

国科大

自然语言处理

---

## 文本分类系统说明文档

---

姓名：张笃振

学号：201918014628019

2019 年 12 月 20 日

2019  
UCAS

# 目录

<b>1</b>	<b>小组成员信息及分工</b>	<b>2</b>
<b>2</b>	<b>任务定义</b>	<b>2</b>
<b>3</b>	<b>方法描述及实现细节</b>	<b>2</b>
3.1	FastText . . . . .	2
3.2	TextCNN . . . . .	3
3.3	RCNN . . . . .	5
3.4	多层双向 LSTM . . . . .	7
3.5	多层双向 LSTM with Attention . . . . .	9
3.6	DPCNN . . . . .	11
<b>4</b>	<b>数据集及预处理</b>	<b>14</b>
4.1	数据集简介 . . . . .	14
4.2	数据预处理 . . . . .	14
4.3	实现细节 . . . . .	15
<b>5</b>	<b>项目组织结构</b>	<b>17</b>
5.1	文件组织结构 . . . . .	17
5.2	数据预处理和加载 . . . . .	18
5.3	模型定义 . . . . .	20
5.4	配置文件 . . . . .	21
5.5	main.py . . . . .	23
<b>6</b>	<b>项目使用方式</b>	<b>29</b>
<b>7</b>	<b>实验结果与分析</b>	<b>29</b>
7.1	实验结果 . . . . .	29
7.2	实验分析 . . . . .	30
<b>8</b>	<b>预测与网页 Demo</b>	<b>30</b>
8.1	通过命令行执行预测 . . . . .	30
8.2	通过网页 Demo 执行预测 . . . . .	31
<b>9</b>	<b>所需环境</b>	<b>32</b>
<b>10</b>	<b>总结</b>	<b>32</b>

## 1 小组成员信息及分工

表 1: 小组成员信息及分工

姓名	学号	分工	备注
张笃振	201918014628019	调研 DPCNN 模型并实现；编写文档；整合各个模型，组织项目结构，训练模型	组长
彭聪	2019E8020261037	调研 FastText 模型并实现；网页 Demo 的搭建	组员
曹萌	201928013229017	调研 TextCNN 模型并实现；调研 RCNN 模型并实现	组员
库鑫	201928013229055	调研 Attention 机制；调研并实现多层双向 LSTM with Attention 模型；	组员
冯春波	201928002829001	数据预处理；调研多层双向 LSTM 模型并实现	组员

## 2 任务定义

在过去的几十年中，文本分类问题在许多实际应用中得到了广泛的研究和解决。文本分类技术在信息检索、信息过滤、推荐系统、情感分析以及推荐系统等多个领域都有广泛的应用。文本分类任务是指根据已经定义好的类别标签对现有的一段文本进行标注的任务。一般来说，文本数据集包含一系列长短不一的文本片段，如  $D = \{X_1, X_2, \dots, X_N\}$ ，其中  $X_i$  代表一个数据点或文本片段，每个  $X_i$  包含  $s$  个句子，每个句子又包含  $w_s$  个单词，每个单词又可以包含  $l_w$  个字符。每个数据点  $X_i$  都用一组  $k$  个不同离散值索引中的类别标签进行标记。

## 3 方法描述及实现细节

我们的文本分类系统实现了一些基于深度学习的分类算法，主要是一些基于 CNN 和 RNN 的模型，包括 Fast-Text、TextCNN、DPCNN、RCNN、多层双向 LSTM/GRU 以及多层双向 LSTM/GRU with Attention 这六个模型。接下来，将逐一介绍上述模型的原理和实现细节。

### 3.1 FastText

#### 1. 分类原理

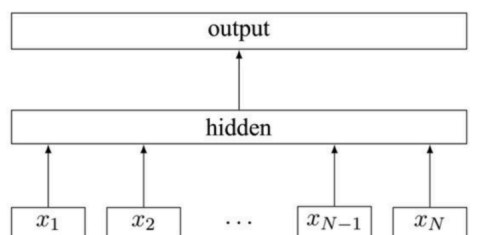


图 1: FastText 分类原理图示

输入样本是一系列整数索引  $(x_1, \dots, x_N)$ ，对应词典中相应的词，通过 embedding 得到每个词对应的词向量。对样本（文本）中每个词对应的词向量求平均，再通过一个全连接层进行分类即可。

## 2. 实现细节

```

1
2     class FastText(BasicModule): #继承自 BasicModule 其中封装了保存加载模型的接口, BasicModule继承自 nn
      .Module
3
4     def __init__(self, vocab_size,opt): #opt是 config类的实例 里面包括所有模型超参数的配置
5         super(FastText, self).__init__()
6
7
8         # 嵌入层
9         self.embedding = nn.Embedding(vocab_size, opt.embed_size) #词嵌入矩阵 每一行代表词典中一
      个词对应的词向量;
10        # 词嵌入矩阵可以随机初始化连同分类任务一起训练, 也可以用预训练词向量初始化 (冻结或微调)
11
12        self.content_fc = nn.Sequential( #可以使用多个全连接层或 batchnorm、dropout等 可以把这些
      模块用 Sequential包装成一个大模块
13        nn.Linear(opt.embed_size, opt.linear_hidden_size),
14        nn.BatchNorm1d(opt.linear_hidden_size),
15        nn.ReLU(inplace=True),
16        #可以再加一个隐层
17        # nn.Linear(opt.linear_hidden_size, opt.linear_hidden_size),
18        # nn.BatchNorm1d(opt.linear_hidden_size),
19        # nn.ReLU(inplace=True),
20        #输出层
21        nn.Linear(opt.linear_hidden_size, opt.classes)
22    )
23
24
25    def forward(self, inputs):
26        #inputs(batch_size, seq_len)
27        embeddings = self.embedding(inputs) # (batch_size, seq_len, embed_size)
28
29        #对 seq_len 维取平均
30        content = torch.mean(embeddings,dim=1) #(batch_size,1,embed_size)
31
32        out = self.content_fc(content.squeeze(1)) #先压缩 seq_len 维 (batch_size, embed_size) 然后
      作为全连接层的输入
33        #输出 (batch_size, classes)
34
35    return out

```

## 3.2 TextCNN

### 1. 分类原理

TextCNN 可以从两个角度来解读, 既可以把它看作但输入通道的 2 维卷积也可以把它看作多输入通道的 1 维卷积 (其中词嵌入维度为通道维), 二者其实是等价的。

如果把它看作一个单输入通道的 2 维卷积的话, 它的分类流程就如图2所示。

1) 把输入文本中的词转换为其对应的词向量, 那么每个输入文本就可以表示为一个  $n \times d$  的矩阵 ( $n$  是输入文本包含的词数,  $d$  为词向量的维数)。

2) 对输入矩阵进行卷积操作。可以使用不同大小的卷积核, 每种类型的卷积核可以有多个。假设卷积

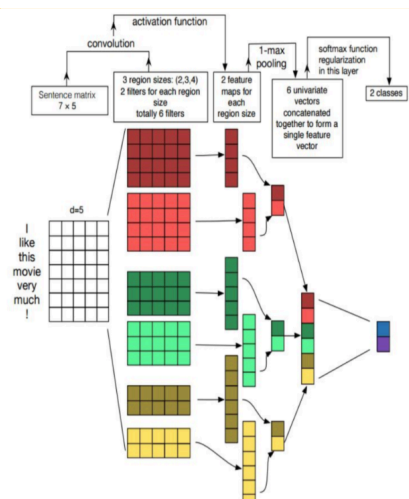


图 2: TextCNN 分类原理图示

核的大小是  $(f,d)$ ,  $f$  可以是不同的取值 (如  $f=2, 3, 4$ ), 而  $d$  是固定的, 是词向量的维度, 并且假设总共使用了  $k$  个卷积核, 步长为 1。经过卷积操作后我们会得到  $k$  个向量, 每个向量的长度是  $n-f+1$ 。我们使用不同大小的卷积核, 从输入文本中提取丰富的特征, 这和  $n$ -gram 特征有点相似 ( $f=2, 3, 4$  分别对应于 2-gram, 3-gram-4-gram)。

3) 对卷积操作的输出进行全局 max-pooling 操作。作用于  $k$  个长度为  $n-f+1$  的向量上, 每个向量整体取最大值, 得到  $k$  个标量。

4) 把  $k$  个标量拼接起来, 组成一个向量表示最后提取的特征。他的长度是固定的, 取决于我们所使用的不同大小的卷积核的总数 ( $k$ )。

5) 最后在接一个全联接层作为输出层, 如果是 2 分类的话使用 sigmoid 激活函数, 多分类则使用 softmax 激活函数, 得到模型的输出。

## 2. 实现细节

```

1  #自定义时序 (全局) 最大池化层
2  class GlobalMaxPool1d(nn.Module):
3  def __init__(self):
4  super(GlobalMaxPool1d, self).__init__()
5  def forward(self, x):
6  # x (batch_size, channel, seq_len)
7  return F.max_pool1d(x, kernel_size=x.shape[2]) # (batch_size, channel, 1)
8
9
10 # 多输入通道的一维卷积和单输入通道的2维卷积等价
11 # 这里按多输入通道的一维卷积来做 也可以用单输入通道的2维卷积来做
12 class TextCNN(BasicModule): #继承自BasicModule 其中封装了保存加载模型的接口, BasicModule继承自nn.Module
13
14     def __init__(self, vocab_size, opt):#opt是config类的实例 里面包括所有模型超参数的配置
15
16         super(TextCNN, self).__init__()
17
18
19         # 嵌入层

```

```

20     self.embedding = nn.Embedding(vocab_size, opt.embed_size) # 词嵌入矩阵 每一行代表词典中一个词对应的
    词向量;
21     # 词嵌入矩阵可以随机初始化连同分类任务一起训练, 也可以用预训练词向量初始化 (冻结或微调)
22
23     # 创建多个一维卷积层
24     self.convs = nn.ModuleList()
25     for c, k in zip(opt.num_channels, opt.kernel_sizes): # num_channels 定义了每种卷积核的个数
        kernel_sizes 定义了每种卷积核的大小
26     self.convs.append(nn.Conv1d(in_channels=opt.embed_size,
    out_channels=c,
27     kernel_size=k))
28     # 定义 dropout 层
29     self.dropout = nn.Dropout(opt.drop_prop)
30     # 定义输出层
31     self.fc = nn.Linear(sum(opt.num_channels), opt.classes)
32     # 时序最大池化层没有权重, 所以可以共用一个实例
33     self.pool = GlobalMaxPool1d()
34
35
36
37     def forward(self, inputs):
38         # inputs (batch_size, seq_len)
39         embeddings = self.embedding(inputs) # (batch_size, seq_len, embed_size)
40
41         # 根据 conv1d 的输入要求 把通道维提前 (这里的通道维是词向量维度)
42         # (batch_size, channel/embed_size, seq_len)
43         embeddings = embeddings.permute(0, 2, 1)
44         # 对于每个一维卷积层, 会得到一个 (batch_size, num_channel (卷积核的个数), seq_len - kernel_size + 1) 大
        小的 tensor
45         # 在时序最大池化后会得到一个形状为 (batch_size, num_channel, 1) 的 tensor
46         # 使用 squeeze 去掉最后一维 并在通道维上连结 得到 (batch_size, sum(num_channels)) 大小的 tensor
47         encoding = torch.cat([self.pool(F.relu(conv(embeddings))).squeeze(-1) for conv in self.convs],
        dim=1)
48
49         # 应用丢弃法后使用全连接层得到输出 (batch_size, classes)
50         outputs = self.fc(self.dropout(encoding))
51         return outputs

```

### 3.3 RCNN

#### 1. 分类原理

图3中中间是输入文本中每个单词的嵌入表示, 左右使用双向 RNN 分别学习当前词  $w_i$  的左上下文表示  $cl(w_i)$  和右上下文表示  $cr(w_i)$ , 与当前词  $w_i$  本身的词向量连接, 构成后续卷积层的输入  $x_i$ 。具体如下:

$$cl(w_i) = f(W^{(l)}cl(w_{i-1}) + W^{(sl)}e(w_{i-1})) \quad (1)$$

$$cr(w_i) = f(W^{(r)}cr(w_{i-1}) + W^{(sr)}e(w_{i-1})) \quad (2)$$

$$x_i = [cl(w_i); e(w_i); cr(w_i)] \quad (3)$$

与 TextCNN 比较类似, 都是把文本表示为一个嵌入矩阵, 再进行卷积操作。不同的是 TextCNN 中的文本嵌入矩阵每一行只是文本中一个词的向量表示, 而在 RCNN 中, 文本嵌入矩阵的每一行是当前词的词向量以及上下文嵌入表示的拼接。

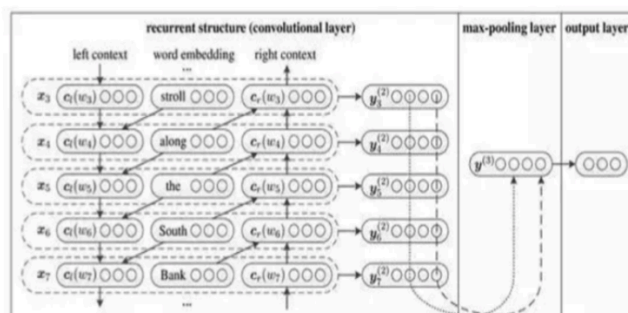


图 3: RCNN 分类原理图示

例如：图3中的 `along` 这一个单词，中间是他的词向量，左边利用一个 RNN 得到 `along` 之前上文所有单词的嵌入表示，右边同样利用一个反向的 RNN 得到 `along` 之后所有单词的嵌入表示。三者共同构成 `along` 的嵌入表示，文本中的其他词同理。最后文本被表示为上图左边的一个嵌入矩阵。

然后  $x_i$  将作为当前词  $w_i$  的嵌入表示，输入到卷积核大小为  $1*d$  ( $d$  为  $x_i$  的长度) 激活函数为  $\tanh$  的卷积层，得到  $w_i$  潜在的语义向量：

$$y_i^{(2)} = \tanh(W^{(2)}x_i + b^{(2)}) \quad (4)$$

在 TextCNN 中我们曾经设置了多个卷积核  $f*d$ , RCNN 中将卷积核大小设置为  $1*d$  的原因是中已经包含了左右上下文的信息，无需再使用窗口大于 1 的卷积核进行特征提取。需要说明的是，实践中依然可以同时使用多个不同大小的卷积核，如  $[1,2,3]$ ，可能会取得更好的实践效果，一种解释是窗口大于 1 的卷积核强化了左右最近的上下文信息。此外实践中一般使用更复杂的 RNN 来捕捉的上下文信息，如 LSTM 和 GRU 等。

在经过卷积层以后，获得了文本中所有词的语义表示  $y_i^{(2)}$ ，然后经过一个 max-pooling 层和 softmax 层（输出层使用 softmax 激活函数）进行分类：

$$y^{(3)} = \max_{i=1}^n y_i^{(2)} \quad (5)$$

$$y^{(4)} = W^{(4)}y^{(3)} + b^{(4)} \quad (6)$$

$$p_i = \frac{\exp(y_i^{(4)})}{\sum_{k=1}^n \exp(y_k^{(4)})} \quad (7)$$

## 2. 实现细节

```
1 class RCNN(BasicModule):#继承自 BasicModule 其中封装了保存加载模型的接口, BasicModule继承自 nn.Module
2
3 def __init__(self, vocab_size, opt):#opt是 config类的实例 里面包括所有模型超参数的配置
4
5     super(RCNN, self).__init__()
6     # 嵌入层
7     self.embedding = nn.Embedding(vocab_size, opt.embed_size)#词嵌入矩阵 每一行代表词典中一个词对应的词向量;
8     # 词嵌入矩阵可以随机初始化连同分类任务一起训练, 也可以用预训练词向量初始化 (冻结或微调)
```

```

9
10     #双向 lstm 由于RCNN中双向lstm一般只有一层 所以opt.drop_prop_rcnn=0.0(丢弃率)
11     self.lstm = nn.LSTM(opt.embed_size,opt.recurrent_hidden_size,num_layers=opt.num_layers_rcnn,
12         bidirectional=True,batch_first=True,dropout=opt.drop_prop_rcnn)
13
14     #全连接层 维度转换 卷积操作可以用全连接层代替
15     self.linear = nn.Linear(2*opt.recurrent_hidden_size+opt.embed_size,opt.recurrent_hidden_size)
16
17     #池化层
18     self.max_pool = nn.MaxPool1d(opt.max_len)
19
20     #全连接层分类
21     self.content_fc = nn.Sequential(
22         nn.Linear(opt.recurrent_hidden_size, opt.linear_hidden_size),
23         nn.BatchNorm1d(opt.linear_hidden_size),
24         nn.ReLU(inplace=True),
25         nn.Dropout(opt.drop_prop),
26         # 可以再加一个隐层
27         # nn.Linear(opt.linear_hidden_size, opt.linear_hidden_size),
28         # nn.BatchNorm1d(opt.linear_hidden_size),
29         # nn.ReLU(inplace=True),
30         # 输出层
31         nn.Linear(opt.linear_hidden_size, opt.classes)
32     )
33
34     def forward(self, inputs):
35         #inputs(batch_size, seq_len)
36         # 由于 batch_first = True 所以inputs不用转换维度
37         embeddings = self.embedding(inputs) # (batch_size, seq_len, embed_size)
38
39         outputs,_ = self.lstm(embeddings) #(batch_size, seq_len, recurrent_hidden_size*2)
40
41         #将前后向隐藏状态和embedding拼接
42         outputs = torch.cat((outputs[:,:,:outputs.size(2)//2],embeddings,outputs[:,:,:outputs.size(2)
43             //2:]),dim=2) #(batch_size, seq_len, embed_size+recurrent_hidden_size*2)
44
45         #做维度转换
46         outputs = self.linear(outputs) #(batch_size, seq_len, recurrent_hidden_size)
47
48         #沿seq_len维做最大池化 (全局池化)
49         #先调整维度 交换recurrent_hidden_size维和seq_len维
50         #即把recurrent_hidden_size作为通道维 符合一维池化的输入
51         outputs = outputs.permute(0,2,1) #(batch_size, recurrent_hidden_size, seq_len)
52         outputs = self.max_pool(outputs).squeeze(2) #(batch_size, recurrent_hidden_size)
53
54         #通过全连接层 分类
55         outputs = self.content_fc(outputs) #(batch_size, classes)
56
57         return outputs

```

### 3.4 多层双向 LSTM

#### 1. 分类原理



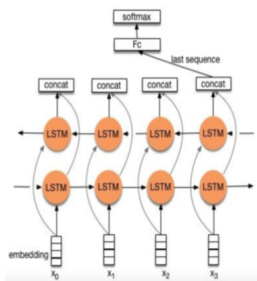


图 4: 多层双向 LSTM 分类原理图示

对于输入文本序列，在 LSTM 的每个时间步输入序列中一个单词的嵌入表示，计算当前时间步的隐藏状态，用于当前时间步的输出以及传递给下一个时间步和下一个单词的词向量一起作为 LSTM 单元输入，然后再计算下一个时间步的 LSTM 隐藏状态，以此重复... 直到处理完输入文本序列中的每一个单词，如果输入文本序列的长度为  $n$ ，那么就要经历  $n$  个时间步。

一般取前向和反向 LSTM 在最后一个时间步的隐藏状态，进行拼接，再接一个全连接层进行分类。由于 LSTM 训练比较困难，层数一般不超过两层。

## 2. 实现细节

```

1 class MulBiLSTM(BasicModule):#继承自 BasicModule          其中封装了保存加载模型的接口, BasicModule继承自 nn
    .Module
2     def __init__(self, vocab_size,opt):#opt是 config 类的实例  里面包括所有模型超参数的配置
3         super(MulBiLSTM, self).__init__()
4
5         #嵌入层
6         self.embedding = nn.Embedding(vocab_size, opt.embed_size)#词嵌入矩阵 每一行代表词典中一个词对应的
            词向量;
7         # 词嵌入矩阵可以随机初始化连同分类任务一起训练, 也可以用预训练词向量初始化 (冻结或微调)
8
9         # bidirectional 设为 True 即得到双向循环神经网络
10        self.encoder = nn.LSTM(input_size=opt.embed_size,
11                                hidden_size=opt.recurrent_hidden_size,
12                                num_layers=opt.num_layers,
13                                bidirectional=True,
14                                dropout=opt.drop_prop
15        )
16        self.fc = nn.Linear(4 * opt.recurrent_hidden_size, opt.classes) # 初始时间步和最终时间步的隐藏
            状态作为全连接层输入
17
18    def forward(self, inputs):
19        # inputs 的形状是(批量大小, 词数), 因为上述定义的 LSTM 没有设置参数 batch_first=True(默认 False), 所以
            需要将序列长度(seq_len)作为第一维, 所以将输入转置后再提取词特征
20        embeddings = self.embedding(inputs.permute(1,0)) # (seq_len, batch_size, embed_size)
21
22        # rnn.LSTM 只传入输入 embeddings (第一层的输入), 因此只返回最后一层的隐藏层在各时间步的隐藏状态。
23        # outputs 形状是(seq_len, batch_size, 2 * recurrent_hidden_size)
24        outputs, _ = self.encoder(embeddings) # output, (h, c)
25        # 连结初始时间步和最终时间步的隐藏状态作为全连接层输入。它的形状为
26        # (batch_size, 4 * recurrent_hidden_size)。
27        encoding = torch.cat((outputs[0], outputs[-1]), -1)
28        outs = self.fc(encoding)

```

```

29     #(batch_size, classes)
30     return outs

```

### 3.5 多层双向 LSTM with Attention

#### 1. 分类原理

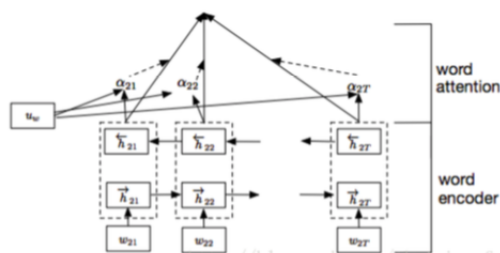


图 5: 多层双向 LSTM with Attention 分类原理图示

多层双向 LSTM with Attention 对多层双向 LSTM 模型做了一些改进，不再单纯地只利用最后时刻的隐藏状态（两个方向最后时刻隐藏状态拼接）进行分类，而是考虑每个时间步的隐藏状态，对每个时间步的隐藏状态（两个方向拼接）进行加权求和，然后对加权求和结果进行分类。其分类流程如图5所示。

Attention 机制的原理：

$$v_t = \tanh(W_w h_t + b_w) \quad (8)$$

$$\alpha_t = \frac{\exp(u_t^T v_t)}{\sum_t (\exp(u_t^T v_t))} \quad (9)$$

$$s = \sum_t \alpha_t h_t \quad (10)$$

其中， $W_w, b_w, u_w$  都是需要自己设置的权重参数（随模型训练）， $u_w$  可以认为是 Attention 机制中的 Query,  $u_t$  可以认为是 Attention 机制中的 Key,  $h_t$  是第  $t$  个时间步的隐藏状态（两个方向拼接），同时也是 Attention 机制中的 Value,  $\alpha_t$  为每个时间步隐藏状态对应的权重， $s$  是加权求和后得到的特征向量，用于最后接全连接层进行分类。

#### 2. 实现细节

```

1 class MulBiLSTM_Atten(BasicModule):#继承自 BasicModule 其中封装了保存加载模型的接口, BasicModule继承自 nn.
    Module
2
3     def __init__(self, vocab_size, opt):#opt是 config 类的实例 里面包括所有模型超参数的配置
4
5         super(MulBiLSTM_Atten, self).__init__()
6         # 嵌入层
7         self.embedding = nn.Embedding(vocab_size, opt.embed_size)#词嵌入矩阵 每一行代表词典中一个词对应的
            词向量;
8         # 词嵌入矩阵可以随机初始化连同分类任务一起训练, 也可以用预训练词向量初始化 (冻结或微调)
9
10        #多层双向LSTM 默认 seq_len 作为第一维 也可以通过 batch_first=True 设置 batch_size 为第一维

```

```

11     self.lstm = nn.LSTM(opt.embed_size, opt.recurrent_hidden_size, opt.num_layers,
12         bidirectional=True, batch_first=True, dropout=opt.drop_prop)
13
14     self.tanh1 = nn.Tanh()
15     self.u = nn.Parameter(torch.Tensor(opt.recurrent_hidden_size * 2, opt.recurrent_hidden_size * 2)
16         )
17     #定义一个参数 (变量) 作为 Attention 的 Query
18     self.w = nn.Parameter(torch.Tensor(opt.recurrent_hidden_size*2))
19
20     #均匀分布 初始化
21     nn.init.uniform_(self.w, -0.1, 0.1)
22     nn.init.uniform_(self.u, -0.1, 0.1)
23
24     #正态分布 初始化
25     #nn.init.normal_(self.w, mean=0, std=0.01)
26     self.tanh2 = nn.Tanh()
27     #最后的全连接层
28     self.content_fc = nn.Sequential(
29         nn.Linear(opt.recurrent_hidden_size*2, opt.linear_hidden_size),
30         nn.BatchNorm1d(opt.linear_hidden_size),
31         nn.ReLU(inplace=True),
32         nn.Dropout(opt.drop_prop),
33         # 可以再加一个隐层
34         # nn.Linear(opt.linear_hidden_size, opt.linear_hidden_size),
35         # nn.BatchNorm1d(opt.linear_hidden_size),
36         # nn.ReLU(inplace=True),
37         # 输出层
38         nn.Linear(opt.linear_hidden_size, opt.classes)
39     )
40
41     def forward(self, inputs):
42         #由于 batch_first = True 所有 inputs 不用转换维度
43         embeddings = self.embedding(inputs) # (batch_size, seq_len, embed_size)
44
45         outputs, _ = self.lstm(embeddings) #(batch_size, seq_len, recurrent_hidden_size*2)
46
47         #M = self.tanh1(outputs) #(batch_size, seq_len, recurrent_hidden_size*2)
48         M = torch.tanh(torch.matmul(outputs, self.u)) #也可以先做一个线性变换 再通过激活函数 作为 Key
49
50         #M (batch_size, seq_len, recurrent_hidden_size*2) self.w (recurrent_hidden_size*2,)
51         #torch.matmul(M, self.w) (batch_size, seq_len) w作为 Query与各个隐藏状态 (Key) 做内积
52         #再对第一维 seq_len 做 softmax 转换为概率分布 (batch_size, seq_len) 得到权重
53         alpha = F.softmax(torch.matmul(M, self.w), dim=1).unsqueeze(-1) # (batch_size, seq_len, 1)
54
55         #对各个隐藏状态和权重 对应相乘
56         out = alpha * outputs #(batch_size, seq_len, recurrent_hidden_size*2)
57
58         #对乘积求和 out为加权求和得到的特征向量
59         out = torch.sum(out, dim=1) #(batch_size, recurrent_hidden_size*2)
60
61         #out = F.relu(out)
62
63         out = self.content_fc(out) #(batch_size, classes)
64
65     return out

```

### 3.6 DPCNN

#### 1. 分类原理

ACL2017 年中，腾讯 AI-lab 提出了 Deep Pyramid Convolutional Neural Networks for Text Categorization(DPCNN)。论文中提出了一种基于 word-level 级别的网络-DPCNN，由于之前介绍的 TextCNN 不能通过卷积获得文本的长距离依赖关系，而论文中 DPCNN 通过不断加深网络，可以抽取长距离的文本依赖关系。实验证明在不增加太多计算成本的情况下，增加网络深度就可以获得最佳的准确率。

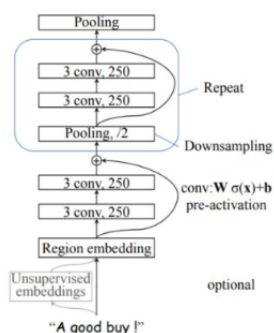


图 6: DPCNN 分类原理图示

#### 1) 等长卷积

首先交代一下卷积的一个基本概念。一般常用的卷积有以下三类：

假设输入的序列长度为  $n$ ，卷积核大小为  $m$ ，步长 (stride) 为  $s$ ，输入序列两端各填补  $p$  个零 (zero padding)，那么该卷积层的输出序列为  $(n-m+2p)/s+1$ 。

- 窄卷积：步长  $s=1$ ，两端不补零，即  $p=0$ ，卷积后输出长度为  $n-m+1$ 。
- 宽卷积：步长  $s=1$ ，两端补零  $p=m-1$ ，卷积后输出长度  $n+m-1$ 。
- 等长卷积：步长  $s=1$ ，两端补零  $p=(m-1)/2$ ，卷积后输出长度为  $n$ 。（不改变序列长度）

#### 2) 池化

那么 DPCNN 是如何捕捉长距离依赖的呢？——Downsampling with the number of feature maps fixed。

作者选择了适当的两层等长卷积来提高词位 embedding 的表示的丰富性。然后接下来就开始 Downsampling (池化)。再每一个卷积块 (两层的等长卷积) 后，使用一个  $\text{size}=3$  和  $\text{stride}=2$  进行 maxpooling 进行池化。序列的长度就被压缩成了原来的一半。其能够感知到的文本片段就比之前长了一倍。

DPCNN 中重复使用多个这种模块，每个模块包含两层等长卷积（不改变序列长度）和一个池化下采样（序列长度减半）。

例如之前是只能感知 3 个词位长度的信息，经过  $1/2$  池化层后就能感知 6 个词位长度的信息啦，这时把  $1/2$  池化层和  $\text{size}=3$  的卷积层组合起来如图7所示。

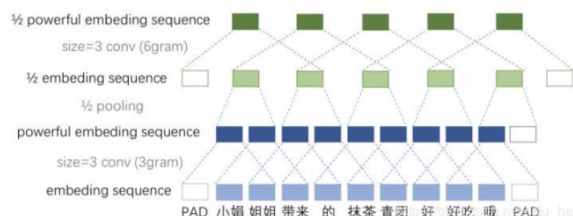


图 7: 感受野

### 3) 固定 feature map (filters/卷积核) 的数量

为什么要固定 feature maps 的数量呢? 许多模型每当执行池化操作时, 增加 feature maps 的数量 (即序列减半, 输出通道数加倍), 导致总计算复杂度是深度的函数。与此相反, 作者对 feature map 的数量进行了修正, 他们实验发现增加 feature map 的数量只会大大增加计算时间, 而没有提高精度。

固定了 feature map 的数量, 每当使用一个 size=3 和 stride=2 进行 max pooling 进行池化时, 每个卷积层的计算时间减半 (数据大小减半/序列长度减半), 从而形成一个金字塔。

剩下的我们就只需要重复的进行等长卷积 + 等长卷积 + 使用一个 size=3 和 stride=2 进行 maxpooling 进行池化就可以啦, DPCNN 就可以捕捉文本的长距离依赖啦! 还有一个问题, 网络加深时, 需要解决梯度消失问题。

#### 4) shortcut connections with pre-activation

a. 初始化 CNN 的时, 往往各层权重都初始化为很小的值, 这导致了最开始的网络中, 后续几乎每层的输入都是接近 0, 这时的网络输出没有意义;

b. 小权重阻碍了梯度的传播, 使得网络的初始训练阶段往往要迭代好久才能启动;

c. 就算网络启动完成, 由于深度网络中仿射矩阵 (每两层间的连接边) 近似连乘, 训练过程中网络也非常容易发生梯度爆炸或弥散问题。

当然, 上述这几点问题本质就是梯度弥散问题。那么如何解决深度 CNN 网络的梯度弥散问题呢? ResNet 中提出的 shortcut-connection/ skip-connection/ residual-connection (残差连接) 就是一种非常简单、合理、有效的解决方案。

类似地, 为了使深度网络的训练成为可能, 作者为了恒等映射, 所以使用加法进行 shortcut connections, 即  $z+f(z)$ , 其中  $f$  用的是两层的等长卷积 ( $z$  和  $f(z)$  的形状相同, 可以直接相加)。这样就可以极大的缓解了梯度消失问题。

另外, 作者也使用了 pre-activation。直观上, 这种“线性”简化了深度网络的训练, 类似于 LSTM 中 constant error carousels 的作用。而且实验证明 pre-activation 优于 post-activation (即把激活函数放在卷积操作前面, 卷积层的顺序为 BatchNorm-ReLU-Conv)。

#### 5) region embedding

同时 DPCNN 的底层保持了跟 TextCNN 一样的结构, 这里作者将 TextCNN 的包含多尺寸卷积滤波器的卷积层的卷积结果称之为 Region embedding, 意思就是对一个文本区域/片段 (比如 3gram, 卷积核大小为 3) 进行一组卷积操作后生成的 embedding。

### 2. 实现细节

```
1 class ResnetBlock(nn.Module):
2     def __init__(self, channel_size):
3         super(ResnetBlock, self).__init__()
```

```

4
5     self.channel_size = channel_size
6     #将序列长度减半
7     self.maxpool = nn.Sequential(
8         nn.ConstantPad1d(padding=(0, 1), value=0), #在每个通道上(一维) 一边填充0个0(不填充) 另一边填充1个0
9         nn.MaxPool1d(kernel_size=3, stride=2) #序列长度减半 height = (height-kernel_size+padding+stride)
10        //stride=height // 2
11    )
12    self.conv = nn.Sequential( #等长卷积 不改变height
13        nn.BatchNorm1d(num_features=self.channel_size),
14        nn.ReLU(),
15        nn.Conv1d(self.channel_size, self.channel_size, kernel_size=3, padding=1),
16
17        nn.BatchNorm1d(num_features=self.channel_size),
18        nn.ReLU(),
19        nn.Conv1d(self.channel_size, self.channel_size, kernel_size=3, padding=1),
20    )
21
22    def forward(self, x): #(batch_size, channel_size, seq_len)
23        x_shortcut = self.maxpool(x) # (batch_size, channel_size, seq_len//2)
24        x = self.conv(x_shortcut)#(batch_size, channel_size, seq_len//2)
25        x = x + x_shortcut#(batch_size, channel_size, seq_len//2) shortcut 残差连结
26
27    return x
28
29 class DPCNN(BasicModule):#继承自 BasicModule 其中封装了保存加载模型的接口, BasicModule继承自 nn.Module
30
31    def __init__(self, vocab_size, opt): #opt是 config类的实例 里面包括所有模型超参数的配置
32
33        super(DPCNN, self).__init__()
34        # 嵌入层
35        self.embedding = nn.Embedding(vocab_size, opt.embed_size)#词嵌入矩阵 每一行代表词典中一个词对应的词向量;
36        # 词嵌入矩阵可以随机初始化连同分类任务一起训练, 也可以用预训练词向量初始化(冻结或微调)
37
38        # region embedding
39        self.region_embedding = nn.Sequential(
40            nn.Conv1d(opt.embed_size, opt.channel_size, kernel_size=3, padding=1), #same卷积 不改变height/序列长度
41            nn.BatchNorm1d(num_features=opt.channel_size),
42            nn.ReLU(),
43            nn.Dropout(opt.drop_prop_dpcnn)
44        )
45
46        #卷积块 same卷积 不改变height
47        self.conv_block = nn.Sequential(
48            nn.BatchNorm1d(num_features=opt.channel_size),
49            nn.ReLU(),
50            nn.Conv1d(opt.channel_size, opt.channel_size, kernel_size=3, padding=1),
51            nn.BatchNorm1d(num_features=opt.channel_size),
52            nn.ReLU(),
53            nn.Conv1d(opt.channel_size, opt.channel_size, kernel_size=3, padding=1),
54        )
55
56        self.seq_len = opt.max_len #序列最大长度

```

```

56     resnet_block_list = [] #存储多个残差块
57
58     while (self.seq_len > 2): #每经过一个残差块 序列长度减半 只要长度>2 就不停地加残差块
59         resnet_block_list.append(ResnetBlock(opt.channel_size))
60         self.seq_len = self.seq_len // 2
61
62     #将残差块 构成残差层 作为一个子模块
63     self.resnet_layer = nn.Sequential(*resnet_block_list)
64
65     #输出层 分类
66     self.linear_out = nn.Linear(self.seq_len * opt.channel_size, opt.classes)
67
68     def forward(self, inputs):
69
70         embeddings = self.embedding(inputs) #(batch_size,max_len,embed_size)
71
72         x = embeddings.permute(0, 2, 1) #(batch_size,embed_size,max_len) 交换维度 作为1维卷积的输入
73         #embed_size 作为通道维
74         x = self.region_embedding(x) #(batch_size,channel_size,max_len)
75
76         x = self.conv_block(x) #(batch_size,channel_size,max_len)
77
78         x = self.resnet_layer(x) #经过多个残差块 每次长度减半 (batch_size,channel_size,self.seq_len)
79
80         x = x.permute(0, 2, 1) #(batch_size,self.seq_len,channel_size)
81
82         x = x.contiguous().view(x.size(0), -1) #(batch_size,self.seq_len*channel_size) 拉伸为向量
83         out = self.linear_out(x) #(batch_size,classes)
84         return out

```

## 4 数据集及预处理

### 4.1 数据集简介

THUCNews 是根据新浪新闻 RSS 订阅频道 2005 2011 年间的历史数据筛选过滤生成, 包含 74 万篇新闻文档 (2.19 GB), 均为 UTF-8 纯文本格式。我们在原始新浪新闻分类体系的基础上, 重新整合划分出 14 个候选分类类别: 财经、彩票、房产、股票、家居、教育、科技、社会、时尚、时政、体育、星座、游戏、娱乐。(完整数据集压缩包下载)

### 4.2 数据预处理

在进行特征提取之前, 需要对原始文本数据进行预处理, 这对于特征提取来说至关重要, 一个好的预处理过程会显著的提高特征提取的质量以及分类算法的性能。文本预处理一般包括以下步骤:

(1) 分词: 首先, 需要把文本切分成单词或短语。对于英文文本, 可以直接按照空格进行切分 (此时句末的标点不会单独切分出来) 或使用一些英文分词工具如 `nltk` 中的分词工具; 对于中文文本, 可以使用分词工具 (如 `jieba` 等) 进行切分。(文本分类算法有基于词和基于字符两种处理方式, 一般来说基于词的文本分类算法效果更好, 本专栏介绍的文本分类算法都是基于词的处理方式)。



(2) 去停止词: 所谓停止词, 就是在文本中大量出现但对分类并没有太多作用的词。如英文里的 'a', 'an', 'the', 'above', 'after', 'and', 'are', 'as', 'at', 'be', 'but', 'by', 'can', 'could', 'do', 'each', 'for', 'from', 'has', 'have', 'he', 'her', 'his', 'hundred', 'if', 'in', 'into', 'is', 'it', 'its', 'me', 'more', 'most', 'my', 'no', 'not', 'of', 'off', 'on', 'or', 'other', 'out', 'over', 'she', 'some', 'that', 'the', 'there', 'these', 'they', 'this', 'to', 'too', 'two', 'us', 'very', 'was', 'we', 'were', 'what', 'when', 'where', 'which', 'who', 'with', 'without', 'would', 'yet', 'you', 'your' 等。中文里的 '的', '这', '那', ... 在这一步, 把文本中的停止词过滤掉。(英文文本去停止词, 可以直接使用 nltk 中封装的函数, 中文文本可以自行下载停用词表(我们使用的是哈工大停用词表), 来进行过滤)。

(3) 小写化: 这一步主要针对英文文本, 大多数情况下需要把英文文本统一转换为小写形式。

(4) 噪声移除: 去除文本中的特殊符号, 如特殊标点等。这些特殊符号对于人类理解文本可能很重要, 但对于分类算法并没有太大意义(可以和停止词一并去除)。

(5) 拼写检查: 数据集可能来自一些社交媒体, 存在一些拼写错误, 可以尝试对拼写问题进行纠正。

(6) 俚语和缩写: 可以尝试把缩写还原为完整形式或将一些较口语化的表示转换为书面语。

(7) 词干提取和词型还原: 这一步主要针对英文文本, 即将单词转换为最基本的形式。如 studying → study, apples → apple 等。

(8) 词频统计与过滤: 对文本训练集进行上述处理后, 统计剩余单词的词频, 并过滤低频词(可以设置一个阈值(如, 5), 保留词频大于阈值的单词; 或者基于词频进行排序, 取前 k 个词。前者词典的大小不确定, 后者词典大小是确定的, 就是 k (或 k+2 包括位置符号 unk 和填充符号 pad))。对剩余单词, 构建词典。

## 4.3 实现细节

### 1. 读取数据

```
1 def read_cnews():
2     #opt.data_root为数据集解压后 文件夹所在路径 定义在config.py中
3     data = []
4     #所有的主题标签
5     labels = [label for label in os.listdir(opt.data_root) if label != '.DS_Store']
6     print(labels)
7     #标签到整数索引的映射
8     labels2index = dict(zip(labels, list(range(len(labels)))))
9     #整数索引到标签的映射
10    index2labels = dict(zip(list(range(len(labels))), labels))
11    print(labels2index)
12    print(index2labels)
13
14    #存储整数索引到标签的映射 以便预测时使用
15    with open('index2labels.json', 'w') as f:
16        json.dump(index2labels, f)
17    #存储类别标签 打印分类报告会用到
18    with open('labels.json', 'w') as f:
19        json.dump(labels, f)
20
21
22    for label in labels:
23        folder_name = os.path.join(opt.data_root, label)
24        datasub = [] #存储某一类的数据 [[string, index], ...]
25        for file in tqdm(os.listdir(folder_name)):
26            with open(os.path.join(folder_name, file), 'rb') as f:
27                #去除文本中空白符
28                review = f.read().decode('utf-8').replace('\n', '').replace('\r', '').replace('\t', '')
29                datasub.append([review, labels2index[label]])
30        data.append(datasub) #存储所有类的数据 [[[string, index], ...], [[string, index], ...], ...]
```



```

31
32     return data

```

## 2. 切分训练集、验证集和测试集

```

1 def split_data(data):
2     #切分数据集 为训练集、验证集和测试集
3     train_data = []
4     val_data = []
5     test_data = []
6
7     #对每一类数据进行打乱
8     #设置验证集和测试集中每一类样本数都为200（样本均衡）
9     for data1 in data: #遍历每一类数据
10        np.random.shuffle(data1) #打乱
11        val_data += data1[:200]
12        test_data += data1[200:400]
13        train_data += data1[400:]
14
15    np.random.shuffle(train_data) #打乱训练集 测试机和验证集不用打乱
16
17    print(len(train_data))
18    print(len(val_data))
19    print(len(test_data))
20
21    return train_data, val_data, test_data

```

## 3. 基于训练集构建词典

```

1 #读取停用词
2 def stopwords(fileroot):
3     #fileroot为下载的停用词表所在的路径
4     with open(fileroot, 'r') as f:
5         stopword = [line.strip() for line in f]
6     print(stopword[:5])
7     return stopword
8
9 #分词 去除停用词
10 def get_tokenized(data, stopword):
11     """
12     data: list of [string, label]
13     """
14     def tokenizer(text):
15         return [tok for tok in jieba.lcut(text) if tok not in stopword]
16     return [tokenizer(review) for review, _ in data]
17
18 #构建词典
19 def get_vocab(data, stopword):
20     tokenized_data = get_tokenized(data, stopword) #分词、去除停用词
21     counter = collections.Counter([tk for st in tokenized_data for tk in st]) #统计词频
22     return Vocab.Vocab(counter, min_freq=5, specials=['<pad>', '<unk>']) #保留词频大于5的词 <pad>对应填充
    项（词典中第0个词） <unk>对应低频词和停止词等未知词（词典中第1个词）

```

## 4. 基于词典对训练集、验证集以及测试集进行处理

```

1 def preprocess_imdb(data, vocab, stopword):
2     #将训练集、验证集、测试集中单词转换为词典中对应的索引

```

```
3     max_l = 500 # 将每条新闻通过截断或者补0,使得长度变成500(所有数据统一成一个长度,方便用矩阵并行计
      算(其实也可以每个 batch填充成一个长度, batch间可以不一样))
4
5     def pad(x):
6         return x[:max_l] if len(x) > max_l else x + [0] * (max_l - len(x))
7
8     tokenized_data = get_tokenized(data,stopword) #分词、去停止词
9     features = torch.tensor([pad([vocab.stoi[word] for word in words]) for words in tokenized_data]) #把
      单词转换为词典中对应的索引,并填充成固定长度,封装为tensor
10    labels = torch.tensor([score for _, score in data])
11    return features, labels
```

## 5. 保存一些中间结果

```
1 with open('word2index.json','w') as f: #保存词到索引的映射,预测和后续加载预训练词向量时会用到
2     json.dump(vocab.stoi,f)
3
4 with open('vocabsize.json','w') as f:
5     #保存词典的大小(因为我们基于词频阈值过滤低频词,词典大小不确定,需要保存,后续模型中会用到)
6     json.dump(len(vocab),f)
7
8 #保存预处理好的训练集、验证集和测试集 以便后续训练时使用
9 torch.save(X_train,'X_train.pt')
10 torch.save(y_train, 'y_train.pt')
11 torch.save(X_val, 'X_val.pt')
12 torch.save(y_val, 'y_val.pt')
13 torch.save(X_test, 'X_test.pt')
14 torch.save(y_test, 'y_test.pt')
```

## 5 项目组织结构

在学习某个深度学习框架时,掌握其基本知识和接口固然重要,但如何合理组织代码,使得代码具有良好的可读性和可扩展性也必不可少。做深度学习实验或项目时,为了得到最优的模型结果,中间往往需要很多次的尝试和修改(也就是所谓地调参)。从事大多数深度学习研究时,程序都需要实现以下几个功能:

- 1) 模型定义
- 2) 数据处理和加载
- 3) 训练模型 (Train&Validate)
- 4) 训练过程的可视化或相关指标的计算
- 5) 测试/预测 (Test/Inference)

另外程序还应该满足以下几个要求:

- 1) 模型需具有高度可配置性,便于修改参数、修改模型,反复实验
- 2) 代码应具有良好的组织结构,使人一目了然
- 3) 代码应具有良好的说明,使其他人能够理解

接下来我将应用这些内容,并结合实际的例子,来讲解如何合理组织我们的文本分类项目。

### 5.1 文件组织结构

首先来看程序文件的组织结构:

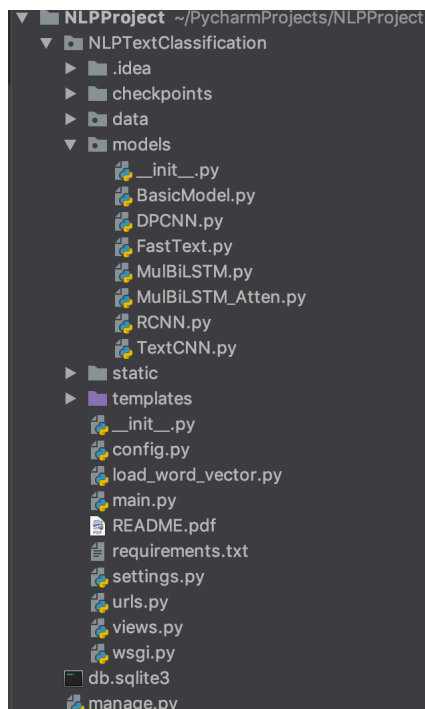


图 8: 程序文件的组织结构

其中:

- 1) checkpoints/: 用于保存训练好的模型, 可使程序在异常退出后仍能重新载入模型, 恢复训练
- 2) data/: 数据相关操作, 包括数据预处理、dataset 实现等
- 3) models/: 模型定义, 可以有多个模型, 例如上面的 FastText、TextCNN 等, 一个模型对应一个文件
- 4) config.py: 配置文件, 所有可配置的变量都集中在此, 并提供默认值
- 5) main.py: 主文件, 训练和预测程序的入口, 可通过不同的命令来指定不同的操作和参数
- 6) load\_word\_vector.py: 定义加载预训练词向量的函数
- 7) .idea/, static/, templates/, settings.py、urls.py、views.py、wsgi.py、db.sqlite3、manage.py: 网页 Demo 运行支撑文件。
- 8) requirements.txt: 程序依赖的第三方库
- 9) README.pdf: 项目说明文档

## 5.2 数据预处理和加载

数据的相关预处理函数主要保存在 data/dataset.py 中。关于数据加载的相关操作, 其基本原理就是使用 Dataset 进行数据集的封装, 再使用 Dataloader 实现数据并行加载。

具体的预处理过程和实现细节在第三部分中已经详细介绍了。

使用时, 我们可通过 dataloader 加载数据:

```
1 #读取之前预处理过程 保存的处理好的训练集、验证集和测试集
2 X_train = torch.load('./data/X_train.pt')
```

```

3 y_train = torch.load('./data/y_train.pt')
4 X_val = torch.load('./data/X_val.pt')
5 y_val = torch.load('./data/y_val.pt')
6 X_test = torch.load('./data/X_test.pt')
7 y_test = torch.load('./data/y_test.pt')
8
9 #封装成DataSet
10 trainset = Data.TensorDataset(X_train,y_train)
11 valset = Data.TensorDataset(X_val,y_val)
12 testset = Data.TensorDataset(X_test,y_test)
13
14 #使用DataLoader并行加载数据
15 train_iter = Data.DataLoader(trainset,opt.batch_size,shuffle=True,num_workers=opt.num_workers)
16 val_iter = Data.DataLoader(valset,opt.batch_size)
17 test_iter = Data.DataLoader(testset,opt.batch_size)

```

加载预训练词向量:

`load_word_vector.py` 定义了加载预训练词向量的函数:

```

1 def read_word_vector(path): #path为 下载的预训练词向量 解压后的文件所在的路径
2     #读取预训练词向量
3     with open(path, 'r') as f:
4         words = set() # 定义一个words集合
5         word_to_vec_map = {} # 定义词到向量的映射字典
6         for line in f: #跳过文件的第一行
7             break
8
9         for line in f: # 遍历f中的每一行
10            line = line.strip().split() # 去掉首尾空格, 每一行以空格切分 返回一个列表 第一项为单词 其
11                余为单词的嵌入表示
12            curr_word = line[0] # 取出单词
13            words.add(curr_word) # 加到集合/词典中
14            # 定义词到其嵌入表示的映射字典
15            word_to_vec_map[curr_word] = np.array(line[1:], dtype=np.float64)
16
17        return words, word_to_vec_map
18
19 def load_pretrained_embedding(word2index, word2vector):
20     #word2index是构建的词典 (单词到索引的映射), word2vector是预训练词向量 (单词到词向量的映射)
21
22     embed = torch.zeros(len(word2index), opt.embed_size) # 初始化词嵌入矩阵为0
23     oov_count = 0 # 找不到预训练词向量的词典中单词的个数
24
25     for word, index in word2index.items(): #遍历词典中的每个单词 及其在词典中的索引
26         try: #如果单词有对应的预训练词向量 则用预训练词向量对词嵌入矩阵的对应行进行赋值
27             embed[index, :] = torch.from_numpy(word2vector[word])
28         except KeyError:
29             oov_count += 1
30
31     if oov_count > 0:
32         print("There are %d oov words."%oov_count)
33     return embed #返回词嵌入矩阵

```

在主程序 `main.py` 的 `train` 函数中调用:

```

1 #加载预训练词向量

```

```

2 if opt.use_pretrained_word_vector:
3     words, word2vec = read_word_vector(opt.word_vector_path) #opt.word_vector_path为下载的预训练词向量解
        压后的文件所在的路径
4     print("预训练词向量读取完毕!")
5     #读取之前预处理过程保存的词典 (词到索引的映射)
6     with open('./data/word2index.json') as f:
7         word2index = json.load(f)
8
9     model.embedding.weight.data.copy_(load_pretrained_embedding(word2index, word2vec)) #使用加载完预训练
        词向量的词嵌入矩阵 对embedding层的词嵌入矩阵赋值
10    print("预训练词向量加载完毕!")
11    if opt.frozen: #冻结还是finetuning
12        model.embedding.weight.requires_grad = False

```

### 5.3 模型定义

各个模型的定义主要保存在 models/目录下，其中 BasicModule 是对 nn.Module 的简易封装，提供快速加载（可以处理 GPU 训练、CPU 加载的情况）和保存模型（提供多 GPU 训练时的模型保存方法）的接口，其他模型都继承自 BasicModule。

```

1
2 class BasicModule(nn.Module):
3     '''
4     封装了nn.Module，主要提供save和load两个方法
5     '''
6     def __init__(self, opt=None):
7         super(BasicModule, self).__init__()
8         self.model_name = str(type(self)) # 模型的默认名字
9
10    def load(self, path):
11        '''
12        加载模型
13        可指定路径
14        '''
15        self.load_state_dict(torch.load(path))
16
17    def load_map(self, path, device): #如果在GPU上训练 在CPU上加载 可以调用这个函数
18        '''
19        加载模型
20        可指定路径
21        '''
22        self.load_state_dict(torch.load(path, map_location=device))
23
24    def save(self, name=None):
25        '''
26        保存模型，默认使用“模型名字_best”作为文件名，
27        '''
28        if name is None:
29            prefix = 'checkpoints/' + self.model_name.split('.')[ -2] + '_best.pth'
30            #name = time.strftime(prefix + '%n%d_%H:%M:%S.pth')
31            torch.save(self.state_dict(), prefix) #只保存模型的参数
32            return name
33

```

```

34 |     def save_multiGPU(self, name=None): #如果使用多GPU训练, 保存模型时, 可以调用这个函数。
35 |         """
36 |         保存模型, 默认使用“模型名字_best”作为文件名,
37 |         """
38 |         if name is None:
39 |             prefix = 'checkpoints/' + self.model_name.split('.')[ -2] + '_best.pth'
40 |             # name = time.strftime(prefix + '%m%d_%H:%M:%S.pth')
41 |             torch.save(self.module.state_dict(), prefix) # 只保存模型的参数
42 |         return name

```

在实际使用中, 直接调用 `model.save()` 及 `model.load(opt.load_path)` 即可, 我们已经对保存和加载做了封装。

其它自定义模型一般继承 `BasicModule`, 然后实现自己的模型。其中 `TextCNN.py` 实现了 `TextCNN`, `FastText.py` 实现了 `FastText` 等。在 `models/__init__.py` 中, 代码如下:

```

1 | #本项目可选择的模型:
2 | from .FastText import FastText
3 | from .TextCNN import TextCNN
4 | from .MulBiLSTM import MulBiLSTM
5 | from .MulBiLSTM_Atten import MulBiLSTM_Atten
6 | from .RCNN import RCNN
7 | from .DPCNN import DPCNN

```

这样在主函数中就可以写成:

```

1 | from models import TextCNN

```

或:

```

1 | import models
2 | model = models.TextCNN()

```

或:

```

1 | import models
2 | model = getattr(models, 'TextCNN')()

```

其中最后一种写法最为关键, 这意味着我们可以通过字符串直接指定使用的模型, 而不必使用判断语句, 也不必在每次新增加模型后都修改代码。新增模型后只需要在 `models/__init__.py` 中加上

```

1 | from .new_module import NewModule

```

各个模型的原理和实现细节在第二部分中已经详细介绍过了。

## 5.4 配置文件

在模型定义、数据处理和训练等过程都有很多变量, 这些变量应提供默认值, 并统一放置在配置文件中, 这样在后期调试、修改代码或迁移程序时会比较方便, 在这里我们将所有可配置项放在 `config.py` 中。

```

1 | class DefaultConfig(object):
2 |
3 |     model = 'FastText' # 使用的模型, 名字必须与models/__init__.py中的名字一致
4 |
5 |     load_model_path = None # 加载预训练的模型的路径, 为None代表不加载
6 |
7 |     batch_size = 256 # batch size

```

```
8     num_workers = 4 # 加载数据使用的线程数
9
10    #下载数据集 解压缩后得到的文件夹所在的路径
11    data_root = '/Users/apple/Downloads/THUCNews-1'
12
13
14    max_epoch = 20
15    lr = 0.01 # initial learning rate
16    weight_decay = 1e-4 # 损失函数 正则化
17    embed_size = 100 #词嵌入维度
18    drop_prop = 0.5 #丢弃率
19    classes = 14 #分类类别数
20    max_len = 500 #序列最大长度
21
22    #学习率衰减相关超参数
23    use_lrdecay = True #是否使用学习率衰减
24    lr_decay = 0.95 # 衰减率
25    n_epoch = 1 #每隔n_epoch个epoch衰减一次 lr = lr * lr_decay
26
27
28    #TextCNN相关的超参数
29    kernel_sizes = [3,4,5] #一维卷积核的大小
30    num_channels = [100,100,100] #一维卷积核的数量
31
32    #FastText相关的超参数
33    linear_hidden_size = 512 #隐层单元数
34
35    #MulBiLSTM/MulBiLSTM_Atten相关超参数
36    recurrent_hidden_size = 128 #循环层 单元数
37    num_layers = 2 #循环层 层数
38
39    #RCNN相关超参数
40    num_layers_rcnn = 1 #循环层 层数
41    drop_prop_rcnn = 0.0 #1个循环层设置为0 丢弃率
42
43    #DPCNN相关超参数
44    channel_size = 250
45    drop_prop_dpcnn = 0.2
46
47    #梯度裁剪相关超参数
48    use_rnn = False
49    norm_type = 1
50    max_norm = 5
51
52    #预训练词向量相关超参数
53    use_pretrained_word_vector = False
54    word_vector_path = '/Users/apple/Downloads/sgns.sogou.word'
55    frozen = False
56
57    #待分类文本
58    text="众所周知，一支球队想要夺冠，超级巨星必不可少，不过得到超级巨星并不简单，方式无非两种，一是自己
        培养，这种方式适用于所有球队，二是交易，这种方式基本只适用于大市场球队——事实就是，30支球队之
        间并非完全公平，超级巨星依然更愿意前往大城市。"
59
60    #预测时是否对文本进行填充或截断
```

```
61 | predict_pad = False
```

可配置的参数主要包括：

- 1) 训练参数（学习率、训练 epoch 等）
- 2) 各个模型相关的参数

这样我们在程序中就可以这样使用：

```
1 | import models
2 | from config import DefaultConfig
3 |
4 | opt = DefaultConfig()
5 | lr = opt.lr
6 | model = getattr(models, opt.model)
```

这些都只是默认参数（如果后续不在命令行指定的话，就是用默认参数），在这里还提供了更新函数（根据命令行中指定的参数进行更新），根据字典更新配置参数。

```
1 | def parse(self, kwargs):
2 |     '''
3 |     根据字典kwargs 更新 默认的 config 参数
4 |     '''
5 |     # 更新配置参数
6 |     for k, v in kwargs.items():
7 |         if not hasattr(self, k):
8 |             # 警告还是报错，取决于个人喜好
9 |             warnings.warn("Warning: opt has not attribut %s" % k)
10 |            setattr(self, k, v)
11 |
12 |    # 打印配置信息
13 |    print('user config:')
14 |    for k, v in self.__class__.__dict__.items(): #python3 中 iteritems() 已经废除了
15 |        if not k.startswith('__'):
16 |            print(k, getattr(self, k))
```

这样我们在实际使用时，并不需要每次都修改 config.py（默认配置），只需要通过命令行传入所需参数，覆盖默认配置即可。

```
1 | opt = DefaultConfig()
2 | new_config = {'lr':0.1,'use_gpu':False}
3 | opt.parse(new_config)
4 | opt.lr == 0.1
```

## 5.5 main.py

在讲解主程序 main.py 之前，我们先来看看 2017 年 3 月谷歌开源的一个命令行工具 fire，通过 pip install fire 即可安装。下面来看看 fire 的基础用法，假设 example.py 文件内容如下：

```
1 | import fire
2 | def add(x, y):
3 |     return x + y
4 |
5 | def mul(**kwargs):
6 |     a = kwargs['a']
7 |     b = kwargs['b']
```



```
8     return a * b
9
10 if __name__ == '__main__':
11     fire.Fire()
```

那么我们可以使用：

```
1 python example.py add 1 2 # 执行add(1, 2)
2 python example.py mul --a=1 --b=2 # 执行mul(a=1, b=2),kwargs={'a':1, 'b':2}
3 python example.py add --x=1 --y=2 # 执行add(x=1, y=2)
```

可见，只要在程序中运行 fire.Fire()，即可使用命令行参数 python file <function> [args,] -kwargs,。

在主程序 main.py 中，主要包含四个函数，其中三个需要命令行执行，main.py 的代码组织结构如下：

```
1 def train(**kwargs):
2     '''
3     训练
4     '''
5     pass
6
7 def evaluate_accuracy(data_iter, net, flag=False, labels=None):
8     '''
9     计算模型在验证集/测试集上的准确率等信息，用以辅助训练
10    '''
11    pass
12
13 def predict(**kwargs):
14    '''
15    对新样本进行预测
16    '''
17    pass
18
19 def help():
20    '''
21    打印帮助的信息
22    '''
23    print('help')
24
25 if __name__ == '__main__':
26     import fire
27     fire.Fire()
```

## 1. 训练

训练的主要步骤如下：

- 1) 定义网络
- 2) 定义数据
- 3) 定义损失函数和优化器
- 4) 计算重要指标
- 5) 开始训练
- 6) 训练网络
- 7) 计算在验证集上的指标

训练函数的代码如下：

```
1 def train(**kwargs):
2
3     # 根据命令行参数更新配置 否则使用默认配置
4     opt.parse(kwargs)
5
6     # step1: 数据
7     #词典大小
8     with open('./data/vocabsize.json') as f:
9         vocab_size = json.load(f)
10    print("词典大小:", vocab_size)
11    #标签
12    with open('./data/labels.json') as f:
13        labels = json.load(f)
14
15    #读取之前预处理过程 保存的处理好的训练集、验证集和测试集
16    X_train = torch.load('./data/X_train.pt')
17    y_train = torch.load('./data/y_train.pt')
18    X_val = torch.load('./data/X_val.pt')
19    y_val = torch.load('./data/y_val.pt')
20    X_test = torch.load('./data/X_test.pt')
21    y_test = torch.load('./data/y_test.pt')
22
23    #封装成DataSet
24    trainset = Data.TensorDataset(X_train, y_train)
25    valset = Data.TensorDataset(X_val, y_val)
26    testset = Data.TensorDataset(X_test, y_test)
27
28    #使用DataLoader并行加载数据
29    train_iter = Data.DataLoader(trainset, opt.batch_size, shuffle=True, num_workers=opt.num_workers)
30    val_iter = Data.DataLoader(valset, opt.batch_size)
31    test_iter = Data.DataLoader(testset, opt.batch_size)
32
33    # step2: 模型
34    model = getattr(models, opt.model)(vocab_size, opt)
35    if opt.load_model_path:
36        model.load(opt.load_model_path)
37
38    #加载预训练词向量
39    if opt.use_pretrained_word_vector:
40        words, word2vec = read_word_vector(opt.word_vector_path) #opt.word_vector_path为下载的预训练词向
41        #量 解压后的文件所在的路径
42        print("预训练词向量读取完毕！")
43        #读取之前预处理过程保存的词典（词到索引的映射）
44        with open('./data/word2index.json') as f:
45            word2index = json.load(f)
46
47        model.embedding.weight.data.copy_(load_pretrained_embedding(word2index, word2vec)) #使用加载完预
48        #训练词向量的词嵌入矩阵 对embedding层的词嵌入矩阵赋值
49        print("预训练词向量加载完毕！")
50        if opt.frozen: #冻结还是finetuning
51            model.embedding.weight.requires_grad = False
52
53    print("使用设备: ", device)
54    if torch.cuda.device_count() > 1: #使用多GPU进行训练
```

```

54     print("Let's use", torch.cuda.device_count(), "GPUs!")
55     model = torch.nn.DataParallel(model)
56
57     model.to(device)
58
59     # step3: 目标函数和优化器
60     criterion = torch.nn.CrossEntropyLoss()
61     optimizer = torch.optim.Adam(model.parameters(),
62     lr = opt.lr,
63     weight_decay = opt.weight_decay)
64     scheduler = lr_scheduler.StepLR(optimizer, opt.n_epoch, opt.lr_decay)
65     # 训练
66     batch_count = 0
67     best_f1_val = 0.0
68
69
70     for epoch in range(opt.max_epoch):
71         train_l_sum, train_acc_sum, n, start = 0.0, 0.0, 0, time.time()
72         if opt.use_lrdecay:
73             scheduler.step()
74         for X, y in train_iter:
75             X = X.to(device)
76             y = y.to(device)
77             y_hat = model(X)
78             loss = criterion(y_hat, y)
79             optimizer.zero_grad()
80             loss.backward()
81             if opt.use_rnn: #梯度裁剪
82                 nn.utils.clip_grad_norm_(model.parameters(), max_norm=opt.max_norm, norm_type=opt.
83                     norm_type)
84             optimizer.step()
85             train_l_sum += loss.cpu().item()
86             train_acc_sum += (y_hat.argmax(dim=1) == y).sum().cpu().item()
87             n += y.shape[0]
88             batch_count += 1
89
90     #一个epoch后在验证集上做一次验证
91     val_f1, val_acc = evaluate_accuracy(val_iter, model)
92     if val_f1 > best_f1_val:
93         best_f1_val = val_f1
94         # 保存在验证集上 weighted average f1最高的参数 (最好的参数)
95         if torch.cuda.device_count() > 1: #多GPU训练时保存参数
96             print("Saving on ", torch.cuda.device_count(), "GPUs!")
97             model.save_multiGPU()
98         else:
99             print("Saving on one GPU!") #单GPU训练时保存参数
100             model.save()
101         #使用当前最好的参数, 在测试集上再跑一遍
102         best_f1_test, best_acc_test = evaluate_accuracy(test_iter, model, True, labels)
103
104     print('epoch %d, lr %.6f, loss %.4f, train acc %.3f, val acc %.3f, val weighted f1 %.3f, val
105           best_weighted f1 %.3f, test best_acc %.3f, test best_weighted f1 %.3f, time %.1f sec'
106           % (epoch + 1, optimizer.state_dict()['param_groups'][0]['lr'], train_l_sum / batch_count,
107             train_acc_sum / n, val_acc, val_f1, best_f1_val, best_acc_test, best_f1_test, time.time() -

```

```
start))
```

## 2. 验证/测试

验证相对来说比较简单，但要注意需将模型置于验证模式 (`model.eval()`)，验证完成后还需要将其置回为训练模式 (`model.train()`)，这两句代码会影响 BatchNorm 和 Dropout 等层的运行模式。

多分类我们使用 weighed average f1-score 作为评估指标，主要使用 sklearn 中的指标计算函数。

代码如下：

```
1 def evaluate_accuracy(data_iter, net, flag=False, labels=None):
2     #计算模型在验证集上的相关指标 多分类我们使用 weighed average f1-score
3
4     acc_sum, n = 0.0, 0
5     net.eval() # 评估模式，这会关闭dropout
6     y_pred_total = []
7     y_total = []
8     with torch.no_grad():
9         for X, y in data_iter:
10             #acc_sum += (net(X.to(device)).argmax(dim=1) == y.to(device)).float().sum().cpu().item()
11             #n += y.shape[0]
12             y_pred = net(X.to(device)).argmax(dim=1).cpu().numpy()
13             y_pred_total.append(y_pred)
14             y_total.append(y.numpy())
15
16     y_pred = np.concatenate(y_pred_total)
17     y_label = np.concatenate(y_total)
18     weighted_f1 = f1_score(y_label, y_pred, average='weighted') #weighed average f1-score
19
20     accuracy = accuracy_score(y_label, y_pred) #准确率
21     if flag: #当在测试集上验证时 flag设置为True 额外打印分类报告和混淆矩阵
22         print(classification_report(y_label, y_pred, digits=4, target_names = labels))
23         cm = confusion_matrix(y_label, y_pred)
24         print(cm)
25     net.train() # 改回训练模式
26
27     return weighted_f1, accuracy
```

## 3. 预测

对于新的输入文本，我们加载训练好的模型进行预测，输出类别标签：

```
1 def predict(**kwargs):
2     # 根据命令行参数更新配置 否则使用默认配置
3     device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
4     print("使用设备:", device)
5     opt.parse(kwargs)
6     text = opt.text #待分类文本
7
8     # 词典大小
9     with open('./data/vocabsize.json') as f:
10         vocab_size = json.load(f)
11     print(vocab_size)
12
13     #创建指定的模型对象
14     model = getattnr(models, opt.model)(vocab_size, opt)
15
```

```

16 #加载训练好的模型参数
17 if device.type=='cpu': #GPU训练 CPU预测 加载参数时需要对参数进行映射
18     model.load_map('./checkpoints/'+opt.model+'_best.pth',device)
19 else:
20     model.load('./checkpoints/' + opt.model + '_best.pth')
21
22 #加载之前预处理过程 保存的词到索引的映射字典
23 with open('./data/word2index.json') as f:
24     word2index = json.load(f)
25
26 #device = list(model.parameters())[0].device
27 if opt.predict_pad: #预测时对文本进行填充 (若文本长度<opt.max_len) 或截断 (若文本长度>
    opt.max_len)
28     sentence = [word2index.get(word, 1) for word in jieba.lcut(text)]
29     sentence = sentence[:opt.max_len] if len(sentence) > opt.max_len else sentence + [0] * (opt.
        max_len - len(sentence))
30     sentence = torch.tensor(sentence,device=device)
31 else:
32     sentence = torch.tensor([word2index.get(word,1) for word in jieba.lcut(text)],device=device)
33 print(sentence)
34 #预测
35 with torch.no_grad():
36     model.eval()
37     label = torch.argmax(model(sentence.view((1,-1))),dim=1)
38
39 # 加载之前预处理过程 保存的索引到类别标签的映射字典
40 with open('./data/index2labels.json') as f:
41     index2labels = json.load(f)
42 #输出新文本的类别标签
43 print(index2labels[str(label.item())])

```

#### 4. 帮助函数

为了方便他人使用, 程序中还应当提供一个帮助函数, 用于说明函数是如何使用。程序的命令行接口中有众多参数, 如果手动用字符串表示不仅复杂, 而且后期修改 config 文件时, 还需要修改对应的帮助信息, 十分不便。这里使用了 Python 标准库中的 inspect 方法, 可以自动获取 config 的源代码。help 的代码如下:

```

1 def help():
2
3     '''
4     打印帮助的信息:  python file.py help
5     '''
6
7     print('')
8     usage : python {0} <function> [--args=value,]
9     <function> := train | test | help
10    example:
11    python {0} train --model='TextCNN' --lr=0.01
12    python {0} test --text='xxxx'
13    python {0} help
14    avai    able
15    args:  ''.format(__file__)
16
17    from inspect import getsource

```

```
18 | source = (getsource(opt.__class__))
19 | print(source)
```

## 6 项目使用方式

1. 下载 THUCnews 数据集 ([完整数据集压缩包下载](#)), 解压缩, 并修改 config.py:

```
1 | #下载数据集 解压缩后得到的文件夹所在的路径
2 | data_root = '/Users/apple/Downloads/THUCNews-1'
```

2. 下载预训练词向量 ([项目所使用预训练词向量](#)), 解压缩, 并修改 config.py:

```
1 | #预训练词向量相关超参数
2 | word_vector_path = '/Users/apple/Downloads/sgns.sogou.word' #下载的预训练词向量 解压后的文件所在的路径
```

3. 进入 data 目录下, 运行 dataset.py, 对数据进行预处理并生成必要的中间文件 (数据集非常大, 此过程需要 3-4 小时):

```
1 | python dataset.py
```

4. 之后便可以训练模型 (在 main.py 所在的目录下运行):

```
1 | 可以在命令行指定新的超参数覆盖默认超参数配置 不然将使用config.py中的默认超参数
2 | # 训练模型 单GPU
3 | CUDA_VISIBLE_DEVICES=5 nohup python -u main.py train
4 | --model='TextCNN'
5 | --lr=0.01
6 | --batch-size=256
7 | --max-epoch = 20 >zdz.log 2>&1 &
8 |
9 | # 训练模型 多GPU
10 | CUDA_VISIBLE_DEVICES=0,1,2,5 nohup python -u main.py train
11 | --model='DPCNN'
12 | --drop_prop_dpcnn=0.2
13 | --batch-size=256
14 | --max-epoch = 10 >zdz.log 2>&1 &
15 |
16 | # 打印帮助信息
17 | python main.py help
```

## 7 实验结果与分析

### 7.1 实验结果

训练各个模型的超参数均采用默认超参数, 具体的配置在 config.py 中。各个模型在测试集上的 weight f1-score 值如表2所示:

表 2: 实验结果

模型	测试集 weighted f1-score	备注
FastText	92.6%	Baseline
TextCNN	93.8%	
RCNN	93.0%	
多层双向 LSTM	94.4%	
多层双向 LSTM with Attention	94.7%	
DPCNN	95.4%	

## 7.2 实验分析

- FastText 模型是我们的 Baseline，在测试集上的 weighted f1-score 为 92.6%
- DPCNN 模型在测试集上取得了最好的性能，其 weighted f1-score 为 95.4%
- Attention 机制在分类问题上的效果不是很明显，多层双向 LSTM 和多层双向 LSTM with Attention 的性能几乎差不多
- 比较意外地是 RCNN 模型 (结合 RNN 与 CNN) 并没有取得预期的性能，相比 FastText 提升不大，可能模型细节还需要细究
- TextCNN 是”性价比”最高的模型，模型比较简单，训练很容易，但效果非常不错。
- 上述实验结果均基于默认的超参数配置，调参可能会取得更好的性能
- 基于 RNN 的模型相对基于 CNN 的模型更难训练 (相同配置下，训练时间更长)，但效果并没有显著优势。基于 RNN 的模型训练一定轮数后可能会出现梯度爆炸，注意使用梯度剪切技巧
- 深层网络 (DPCNN) 相对于浅层网络 (TextCNN) 更难训练，但效果提升比较明显。深层网络要注意缓解梯度消失现象 (比如，使用残差连结)。

## 8 预测与网页 Demo

### 8.1 通过命令行执行预测

模型训练完成后，便可以对新闻文本执行预测过程，预测其对应的主题标签 (在 main.py 所在的目录下运行)：

```

1 # 通过命令行运行预测过程
2 python main.py predict
3     --model='RCNN' #指定预测所使用的model
4     --text='众所周知，一支球队想要夺冠，超级巨星必不可少，不过得到超级巨星并不简单，方式无非两种，一是自己培养，这种方式适用于所有球队，二是交易。' #待分类文本

```

## 8.2 通过网页 Demo 执行预测

1) 进入项目目录，运行 manage.py（在 manage.py 所在的目录下运行）：

```
1 | python manage.py runserver
```

```
(base) localhost:NLPProject apple$ python manage.py runserver
Watching for file changes with StatReloader
Performing system checks...

System check identified no issues (0 silenced).

You have 17 unapplied migration(s). Your project may not work properly until
you apply the migrations for app(s): admin, auth, contenttypes, sessions.
Run 'python manage.py migrate' to apply them.

December 20, 2019 - 00:46:57
Django version 3.0.1, using settings 'NLPTextClassification.settings'
Starting development server at http://127.0.0.1:8000/
Quit the server with CONTROL-C.
```

图 9: 运行网页 Demo

2) 在浏览器打开网址：

```
1 | http://127.0.0.1:8000/index #不要忘了index
```

3) 界面如下所示：



图 10: 网页 Demo 界面

把待分类的文本粘贴到上图中的红框中，在蓝框的下拉列表中选择模型，点击上传分析后，在下方的绿框中便可显示分类结果。

4) 支持的新闻类别有 14 个：财经、彩票、房产、股票、家居、教育、科技、社会、时尚、时政、体育、星座、游戏、娱乐。

以下是一些示例（注意：我们的系统是对篇章进行分类，最长不要超过 500 个词。）：

a. 体育：“众所周知，一支球队想要夺冠，超级巨星必不可少，不过得到超级巨星并不简单，方式无非两种，一是自己培养，这种方式适用于所有球队，二是交易，这种方式基本只适用于大市场球队——事实就是，30 支球队之间并非完全公平，超级巨星依然更愿意前往大城市。”

b. 体育：“德国球星萨内告知拜仁慕尼黑，他希望在二月份转会加盟，但据英国媒体报道，曼城方面对他的要价高达 1 亿英镑。尽管萨内自从 8 月份以来就一直伤停，但拜仁仍对他感兴趣。曼城方面也清楚萨内想走，但他们的立场是，买家出价合适才会放人，而瓜迪奥拉对他的要价是 1 亿英镑，这笔钱将用来买进新球员补充阵容。”



c. 娱乐：“有网友晒出了范冰冰现身好莱坞华裔导演温子仁新作《恶毒》的杀青晚宴现场的照片。照片中范冰冰戴着暗绿色帽子，穿着黑色皮衣，留着大波浪长发，五官精致笑容甜美，与众主创合照站 C 位，很有排面。”

d. 娱乐：“北京时间 12 月 20 日消息，据香港媒体报道，电影《急先锋》日前在北京举行发布会，导演唐季礼联同成龙、杨洋、艾伦、朱正廷、母其弥雅（MIYA）等主角齐齐亮相，并分享台前幕后的故事。其中导演唐季礼特别提到，在拍摄过程中快艇不慎被石头掀翻，一下将成龙扣在水下，把自己吓哭了！今年成龙大哥又有新作，与老拍档唐季礼导演联手打造新片《急先锋》，并且找来一班新血组成中国版“复仇者”，在大年初一与观众贺岁，让影迷万分期待，《急先锋》日前在北京举行发布会。”

e. 游戏：“尽管新英雄厄斐琉斯还没有正式登陆各大服务器，但拳头已经通过邮件向国外网友发送了另一名新英雄的神秘技能卡片，卡片描绘出其他的英雄被某个技能打中的效果，但是这个打击效果的来源却是未知的。之前拳头曾在英雄制作大纲写道：“在厄斐琉斯之后的英雄会是一位来自艾欧尼亚的斗士。这名英雄在打斗中茁壮成长，在受到过对方强烈的击打后，他（她）会狂笑，并且将所有受到的挑衅全部释放到对方的脸上。如果你喜欢用拳头说话，喜欢致命搏击，或者喜欢在激烈的战斗中对对手的头打得粉碎，他（她）可能是你的本命英雄。”这是铁拳要进入联盟了吗”

## 9 所需环境

```
1 | numpy      >=1.16.2
2 | json       >=2.0.9
3 | jieba      >=0.39
4 | torch      >=1.1.0
5 | torchtext  >=0.4.0
6 | sklearn    >=0.20.3
7 | django     >=3.0.1
8 | fire       >=0.2.1
```

## 10 总结

1) 本项目并没有解决训练集样本类别分布不均衡的问题，之后会考虑解决这个问题。

2) 尽管各个模型在测试集上的准确率都在 90% 以上，但在预测时，准确率会有一定程度的下降，原因在于真实（预测）数据的分布和训练集数据的分布不同。

3) 本项目主要实现了一些基于 CNN、RNN（Attention）的文本分类模型，并没有实现一些基于预训练语言模型（如 Bert、XLNet 等）的分类模型，之后会逐步完善。