

Multi-class Sentiment Analysis using Deep Learning

Parth Valani
Computer Science
Lakehead University 1116576
Thunder Bay, Canada
Valaniparth13@gmail.com

Abstract—This paper gives idea about how to use simple convolution neural network on the multi-class dataset and doing the multi-class sentiment classification get accuracy, recall, Precision and figure of merits using keras library.

Index Terms—Convolution Neural Network, Overfitting, Kernel, pytorch, pooling, activation, epoch, batch size, learning rate.

I. INTRODUCTION

In this report, I will going to discuss how we can use 1D convolution neural network to train the model with multi-class values. I use roton tomatoes dataset [1] which has 5 classes like negative,somewhat negative,nuaral,somewhat positive,positive.

The goal of this report is to explore and apply the most recent technique's used in this task involving Natural language processing and to see and quantify the benefit of using those techniques to this kind of problem. [2]

II. LITERATURE REVIEW

With the rise of social media, Sentiment Analysis, which is one of the most well-known NLP tasks, gained a lot of importance over the years. The recent advances made in Machine Learning and Deep Learning made it an even more active task where a lot of work and research is still done.

CNN, which has been widely used on image datasets, extracts the significant features of the image, as the “convolutional” filter (i.e., kernel) moves through the image. If the input data are given as one-dimensional, the same function of CNN could be used in the text as well. In the text area, while the filter moves, local information of texts is stored, and important features are extracted. Therefore, using CNN for text classification is effective. [3]

There are many other model that have been used in the past like RNN, Embedding and LSTM and other model can also give far more accuracy than CNN network.

III. PROPOSED MODEL

Here is the step by step process:

- First, import all the necceesary libraries in the python. for example pandas for data visulization, numpy for array operations, sklearn for machine learning operations.
- Now, load the roton tomatos dataset from the given github link.
- Then, store it into a tabular format using pandas and get the full sentence from the tabular data.

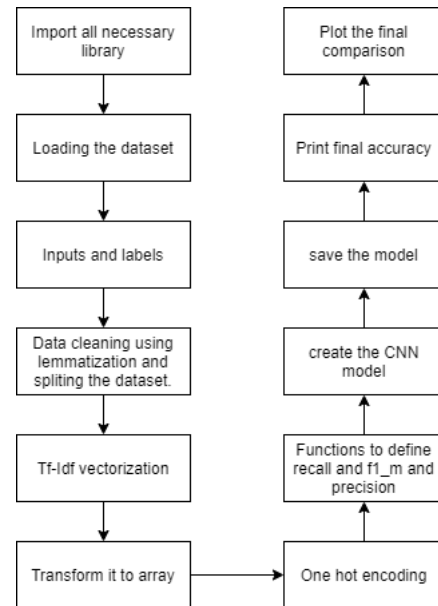


Fig. 1. Block diagram of process

- Divide the training data into inputs(phrase) and labels(sentiment) for further processing.
- remove stop words and use lemmatazation.
- Now, split the data in to training and testing data in the ratio of 70 to 30 percentage.
- In the next step, initiate the Tfidf vectorizer and fit the input phrase into into vectorizer.
- Converting train data into the 3 dimensional vector from 2 dimension. so that, we can feed it to the network.
- Then, convert the value to binary format using one hot encoding..
- Now, define all the necessary functions like recall, precision, figure-of metric.
- Now, define the convolution network of 4 layer and at the end connect fully connected network and compile the model.
- Finally, In the training step,feed the training data into 3-D shape and classify it.
- Now, save our model and get metrics including accuracy,recall,precision,f1_score.
- plot the training and validation accuracy vs epoch.

(A) **Import necessary libraries:** As per appendix-1, we

can call numpy for array operations, pandas for data visualization, sklearn for data analysis.

- (B) **Load 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. I used pandas's 'read_csv' function file to read the data from the given link into tabular format which is separated by tab. [4]
- (C) **Data-preprocessing:** As per appendix-2, get full sentence and convert it to the tokens using word_tokenize function.
- (D) **Removing stop words and use lammetization:** As per appendix-3, I have removed all stop words and punctuation and use wordnet_lemmatizer and store it in the document list. and then take 'text' column of the all_data as an input and 'sentiment' column as an output.
- (E) **Splitting dataset:** As per appendix-4, 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.
- (F) **Vectorization:** As per appendix-5, I used Tfidf vectorizer which convert text into the vector form. It removes the stop words in English languages as well and make a list of words.
- (G) **Convert train and test data into vectorize array:** As per appendix-6, I convert all the training and testing data into the array.
- (H) **Convert 2-D to 3-D:** As per appendix-7, I convert inputs from two dimension to three dimensional space.
- (I) **One-hot encoding** As per appendix-8, 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].
- (J) **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.
- (K) **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.
- (L) **Compile the model:** As per appendix-11, we compile our cnn model using Nadam optimizer. we set the learning rate 0.001. number of classes are five.
- (M) **Train model:** As per appendix-12, 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.
- (N) **Save the model:** As per appendix-13, I saved the model

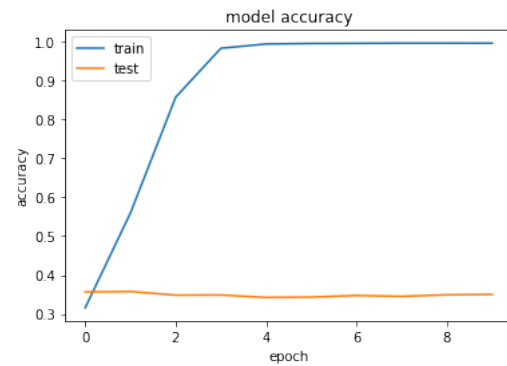


Fig. 2. epoch vs accuracy

as '1116576_1dconv_reg.h5' using model.save() function.

- (O) **Test model:** As per appendix-14, 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.
- (P) **Plot the epoch vs accuracy graph:** As per appendix-15, define a Plotlist() function which plot the graph for training and testing accuracy per epochs.

By that way, The model's accuracy is 35% , precision is 36% , recall is 34% and f1-score is also 35%. My training accuracy is quite high but testing accuracy is quite low because i used less number of data. The reason behind that is when I convert the data into corpus of words it will go out of memory utilization.

IV. CONCLUSION

To conclude, I have received 35% accuracy using just 3 convolution layer and 1 dense layer. I have also tried using embedding layer and I got little high accuracy.

From the above results, it has been clear that proper pre-processing of data based on natural language processing approaches as well as incorporating already existing models in the domain of sentiment analysis altogether with appropriate classification process can improve the performance of the model for multi-class classification.

V. APPENDIX

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 import numpy as np
```

Listing 1. Import all necessary libraries

```
1 data = pd.read_csv('https://raw.githubusercontent.com/cacoderquan/Sentiment-Analysis-on-the-Rotten-
2 com/cacoderquan/Sentiment-Analysis-on-the-Rotten-
3 -Tomatoes-movie-review-dataset/master/train.tsv')
4 data.head(20)
5 numSentences = data['SentenceId'].max()
6 # extract full sentences only from the dataset
7 fullSentences = []
8 curSentence = 0
```

```

7 for i in range(data.shape[0]):
8     if data['SentenceId'][i]> curSentence:
9         fullSentences.append((data['Phrase'][i], data['
10             Sentiment'][i]))
11         curSentence = curSentence + 1
12 # put data into a df
13 fullSentDf = pd.DataFrame(fullSentences,
14                             columns=['Phrase', '
15             Sentiment'])
16 documents = []
17 # Use only complete sentences
18 for i in range(fullSentDf.shape[0]):
19     tmpWords = word_tokenize(fullSentDf['Phrase'][i])
20     documents.append((tmpWords, fullSentDf['Sentiment'
21         ][i]))

```

Listing 2. Data-preprocessing

```

1 from nltk.corpus import stopwords
2 from nltk.stem import WordNetLemmatizer,
3   PorterStemmer, LancasterStemmer
4 porter = PorterStemmer()
5 lancaster=LancasterStemmer()
6 wordnet_lemmatizer = WordNetLemmatizer()
7 stopwords_en = stopwords.words("english")
8 punctuations="?:!.,;'\\"-()"
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])
32 all_data = pd.DataFrame(documents,
33                           columns=['text', '
34           sentiment'])

```

Listing 3. Removing stopwords and use lammetizer

```

1 # Splits the dataset so 70% is used for training and
2   30% for testing
3 x_train_raw, x_test_raw, y_train_raw, y_test_raw =
4   train_test_split(all_data['text'], all_data['
5       sentiment'], test_size=0.3)

```

Listing 4. Splitting data into 70:30 ratio

```

1 from sklearn.feature_extraction.text import
2   TfidfVectorizer
3 vectorizer = TfidfVectorizer(stop_words="english",
4                               ngram_range=(1, 1))
5 x_train = vectorizer.fit_transform(x_train_raw)
6 y_train = y_train_raw
7 x_test = vectorizer.transform(x_test_raw)

```

```
y_test = y_test_raw
```

Listing 5. Tfidf vectorizer

```

1 # Converts the datasets to numpy arrays to work with
2   our keras model
3 x_train_np = x_train.toarray()
4 y_train_np = np.array(y_train)
5
6 # Convert the testing data
7 x_test_np = x_test.toarray()
8 y_test_np = np.array(y_test)

```

Listing 6. Converts the datasets to numpy arrays

```

1 # Converting 2-D to 3-D
2 x_train_np = np.expand_dims(x_train_np, 2)
3 x_test_np = np.expand_dims(x_test_np, 2)
4 print(x_train_np.shape)

```

Listing 7. Convert 2-dimension data to 3-dimension for convolution

```

1 from keras.utils import to_categorical
2 y_train_np = to_categorical(y_train_np)
3 y_test_np = to_categorical(y_test_np)
4 print(y_train_np.shape)

```

Listing 8. One-hot encoding

```

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

```

Listing 9. Define necessary functions

```

1 # import the neccesary libraries for the training
2   the model
3 from keras.models import Sequential
4 from keras.layers import Conv1D, MaxPooling1D, Dense
5   , Dropout, Flatten
6 from keras.layers import Activation, Conv1D,
7   GlobalMaxPooling1D
8 from keras import optimizers
9 from keras import layers
10 # define convolution neural network
11 def cnn_model(data, num_class, opt):
12     model= Sequential()

```

REFERENCES

- [1] <https://raw.githubusercontent.com/cacoderquan/Sentiment-Analysis-on-the-Rotten-Tomatoes-movie-review-dataset/master/train.tsv>.
- [2] <https://www.kaggle.com/ynouri/rotten-tomatoes-sentiment-analysis>.
- [3] H. Kim and Y.-S. Jeong, "Sentiment classification using convolutional neural networks," *Applied Sciences*, vol. 9, no. 11, p. 2347, 2019.
- [4] <https://github.com/pratikmjoshi/sentiment-analysis-movie-reviews/blob/master/SentimentAnalysisofMovieReviews.ipynb>.

```

10 model.add(Conv1D(filters=32, kernel_size=2,
    activation= 'relu', input_shape= (data.shape[1],
    data.shape[2])))
11 model.add(MaxPooling1D(pool_size=2))
12 model.add(Conv1D(filters=64, kernel_size=2,
    activation= 'relu'))
13 model.add(MaxPooling1D(pool_size=2))
14 model.add(Conv1D(filters=128, kernel_size=2,
    activation= 'relu'))
15 model.add(MaxPooling1D(pool_size=2))
16 model.add(Flatten())
17 model.add(Activation('relu'))
18 model.add(Dense(num_class))
19
20 model.add(Activation('softmax'))
21 model.compile(optimizer= opt, loss='
    categorical_crossentropy', metrics=['acc', f1_m,
    precision_m, recall_m])
22
23 return model

```

Listing 10. 1D-CNN

```

1 # assign values to variables
2 batch_size= 64
3 num_epochs= 10
4 n_class = 5
5 Nadam = optimizers.Nadam(lr = 1e-3 , beta_1=0.9,
    beta_2=0.999, epsilon=1e-08)
6 # compile the model
7 model=cnn_model(x_train_np, n_class, Nadam)

```

Listing 11. Compile the model

```

1 # run the model
2 x = model.fit(x_train_np, y_train_np, batch_size=
    batch_size, epochs=num_epochs, verbose=1,
    validation_split=0.3)

```

Listing 12. Train the model

```

1 # save model with 50 epoch
2 model.save('1116576_1dconv_reg.h5')

```

Listing 13. Save the model

```

1 def metrics(accuracy, f1_score, precision, recall):
2     print('CNN model performance')
3     print('Accuracy:', np.round(accuracy,2))
4     print('Precision: ', np.round(precision,2))
5     print("recall:", np.round(recall,2))
6     print("f1score:", np.round(f1_score,2))
7 loss, accuracy, f1_score, precision, recall = model.
    evaluate(x_test_np, y_test_np, verbose=False)
8 metrics(accuracy, f1_score, precision, recall)

```

Listing 14. Get accuracy&f1-score&precision&recall and print the values

```

1 import matplotlib.pyplot as plt
2 def plothist(hist):
3     plt.plot(hist.history['acc'])
4     plt.plot(hist.history['val_acc'])
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 15. Plot the epoch vs accuracy