(documents[1][0])
55
56 len(documents)
57
58 from nltk.corpus import stopwords
59 from nltk.stem import WordNetLemmatizer,
60     PorterStemmer, LancasterStemmer
61 porter = PorterStemmer()
62 lancaster=LancasterStemmer()
63 wordnet_lemmatizer = WordNetLemmatizer()
64 stopwords_en = stopwords.words("english")
65 punctuations="?!,.,;'\\"-()"
66
67 #parameters to adjust to see the impact on outcome
68 remove_stopwords = True
69 useStemming = False
70 useLemma = True
71 removePuncs = True
72
73 for l in range(len(documents)):
74     label = documents[l][1]
75     tmpReview = []
76     for w in documents[l][0]:
77         newWord = w
78         if remove_stopwords and (w in stopwords_en):
79             continue
80         if removePuncs and (w in punctuations):
81             continue
82         if useStemming:
83             newWord = porter.stem(newWord)
84             newWord = lancaster.stem(newWord)
85         if useLemma:
86             newWord = wordnet_lemmatizer.lemmatize(newWord)
87
88     tmpReview.append(newWord)
89 documents[l] = (' '.join(tmpReview), label)
90
91 print(documents[2])
92
93 all_data = pd.DataFrame(documents,
94                         columns=['text', '
95                             sentiment'])
96 # Splits the dataset so 70% is used for training and
97     30% for testing
98 x_train_raw, x_test_raw, y_train_raw, y_test_raw =
99     train_test_split(all_data['text'], all_data['
100         sentiment'], test_size=0.3, random_state=2003)
101
102 len(x_train_raw)
103
104 from sklearn.feature_extraction.text import
105     TfidfVectorizer
106 from sklearn.feature_extraction.text import
107     CountVectorizer
108
109 vectorizer = TfidfVectorizer(stop_words="english",
110                             ngram_range=(1, 1))
111
112 x_train = vectorizer.fit_transform(x_train_raw)
113 y_train = y_train_raw
114 x_test = vectorizer.transform(x_test_raw)
115 y_test = y_test_raw
116
117 # Converts the datasets to numpy arrays to work with
118     our PyTorch model
119 x_train_np = x_train.toarray()
120 y_train_np = np.array(y_train)
121
122 # Convert the testing data
123 x_test_np = x_test.toarray()
124 y_test_np = np.array(y_test)
125
126 x_train_np.shape
127 x_train_np = np.expand_dims(x_train_np, -1)
128
129 x_test_np = np.expand_dims(x_test_np, -1)
130 x_train_np.shape

```

Listing 1. Python code

The model and evaluation part is given below.

```

1
2
3 from keras import backend as K
4
5 def recall_m(y_true, y_pred):
6     true_positives = K.sum(K.round(K.clip(y_true*
7         y_pred, 0, 1)))
8     possible_positives = K.sum(K.round(K.clip(y_true
9         , 0, 1)))
10    recall = true_positives / (possible_positives + K.
11        epsilon())
12    return recall
13
14 def precision_m(y_true, y_pred):
15     true_positives = K.sum(K.round(K.clip(y_true*
16         y_pred, 0, 1)))
17     predicted_positives = K.sum(K.round(K.clip(y_pred
18         , 0, 1)))
19     precision = true_positives / (predicted_positives
20         + K.epsilon())
21    return precision
22
23 def f1_m(y_true, y_pred):
24     precision = precision_m(y_true, y_pred)
25     recall = recall_m(y_true, y_pred)
26    return 2*((precision*recall)/(precision+recall+K.
27        epsilon()))
28
29 from keras.models import Sequential

```

```

23 from keras.layers import Conv1D, MaxPooling1D, Dense
    , Flatten, Activation
24 from keras import optimizers
25 from keras import layers
26 from keras.utils import to_categorical
27
28 def cnn_model(fea_matrix, compiler):
29     model = Sequential()
30     model.add(Conv1D(filters = 16, kernel_size=2,
        activation='relu', input_shape = (fea_matrix.
        shape[1], fea_matrix.shape[2])))
31     model.add(MaxPooling1D(pool_size=2))
32     model.add(Conv1D(filters = 32, kernel_size=2,
        activation='relu'))
33     model.add(MaxPooling1D(pool_size=2))
34     model.add(Conv1D(filters = 64, kernel_size=2,
        activation='relu'))
35     model.add(MaxPooling1D(pool_size=2))
36     model.add(Flatten())
37     model.add(Activation('relu'))
38     model.add(Dense(5))
39     model.add(Activation('softmax'))
40     model.compile(optimizer = compiler, loss = '
        categorical_crossentropy', metrics = ['acc',
        f1_m, precision_m, recall_m])
41     return model
42
43 model = cnn_model(x_train_np, optimizers.Nadam(lr =
    1e-3))
44
45 y_train_np = to_categorical(y_train_np)
46 y_test_np = to_categorical(y_test_np)
47 model.fit(x_train_np, y_train_np, batch_size=64,
    epochs = 30, verbose = 1, validation_split =
    0.3)
48 model.save("1117675_ldconv_reg")
49
50 loss, accuracy, f1_score, precision, recall = model.
    evaluate(x_test_np, y_test_np, verbose = False)
51
52 #get_metrics(accuracy, f1_score, precision, recall)
53 print('Accuracy:', np.round(accuracy, 4))
54 print('Precision:', np.round(precision, 4))
55 print('Recall:', np.round(recall, 4))
56 print('F1 Score:', np.round(f1_score, 4))

```

Listing 2. Python code

IV. EXPERIMENTAL ANALYSIS

Selecting a perfect model for text classification always require some trial and error. The number of convolutional layers, the batch size, learning rate, number of epoch and the number of filters decide the efficiency of the model. After multiple combinations of these deciding parameters, the final combination of model is three convolutional layers with a filter size of 16,32,64 respectively. The activation function used everywhere is relu and softmax in the end before output layer because our problem is multiple class. The table below indicates the results for different number of epochs.

Epochs	Accuracy	Precision	F1 Score	Recall
100	0.3077	0.3092	0.3054	0.3019
30	0.3147	0.3176	0.3145	0.3116
10	0.3186	0.3167	0.3089	0.3127

Fig. 1. Block diagram of proposed model

V. CONCLUSION

The final accuracy, precision, recall and f1-score for our Convolutional Neural Network (CNN) model are 0.3077,0.3092,0.3019 and 0.3054 respectively. SO preprocessing helped in improving the performance.

REFERENCES

- [1] Sorostinean, Mihaela, Katia Sana, Mohamed Mohamed, and Amal Targhi. "Sentiment analysis on movie reviews." In Journal Agroparis-tech. 2017.
- [2] Camacho-Collados, Jose, and Mohammad Taher Pilehvar. "On the role of text preprocessing in neural network architectures: An evaluation study on text categorization and sentiment analysis." arXiv preprint arXiv:1707.01780 (2017).