

Python and Deep Learning Programming Lab 4

Lab ID:14

Student ID: 20

Student Name : Raju Nekadi .

Student ID: 19

Student Name: Sneha Mishra

Github Link:

https://github.com/rnekadi/CSEE5590_PYTHON_DEEPLARNING_FALL2018/tree/master/Lab4

Video Link:

<https://youtu.be/KUsG7yV8oz0>

Introduction:

To apply CNN, LSTM and RNN models on text and image datasets for Text and Image Classification.

Objective:

Implementing the Text Classification using the CNN Model on text dataset.

Approaches:

For Text Classification using CNN we have used the Positive and Negative Polarity datasets. First of all we have loaded preprocessed polarity data from files, split the data into words and generate labels and returns vectors, labels, vocabulary, and inverse vocabulary.

Split the processed data into train and test for our model.

The setup the CNN Model using Embedding of Size 256 using the vocabulary we created. We have set 3 convolutional layers of filter size 3 4 5 ,3 MaxPool layer, Relu activation function and then compile our model using Adam optimizer. Perform Model fitting on train data and evaluated the model score and accuracy.

Also plotted the Loss and Accuracy on graph. The detailed explanation given below in workflow section.

Workflow:

First of all loaded the preprocessed file and extracted the input vectors, labels, vocabulary, and inverse vocabulary using the data_helpers.py file.

```

data_helpers.py x
61     return padded_sentences
62
63
64 def build_vocab(sentences):
65     """
66     Builds a vocabulary mapping from word to index based on the sentences.
67     Returns vocabulary mapping and inverse vocabulary mapping.
68     """
69     # Build vocabulary
70     word_counts = Counter(itertools.chain(*sentences))
71     # Mapping from index to word
72     vocabulary_inv = [x[0] for x in word_counts.most_common()]
73     vocabulary_inv = list(sorted(vocabulary_inv))
74     # Mapping from word to index
75     vocabulary = {x: i for i, x in enumerate(vocabulary_inv)}
76     return [vocabulary, vocabulary_inv]
77
78
79 def build_input_data(sentences, labels, vocabulary):
80     """
81     Maps sentences and labels to vectors based on a vocabulary.
82     """
83     x = np.array([[vocabulary[word] for word in sentence] for sentence in sentences])
84     y = np.array(labels)
85     return [x, y]
86
87
88 def load_data():
89     """
90     Loads and preprocesses data for the dataset.
91     Returns input vectors, labels, vocabulary, and inverse vocabulary.
92     """
93     # Load and preprocess data
94     sentences, labels = load_data_and_labels()
95     sentences_padded = pad_sentences(sentences)
96     vocabulary, vocabulary_inv = build_vocab(sentences_padded)
97     x, y = build_input_data(sentences_padded, labels, vocabulary)
98     return [x, y, vocabulary, vocabulary_inv]
99

```

Train and Test Data Creation

```

cnn_model.py x
2   from keras.layers import Reshape, Flatten, Dropout, Concatenate
3   from keras.callbacks import ModelCheckpoint
4   from keras.optimizers import Adam
5   from keras.models import Model
6   from sklearn.model_selection import train_test_split
7   from data_helpers import load_data
8   import matplotlib.pyplot as plt
9
10  print('Loading data')
11  x, y, vocabulary, vocabulary_inv = load_data()
12
13  # x.shape -> (10662, 56)
14  # y.shape -> (10662, 2)
15  # len(vocabulary) -> 18765
16  # len(vocabulary_inv) -> 18765
17
18  X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)
19

```

Creating Model , adding the Embedding CONV2D,MaxPool and Compile the model.

```

cnn_model.py x
26  sequence_length = x.shape[1]_# 56
27  vocabulary_size = len(vocabulary_inv)_# 18765
28  embedding_dim = 256
29  filter_sizes = [3, 4, 5]
30  num_filters = 512
31  drop = 0.5
32
33  epochs = 20
34  batch_size = 80
35
36  # this returns a tensor
37  print("Creating Model...")
38  inputs = Input(shape=(sequence_length,), dtype='int32')
39  embedding = Embedding(input_dim=vocabulary_size, output_dim=embedding_dim, input_length=sequence_length)(inputs)
40  reshape = Reshape((sequence_length, embedding_dim, 1))(embedding)
41
42  conv_0 = Conv2D(num_filters, kernel_size=(filter_sizes[0], embedding_dim), padding='valid', kernel_initializer='normal', activation='relu')
43  conv_1 = Conv2D(num_filters, kernel_size=(filter_sizes[1], embedding_dim), padding='valid', kernel_initializer='normal', activation='relu')
44  conv_2 = Conv2D(num_filters, kernel_size=(filter_sizes[2], embedding_dim), padding='valid', kernel_initializer='normal', activation='relu')
45
46  maxpool_0 = MaxPool2D(pool_size=(sequence_length - filter_sizes[0] + 1, 1), strides=(1, 1), padding='valid')(conv_0)
47  maxpool_1 = MaxPool2D(pool_size=(sequence_length - filter_sizes[1] + 1, 1), strides=(1, 1), padding='valid')(conv_1)
48  maxpool_2 = MaxPool2D(pool_size=(sequence_length - filter_sizes[2] + 1, 1), strides=(1, 1), padding='valid')(conv_2)
49
50  concatenated_tensor = Concatenate(axis=1)([maxpool_0, maxpool_1, maxpool_2])
51  flatten = Flatten()(concatenated_tensor)
52  dropout = Dropout(drop)(flatten)
53  output = Dense(units=2, activation='softmax')(dropout)
54
55  # this creates a model that includes
56  model = Model(inputs=inputs, outputs=output)
57
58  checkpoint = ModelCheckpoint('weights.{epoch:03d}-{val_acc:.4f}.hdf5', monitor='val_acc', verbose=1, save_best_only=True, mode='auto')
59  adam = Adam(lr=1e-4, beta_1=0.9, beta_2=0.999, epsilon=1e-08, decay=0.0)
60  model.compile(optimizer=adam, loss='binary_crossentropy', metrics=['accuracy'])
61
62

```

Fitting the Model on train data and train it . Calculated the accuracy and score and plot loss and accuracy on graph.

```

61 model.compile(optimizer=adam, loss='binary_crossentropy', metrics=['accuracy'])
62
63
64 print("Simplified convolutional neural network")
65 model.summary()
66
67 print("Traning Model...")
68 history = model.fit(X_train, y_train, batch_size=batch_size, epochs=epochs, verbose=1, callbacks=[checkpoint],
69                     validation_data=(X_test, y_test)) # starts training
70
71 print("Evaluating the Model")
72
73 score, accuracy = model.evaluate(X_test, y_test, batch_size=batch_size, verbose=1)
74
75 print('Score: %.2f' %(score))
76 print('Validation Accuracy: %.2f' %(accuracy))
77
78 # Plot The CNN LOSS
79 fig1 = plt.figure()
80 plt.plot(history.history['loss'], 'r', linewidth=3.0)
81 plt.plot(history.history['val_loss'], 'b', linewidth=3.0)
82 plt.legend(['Training loss', 'Validation Loss'], fontsize=18)
83 plt.xlabel('Epochs ', fontsize=16)
84 plt.ylabel('Loss', fontsize=16)
85 plt.title('Loss Curves :CNN', fontsize=16)
86 fig1.savefig('loss_cnn.png')
87
88
89 # Plot the CNN Accuracy
90
91 fig2=plt.figure()
92 plt.plot(history.history['acc'], 'r', linewidth=3.0)
93 plt.plot(history.history['val_acc'], 'b', linewidth=3.0)
94 plt.legend(['Training Accuracy', 'Validation Accuracy'], fontsize=18)
95 plt.xlabel('Epochs ', fontsize=16)
96 plt.ylabel('Accuracy', fontsize=16)
97 plt.title('Accuracy Curves : CNN', fontsize=16)
98 fig2.savefig('accuracy_cnn.png')
99 plt.show()

```

This concludes our workflow.

Dataset:

<http://www.cs.cornell.edu/people/pabo/movie-review-data/>

Parameters:

None

Evaluation:

After running the model for the for 20 Epochs, batch size of 80 and dropout of 50% we got the got the score and accuracy of 72% and 75% respectively.

```
cnn_model x
80/2133 [>.....] - ETA: 4s
160/2133 [=>.....] - ETA: 4s
240/2133 [==>.....] - ETA: 4s
320/2133 [===>.....] - ETA: 4s
400/2133 [====>.....] - ETA: 4s
480/2133 [=====>.....] - ETA: 3s
560/2133 [=====>.....] - ETA: 3s
640/2133 [======>.....] - ETA: 3s
720/2133 [======>.....] - ETA: 3s
800/2133 [======>.....] - ETA: 3s
880/2133 [======>.....] - ETA: 2s
960/2133 [======>.....] - ETA: 2s
1040/2133 [======>.....] - ETA: 2s
1120/2133 [======>.....] - ETA: 2s
1200/2133 [======>.....] - ETA: 2s
1280/2133 [======>.....] - ETA: 1s
1360/2133 [======>.....] - ETA: 1s
1440/2133 [======>.....] - ETA: 1s
1520/2133 [======>.....] - ETA: 1s
1600/2133 [======>.....] - ETA: 1s
1680/2133 [======>.....] - ETA: 1s
1760/2133 [======>.....] - ETA: 0s
1840/2133 [======>.....] - ETA: 0s
1920/2133 [======>...] - ETA: 0s
2000/2133 [======>..] - ETA: 0s
2080/2133 [======>.] - ETA: 0s
2133/2133 [=====] - 5s 2ms/step
Score: 0.72
Validation Accuracy: 0.75
```

Conclusion:

CNN Model is very good for Text Classification.

Objective:

Implementing the Text Classification using the LSTM Model on text dataset.

Approaches:

For Text Classification using CNN we have used the Positive and Negative Polarity datasets. First of all we have loaded preprocessed polarity data from files, split the data into words and generate labels and returns vectors, labels, vocabulary, and inverse vocabulary.

Split the processed data into train and test for our model.

The setup the LSTM Model using Embedding. Add the LSTM layer of 100 with dropout 20% .Added the Dense of size 2 and activation used is sigmoid . Compile model using adam optimizer Perform Model fitting on train data and evaluated the model score and accuracy.

Also plotted the Loss and Accuracy on graph. The detailed explanation given below in workflow section.

Workflow:

Fist of all loaded the preprocessed file and extracted the input vectors, labels, vocabulary, and inverse vocabulary using the data_helpers.py file.


```

data_helpers.py x
61     return padded_sentences
62
63
64 def build_vocab(sentences):
65     """
66     Builds a vocabulary mapping from word to index based on the sentences.
67     Returns vocabulary mapping and inverse vocabulary mapping.
68     """
69     # Build vocabulary
70     word_counts = Counter(itertools.chain(*sentences))
71     # Mapping from index to word
72     vocabulary_inv = [x[0] for x in word_counts.most_common()]
73     vocabulary_inv = list(sorted(vocabulary_inv))
74     # Mapping from word to index
75     vocabulary = {x: i for i, x in enumerate(vocabulary_inv)}
76     return [vocabulary, vocabulary_inv]
77
78
79 def build_input_data(sentences, labels, vocabulary):
80     """
81     Maps sentences and labels to vectors based on a vocabulary.
82     """
83     x = np.array([[vocabulary[word] for word in sentence] for sentence in sentences])
84     y = np.array(labels)
85     return [x, y]
86
87
88 def load_data():
89     """
90     Loads and preprocesses data for the dataset.
91     Returns input vectors, labels, vocabulary, and inverse vocabulary.
92     """
93     # Load and preprocess data
94     sentences, labels = load_data_and_labels()
95     sentences_padded = pad_sentences(sentences)
96     vocabulary, vocabulary_inv = build_vocab(sentences_padded)
97     x, y = build_input_data(sentences_padded, labels, vocabulary)
98     return [x, y, vocabulary, vocabulary_inv]
99

```

Train and Test Data Creation


```

lstm_model.py x
1 from keras.layers import Dense, Embedding, LSTM
2 from keras.optimizers import Adam
3 from keras.models import Model
4 from keras.models import Sequential
5 from sklearn.model_selection import train_test_split
6 from data_helpers import load_data
7 import matplotlib.pyplot as plt
8
9 print('Loading data')
10 x, y, vocabulary, vocabulary_inv = load_data()
11
12 # x.shape -> (10662, 56)
13 # y.shape -> (10662, 2)
14 # len(vocabulary) -> 18765
15 # len(vocabulary_inv) -> 18765
16
17 X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)

```

Creating model, Embedding, Adding LSTM Layer Dense layer and compiling using adam optimizer.

```

lstm_model.py x
31
32 epochs = 20
33 batch_size = 80
34
35 # Creating the Model
36 print("Creating Model...")
37
38
39 model = Sequential()
40 model.add(Embedding(20000, 100, input_length=56))
41 model.add(LSTM(100, dropout=0.2, recurrent_dropout=0.2))
42 model.add(Dense(2, activation='sigmoid'))
43 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
44
45 model.summary()
46
47

```

We fitted the model on train data and train it . Calculated the accuracy and score and plot loss and accuracy on graph.

```

46
47
48     ## Fit the model
49
50     history = model.fit(X_train, y_train, validation_data=(X_test, y_test), batch_size=batch_size, epochs=epochs, verbose=1)
51
52
53     # Evaluate the model
54
55     score, accuracy = model.evaluate(X_test, y_test, batch_size=batch_size, verbose=1)
56
57     print('LSTM Score: %.2f' %(score))
58     print('LSTM Validation Accuracy: %.2f' % (accuracy))
59
60
61     # Plot The LSTM LOSS
62     fig1 = plt.figure()
63     plt.plot(history.history['loss'], 'r', linewidth=3.0)
64     plt.plot(history.history['val_loss'], 'b', linewidth=3.0)
65     plt.legend(['Training loss', 'Validation Loss'], fontsize=18)
66     plt.xlabel('Epochs ', fontsize=16)
67     plt.ylabel('Loss', fontsize=16)
68     plt.title('Loss Curves :LSTM', fontsize=16)
69     fig1.savefig('loss_lstm.png')
70
71
72     # Plot the LSTM Accuracy
73
74     fig2=plt.figure()
75     plt.plot(history.history['acc'], 'r', linewidth=3.0)
76     plt.plot(history.history['val_acc'], 'b', linewidth=3.0)
77     plt.legend(['Training Accuracy', 'Validation Accuracy'], fontsize=18)
78     plt.xlabel('Epochs ', fontsize=16)
79     plt.ylabel('Accuracy', fontsize=16)
80     plt.title('Accuracy Curves : LSTM', fontsize=16)
81     fig2.savefig('accuracy_lstm.png')
82     plt.show()

```

This Concludes our workflow.

Dataset:

<http://www.cs.cornell.edu/people/pabo/movie-review-data/>

Parameters:

None

Evaluation:

After running the model for the for 20 Epochs, batch size of 80 and dropout of 50% we got the got the score and accuracy of 69% and 50% respectively.

```
rnn_model x
8529/8529 [=====] - 28s 3ms/step - loss: 0.6932 - acc: 0.4991 - val_loss: 0.6933 - val_acc: 0.4979

 80/2133 [>.....] - ETA: 1s
160/2133 [=>.....] - ETA: 1s
240/2133 [==>.....] - ETA: 1s
320/2133 [===>.....] - ETA: 1s
400/2133 [====>.....] - ETA: 1s
480/2133 [=====] - ETA: 1s
560/2133 [=====] - ETA: 1s
640/2133 [=====] - ETA: 1s
720/2133 [=====] - ETA: 1s
800/2133 [=====] - ETA: 1s
880/2133 [=====] - ETA: 1s
960/2133 [=====] - ETA: 0s
1040/2133 [=====] - ETA: 0s
1120/2133 [=====] - ETA: 0s
1200/2133 [=====] - ETA: 0s
1280/2133 [=====] - ETA: 0s
1360/2133 [=====] - ETA: 0s
1440/2133 [=====] - ETA: 0s
1520/2133 [=====] - ETA: 0s
1600/2133 [=====] - ETA: 0s
1680/2133 [=====] - ETA: 0s
1760/2133 [=====] - ETA: 0s
1840/2133 [=====] - ETA: 0s
1920/2133 [=====] - ETA: 0s
2000/2133 [=====] - ETA: 0s
2080/2133 [=====] - ETA: 0s
2133/2133 [=====] - 2s 852us/step
LSTM Score: 0.69
LSTM Validation Accuracy: 0.50
Process finished with exit code 0
```

Conclusion:

LSTM Model gives average performance for Text Classification as compared to CNN.

Objective:

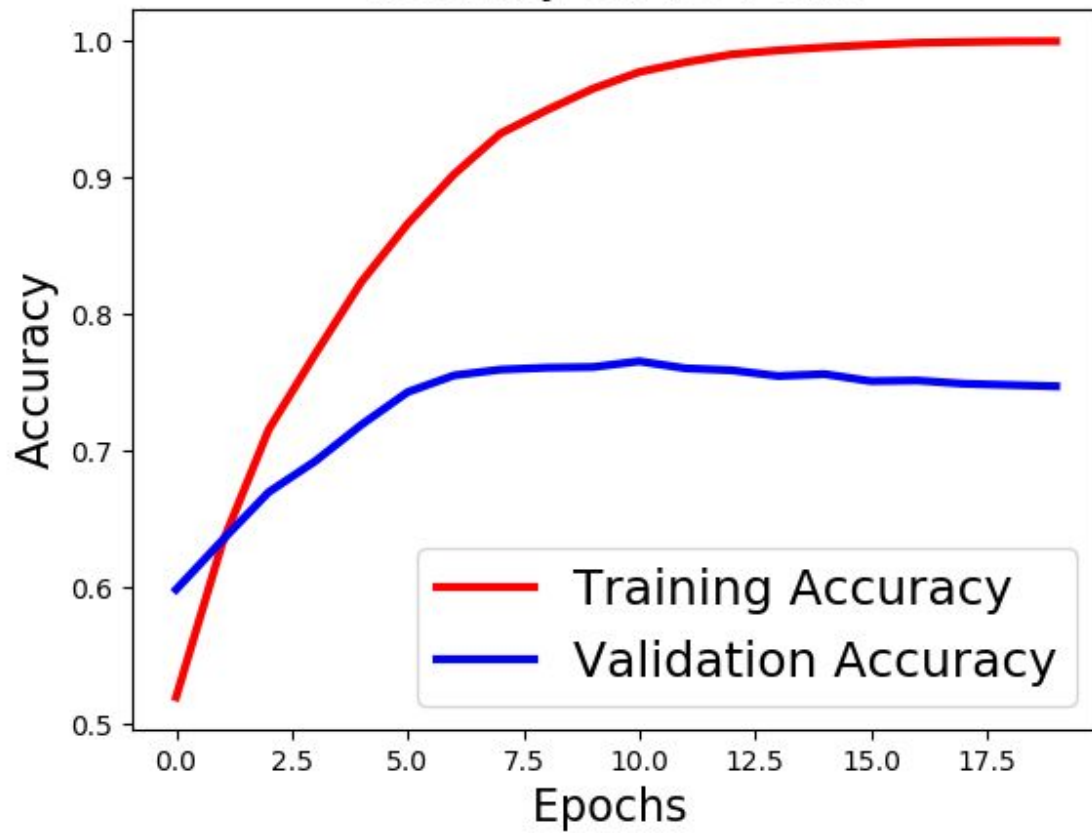
Comparing the LSTM and CNN Model Using above models.

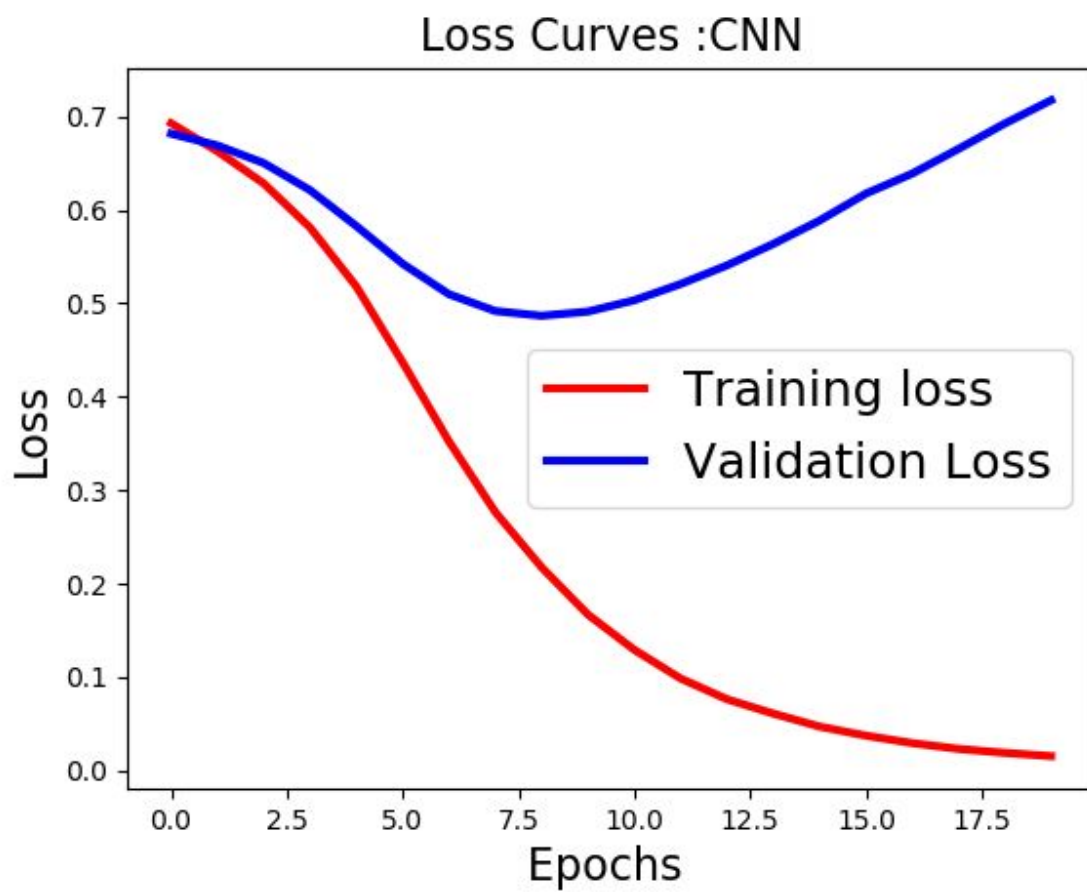
Comparison Results:

From the above models we got below Accuracy and Loss plots.

CNN LOSS and ACCURACY

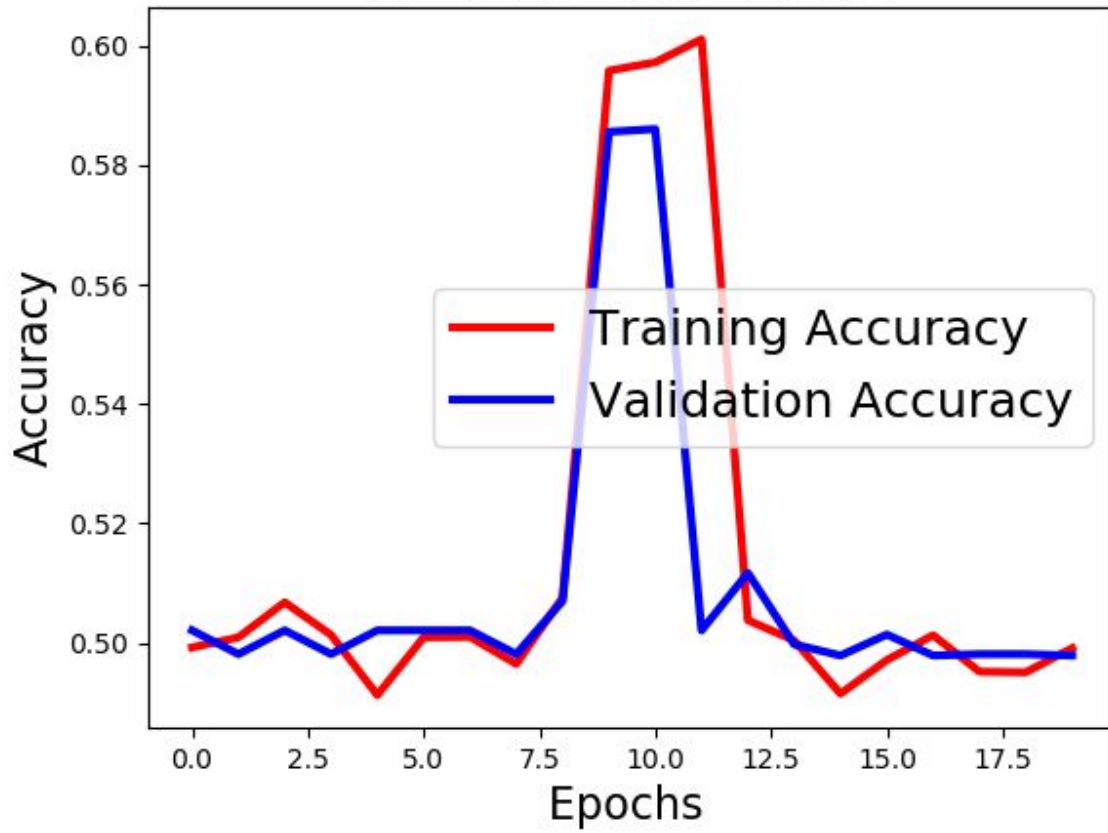
Accuracy Curves : CNN

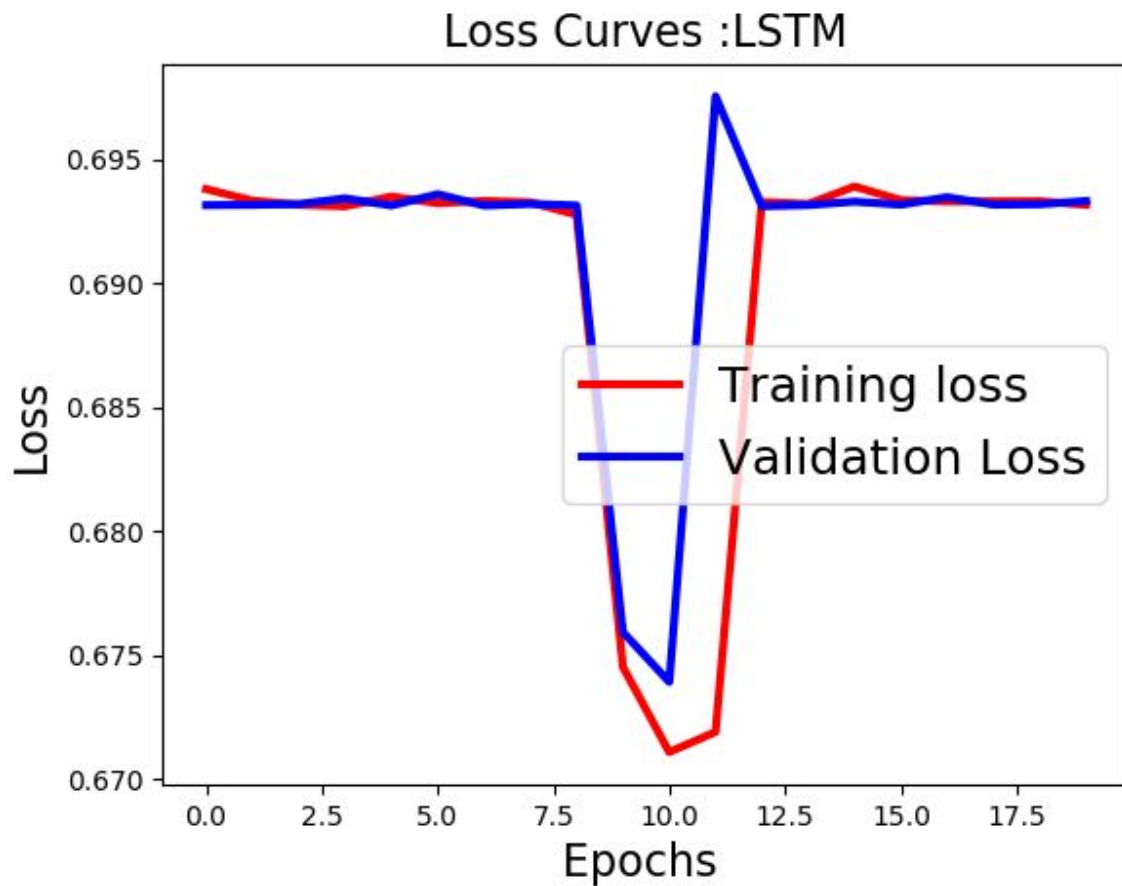




LSTM LOSS AND ACCURACY

Accuracy Curves : LSTM





Observations:

- CNN has achieved good validation accuracy with high consistency , LSTM gave average performance but not consistent as CNN.
- LSTM model architecture is not for deployment in production .
- CNN outperformed LSTM in training time.

Objective:

To apply CNN Model on Image Dataset.

Approaches:

For CNN model training and testing. I had used the small CFIR-10 dataset. CNN Model take input image for classification and follow it definitions.

We used the keras and Tensorflow as backend . Loaded the inbuilt CIFR-10 and created the CNN model with multiple Conv2D layer, Maxpooling and dropout layer.

Compile the model with loss categorical_crossentropy and optimizer rmsprop. Fitted the model on tain and test data for training. Calculated score and accuracy plotted same on graph.

Workflow:

First of all we have loaded the cifr-10 .The data, split between train and test sets,converts the vectors into binary class matrices.

```
image_classification_cnn.py x
3 from __future__ import print_function
4 import keras
5 from keras.datasets import cifar10
6
7 from keras.models import Sequential
8 from keras.layers import Dense, Dropout, Activation, Flatten
9 from keras.layers import Conv2D, MaxPooling2D
10 import matplotlib.pyplot as plt
11
12 batch_size = 32
13 num_classes = 10
14 epochs = 15
15 num_predictions = 20
16
17
18 # The data, split between train and test sets:
19 (x_train, y_train), (x_test, y_test) = cifar10.load_data()
20 print('x_train shape:', x_train.shape)
21 print(x_train.shape[0], 'train samples')
22 print(x_test.shape[0], 'test samples')
23
24 # Convert class vectors to binary class matrices.
25 y_train = keras.utils.to_categorical(y_train, num_classes)
26 y_test = keras.utils.to_categorical(y_test, num_classes)
27
```

Created and Compiled the model.

image_classification_cnn.py x

```
27
28 model = Sequential()
29 model.add(Conv2D(32, (3, 3), padding='same',
30                 input_shape=x_train.shape[1:]))
31 model.add(Activation('relu'))
32 model.add(Conv2D(32, (3, 3)))
33 model.add(Activation('relu'))
34 model.add(MaxPooling2D(pool_size=(2, 2)))
35 model.add(Dropout(0.25))
36
37 model.add(Conv2D(64, (3, 3), padding='same'))
38 model.add(Activation('relu'))
39 model.add(Conv2D(64, (3, 3)))
40 model.add(Activation('relu'))
41 model.add(MaxPooling2D(pool_size=(2, 2)))
42 model.add(Dropout(0.25))
43
44 model.add(Flatten())
45 model.add(Dense(512))
46 model.add(Activation('relu'))
47 model.add(Dropout(0.5))
48 model.add(Dense(num_classes))
49 model.add(Activation('softmax'))
50
51 # initiate RMSprop optimizer
52 opt = keras.optimizers.rmsprop(lr=0.0001, decay=1e-6)
53
54 # Let's train the model using RMSprop
55 model.compile(loss='categorical_crossentropy',
56              optimizer=opt,
57              metrics=['accuracy'])
58
59 x_train = x_train.astype('float32')
60 x_test = x_test.astype('float32')
61 x_train /= 255
62 x_test /= 255
63
64 # Model Summary
65 model.summary()
```

Fitted the model on train and test. Calculated the score, accuracy and plot the graph for loss and accuracy.

```
# Model Fitting
history = model.fit(x_train, y_train, batch_size=batch_size, epochs=epochs, validation_data=(x_test, y_test), shuffle=True)

# Score trained model.
scores = model.evaluate(x_test, y_test, verbose=1)
print('Test loss:', scores[0])
print('Test accuracy:', scores[1])

# Plot The CIFR-10 CNN LOSS
fig1 = plt.figure()
plt.plot(history.history['loss'], 'r', linewidth=3.0)
plt.plot(history.history['val_loss'], 'b', linewidth=3.0)
plt.legend(['Training loss', 'Validation Loss'], fontsize=18)
plt.xlabel('Epochs ', fontsize=16)
plt.ylabel('Loss', fontsize=16)
plt.title('Loss Curves : CIFR-10 CNN', fontsize=16)
fig1.savefig('cifr_loss_lstm.png')

# Plot the CIFR-10 CNN Accuracy
fig2=plt.figure()
plt.plot(history.history['acc'], 'r', linewidth=3.0)
plt.plot(history.history['val_acc'], 'b', linewidth=3.0)
plt.legend(['Training Accuracy', 'Validation Accuracy'], fontsize=18)
plt.xlabel('Epochs ', fontsize=16)
plt.ylabel('Accuracy', fontsize=16)
plt.title('Accuracy Curves : CIFR-10 CNN', fontsize=16)
fig2.savefig('cifr_accuracy_lstm.png')
plt.show()
```

This Concludes the workflow.

Dataset:

Used inbuilt dataset from

<https://www.cs.toronto.edu/~kriz/cifar.html>

Parameters:

None

Evaluation:

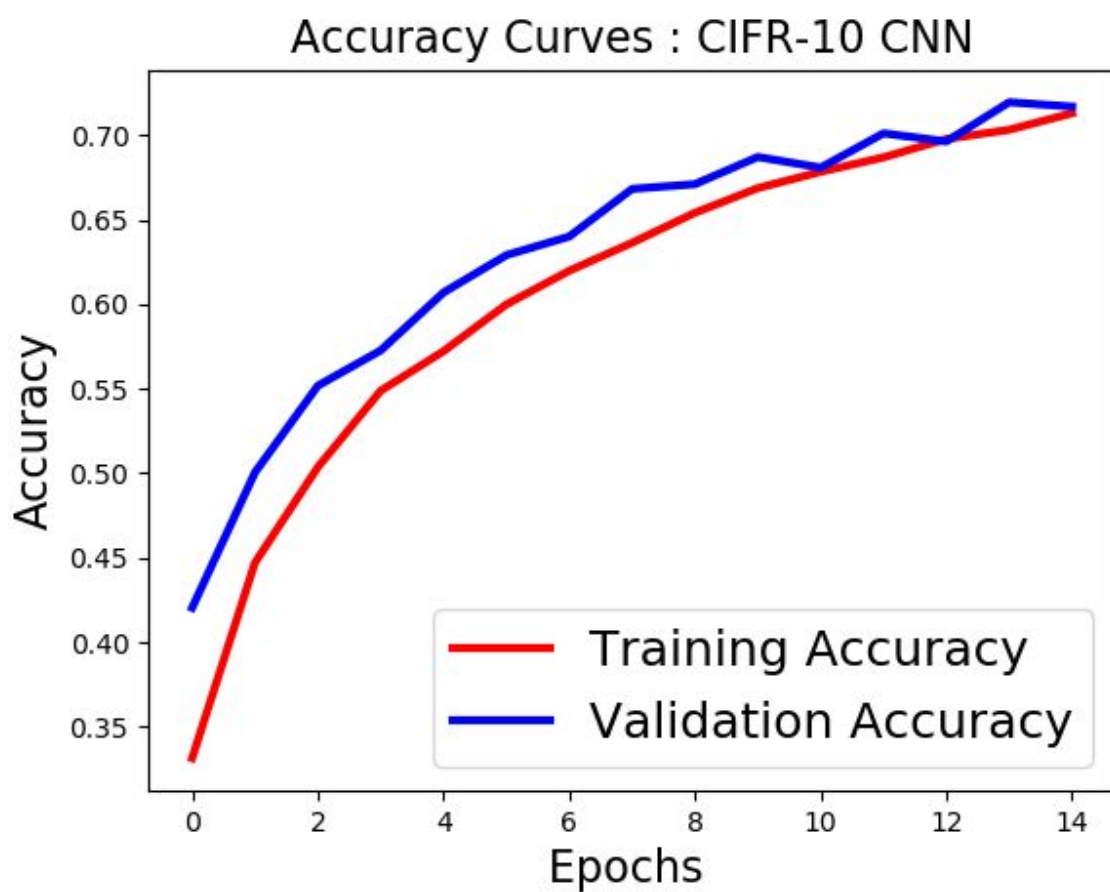
With 15 Epoch and batch size 32 we will give accuracy of 71 percent. We can improve the accuracy with more epochs.



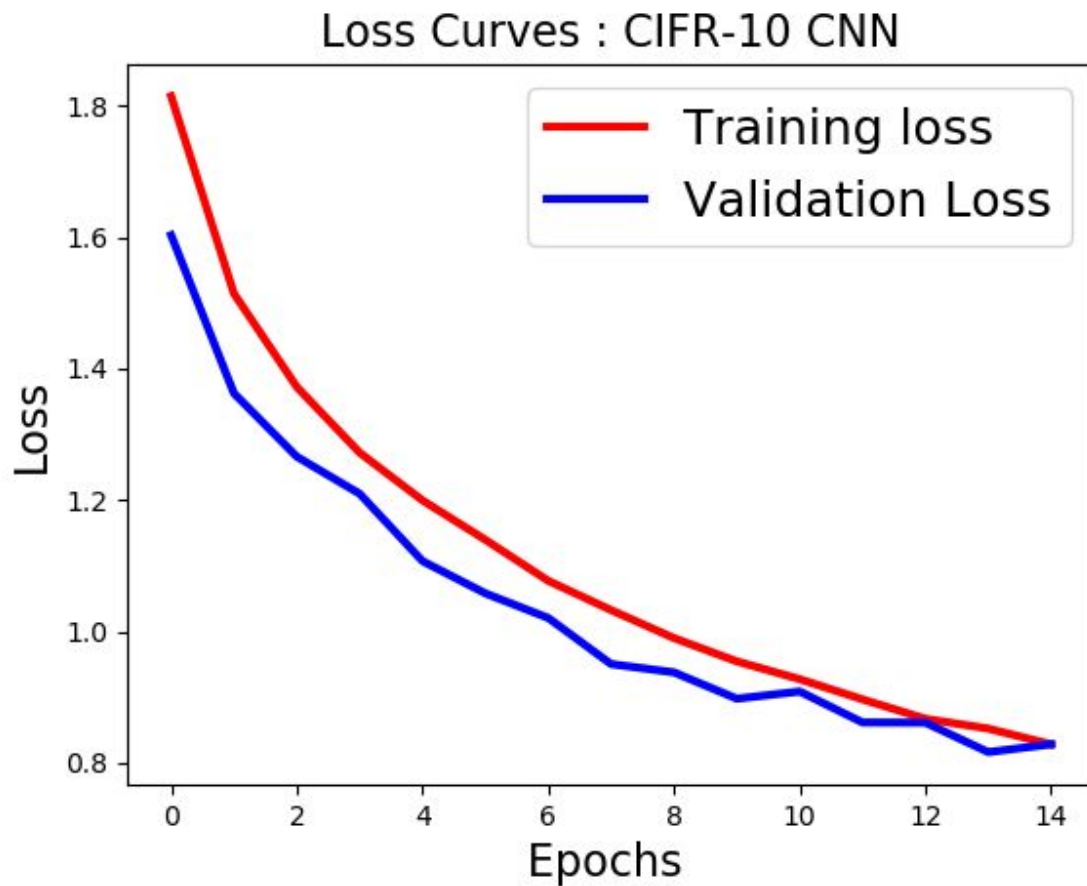
Conclusion:

We can see from below loss and Accuracy graph CNN performs very well on image datasets.

Accuracy:



Loss



References:

<https://medium.com/@ksusorokina/image-classification-with-convolutional-neural-networks-496815db12a8>

<https://towardsdatascience.com/recurrent-neural-networks-and-lstm-4b601dd822a5>

<https://medium.com/@sabber/classifying-yelp-review-comments-using-lstm-and-word-embeddings-part-1-eb2275e4066b>