

NLP - 15 分钟搭建中文文本分类模型

原创 艾力亚尔 AINLP 2019-01-29

本文系作者投稿，作者：艾力亚尔（微博 @艾力亚尔），暴风大脑研究院研发工程师，现负责电视端的语音助手相关工作。

原文地址，可点击文末“阅读原文”直达：

https://eliyar.biz/nlp_chinese_text_classification_in_15mins/

欢迎大家投稿，AI、NLP相关即可。

文本分类是自然语言处理核心任务之一，常见用文本审核、广告过滤、情感分析、语音控制和反黄识别等NLP领域。本文主要是介绍我最近开源的极简文本分类和序列标注框架 Kashgari (<https://github.com/BrikerMan/Kashgari>)

搭建环境和数据准备

准备工作，先准备 python 环境和数据集。

- Python 3.6 环境
- THUCNews 数据集子集,链接: <https://pan.baidu.com/s/1hugrfRu> 密码: qfud

如果需要完整数据集请自行到 THUCTC：一个高效的中文文本分类工具包 下载，请遵循数据提供方的开源协议。上面的子数据集包括一下 10 个分类。

1 体育，财经，房产，家居，教育，科技，时尚，时政，游戏，娱乐



每个分类 6500 条数据。感谢 @gaussic 在使用卷积神经网络以及循环神经网络进行中文文本分类 分享。

虚拟环境中安装所有需要的依赖

```
1 pip install jieba
2 pip install kashgari
3 pip install tensorflow
```

```
4 # 如果有 GPU 则可以安装
5 pip install tensorflow-gpu
```

数据分别为格式为一样一条新闻，每一行是 `分类\t新闻内容` 格式。我们需要把新闻内容分词后作为输入喂给模型。

```
1 import tqdm
2 import jieba
3
4 def read_data_file(path):
5     lines = open(path, 'r', encoding='utf-8').read().splitline
6     x_list = []
7     y_list = []
8     for line in tqdm.tqdm(lines):
9         rows = line.split('\t')
10        if len(rows) >= 2:
11            y_list.append(rows[0])
12            x_list.append(list(jieba.cut('\t'.join(rows[1:])))
13        else:
14            print(rows)
15    return x_list, y_list
16
17 test_x, test_y = read_file('cnews/cnews.test.txt')
18 train_x, train_y = read_file('cnews/cnews.train.txt')
19 val_x, val_y = read_file('cnews/cnews.val.txt')
```

训练与验证模型

Kashgari 目前提供了三种分类模型结构 `CNNModel` `CNNLSTMModel` 和 `BLSTMModel`。我们先使用 `CNNModel`

```
1 from kashgari.tasks.classification import CNNModel
2
3 model = CNNModel()
4 model.fit(train_x, train_y, val_x, val_y, batch_size=128)
```

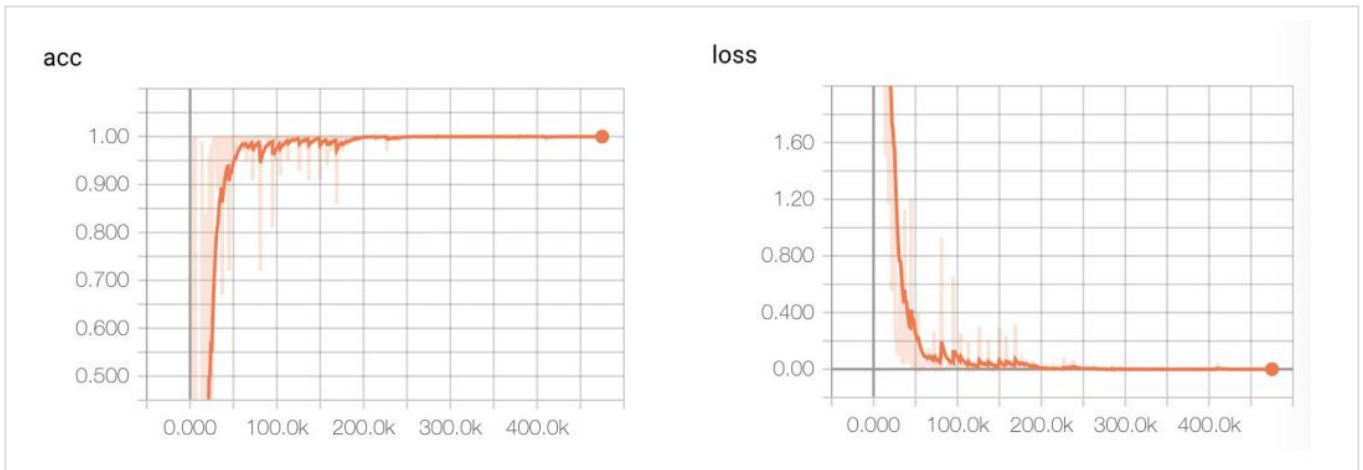
运行结果

1	Layer (type)	Output Shape	Param #
2	=====		

```

3  input_17 (InputLayer)          (None, 1462)          0
4  -----
5  embedding_17 (Embedding)       (None, 1462, 100)     19139400
6  -----
7  conv1d_9 (Conv1D)              (None, 1458, 128)     64128
8  -----
9  global_max_pooling1d_8 (Glob (None, 128)          0
10 -----
11 dense_16 (Dense)                (None, 64)            8256
12 -----
13 dense_17 (Dense)                (None, 11)            715
14 =====
15 Total params: 19,212,499
16 Trainable params: 19,212,499
17 Non-trainable params: 0
18 -----
19 Epoch 1/5
20 950/950 [=====] - 28s 30ms/step - loss: 0.6581 - acc
21 Epoch 2/5
22 950/950 [=====] - 28s 29ms/step - loss: 0.0403 - acc
23 Epoch 3/5
24 950/950 [=====] - 27s 29ms/step - loss: 0.0068 - acc
25 ...

```



由于数据特征比较明显，几轮就达到了 0.9999，val_acc 都 1.0 了。

再拿验证机验证一下模型，测试集上 F1 达到了 0.98，相当不错的成绩了。

```

1  model.evaluate(test_x, test_y)
2  """

```

		precision	recall	f1-score	support
3					
4					
5	体育	1.00	1.00	1.00	10
6	娱乐	0.99	0.99	0.99	10
7	家居	0.99	0.94	0.97	10
8	房产	1.00	1.00	1.00	10
9	教育	0.98	0.95	0.97	10
10	时尚	0.99	0.99	0.99	10
11	时政	0.96	0.97	0.97	10
12	游戏	0.99	0.99	0.99	10
13	科技	0.97	0.99	0.98	10
14	财经	0.96	1.00	0.98	10
15					
16	micro avg	0.98	0.98	0.98	10000
17	macro avg	0.98	0.98	0.98	10000
18	weighted avg	0.98	0.98	0.98	10000
19	"""				

保存模型和加载保存模型

模型的保存和重新加载都非常简单

```

1 model.save('./model')
2
3 new_model = CNNModel.load_model('./model')
4 news = """ [DeepMind 击败人类职业玩家的方式与他们声称的 AI 使命，以及所声称的『正确』
5 DeepMind 的人工智能 AlphaStar 一战成名，击败两名人类职业选手。掌声和欢呼之余，它也引起
6 """
7 x = list(jieba.cut(news))
8 new_model.predict(x)
9 '游戏'

```

使用 tensorboard 可视化训练过程

Kashgari 是基于 Keras 封装，所以可以很方便的使用 keras 的各种回调函数来记录训练过程，比如我们可以使用 `keras.callbacks.TensorBoard` 来可视化训练过程。

```

1 import keras
2

```

```
3 tf_board_callback = keras.callbacks.TensorBoard(log_dir='./logs', update_freq=
4
5
6 model = CNNModel()
7 model.fit(train_x,
8           train_y,
9           val_x,
10          val_y,
11          batch_size=100,
12          fit_kwargs={'callbacks': [tf_board_callback]})
```

在项目目录运行下面代码即可启动 tensorboard 查看可视化效果

```
1 $ tensorboard --log-dir log:
```

使用预训练词向量

由于长新闻特征比较明显，语料量也比较大，很容易取得比较不错的结果。但是如果我们的语料比较少，特征不是很明显时候直接训练可能会导致模型过拟合，泛化能力很差，此时我们可以使用预训练的词 Embedding 层来提高模型的泛化能力。

```
1 # 初始化 word2vec embedding
2 from kashgari.embeddings import WordEmbeddings
3 embedding = WordEmbeddings('<embedding-file-path>', sequence_length=6
4
5 # 初始化 BERT embedding
6 from kashgari.embeddings import BERTEmbedding
7 embedding = BERTEmbedding('bert-base-chinese', sequence_length=600)
8
9 # 使用 embedding 初始化模型
10 from kashgari.tasks.classification import CNNModel
11 model = CNNModel(embedding)
```

参考

- 使用卷积神经网络以及循环神经网络进行中文文本分类
- 中文文本分类对比（经典方法和CNN）