ECE471, Selected Topics in Machine Learning – Assignment 5
AG News Dataset Classifier
Nikola Janjušević
October 9, 2018

# Result

Test Set Accuracy: 90.32%

# Program:

```python
1   import os
2   import pandas
3   import collections
4   import numpy as np
5   from tqdm import tqdm
6   import tensorflow as tf
7
8   # for tensorboard
9   logs_path = "./tf_logs/"
10  # disable tensorflow warnings
11  os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2'
12
13  num_classes = 4
14  embedding_dim = 32
15  VOCAB_SIZE = 10000
16
17  # hyper-parameters
18  BATCH_SIZE = 128
19  NUM_EPOCHS = 2
20  display_epoch = 1
21  LEARNING_RATE = .01
22  # regularization parameters
23  drop_prob = 0.25
24  reg_scale = 1e-6
25
26  def main():
27
28      class Data(object):
29          def __init__(self):
30              self.data_train, self.y_train, all_words \
31                  = get_data_csv('ag_news_csv/train.csv')
32              data_test, y_test, _ = get_data_csv('ag_news_csv/test.csv')
33              # split test into test and validation set
34              [(self.data_test, self.data_val), (self.y_test, self.y_val)] \
35                  = [np.split(xy,2) for xy in [data_test, y_test]]
36
37              # build dictionaries of test words, num limited to vocab_size
38              self.dictionary, self.reversed_dictionary \
39                  = build_dataset(all_words, VOCAB_SIZE)
40
41              self.train_size = self.y_train.shape[0]
42
43              # encodes strings by dictionary number
44              [(self.x_train,longest_str), (self.x_val,_), (self.x_test,_)] \
45                  = [code_data(x, self.dictionary, max_len=200) \
46                          for x in  \
47                              [self.data_train, self.data_val, self.data_test] ]
48              # maximum sentence length
49              self.max_len = longest_str
50
51          def get_batch(self):
52              choices = np.random.choice(self.train_size, size=BATCH_SIZE)
53              return self.x_train[choices], self.y_train[choices,:]
54
55      print("constructing dataset...")
56      data = Data()
57      print("done.")
58
```

```
59      x = tf.placeholder(tf.int32, shape=[None, data.max_len])
60      y = tf.placeholder(tf.int32, shape=[None, num_classes])
61      embeddings = tf.Variable(
62          tf.random_uniform([VOCAB_SIZE, embedding_dim], -1.0, 1.0)
63      )
64      phase = tf.placeholder(tf.bool) # is_training
65      lr = tf.placeholder(tf.float32) # learning rate [not used in this program]
66
67      def f(x):
68          x = tf.nn.embedding_lookup(embeddings, x)
69          x = tf.layers.flatten(x)
70          x = fc_layer(x, 8, is_training=phase)
71          x = tf.layers.dense(x, num_classes)
72          return x
73
74      # models
75      logits = f(x)
76      prediction = tf.nn.softmax(logits)
77
78      correct_pred = tf.equal(tf.argmax(prediction, 1), tf.argmax(y, 1))
79      accuracy = tf.reduce_mean(tf.cast(correct_pred, tf.float32))
80
81      # LOSS
82      update_ops = tf.get_collection(tf.GraphKeys.UPDATE_OPS)
83      with tf.control_dependencies(update_ops):
84          loss = tf.reduce_mean( tf.losses.softmax_cross_entropy(y, logits) ) \
85              + tf.reduce_mean( tf.losses.get_regularization_loss() )
86          optim = tf.train.AdamOptimizer(learning_rate=lr).minimize(loss)
87
88      init = tf.global_variables_initializer()
89
90      # Create a summary to monitor cost tensor
91      tf.summary.scalar("loss", loss)
92      # Create a summary to monitor accuracy tensor
93      tf.summary.scalar("accuracy", accuracy)
94      # Merge all summaries into a single op
95      merged_summary_op = tf.summary.merge_all()
96
97      with tf.Session() as sess:
98          sess.run(init)
99          # op to write logs to Tensorboard
100         summary_writer = tf.summary.FileWriter(logs_path,
101             graph=tf.get_default_graph())
102         learning_rate = LEARNING_RATE
103         # training
104         for epoch in range(NUM_EPOCHS):
105             avg_cost = 0.
106             num_batches = int( np.ceil( data.train_size / BATCH_SIZE ) )
107
108             for i in tqdm(range(num_batches)):
109                 xb, yb = data.get_batch()
110                 fd = {x: xb, y: yb, phase: True, lr: learning_rate}
111                 loss_np, _, summary \
112                     = sess.run([loss, optim, merged_summary_op], feed_dict=fd)
113                 # logs every batch
114                 summary_writer.add_summary(summary, epoch * num_batches + i)
115                 avg_cost  += loss_np/num_batches
116
117             # Display logs per epoch step
118             if (epoch+1) % display_epoch == 0:
119                 print("Epoch:", '%02d' % (epoch+1),
120                     "cost=", "{:.6f}".format(avg_cost))
121
122             print('Validation Set Accuracy:',
123                 accuracy.eval({x: data.x_val, y: data.y_val, phase: False}))
124
125         # Test the model on separate data
126         print('Test Set Accuracy:',
127             accuracy.eval({x: data.x_test, y: data.y_test, phase: False}))
128
129         print("Run the command line:\n--> tensorboard --logdir=./tf_logs ")
130
131  # ------------------------------------------------
```

```
132  # -------------- MODEL FUNCTIONS -----------------
133  # ------------------------------------------------
134
135  # --- CONV LAYER WRAPPER --- w/ L2 regularization
136  def conv_layer(input, filters, kernel_size, strides=2, is_training=True):
137      x = tf.layers.conv2d(
138          input, filters, kernel_size, strides=strides, padding='same',
139          kernel_regularizer=tf.contrib.layers.l2_regularizer(scale=reg_scale)
140      )
141      return x
142
143  # --- FULLY CONNECTED LAYER WRAPPER ---
144  # matmul -> BN -> relu -> dropout
145  def fc_layer(input, units, is_training=True):
146      x = tf.layers.dense(
147          input, units,
148          kernel_regularizer=tf.contrib.layers.l2_regularizer(scale=reg_scale)
149      )
150      x = tf.layers.batch_normalization(x, training=is_training, renorm=True)
151      x = tf.nn.relu6(x)
152      x = tf.layers.dropout(x, rate=drop_prob, training=is_training)
153      return x
154
155  # --------------------------------------------
156  # ---------- DATA BUILDING FUNCTIONS ----------
157  # --------------------------------------------
158
159  # https://stackoverflow.com/questions/38592324/one-hot-encoding-using-numpy/38592416
160  def get_one_hot(targets, nb_classes):
161      res = np.eye(nb_classes)[np.array(targets).reshape(-1)]
162      return res.reshape(list(targets.shape)+[nb_classes])
163
164  def get_data_csv(pathname):
165      # data from csv
166      df = pandas.read_csv(pathname, index_col=False, \
167          header=None, names=['label', 'headline', 'description'],
168          quotechar='"', doublequote=True, lineterminator='\n')
169      # joining headline and description into one string
170      # https://stackoverflow.com/questions/39571832/how-to-row-wise-concatenate-several-
                  columns-containing-strings
171      df['cat'] = df[df.columns[1:]].apply(tuple, axis=1).str.join(' ')
172      # puts all the data into one list
173      all_words = "".join(df['cat'].tolist()).lower().split()
174      # shuffling
175      s = np.arange(len(df['label']))
176      np.random.shuffle(s)
177      x = np.array(df['cat'])[s]
178      y = get_one_hot(np.array(df['label']) - 1, num_classes)[s,:]
179      return x, y, all_words
180
181  # makes dictionary mapping of words to unique integer id
182  # http://adventuresinmachinelearning.com/word2vec-tutorial-tensorflow/
183  def build_dataset(words, n_words):
184      """Process raw inputs into a dataset."""
185      count = [['UNK', -1]]
186      count.extend(collections.Counter(words).most_common(n_words - 1))
187      dictionary = dict()
188      for word, _ in count:
189          dictionary[word] = len(dictionary)
190      data = list()
191      unk_count = 0
192      for word in words:
193          if word in dictionary:
194              index = dictionary[word]
195          else:
196              index = 0   # dictionary['UNK']
197              unk_count += 1
198          data.append(index)
199      count[0][1] = unk_count
200      reversed_dictionary = dict(zip(dictionary.values(), dictionary.keys()))
201      return dictionary, reversed_dictionary
202
203  # codes any string into an array of integers corresponding to the mapping
```

```
204  # provided by the dictionary. Pads / clips strings to max_len number of words
205  # if provided.
206  def code_data(data, dictionary, max_len=None):
207      coded_data = []
208      longest_str = 0
209      for i in range(data.shape[0]):
210          word_list = data[i].lower().split()
211          for j in range(len(word_list)):
212              try:
213                  key = dictionary[word_list[j]]
214              except KeyError:
215                  key = 0
216              word_list[j] = key
217          if len(word_list)>longest_str:
218              longest_str = len(word_list)
219          coded_data.append(word_list)
220
221      if max_len is not None:
222          longest_str = max_len
223
224      cd = np.zeros((data.shape[0], longest_str))
225
226      for i in range(len(coded_data)):
227          cd[i,:] = np.pad( np.asarray(coded_data[i][:longest_str]), \
228              ( 0, longest_str-len(coded_data[i][:longest_str]) ), 'constant')
229      return cd, longest_str
230
231  if __name__ == "__main__":
232      main()
```