

Multi-class Sentiment Analysis using Deep Learning

BIRVA PATEL

Department of Computer Science

Lakehead University 1111092

Thunder Bay, Ontario

bpatel24@lakeheadu.ca

Abstract—The aim of the project is to experiment with different Deep learning algorithms(mainly by CNN) to find the solution for the problem of text-based movie review multi-class sentiment analysis And returns the accuracy, recall, Precision and figure of merits(f-1) score.

Index Terms—Convolutional neural network, Tf-Idf, tokenization, vectorisation,n-gram,Deep Learning, Neural Network

I. INTRODUCTION

The human language is complex therefore teaching a machine to analyze the various grammatical nuances, cultural variations, slang and misspellings that occur in reviews provided by users is a difficult process. Teaching a machine to understand how context can affect tone is even more difficult. This method of sentiment analysis is useful in a wide range of domains, such as business or politics. [1] The purpose is to experiment different machine learning models for the task of sentiment analysis, using the Rotten tomatoes movie review corpus. More specifically, phrases from the dataset have to be classified in 5 categories, according to the intended tone of the user: negative, somewhat negative, neutral, somewhat positive, positive.

In this paper i will mainly focus on how you can predict the sentiment of movie reviews.in Python using the Keras deep learning library. This is a multi-class classification problem, which simply means the data set have more than 2 classes(binary classifier) that is 5.I am using raw train data of Rotten Tomatoes movie reviews.which is showed in fig-1, The dataset used is Pang and Lee’s movie review sentiment polarity dataset, which consists of 5,331 positive and 5,331 negative sentences, formed from a vocabulary size of 20,000. For all the models, the data is shuffled and 30% of the dataset is used as the test set [2]. For that I can use several sequence handling

PhraseId	SentenceId	Phrase	Sentiment
0	1	1 A series of escapades demonstrating the adage ...	1
1	2	1 A series of escapades demonstrating the adage ...	2
2	3	1 A series	2
3	4	1 A	2
4	5	1 series	2

Fig. 1. First five row of dataset.

components like GRU, RNN, LSTM but here i implemented it using convolution neural network.

TF-IDF (term frequency-inverse document frequency) is a statistical measure that evaluates how relevant a word is to a document in a collection of documents. This is done by multiplying two metrics: how many times a word appears in a document, and the inverse document frequency of the word across a set of documents.

II. LITERATURE REVIEW

Early works of sentiment classification mainly focus on polarity detection of English product reviews or movie reviews.

Kennedy and Inkpen (2006) [3] present two methods for determining the sentiment expressed by a movie review. The first method uses GI to classify customers’ reviews based on the number of positive and negative terms they contain, as well as negations, overstatements and understatements. Negations are used to reverse the semantic polarity of a particular term, while intensifiers and diminishers are used to increase and decrease, respectively, the degree to which a term is positive or negative. The second method uses a ML algorithm, SVM. Authors start with unigram features and then add bigrams that consist of a valence shifter and another word. The accuracy of classification is very high, and the valence shifter bigrams slightly improve it. They also demonstrate that combining the two methods achieves better results than either method alone.

The work in Prabowo and Thelwall (2009) [4] combines rule-based classification, supervised learning and machine learning into a new method. This method is tested on movie reviews, product reviews and MySpace comments. The procedure is that if one classifier fails to classify a document, the classifier will pass the document onto the next classifier, until the document is classified or no other classifier exists. A number of approaches that focus on acquiring and defining a set of rules are used along with SVM learning: General Inquirer Based Classifier (GIBC), Rule-Based Classifier (RBC), Statistics Based Classifier (SBC) and Induction Rule Based Classifier (IRBC). Experiments were carried out by applying different sequences of previous classifiers. Results showed that the use of multiple classifiers in a hybrid manner can improve the effectiveness of sentiment analysis. The sequence RBC then SBC then GIBC then SVM of classifiers yielded the highest accuracy in comparison to other classifiers’ sequences.

As you can notice that most of the researchers tries to use the time sequence algorithms which are easy to handle the large data where in this project I am experimenting to use the deep learning algorithms like CNN and trying to get as much

high accuracy as I can. The data set is quite challenging as the dataset is split into separate phrases and each one of them have a sentiment label. Moreover, obstacles like sentence negation, sarcasm, terseness, language ambiguity, and many others make this task very interesting for Natural Language Processing.

here I have mentioned the complete path to my model that i have used fig-2 is the flow chart that exactly shows my workflow.

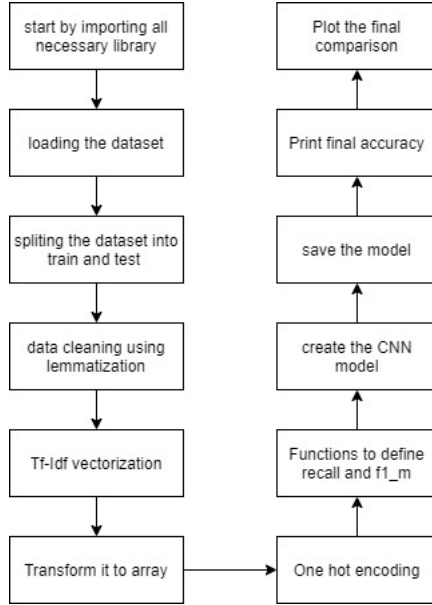


Fig. 2. Flow of the Proposed Model.

III. PROPOSED MODEL

- Import necessary libraries: As per appendix-1, we can call numpy for array operations, pandas for data visualization, sklearn for data analysis.
- Dataset: As per appendix-2, I load the rotten tomato's movie review analysis data from the GitHub link which is provided in the lab. Then I used pandas's 'read csv' function file to read the data from the given link into tabular format which is separated by tab and store it in to fullSent variable as shown in appendix3. Now, I divide all the phrases column as an inputs and its sentiment column as an outputs. Only full sentences are considered and stored in dataframe.
- Data preprocessing: As per appendix5, I removed all punctuations and stop words using Lemmatizer and stemmer
- Splitting dataset: As per appendix-6, I have split the dataset in training and testing data. training dataset is used, while we train the model and testing dataset is used while predicting the data. I use 70% for the training and 30% for the testing of whole dataset. train test split() function from sklearn used for that and as a parameter I set test size= 0.3.
- Vectorization: As per appendix 7, I used TfIdf vectorizer which convert phrases into bigrams. It removes the stop

words in english languages as well and make a list of words. Also I Convert train and test data into vectorize array and I transform the x train and x test into train idf and test idf respectively using vectorizer and then convert it into the array and reshape it into three dimension.

- One-hot encoding: As per appendix8, I convert all the labels into binary array using to categorical from the keras. for example, if an output is 6, it will convert it to the [0 0 0 0 0 1 0 0 0].
- Define necessary functions: As per appendix-9, define the function for recall which means true positive out of possible positives value, precision is true positives out of predicted positives and f1-score can be interpreted as a weighted average of the precision and recall.
- Convolution neural network(CNN) for multi-class classification: As per appendix-10, CNN class is created for the convolution operations. I initialize all the variables. I used convolution layers like conv1d and relu activation function for converting all the negative values into zero. After that, max-pooling is used for sub-sampling and then used flatten() to convert the output from CNN to feed forward network. then add dense layer for the five classes as the dataset has five classes.
- Compile the model: As per appendix-10, we compile our cnn model using adam optimizer and using 'reg' mode. we set the learning rate 0.001. number of classes are five.
- Train model: As per appendix-11, we train and test model using 10 epochs and batch size is 64. If number of epoch is more network is deep and accuracy will increase for certain number of epochs.
- Save the model: As per appendix-12, I saved the model as '1116576 1dconv reg.h5' using model.save() function.
- Test model: As per appendix-13, I test the model on the training data and get accuracy, precision, recall and f1-score and print these values up to two decimal points.
- Plot the accuracy graph: As per appendix14, define a plotist() function which plot the graph for training and testing accuracy per epochs. fig-3 denotes the accuracy graph for both training and testing dataset.

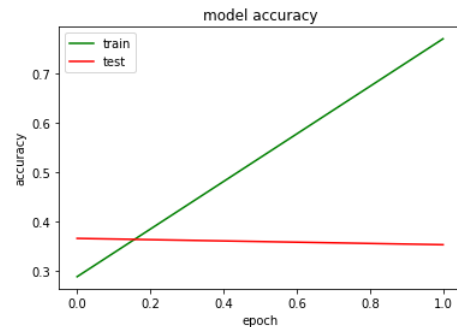


Fig. 3. Accuracy graph.

- plot the metrics graph: to show the training and testing value difference between the precision, recall and f1-score. I plotted this graph using same function. fig-4

indicates the difference of training and testing values of the precision, recall and f1-score

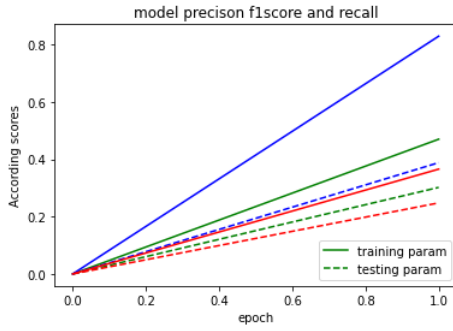


Fig. 4. precision, recall, f1-score.

By that way, The model's accuracy is 34% , precision is 35% , recall is 32% and f1-score is also 33%. My traing accuracy is quite high but testing accuracy is quite low because of the small dataset that I have used. The reason behind that is when i convert the data into corpus of words it will go out of memory utilization.

IV. EXPERIMENTAL ANALYSIS

I have experimented on different models listed below in the table and it indicates the experimental results of the different models.

Layers	Accuracy	Precision	Recall	F1-score
CNN	34	35	32	33
Embedding	59	62	59	58

V. APPENDIX

```

1 # Import the pandas library to read our dataset
2 import pandas as pd
3
4 # Get the train/test split package from sklearn for
  preparing our dataset to
5 # train and test the model with
6 from sklearn.model_selection import train_test_split
7
8 # Import the numpy library to work with and
  manipulate the data
9 import numpy as np

```

Listing 1. importing library

```

1 train_data = pd.read_csv('https://raw.
  githubusercontent.com/cacoderquan/Sentiment-
  Analysis-on-the-Rotten-Tomatoes-movie-review-
  dataset/master/train.tsv', sep='\t')
2
3 train_data.head(5)

```

Listing 2. load dataset

```

1 # Get number of unique sentences
2 numSentences = data['SentenceId'].max()
3 # extract full sentences only from the dataset
4 fullSentences = []
5 curSentence = 0
6 for i in range(data.shape[0]):
7     if data['SentenceId'][i] > curSentence:
8         fullSentences.append((data['Phrase'][i], data['
          Sentiment'][i]))

```

```

9     curSentence = curSentence + 1
10
11 len(fullSentences)
12 # put data into a df
13 fullSentDf = pd.DataFrame(fullSentences,
14                             columns=['Phrase', '
          Sentiment'])
15
16 # Check class imbalance in tokenized sentences
17 data['Sentiment'].value_counts()

```

Listing 3. get the full sentences

```

1 import nltk
2 import random
3 nltk.download('punkt')
4 nltk.download('stopwords')
5 nltk.download('wordnet')
6 from nltk.tokenize import word_tokenize
7
8 documents = []
9 # Use only complete sentences
10 for i in range(fullSentDf.shape[0]):
11     tmpWords = word_tokenize(fullSentDf['Phrase'][i])
12     documents.append((tmpWords, fullSentDf['Sentiment'
          ][i]))
13
14 random.seed(9001)
15 random.shuffle(documents)
16 print(documents[1][0])

```

Listing 4. separate the dataset

```

1 from nltk.corpus import stopwords
2 from nltk.stem import WordNetLemmatizer,
  PorterStemmer, LancasterStemmer
3 porter = PorterStemmer()
4 lancaster = LancasterStemmer()
5 wordnet_lemmatizer = WordNetLemmatizer()
6 stopwords_en = stopwords.words("english")
7 punctuations = "?!.,;'\\"-()"
8
9 #parameters to adjust to see the impact on outcome
10 remove_stopwords = True
11 useStemming = True
12 useLemma = False
13 removePuncs = True
14
15 for l in range(len(documents)):
16     label = documents[l][1]
17     tmpReview = []
18     for w in documents[l][0]:
19         newWord = w
20         if remove_stopwords and (w in stopwords_en):
21             continue
22         if removePuncs and (w in punctuations):
23             continue
24         if useStemming:
25             #newWord = porter.stem(newWord)
26             newWord = lancaster.stem(newWord)
27         if useLemma:
28             newWord = wordnet_lemmatizer.lemmatize(newWord)
29         tmpReview.append(newWord)
30     documents[l] = (' '.join(tmpReview), label)
31 print(documents[2])

```

Listing 5. Data Preprocessing

```

1 all_data = pd.DataFrame(documents,
2                           columns=['text', '
          sentiment'])
3 # Splits the dataset so 70% is used for training and
  30% for testing

```

```

4 x_train_raw, x_test_raw, y_train_raw, y_test_raw =
    train_test_split(all_data['text'], all_data['
      sentiment'], test_size=0.3)
5
6 len(x_train_raw)

```

Listing 6. Splitting the data

```

1 from sklearn.feature_extraction.text import
    TfidfVectorizer
2 from sklearn.feature_extraction.text import
    CountVectorizer
3
4 # Transform each text into a vector of word counts
5
6 vectorizer = TfidfVectorizer(stop_words="english",
7                               ngram_range=(1, 2))
8 train_x = vectorizer.fit_transform(x_train_raw)
9
10
11 train_vectorized = vectorizer.transform(x_train_raw)
12 train_vectorized= np.array(train_vectorized.toarray
13                             ())
14 train_vectorized = train_vectorized.reshape(
15     train_vectorized.shape[0],train_vectorized.shape
16     [1],1)
17 print(train_vectorized.shape)
18
19 test_vectorized = vectorizer.transform(x_test_raw)
20 test_vectorized= np.array(test_vectorized.toarray())
21 test_vectorized = test_vectorized.reshape(
22     test_vectorized.shape[0],test_vectorized.shape
23     [1],1)
24 print(test_vectorized.shape)

```

Listing 7. Vectorization

```

1 from keras.utils import to_categorical
2 Y_train = to_categorical(y_train_raw)
3 Y_test = to_categorical(y_test_raw)

```

Listing 8. One Hot Vectorization

```

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

```

Listing 9. Define necessary functions

```

1 import keras
2 from keras.models import Sequential
3 from keras.layers import Dense, Activation, Dropout,
    Flatten, Conv1D, MaxPooling1D

```

```

4 from keras.layers.normalization import
    BatchNormalization
5 import numpy as np
6 from keras import optimizers
7 np.random.seed(1000)
8 def CNN_model(fea_matrix, n_class, mode, compiler):
9     model = Sequential()
10    model.add(Conv1D(filters=64, kernel_size=1,
11        activation='relu',
12        input_shape=(fea_matrix.shape[1],
13        fea_matrix.shape[2])))
14    model.add(MaxPooling1D(pool_size=2))
15    model.add(Conv1D(filters=128, kernel_size=1,
16        activation='relu'))
17    model.add(MaxPooling1D(pool_size=2))
18    model.add(Flatten())
19    model.add(Activation('relu'))
20    model.add(Dense(n_class))
21
22    model.add(Activation('softmax'))
23    model.compile(optimizer=compiler, loss='
24        categorical_crossentropy',
25        metrics=['acc', f1_m, precision_m,
26        recall_m])
27    return model
28
29 lr = 1e-3
30 decay=1e-4
31 mode="reg"
32 n_class = 5
33
34 adm = optimizers.Adam(lr=lr, decay=decay)
35 model = CNN_model(train_vectorized, n_class, mode, adm)

```

Listing 10. Model Layout

```

1 #fit the model
2 x = model.fit(train_vectorized, Y_train, batch_size
    =64, epochs=10, verbose=1, validation_split=0.3)

```

Listing 11. Fit the Model

```

1 #save the model
2 model.save('1111092_1dconv_reg.h5')

```

Listing 12. Save the Model

```

1 #function call for accuracy f1-score, precision and
    recall
2 def metrics(accuracy, f1_score, precision, recall):
3     print('CNN model performance')
4     print('Accuracy:', np.round(accuracy, 2))
5     print('Precision: ', np.round(precision, 2))
6     print("recall:", np.round(recall, 2))
7     print("f1score:", np.round(f1_score, 2))
8     loss, accuracy, f1_score, precision, recall = model.
9     evaluate(test_vectorized, Y_test, verbose=False)
10
11 # get the accuracy, f1_score, presion, recall
12 metrics(accuracy, f1_score, precision, recall)

```

Listing 13. Test model

```

1 import matplotlib.pyplot as plt
2 def plothist(hist):
3     plt.plot(hist.history['acc'], color='green')
4     plt.plot(hist.history['val_acc'], color='red')
5     plt.title('model accuracy')
6     plt.ylabel('accuracy')
7     plt.xlabel('epoch')
8     plt.legend(['train', 'test'], loc='upper left')
9     plt.show()
10 plothist(x)

```

Listing 14. plot the accuracy graph

```

1 import matplotlib.pyplot as plt
2 def plothist(hist):
3     plt.plot(hist.history['f1_m'],color='green')
4     plt.plot(hist.history['val_f1_m'],color='green',
5             linestyle='dashed')
6     plt.plot(hist.history['precision_m'],color='blue')
7     plt.plot(hist.history['val_precision_m'],color='blue',
8             linestyle='dashed')
9     plt.plot(hist.history['recall_m'],color='red')
10    plt.plot(hist.history['val_recall_m'],color='red',
11            linestyle='dashed')
12    plt.title('model precison f1score and recall')
13    plt.legend(['training param','testing param'],
14              loc='lower right')
15    plt.ylabel('According scores')
16    plt.xlabel('epoch')
17    plt.show()
18 plothist(x)

```

Listing 15. plot the functional graph

VI. CONCLUSION

To conclude,1D CNN works well with the text datasets and in proposed method you can identify that it was not much complex problem where as sequencing models works much better.To find the more accurate results dataset should have to be in large amount.still accuracy, Precision, f1-score, recall are measured from the available data. Also I explained the important method used for text summarization which is tf-idf.

REFERENCES

- [1] M. Sorostinean, K. Sana, M. Mohamed, and A. Targhi, "Sentiment analysis on movie reviews," in *Journal Agroparistech*, 2017.
- [2] J. Hong, A. Nam, and A. Cai, "Multi-class text sentiment analysis," 2019.
- [3] A. Kennedy and D. Inkpen, "Sentiment classification of movie reviews using contextual valence shifters," *Computational intelligence*, vol. 22, no. 2, pp. 110–125, 2006.
- [4] R. Prabowo and M. Thelwall, "Sentiment analysis: A combined approach," *Journal of Informetrics*, vol. 3, no. 2, pp. 143–157, 2009.