

WTF Daily Blog

斗大的熊猫

TensorFlow练习15: 中文语音识别

语音识别的应用领域非常广泛，洋文名Speech Recognition。它所解决的问题是让计算机能够“听懂”人类的语音，将语音中包含的文字信息“提取”出来。

语音识别是前文《[聊天机器人](#)》必不可少的一个组件，本帖就使用TensorFlow做一个中文语音识别。

使用的数据集

[THCHS30](#)是Dong Wang, Xuwei Zhang, Zhiyong Zhang这几位大神发布的开放语音数据集，可用于开发中文语音识别系统。

为了感谢这几位大神，我是跪在电脑前写的本帖代码。

下载中文语音数据集（5G+）：

```
1 $ wget http://data.csalt.org/thchs30/zip/wav.tgz
2 $ wget http://data.csalt.org/thchs30/zip/doc.tgz
3 $ wget http://data.csalt.org/thchs30/zip/lm.tgz
4 # 解压
```

```
5 $ tar xvf wav.tgz
6 $ tar xvf doc.tgz
7 $ tar xvf lm.tgz
```

在开始之前，先好好检视一下数据集。

训练

```
1 import tensorflow as tf # 0.12
2 import numpy as np
3 import os
4 from collections import Counter
5 import librosa # https://github.com/librosa/librosa
6
7 # 训练样本路径
8 wav_path = 'data/wav/train'
9 label_file = 'data/doc/trans/train.word.txt'
10
11 # 获得训练用的wav文件路径列表
12 def get_wav_files(wav_path=wav_path):
13     wav_files = []
14     for (dirpath, dirnames, filenames) in os.walk(wav_path):
15         for filename in filenames:
16             if filename.endswith('.wav') or filename.endswith('.WAV'):
17                 filename_path = os.sep.join([dirpath, filename])
18                 if os.stat(filename_path).st_size < 240000: # 剔除掉一些小文件
19                     continue
20                 wav_files.append(filename_path)
21     return wav_files
22
23 wav_files = get_wav_files()
24
25 # 读取wav文件对应的label
26 def get_wav_label(wav_files=wav_files, label_file=label_file):
27     labels_dict = {}
28     with open(label_file, 'r') as f:
29         for label in f:
30             label = label.strip('\n')
31             label_id = label.split(' ', 1)[0]
```

```
32         label_text = label.split(' ', 1)[1]
33         labels_dict[label_id] = label_text
34
35     labels = []
36     new_wav_files = []
37     for wav_file in wav_files:
38         wav_id = os.path.basename(wav_file).split('.')[0]
39         if wav_id in labels_dict:
40             labels.append(labels_dict[wav_id])
41             new_wav_files.append(wav_file)
42
43     return new_wav_files, labels
44
45 wav_files, labels = get_wav_label()
46 print("样本数:", len(wav_files)) # 8911
47 #print(wav_files[0], labels[0])
48 # wav/train/A11/A11_0.WAV -> 绿 是 阳春 烟 景 大块 文章 的 底色 四月 的 林 峦 更是 绿 得 鲜活 秀媚 诗意 盎然
49
50 # 词汇表(参看练习1和7)
51 all_words = []
52 for label in labels:
53     all_words += [word for word in label]
54 counter = Counter(all_words)
55 count_pairs = sorted(counter.items(), key=lambda x: -x[1])
56
57 words, _ = zip(*count_pairs)
58 words_size = len(words)
59 print('词汇表大小:', words_size)
60
61 word_num_map = dict(zip(words, range(len(words))))
62 to_num = lambda word: word_num_map.get(word, len(words))
63 labels_vector = [list(map(to_num, label)) for label in labels]
64 #print(wav_files[0], labels_vector[0])
65 #wav/train/A11/A11_0.WAV -> [479, 0, 7, 0, 138, 268, 0, 222, 0, 714, 0, 23, 261, 0, 28, 1191, 0, 1, 0, 442,
66 #print(words[479]) #绿
67 label_max_len = np.max([len(label) for label in labels_vector])
68 print('最长句子的字数:', label_max_len)
69
70 wav_max_len = 0 # 673
71 for wav in wav_files:
72     wav, sr = librosa.load(wav, mono=True)
73     mfcc = np.transpose(librosa.feature.mfcc(wav, sr), [1,0])
```

```

74     if len(mfcc) > wav_max_len:
75         wav_max_len = len(mfcc)
76     print("最长的语音:", wav_max_len)
77
78     batch_size = 16
79     n_batch = len(wav_files) // batch_size
80
81     # 获得一个batch
82     pointer = 0
83     def get_next_batches(batch_size):
84         global pointer
85         batches_wavs = []
86         batches_labels = []
87         for i in range(batch_size):
88             wav, sr = librosa.load(wav_files[pointer], mono=True)
89             mfcc = np.transpose(librosa.feature.mfcc(wav, sr), [1,0])
90             batches_wavs.append(mfcc.tolist())
91             batches_labels.append(labels_vector[pointer])
92             pointer += 1
93
94         # 补零对齐
95         for mfcc in batches_wavs:
96             while len(mfcc) < wav_max_len:
97                 mfcc.append([0]*20)
98         for label in batches_labels:
99             while len(label) < label_max_len:
100                 label.append(0)
101         return batches_wavs, batches_labels
102
103     X = tf.placeholder(dtype=tf.float32, shape=[batch_size, None, 20])
104     sequence_len = tf.reduce_sum(tf.cast(tf.not_equal(tf.reduce_sum(X, reduction_indices=2), 0.), tf.int32), red
105     Y = tf.placeholder(dtype=tf.int32, shape=[batch_size, None])
106
107     # conv1d_layer
108     conv1d_index = 0
109     def conv1d_layer(input_tensor, size, dim, activation, scale, bias):
110         global conv1d_index
111         with tf.variable_scope('conv1d_' + str(conv1d_index)):
112             W = tf.get_variable('W', (size, input_tensor.get_shape().as_list()[-1], dim), dtype=tf.float32, init
113             if bias:
114                 b = tf.get_variable('b', [dim], dtype=tf.float32, initializer=tf.constant_initializer(0))
115             out = tf.nn.conv1d(input_tensor, W, stride=1, padding='SAME') + (b if bias else 0)

```

```

116     if not bias:
117         beta = tf.get_variable('beta', dim, dtype=tf.float32, initializer=tf.constant_initializer(0))
118         gamma = tf.get_variable('gamma', dim, dtype=tf.float32, initializer=tf.constant_initializer(1))
119         mean_running = tf.get_variable('mean', dim, dtype=tf.float32, initializer=tf.constant_initializer(0))
120         variance_running = tf.get_variable('variance', dim, dtype=tf.float32, initializer=tf.constant_initializer(1))
121         mean, variance = tf.nn.moments(out, axes=range(len(out.get_shape()) - 1))
122         def update_running_stat():
123             decay = 0.99
124             update_op = [mean_running.assign(mean_running * decay + mean * (1 - decay)), variance_running.assign(variance_running * decay + variance * (1 - decay))]
125             with tf.control_dependencies(update_op):
126                 return tf.identity(mean), tf.identity(variance)
127         m, v = tf.cond(tf.Variable(False, trainable=False, collections=[tf.GraphKeys.LOCAL_VARIABLES]), lambda: update_running_stat(), lambda: (mean, variance))
128         out = tf.nn.batch_normalization(out, m, v, beta, gamma, 1e-8)
129     if activation == 'tanh':
130         out = tf.nn.tanh(out)
131     if activation == 'sigmoid':
132         out = tf.nn.sigmoid(out)
133
134     conv1d_index += 1
135     return out
136 # aconv1d_layer
137 aconv1d_index = 0
138 def aconv1d_layer(input_tensor, size, rate, activation, scale, bias):
139     global aconv1d_index
140     with tf.variable_scope('aconv1d_' + str(aconv1d_index)):
141         shape = input_tensor.get_shape().as_list()
142         W = tf.get_variable('W', (1, size, shape[-1], shape[-1]), dtype=tf.float32, initializer=tf.random_uniform_initializer(-0.1, 0.1))
143         if bias:
144             b = tf.get_variable('b', [shape[-1]], dtype=tf.float32, initializer=tf.constant_initializer(0))
145         out = tf.nn.atrous_conv2d(tf.expand_dims(input_tensor, dim=1), W, rate=rate, padding='SAME')
146         out = tf.squeeze(out, [1])
147         if not bias:
148             beta = tf.get_variable('beta', shape[-1], dtype=tf.float32, initializer=tf.constant_initializer(0))
149             gamma = tf.get_variable('gamma', shape[-1], dtype=tf.float32, initializer=tf.constant_initializer(1))
150             mean_running = tf.get_variable('mean', shape[-1], dtype=tf.float32, initializer=tf.constant_initializer(0))
151             variance_running = tf.get_variable('variance', shape[-1], dtype=tf.float32, initializer=tf.constant_initializer(1))
152             mean, variance = tf.nn.moments(out, axes=range(len(out.get_shape()) - 1))
153             def update_running_stat():
154                 decay = 0.99
155                 update_op = [mean_running.assign(mean_running * decay + mean * (1 - decay)), variance_running.assign(variance_running * decay + variance * (1 - decay))]
156                 with tf.control_dependencies(update_op):
157                     return tf.identity(mean), tf.identity(variance)

```

```

158         m, v = tf.cond(tf.Variable(False, trainable=False, collections=[tf.GraphKeys.LOCAL_VARIABLES],
159         out = tf.nn.batch_normalization(out, m, v, beta, gamma, 1e-8)
160         if activation == 'tanh':
161             out = tf.nn.tanh(out)
162         if activation == 'sigmoid':
163             out = tf.nn.sigmoid(out)
164
165         aconv1d_index += 1
166         return out
167 # 定义神经网络
168 def speech_to_text_network(n_dim=128, n_blocks=3):
169     out = conv1d_layer(input_tensor=X, size=1, dim=n_dim, activation='tanh', scale=0.14, bias=False)
170     # skip connections
171     def residual_block(input_sensor, size, rate):
172         conv_filter = aconv1d_layer(input_sensor, size=size, rate=rate, activation='tanh', scale=0.03, b
173         conv_gate = aconv1d_layer(input_sensor, size=size, rate=rate, activation='sigmoid', scale=0.03,
174         out = conv_filter * conv_gate
175         out = conv1d_layer(out, size=1, dim=n_dim, activation='tanh', scale=0.08, bias=False)
176         return out + input_sensor, out
177     skip = 0
178     for _ in range(n_blocks):
179         for r in [1, 2, 4, 8, 16]:
180             out, s = residual_block(out, size=7, rate=r)
181             skip += s
182
183     logit = conv1d_layer(skip, size=1, dim=skip.get_shape().as_list()[-1], activation='tanh', scale=0.08, bi
184     logit = conv1d_layer(logit, size=1, dim=words_size, activation=None, scale=0.04, bias=True)
185
186     return logit
187
188 class MaxPropOptimizer(tf.train.Optimizer):
189     def __init__(self, learning_rate=0.001, beta2=0.999, use_locking=False, name="MaxProp"):
190         super(MaxPropOptimizer, self).__init__(use_locking, name)
191         self._lr = learning_rate
192         self._beta2 = beta2
193         self._lr_t = None
194         self._beta2_t = None
195     def _prepare(self):
196         self._lr_t = tf.convert_to_tensor(self._lr, name="learning_rate")
197         self._beta2_t = tf.convert_to_tensor(self._beta2, name="beta2")
198     def _create_slots(self, var_list):
199         for v in var_list:

```

```

200         self._zeros_slot(v, "m", self._name)
201     def _apply_dense(self, grad, var):
202         lr_t = tf.cast(self._lr_t, var.dtype.base_dtype)
203         beta2_t = tf.cast(self._beta2_t, var.dtype.base_dtype)
204         if var.dtype.base_dtype == tf.float16:
205             eps = 1e-7
206         else:
207             eps = 1e-8
208         m = self.get_slot(var, "m")
209         m_t = m.assign(tf.maximum(beta2_t * m + eps, tf.abs(grad)))
210         g_t = grad / m_t
211         var_update = tf.assign_sub(var, lr_t * g_t)
212         return tf.group(*[var_update, m_t])
213     def _apply_sparse(self, grad, var):
214         return self._apply_dense(grad, var)
215
216 def train_speech_to_text_network():
217     logit = speech_to_text_network()
218
219     # CTC loss
220     indices = tf.where(tf.not_equal(tf.cast(Y, tf.float32), 0.))
221     target = tf.SparseTensor(indices=indices, values=tf.gather_nd(Y, indices) - 1, shape=tf.cast(tf.shape(Y), tf.int32))
222     loss = tf.nn.ctc_loss(logit, target, sequence_len, time_major=False)
223     # optimizer
224     lr = tf.Variable(0.001, dtype=tf.float32, trainable=False)
225     optimizer = MaxPropOptimizer(learning_rate=lr, beta2=0.99)
226     var_list = [t for t in tf.trainable_variables()]
227     gradient = optimizer.compute_gradients(loss, var_list=var_list)
228     optimizer_op = optimizer.apply_gradients(gradient)
229
230     with tf.Session() as sess:
231         sess.run(tf.global_variables_initializer())
232
233         saver = tf.train.Saver(tf.global_variables())
234
235         for epoch in range(16):
236             sess.run(tf.assign(lr, 0.001 * (0.97 ** epoch)))
237
238             global pointer
239             pointer = 0
240             for batch in range(n_batch):
241                 batches_wavs, batches_labels = get_next_batches(batch_size)

```

```
242         train_loss, _ = sess.run([loss, optimizer_op], feed_dict={X: batches_wavs, Y: batches_labels})
243         print(epoch, batch, train_loss)
244     if epoch % 5 == 0:
245         saver.save(sess, 'speech.module', global_step=epoch)
246
247 # 训练
248 train_speech_to_text_network()
249
250 # 语音识别
251 # 把batch_size改为1
252 def speech_to_text(wav_file):
253     wav, sr = librosa.load(wav_file, mono=True)
254     mfcc = np.transpose(np.expand_dims(librosa.feature.mfcc(wav, sr), axis=0), [0,2,1])
255
256     logit = speech_to_text_network()
257
258     saver = tf.train.Saver()
259     with tf.Session() as sess:
260         saver.restore(sess, tf.train.latest_checkpoint('.'))
261
262         decoded = tf.transpose(logit, perm=[1, 0, 2])
263         decoded, _ = tf.nn.ctc_beam_search_decoder(decoded, sequence_len, merge_repeated=False)
264         predict = tf.sparse_to_dense(decoded[0].indices, decoded[0].shape, decoded[0].values) + 1
265         output = sess.run(decoded, feed_dict={X: mfcc})
266         #print(output)
```

后续：从麦克风获得语音输入，使用上面的模型进行识别。

相关资源：

- [TensorFlow练习8: 基于RNN生成音乐](#)
- [Machine Learning is Fun Part 6: How to do Speech Recognition with Deep Learning](#)
- 深度学习大牛Andrew Ng：[Speech Recognition and Beyond](#)
- <https://github.com/kaldi-asr/kaldi>
- <http://cmusphinx.sourceforge.net>
- <https://pypi.python.org/pypi/SpeechRecognition>

 Facebook  Google+  Twitter  Weibo  Email

相关文章





[Ubuntu 16.04 安装 Tensorflow\(GPU支持\)](#)

[使用Python实现神经网络](#)


[TensorFlow练习1: 对评论进行分类](#)

[TensorFlow练习2: 对评论进行分类](#)

[TensorFlow练习4: CNN, Convolutional Neural Networks...](#)

 2016年12月10日  wtf  ML、coding  TensorFlow、中文语音识别

《TensorFlow练习15: 中文语音识别》有20个想法

 zerozzl

2017年6月8日 下午3:27

大神你好，解码的时候，得到的logit是有数值的，但是解码的结果，是一个空值，请问是哪里出错了吗？

[SparseTensorValue(indices=array([], shape=(0, 2), dtype=int64), values=array([], dtype=int64), dense_shape=array([1, 0]))]



hc

2017年6月9日 上午10:37

我和你同样的问题，我的words大小是2269，python2.7,tensorflow0.12



catmonkey

2017年4月20日 下午12:31

训练结束后，预测的时候报错，可否指点下？

Failed precondition: Attempting to use uninitialized value conv1d_18/W



connor

2017年4月18日 下午4:43

我自己笔记本连着训练了几天，还没训练完，结果笔记本死机了，估计前面那些是白训练了。后续还有的出吗？



李勇

2017年3月5日 下午7:47


all_words += [word for word in label] 这里是不有bug label需要split



hc

2017年6月9日 上午10:37

python版本的原因，你用的是python2.7吧


 **acans**

2017年1月17日 下午6:02

后续：从麦克风获得语音输入，使用上面的模型进行识别。

这个后续什么时候能出？

楼主能给个如何使用这个模型来输出的例子吗？

 **wtf** 

2017年1月17日 下午6:11

后续纯粹唬人，我这水平太水，不要被我坑了，速速弃坑

 **hc**

2017年6月9日 上午10:38

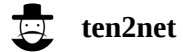
大神你好，解码的时候，得到的logit是有数值的，但是解码的结果，是一个空值，请问是哪里出错了么？

```
[SparseTensorValue(indices=array([], shape=(0, 2), dtype=int64), values=array([], dtype=int64), dense_shape=array([1, 0]))]
```

 **zhanglihai**

2017年1月10日 上午10:10

遇到了wav.zip 文件下载失败的问题。尝试5次每次都在36%失败。



2017年1月10日 下午12:31

用迅雷下载后传到服务器上解压。我就是这么做的



2017年1月9日 下午4:27

我用wav/train/A6/A6_111.wav来预测。结果返回.....

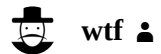
...

```
values=array([ 0, 0, 38, 0, 5, 38, 0, 5, 38, 0, 5, 38, 0, 5, 38, 0, 5,
 9, 0, 5, 9, 0, 5, 15, 0, 5, 9, 0, 16, 0, 6, 0, 5, 6,
 0, 5, 0, 5, 36, 0, 14, 15, 0, 14, 15, 0, 5, 9, 0, 5, 35,
 0, 5, 0, 5, 6, 0, 6, 0, 5, 9, 0, 9, 0, 0, 6, 0, 6,
 15, 0, 6, 15, 0, 5, 15, 0, 14, 9, 0, 14, 15, 0, 0])
```

...

使用如下代码转换，结果却是乱码，请指教哪儿不对：

```
python
msg="".join([words[n] for n in output[0][1]])
print msg
```



2017年1月9日 下午5:48

调整一下坐姿，再来一次



zjun

2017年6月5日 下午1:09

为啥我返回的是空啊。。。您遇到过这种情况么，



meng_xi

2017年1月18日 下午9:14

能大概解释下吗，请问这个是什么意思呢？您是怎么做的呢？



cdjauto

2017年3月14日 上午11:45

我识别的时候都是得到这个结果：

```
[SparseTensorValue(indices=array([[0, 0],
```

```
[0, 1],
```

```
[0, 2],
```

```
[0, 3],
```

```
[0, 4],
```

```
[0, 5],
```

```
[0, 6],
```

```
[0, 7],
```

```
[0, 8],
```


```
[0, 9]]), values=array([7113, 7113, 7113, 7113, 7113, 7113, 7113, 7113, 7113, 7113]), shape=array([ 1, 10]))]
```

是我打开的姿势不对吗？求指教

 **garbo**



2017年7月7日 上午9:38

请问这套代码训练完了之后，做预测应该如何实现？

 **ten2net**


2017年1月9日 下午4:16

期待博主的后续中.....

 **wtf** 

2017年1月9日 下午4:25

何年何月才有后续...

 **ervin**

2016年12月18日 上午10:45

楼主你好！非常感谢你分享的实践经验！我在调试运行你的代码，但是发现训练时间实在太长，需要几天的时间。

- 1.请问楼主你是用什么硬件资源训练的网络？花费的时间怎样？
 - 2.还有，我的笔记本配置有限，你可否方便把你训练好的网络发给我研究一下呢？真诚感谢！
-

Copyright © 2013-2017 WTF Daily Blog | Powered by DigitalOcean