

# DL-Lab 2 Report

VINAYCHANDRA VASAMSETTI, Class ID: 52

## Introduction

In class, we learned about the properties of convolutional neural networks (CNN) and how one can be used to create a text classifier that can determine which of two classes a given line of text belongs to. This binary classification of text can be very useful in the real world, such as when companies wish to classify tweets regarding the company as being either positive or negative. Knowing how useful binary text classification can be, the next question must be whether classifying text into one of many categories can be just as useful. The purpose of this lab is to explore the idea of multi-class classification using a CNN and its implementation in Tensorflow.

## Objectives

The objective of this lab is to explore multi-class classification using convolutional neural networks in Tensorflow to determine which of several categories a given string of text belongs to. In order to meet our objective we will use a data consisting of text files grouped together by subject to train our network. Once the network has been trained we will see whether it is able to correctly predict which subject a never before seen text file belongs to, based on its features.

## Approaches

To achieve our objective, the in-class CNN example will be improved upon to accommodate not only multiple classes, but a corpus of files for each class as well. Allowing more than one file to be used for training a class will let us increase the data set size more easily. The design of the CNN used will remain the same as that of the in-class example so that we may compare the results of the two approaches. In future attempts we may wish to tune some parameters to achieve higher accuracy.

## Workflow

### Process Data

A data set consisting of 5 or more classes needs to be obtained.

The in-class example must be altered to accept more than two classes and more than one file per class.

Each file must be cleaned and labeled before being fed to the CNN.

### Train Network

The text input and labels are to be fed to the network in batches.

The words from the input are embedded to create feature vectors which will get applied to weights in the network and the weights will be updated over many iterations to determine the most important features.

## Evaluate

A portion of the data set is used to test the accuracy of the network after a specified number of steps.

The network predicts the correct class of given feature vectors which is then compared to the correct label.

## Data Set

[20 newsgroups](#) was used as the data set and originally contained 20 categories with 1000 articles within each category. Related categories were condensed down to 7 total categories and only approximately 10 articles were kept for each category to ensure the model could be trained in a reasonable amount of time.

## Parameters

Parameters were kept the same as the in-class example:

Batch Size: 64 lines of text processed at a time

Epochs: 200 iterations over the training data

Steps Before Evaluation: 100 batches processed before calculating accuracy and loss

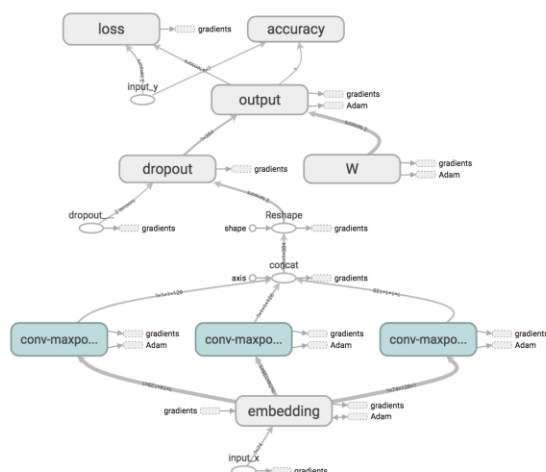
Steps Before Checkpoint: 100 batches processed before writing to file

Stored Checkpoints: 5 of the last checkpoints retained

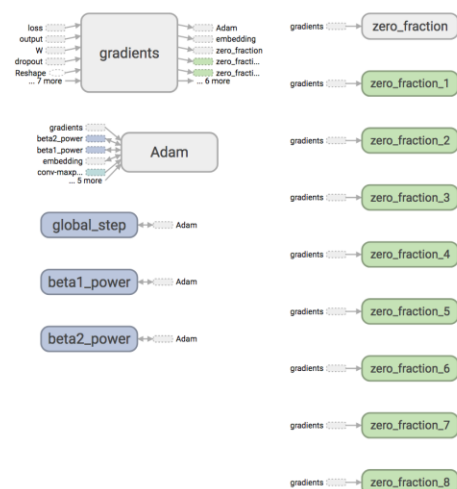
## Evaluation & Discussion

### Graph

Main Graph



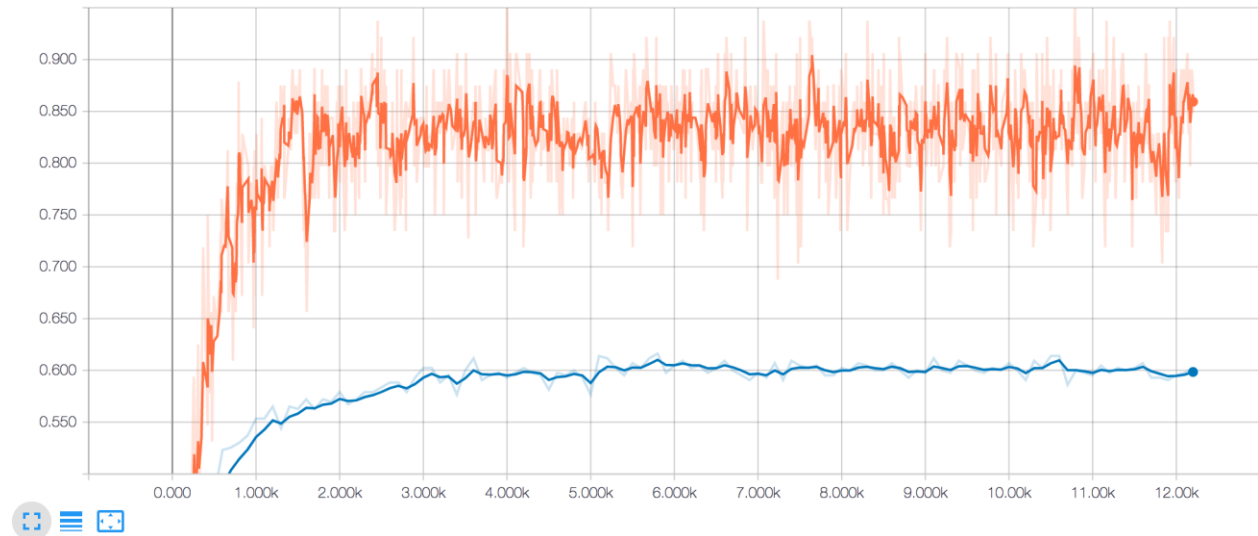
Auxiliary Nodes



*The structure of the CNN is reliable with the in-class case since our goal was to watch how well the same CNN would perform while ordering between different unmistakable classes.*

## Accuracy

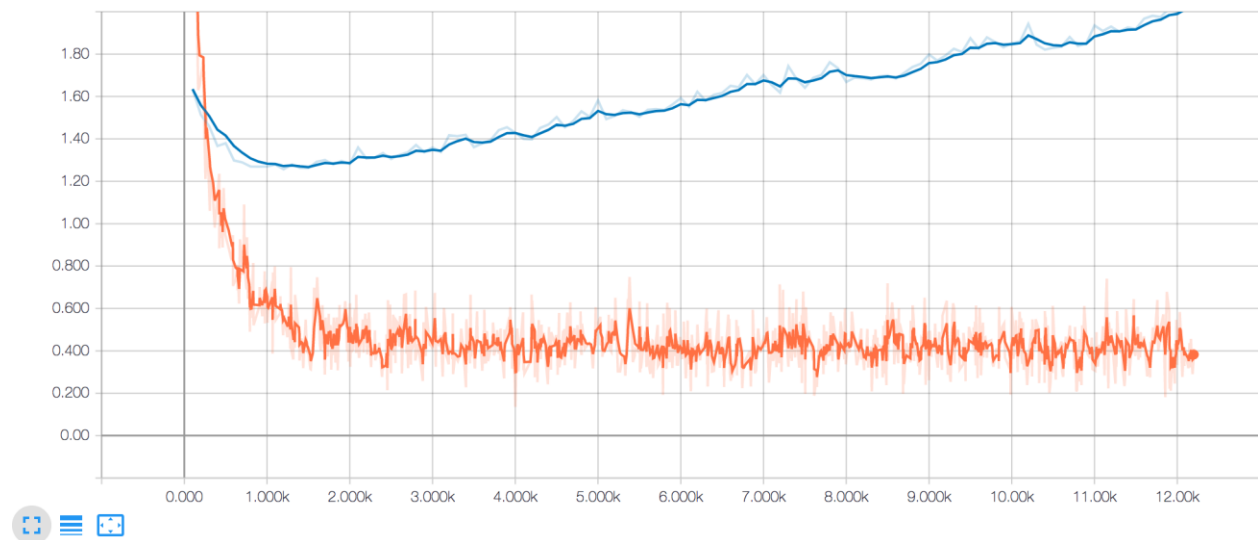
accuracy\_1



*It is normal that the test accuracy will be much lower than the training precision however our outcomes demonstrate a critical contrast between the two. The training accuracy accomplishes satisfactory rates while the test accuracy is no place close worthy.*

## Loss

loss\_1



*The training loss seems to demonstrate an alluring impact by beginning high, of course, and rapidly diminishing before leveling off to a relatively static rate. Nonetheless, the test loss demonstrates a significantly less attractive outcome. While the test loss decreases at the outset, it rapidly achieves its least point, which is still high, before starting to increment. This demonstrates our system isn't enhancing with an ever increasing number of emphasis and that a few parameters should be tuned to accomplish better outcomes.*

## Conclusion

Here we provide a detailed summary /conclusion for the assignment theory uploaded till now.

```
Train.py x DataHelper.py x TextCNN.py x
57 self.h_pool_flat = tf.reshape(self.h_pool, [-1, num_filters_total])
58
59 # Add dropout
60 with tf.name_scope("dropout"):
61     self.h_drop = tf.nn.dropout(self.h_pool_flat, self.dropout_keep_prob)
62
63 # Final (unnormalized) scores and predictions
64 with tf.name_scope("output"):
65     W = tf.get_variable(
66         "W",
67         shape=[num_filters_total, num_classes],
68         initializer=tf.contrib.layers.xavier_initializer())
69     b = tf.Variable(tf.constant(0.1, shape=[num_classes]), name="b")
70     l2_loss += tf.nn.l2_loss(W)
71     l2_loss += tf.nn.l2_loss(b)
72     self.scores = tf.nn.xw_plus_b(self.h_drop, W, b, name="scores")
73     self.predictions = tf.argmax(self.scores, 1, name="predictions")
74
75 # CalculateMean cross-entropy loss
76 with tf.name_scope("loss"):
77     losses = tf.nn.softmax_cross_entropy_with_logits(logits=self.scores, labels=self.input_y)
78     self.loss = tf.reduce_mean(losses) + l2_reg_lambda * l2_loss
79
80 # Accuracy
81 with tf.name_scope("accuracy"):
82     correct_predictions = tf.equal(self.predictions, tf.argmax(self.input_y, 1))
83     self.accuracy = tf.reduce_mean(tf.cast(correct_predictions, "float"), name="accuracy")
84
```

While the in-class example was successfully altered to accommodate multi-class classification, using the previous parameter and hyper-parameter values, and a relatively small data set, seems to result in a text classifier which is only slightly better than random guessing (~ 60% accuracy). Because there are more files available from the original data set that were not used in order to keep training time to a minimum, the easiest way to improve accuracy would be to add those files back into the data set used. However, before this is done, it would be wise to run several more experiments with different variations of parameter values on the reduced data set, to conclude the optimal tuning without wasting effort training the larger data set. In sum, our objective was met by showing that a binary text classifier could be used as a base for a multi-class text classifier. In future work, the focus will shift to reducing the loss and improving the accuracy of the classifier.

## CODE FOR IMPLENTATION :

```
Train.py x
1 import datetime
2 import os
3 import time
4
5 import numpy as np
6 import tensorflow as tf
7 from tensorflow.contrib import learn
8
9 import DataHelper
10 from TextCNN import TextCNN
11
12 # Parameters
13 # =====
14
15 # Data loading params
16 tf.flags.DEFINE_float("dev_sample_percentage", .1, "Percentage of the training data to use for validation")
17 tf.flags.DEFINE_string("alt_data", "./newsgroups/alt", "Data source for articles about atheism.")
18 tf.flags.DEFINE_string("comp_data", "./newsgroups/comp", "Data source for computer articles.")
19 tf.flags.DEFINE_string("misc_data", "./newsgroups/misc", "Data source for sales ads.")
20 tf.flags.DEFINE_string("rec_data", "./newsgroups/rec", "Data source for sports articles.")
21 tf.flags.DEFINE_string("sci_data", "./newsgroups/sci", "Data source for scientific articles.")
22 tf.flags.DEFINE_string("soc_data", "./newsgroups/soc", "Data source for religious articles.")
23 tf.flags.DEFINE_string("talk_data", "./newsgroups/talk", "Data source for open forums.")
24
25 # Model Hyperparameters
26 tf.flags.DEFINE_integer("embedding_dim", 128, "Dimensionality of character embedding (default: 128)")
27 tf.flags.DEFINE_string("filter_sizes", "3,4,5", "Comma-separated filter sizes (default: '3,4,5')")
28 tf.flags.DEFINE_integer("num_filters", 128, "Number of filters per filter size (default: 128)")
29 tf.flags.DEFINE_float("dropout_keep_prob", 0.5, "Dropout keep probability (default: 0.5)")
30 tf.flags.DEFINE_float("l2_reg_lambda", 0.0, "L2 regularization lambda (default: 0.0)")
31
32 # Training parameters
33
34 # Training parameters
35 tf.flags.DEFINE_integer("batch_size", 64, "Batch Size (default: 64)")
36 tf.flags.DEFINE_integer("num_epochs", 200, "Number of training epochs (default: 200)")
37 tf.flags.DEFINE_integer("evaluate_every", 100, "Evaluate model on dev set after this many steps (default: 100)")
38 tf.flags.DEFINE_integer("checkpoint_every", 100, "Save model after this many steps (default: 100)")
39 tf.flags.DEFINE_integer("num_checkpoints", 5, "Number of checkpoints to store (default: 5)")
40
41 # Misc Parameters
42 tf.flags.DEFINE_boolean("allow_soft_placement", True, "Allow device soft device placement")
43 tf.flags.DEFINE_boolean("log_device_placement", False, "Log placement of ops on devices")
44
45 FLAGS = tf.flags.FLAGS
46 FLAGS.parse_flags()
47 print("\nParameters:")
48 for attr, value in sorted(FLAGS.__flags.items()):
49     print("{}={}".format(attr.upper(), value))
50 print("")
51
52 # Data Preparation
53
54 # Load data
55 print("Loading data...")
56 data_flags = [FLAGS.alt_data, FLAGS.comp_data, FLAGS.misc_data, FLAGS.rec_data,
57               FLAGS.sci_data, FLAGS.soc_data, FLAGS.talk_data]
58 x_text, y = DataHelper.load_data_and_labels(data_flags)
59
60 # Build vocabulary
61 max_document_length = max([len(x.split(" ")) for x in x_text])
62 vocab_processor = learn.preprocessing.VocabularyProcessor(max_document_length)
63 x = np.array(list(vocab_processor.fit_transform(x_text)))
```

```

Train.py x
64
65     # Randomly shuffle data
66     np.random.seed(10)
67     shuffle_indices = np.random.permutation(np.arange(len(y)))
68     x_shuffled = x[shuffle_indices]
69     y_shuffled = y[shuffle_indices]
70
71     # Split train/test set
72     # TODO: This is very crude, should use cross-validation
73     dev_sample_index = -1 * int(FLAGS.dev_sample_percentage * float(len(y)))
74     x_train, x_dev = x_shuffled[:dev_sample_index], x_shuffled[dev_sample_index:]
75     y_train, y_dev = y_shuffled[:dev_sample_index], y_shuffled[dev_sample_index:]
76     print(x_train.shape)
77     print(y_train.shape)
78     print("Vocabulary Size: {:d}".format(len(vocab_processor.vocabulary_)))
79     print("Train/Dev split: {:d}/{:d}".format(len(y_train), len(y_dev)))
80
81     # Training
82     # =====
83
84     with tf.Graph().as_default():
85         session_conf = tf.ConfigProto(
86             allow_soft_placement=FLAGS.allow_soft_placement,
87             log_device_placement=FLAGS.log_device_placement)
88         sess = tf.Session(config=session_conf)
89         with sess.as_default():
90             cnn = TextCNN(
91                 sequence_length=x_train.shape[1],
92                 num_classes=y_train.shape[1],
93                 vocab_size=len(vocab_processor.vocabulary_),
94                 embedding_size=FLAGS.embedding_dim,
95                 filter_sizes=list(map(int, FLAGS.filter_sizes.split(","))),
96                 num_filters=FLAGS.num_filters,

```

```

Train.py x
99     # Define Training procedure
100     global_step = tf.Variable(0, name="global_step", trainable=False)
101     optimizer = tf.train.AdamOptimizer(1e-1)
102     grads_and_vars = optimizer.compute_gradients(cnn.loss)
103     train_op = optimizer.apply_gradients(grads_and_vars, global_step=global_step)
104
105     # Keep track of gradient values and sparsity (optional)
106     grad_summaries = []
107     for g, v in grads_and_vars:
108         if g is not None:
109             grad_hist_summary = tf.summary.histogram("{}grad/hist".format(v.name), g)
110             sparsity_summary = tf.summary.scalar("{}grad/sparsity".format(v.name), tf.nn.zero_fraction(g))
111             grad_summaries.append(grad_hist_summary)
112             grad_summaries.append(sparsity_summary)
113     grad_summaries_merged = tf.summary.merge(grad_summaries)
114
115     # Output directory for models and summaries
116     timestamp = str(int(time.time()))
117     out_dir = os.path.abspath(os.path.join(os.path.curdir, "runs", timestamp))
118     print("Writing to {}".format(out_dir))
119
120     # Summaries for loss and accuracy
121     loss_summary = tf.summary.scalar("loss", cnn.loss)
122     acc_summary = tf.summary.scalar("accuracy", cnn.accuracy)
123
124     # Train Summaries
125     train_summary_op = tf.summary.merge([loss_summary, acc_summary, grad_summaries_merged])
126     train_summary_dir = os.path.join(out_dir, "summaries", "train")
127     train_summary_writer = tf.summary.FileWriter(train_summary_dir, sess.graph)
128

```



```

Train.py x
129     # Dev summaries
130     dev_summary_op = tf.summary.merge([loss_summary, acc_summary])
131     dev_summary_dir = os.path.join(out_dir, "summaries", "dev")
132     dev_summary_writer = tf.summary.FileWriter(dev_summary_dir, sess.graph)
133
134     # Checkpoint directory. Tensorflow assumes this directory already exists so we need to create it
135     checkpoint_dir = os.path.abspath(os.path.join(out_dir, "checkpoints"))
136     checkpoint_prefix = os.path.join(checkpoint_dir, "model")
137     if not os.path.exists(checkpoint_dir):
138         os.makedirs(checkpoint_dir)
139     saver = tf.train.Saver(tf.global_variables(), max_to_keep=FLAGS.num_checkpoints)
140
141     # Write vocabulary
142     vocab_processor.save(os.path.join(out_dir, "vocab"))
143
144     # Initialize all variables
145     sess.run(tf.global_variables_initializer())
146
147     def train_step(x_batch, y_batch):
148         """
149         A single training step
150         """
151         feed_dict = {
152             cnn.input_x: x_batch,
153             cnn.input_y: y_batch,
154             cnn.dropout_keep_prob: FLAGS.dropout_keep_prob
155         }
156         _, step, summaries, loss, accuracy = sess.run(
157             [train_op, global_step, train_summary_op, cnn.loss, cnn.accuracy],
158             feed_dict)
159         time_str = datetime.datetime.now().isoformat()

```

```

Train.py x
160         print("{}: step {}, loss {:g}, acc {:g}".format(time_str, step, loss, accuracy))
161         train_summary_writer.add_summary(summaries, step)
162
163     def dev_step(x_batch, y_batch, writer=None):
164         """
165         Evaluates model on a dev set
166         """
167         feed_dict = {
168             cnn.input_x: x_batch,
169             cnn.input_y: y_batch,
170             cnn.dropout_keep_prob: 1.0
171         }
172         step, summaries, loss, accuracy = sess.run(
173             [global_step, dev_summary_op, cnn.loss, cnn.accuracy],
174             feed_dict)
175         time_str = datetime.datetime.now().isoformat()
176         print("{}: step {}, loss {:g}, acc {:g}".format(time_str, step, loss, accuracy))
177         if writer:
178             writer.add_summary(summaries, step)
179
180     # Generate batches
181     batches = DataHelper.batch_iter(
182         list(zip(x_train, y_train)), FLAGS.batch_size, FLAGS.num_epochs)

```

## TEXT CNN CODE :

```
1 import tensorflow as tf
2
3
4 class TextCNN(object):
5     """
6     A CNN for text classification.
7     Uses an embedding layer, followed by a convolutional, max-pooling and softmax layer.
8     """
9     def __init__(self, sequence_length, num_classes, vocab_size,
10                  embedding_size, filter_sizes, num_filters, l2_reg_lambda=0.0):
11
12         # Placeholders for input, output and dropout
13         self.input_x = tf.placeholder(tf.int32, [None, sequence_length], name="input_x")
14         self.input_y = tf.placeholder(tf.float32, [None, num_classes], name="input_y")
15         self.dropout_keep_prob = tf.placeholder(tf.float32, name="dropout_keep_prob")
16
17         # Keeping track of l2 regularization loss (optional)
18         l2_loss = tf.constant(0.0)
19
20         # Embedding layer
21         with tf.device('/cpu:0'), tf.name_scope("embedding"):
22             self.W = tf.Variable(
23                 tf.random_uniform([vocab_size, embedding_size], -1.0, 1.0),
24                 name="W")
25             self.embedded_chars = tf.nn.embedding_lookup(self.W, self.input_x)
26             self.embedded_chars_expanded = tf.expand_dims(self.embedded_chars, -1)
27
28         # Create a convolution + maxpool layer for each filter size
29         pooled_outputs = []
30         for i, filter_size in enumerate(filter_sizes):
31             with tf.name_scope("conv-maxpool-%s" % filter_size):
32                 # Convolution Layer
33                 filter_shape = [filter_size, embedding_size, 1, 1]
```



```
Train.py x DataHelper.py x TextCNN.py x
33 # Convolution Layer
34 filter_shape = [filter_size, embedding_size, 1, num_filters]
35 W = tf.Variable(tf.truncated_normal(filter_shape, stddev=0.1), name="W")
36 b = tf.Variable(tf.constant(0.1, shape=[num_filters]), name="b")
37 conv = tf.nn.conv2d(
38     self.embedded_chars_expanded,
39     W,
40     strides=[1, 1, 1, 1],
41     padding="VALID",
42     name="conv")
43 # Apply nonlinearity
44 h = tf.nn.relu(tf.nn.bias_add(conv, b), name="relu")
45 # Max pooling over the outputs
46 pooled = tf.nn.max_pool(
47     h,
48     ksize=[1, sequence_length - filter_size + 1, 1, 1],
49     strides=[1, 1, 1, 1],
50     padding='VALID',
51     name="pool")
52 pooled_outputs.append(pooled)
53
54 # Combine all the pooled features
55 num_filters_total = num_filters * len(filter_sizes)
56 self.h_pool = tf.concat(pooled_outputs, 3)
57 self.h_pool_flat = tf.reshape(self.h_pool, [-1, num_filters_total])
58
59 # Add dropout
60 with tf.name_scope("dropout"):
61     self.h_drop = tf.nn.dropout(self.h_pool_flat, self.dropout_keep_prob)
62
63 self.h_pool_flat = tf.reshape(self.h_pool, [-1, num_filters_total])
64
65 # Add dropout
66 with tf.name_scope("dropout"):
67     self.h_drop = tf.nn.dropout(self.h_pool_flat, self.dropout_keep_prob)
68
69 # Final (unnormalized) scores and predictions
70 with tf.name_scope("output"):
71     W = tf.get_variable(
72         "W",
73         shape=[num_filters_total, num_classes],
74         initializer=tf.contrib.layers.xavier_initializer())
75     b = tf.Variable(tf.constant(0.1, shape=[num_classes]), name="b")
76     l2_loss += tf.nn.l2_loss(W)
77     l2_loss += tf.nn.l2_loss(b)
78     self.scores = tf.nn.xw_plus_b(self.h_drop, W, b, name="scores")
79     self.predictions = tf.argmax(self.scores, 1, name="predictions")
80
81 # Calculate Mean cross-entropy loss
82 with tf.name_scope("loss"):
83     losses = tf.nn.softmax_cross_entropy_with_logits(logits=self.scores, labels=self.input_y)
84     self.loss = tf.reduce_mean(losses) + l2_reg_lambda * l2_loss
85
86 # Accuracy
87 with tf.name_scope("accuracy"):
88     correct_predictions = tf.equal(self.predictions, tf.argmax(self.input_y, 1))
89     self.accuracy = tf.reduce_mean(tf.cast(correct_predictions, "float"), name="accuracy")
```