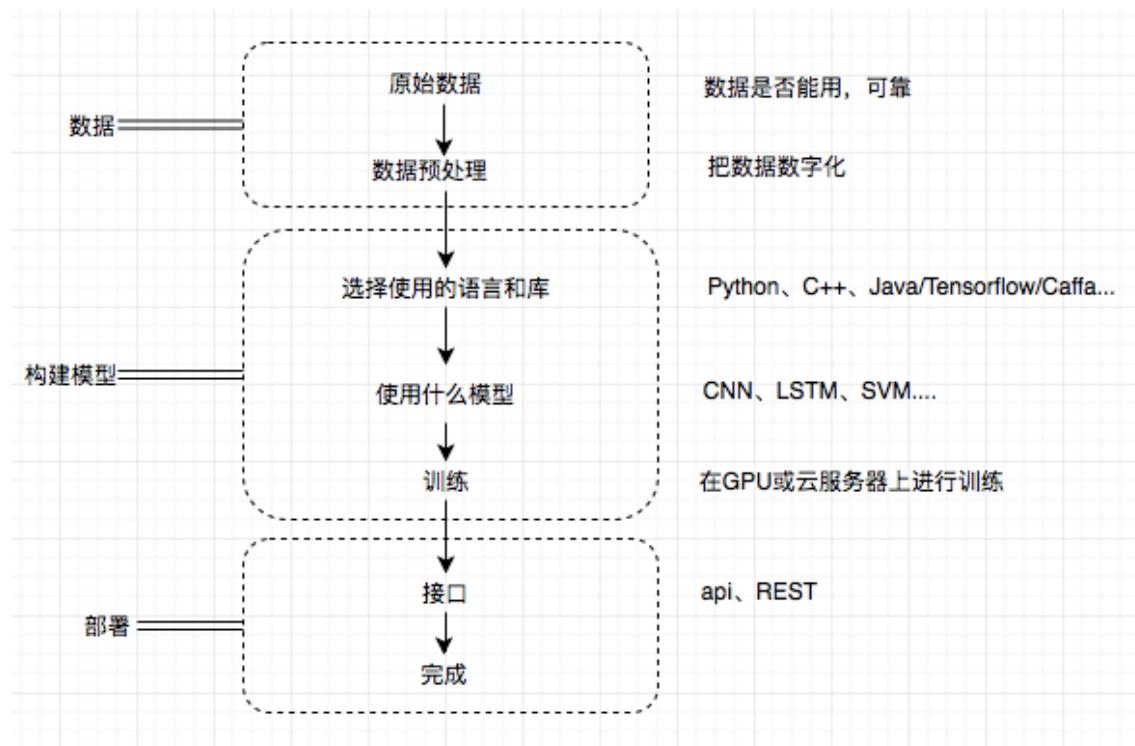# WTF Daily Blog

斗大的熊猫

# TensorFlow练习2: 对评论进行分类

本帖是前一贴的补充：

1. 使用大数据，了解怎么处理数据不能一次全部加载到内存的情况。如果你内存充足，当我没说
2. 训练好的模型的保存和使用
3. 使用的模型没变，还是简单的feedforward神经网络
4. 如果你要运行本帖代码，推荐使用GPU版本或强大的VPS，我使用小笔记本差点等吐血

在正文开始之前，我画了一个机器学习模型的基本开发流程图：

使用的数据集

使用的数据集：http://help.sentiment140.com/for-students/ (情绪分析)

数据集包含1百60万条推特，包含消极、中性和积极tweet。不知道有没有现成的微博数据集。

数据格式：移除表情符号的CSV文件，字段如下：

- 0 – the polarity of the tweet (0 = negative, 2 = neutral, 4 = positive)
- 1 – the id of the tweet (2087)
- 2 – the date of the tweet (Sat May 16 23:58:44 UTC 2009)
- 3 – the query (lyx). If there is no query, then this value is NO_QUERY.
- 4 – the user that tweeted (robotickilldozr)
- 5 – the text of the tweet (Lyx is cool)

training.1600000.processed.noemoticon.csv（238M）

testdata.manual.2009.06.14.csv（74K）

数据预处理

```
1   import nltk
2   from nltk.tokenize import word_tokenize
3   from nltk.stem import WordNetLemmatizer
4
5   import pickle
6   import numpy as np
7   import pandas as pd
8   from collections import OrderedDict
9
10
11  org_train_file = 'training.1600000.processed.noemoticon.csv'
12  org_test_file = 'testdata.manual.2009.06.14.csv'
13
14  # 提取文件中有用的字段
15  def usefull_filed(org_file, output_file):
16      output = open(output_file, 'w')
17      with open(org_file, buffering=10000, encoding='latin-1') as f:
18          try:
19              for line in f:                      # "4","2193601966","Tue Jun
20                  line = line.replace('"', '')
21                  clf = line.split(',')[0]    # 4
```

```python
22                    if clf == '0':
23                        clf = [0, 0, 1]   # 消极评论
24                    elif clf == '2':
25                        clf = [0, 1, 0]   # 中性评论
26                    elif clf == '4':
27                        clf = [1, 0, 0]   # 积极评论
28
29                    tweet = line.split(',')[-1]
30                    outputline = str(clf) + ':%:%:%:' + tweet
31                    output.write(outputline)  # [0, 0, 1]:%:%:%: that's a
32            except Exception as e:
33                print(e)
34    output.close()   # 处理完成，处理后文件大小127.5M
35
36 usefull_filed(org_train_file, 'training.csv')
37 usefull_filed(org_test_file, 'tesing.csv')
38
39 # 创建词汇表
40 def create_lexicon(train_file):
41     lex = []
42     lemmatizer = WordNetLemmatizer()
43     with open(train_file, buffering=10000, encoding='latin-1') as f:
44         try:
45             count_word = {}   # 统计单词出现次数
46             for line in f:
47                 tweet = line.split(':%:%:%:')[1]
48                 words = word_tokenize(line.lower())
49                 for word in words:
50                     word = lemmatizer.lemmatize(word)
51                     if word not in count_word:
52                         count_word[word] = 1
53                     else:
54                         count_word[word] += 1
55
56             count_word = OrderedDict(sorted(count_word.items(), key=l
57             for word in count_word:
58                 if count_word[word] < 100000 and count_word[word] > 1
59                     lex.append(word)
60         except Exception as e:
61             print(e)
62     return lex
63
64 lex = create_lexicon('training.csv')
65
66 with open('lexcion.pickle', 'wb') as f:
67     pickle.dump(lex, f)
68
69
70 """
71 # 把字符串转为向量
72 def string_to_vector(input_file, output_file, lex):
73     output_f = open(output_file, 'w')
```

```
74      lemmatizer = WordNetLemmatizer()
75      with open(input_file, buffering=10000, encoding='latin-1') as f:
76          for line in f:
77              label = line.split(':%:%:%:')[0]
78              tweet = line.split(':%:%:%:')[1]
79              words = word_tokenize(tweet.lower())
80              words = [lemmatizer.lemmatize(word) for word in words]
81
82              features = np.zeros(len(lex))
83              for word in words:
84                  if word in lex:
85                      features[lex.index(word)] = 1   # 一个句子中某个词可能
86
87              features = list(features)
88              output_f.write(str(label) + ":" + str(features) + '\n')
89      output_f.close()
90
91
92  f = open('lexcion.pickle', 'rb')
93  lex = pickle.load(f)
94  f.close()
95
96  # lexcion词汇表大小112k,training.vec大约112k*1600000  170G  太大，只能边转
97  # string_to_vector('training.csv', 'training.vec', lex)
98  # string_to_vector('tesing.csv', 'tesing.vec', lex)
99  """
```

上面代码把原始数据转为training.csv、和tesing.csv，里面只包含label和tweet。
lexcion.pickle文件保存了词汇表。

> 如果数据文件太大，不能一次加载到内存，可以把数据导入数据库
> *Dask*可处理大csv文件

开始漫长的训练

```
1   import os
2   import random
3   import tensorflow as tf
4   import pickle
5   import numpy as np
6   from nltk.tokenize import word_tokenize
7   from nltk.stem import WordNetLemmatizer
8
9   f = open('lexcion.pickle', 'rb')
10  lex = pickle.load(f)
11  f.close()
12
```

```
13
14  def get_random_line(file, point):
15      file.seek(point)
16      file.readline()
17      return file.readline()
18  # 从文件中随机选择n条记录
19  def get_n_random_line(file_name, n=150):
20      lines = []
21      file = open(file_name, encoding='latin-1')
22      total_bytes = os.stat(file_name).st_size
23      for i in range(n):
24          random_point = random.randint(0, total_bytes)
25          lines.append(get_random_line(file, random_point))
26      file.close()
27      return lines
28
29
30  def get_test_dataset(test_file):
31      with open(test_file, encoding='latin-1') as f:
32          test_x = []
33          test_y = []
34          lemmatizer = WordNetLemmatizer()
35          for line in f:
36              label = line.split(':%:%:%:')[0]
37              tweet = line.split(':%:%:%:')[1]
38              words = word_tokenize(tweet.lower())
39              words = [lemmatizer.lemmatize(word) for word in words]
40              features = np.zeros(len(lex))
41              for word in words:
42                  if word in lex:
43                      features[lex.index(word)] = 1
44
45              test_x.append(list(features))
46              test_y.append(eval(label))
47      return test_x, test_y
48
49  test_x, test_y = get_test_dataset('tesing.csv')
50
51
52  ####################################################################
53
54  n_input_layer = len(lex)   # 输入层
55
56  n_layer_1 = 2000      # hide layer
57  n_layer_2 = 2000      # hide layer(隐藏层)听着很神秘，其实就是除输入输出层外的
58
59  n_output_layer = 3        # 输出层
60
61
62  def neural_network(data):
63      # 定义第一层"神经元"的权重和biases
64      layer_1_w_b = {'w_':tf.Variable(tf.random_normal([n_input_layer,
```

```
65        # 定义第二层"神经元"的权重和biases
66        layer_2_w_b = {'w_':tf.Variable(tf.random_normal([n_layer_1, n_l
67        # 定义输出层"神经元"的权重和biases
68        layer_output_w_b = {'w_':tf.Variable(tf.random_normal([n_layer_2
69
70        # w·x+b
71        layer_1 = tf.add(tf.matmul(data, layer_1_w_b['w_']), layer_1_w_b
72        layer_1 = tf.nn.relu(layer_1)  # 激活函数
73        layer_2 = tf.add(tf.matmul(layer_1, layer_2_w_b['w_']), layer_2_
74        layer_2 = tf.nn.relu(layer_2 ) # 激活函数
75        layer_output = tf.add(tf.matmul(layer_2, layer_output_w_b['w_'])
76
77        return layer_output
78
79
80 X = tf.placeholder('float')
81 Y = tf.placeholder('float')
82 batch_size = 90
83
84 def train_neural_network(X, Y):
85        predict = neural_network(X)
86        cost_func = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logi
87        optimizer = tf.train.AdamOptimizer().minimize(cost_func)
88
89        with tf.Session() as session:
90            session.run(tf.initialize_all_variables())
91
92            lemmatizer = WordNetLemmatizer()
93            saver = tf.train.Saver()
94            i = 0
95            pre_accuracy = 0
96            while True:   # 一直训练
97                batch_x = []
98                batch_y = []
99
100               #if model.ckpt文件已存在:
101               #    saver.restore(session, 'model.ckpt')  恢复保存的sessic
102
103               try:
104                   lines = get_n_random_line('training.csv', batch_size
105                   for line in lines:
106                       label = line.split(':%:%:%:')[0]
107                       tweet = line.split(':%:%:%:')[1]
108                       words = word_tokenize(tweet.lower())
109                       words = [lemmatizer.lemmatize(word) for word in
110
111                       features = np.zeros(len(lex))
112                       for word in words:
113                           if word in lex:
114                               features[lex.index(word)] = 1  # 一个句子
115
116                       batch_x.append(list(features))
```

```
117                        batch_y.append(eval(label))
118
119                    session.run([optimizer, cost_func], feed_dict={X:bat
120                except Exception as e:
121                    print(e)
122
123                # 准确率
124                if i > 100:
125                    correct = tf.equal(tf.argmax(predict,1), tf.argmax(Y
126                    accuracy = tf.reduce_mean(tf.cast(correct,'float'))
127                    accuracy = accuracy.eval({X:test_x, Y:test_y})
128                    if accuracy > pre_accuracy:  # 保存准确率最高的训练模型
129                        print('准确率: ', accuracy)
130                        pre_accuracy = accuracy
131                        saver.save(session, 'model.ckpt')  # 保存session
132                    i = 0
133            i += 1
134
135
136 train_neural_network(X,Y)
```

上面程序占用内存600M，峰值1G。


运行：

```
(env) tianshuais-MacBook-Air:tf tianshuai$ python training.py
准确率:   0.459839
准确率:   0.477912
准确率:   0.493976
准确率:   0.495984
准确率:   0.508032
准确率:   0.51004
准确率:   0.516064
准确率:   0.522088
准确率:   0.526104
准确率:   0.528112
准确率:   0.532129
```

训练模型保存为model.ckpt。


使用训练好的模型

```
1  import tensorflow as tf
2  import pickle
3  from nltk.tokenize import word_tokenize
4  from nltk.stem import WordNetLemmatizer
5  import numpy as np
6
```

```
 7  f = open('lexcion.pickle', 'rb')
 8  lex = pickle.load(f)
 9  f.close()
10
11
12  n_input_layer = len(lex)  # 输入层
13
14  n_layer_1 = 2000      # hide layer
15  n_layer_2 = 2000      # hide layer(隐藏层)听着很神秘，其实就是除输入输出层外的
16
17  n_output_layer = 3        # 输出层
18  def neural_network(data):
19      # 定义第一层"神经元"的权重和biases
20      layer_1_w_b = {'w_':tf.Variable(tf.random_normal([n_input_layer,
21      # 定义第二层"神经元"的权重和biases
22      layer_2_w_b = {'w_':tf.Variable(tf.random_normal([n_layer_1, n_la
23      # 定义输出层"神经元"的权重和biases
24      layer_output_w_b = {'w_':tf.Variable(tf.random_normal([n_layer_2,
25
26      # w·x+b
27      layer_1 = tf.add(tf.matmul(data, layer_1_w_b['w_']), layer_1_w_b[
28      layer_1 = tf.nn.relu(layer_1)  # 激活函数
29      layer_2 = tf.add(tf.matmul(layer_1, layer_2_w_b['w_']), layer_2_w
30      layer_2 = tf.nn.relu(layer_2 ) # 激活函数
31      layer_output = tf.add(tf.matmul(layer_2, layer_output_w_b['w_']),
32
33      return layer_output
34
35  X = tf.placeholder('float')
36  def prediction(tweet_text):
37      predict = neural_network(X)
38
39      with tf.Session() as session:
40          session.run(tf.initialize_all_variables())
41          saver = tf.train.Saver()
42          saver.restore(session, 'model.ckpt')
43
44          lemmatizer = WordNetLemmatizer()
45          words = word_tokenize(tweet_text.lower())
46          words = [lemmatizer.lemmatize(word) for word in words]
47
48          features = np.zeros(len(lex))
49          for word in words:
50              if word in lex:
51                  features[lex.index(word)] = 1
52
53          #print(predict.eval(feed_dict={X:[features]})) [[val1,val2,va
54          res = session.run(tf.argmax(predict.eval(feed_dict={X:[featur
55          return res
56
57
58  prediction("I am very happe")
```

f  Facebook      G+  Google+      🐦  Twitter      🔴  Weibo      ✉  Email

## 相关文章

Ubuntu 16.04 安装 Tensorflow(GPU支持)

TensorFlow练习1: 对评论进行分类

TensorFlow练习3: RNN, Recurrent Neural Networks

TensorFlow练习4: CNN, Convolutional Neural Networks...

TensorFlow练习5: 训练一个简单的游戏AI（Deep Q Network）...

用微信 OR 支付宝 扫描二维码

为本文作者 打个赏

微信扫一扫转账

用支付宝扫二维码，加我好友

📅 2016年11月16日    👤 wtf    📁 ML、coding    🏷 TensorFlow、教程