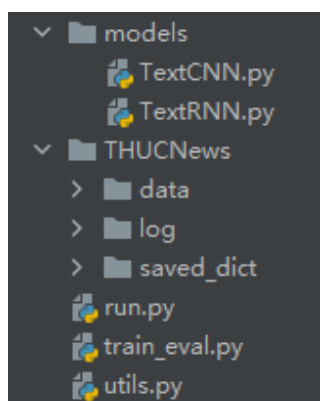


卷积神经网络实现THUCNews新闻文本分类

代码结构

整体代码结构如下图所示：



点击run.py文件，直接运行。可以手动调节参数以及更换模型

1数据集

本文采用的数据集属于清华NLP组提供的THUCNews新闻文本分类数据集的一个子集（原始的数据集大约74万篇文档，训练起来需要花较长的时间）。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。

下载的数据放入THUCNews/data目录中。本次训练使用了其中的10个分类，每个分类6500条，总共65000条新闻数据。

类别如下：

体育, 财经, 房产, 家居, 教育, 科技, 时尚, 时政, 游戏, 娱乐

数据集划分如下：

- 训练集：5000*10
- 验证集：500*10
- 测试集：1000*10

从原始数据集生成子集的过程请参看helper下的两个脚本。其中copy_data.sh用于从每个分类拷贝6500个文件，cnews_group.py用于将多个文件整合到一个文件中。执行该文件后，得到三个数据文件：

- train.txt: 训练集(50000条)

- dev.txt: 验证集(5000条)
- test.txt: 测试集(10000条)

测试集示例：

中华女子学院：本科层次仅1专业招男生 3
 两天价网站背后重重迷雾：做个网站究竟要多少钱 4
 东5环海棠公社230-290平2居准现房98折优惠 1
 卡佩罗：告诉你德国脚生猛的原因 不希望英德战踢点球
 82岁老太为学生做饭扫地44年获授港大荣誉院士 5
 记者回访地震中可乐男孩：将受邀赴美国参观 5
 冯德伦徐若 隔空传情 默认其是女友 9
 传郭晶晶欲落户香港战伦敦奥运 装修别墅当婚房 1
 《赤壁OL》攻城战诸侯战硝烟又起 8
 “手机钱包”亮相科博会 4
 上海2010上半年四六级考试报名4月8日前完成 3
 李永波称李宗伟难阻林丹取胜 透露谢杏芳有望出战 7
 3岁女童下体红肿 自称被幼儿园老师用尺子捅伤 5
 金证顾问：过山车行情意味着什么 2
 谁料地王如此虚 1
 《光环5》Logo泄露 Kinect版几无悬念 8
 海淀区领秀新硅谷宽景大宅预计10月底开盘 1
 柴志坤：土地供应里不断从紧 地价难现07水平(图) 1
 伊达传说EDDA Online 8
 三联书店建起书香巷 4
 宇航员尿液堵塞国际空间站水循环系统 4
 研究发现开车技术差或与基因相关 6

2预处理

调用加载数据的函数返回预处理的数据

```
Python

1 def build_vocab(file_path, tokenizer, max_size, min_freq):
2     vocab_dic = {}
3     with open(file_path, 'r', encoding='UTF-8') as f:
4         for line in tqdm(f):
5             lin = line.strip()
6             if not lin:
7                 continue
8             content = lin.split('\t')[0]
9             for word in tokenizer(content):
10                 vocab_dic[word] = vocab_dic.get(word, 0) + 1
```

```

11     vocab_list = sorted([_ for _ in vocab_dic.items() if _[1] >=
min_freq], key=lambda x: x[1], reverse=True)[:max_size]
12     vocab_dic = {word_count[0]: idx for idx, word_count in
enumerate(vocab_list)}
13     vocab_dic.update({UNK: len(vocab_dic), PAD: len(vocab_dic) + 1})
14     return vocab_dic
15 def build_dataset(config, ues_word):
16     if ues_word:
17         tokenizer = lambda x: x.split(' ') # 以空格隔开, word-level
18     else:
19         tokenizer = lambda x: [y for y in x] # char-level
20     if os.path.exists(config.vocab_path):
21         vocab = pickle.load(open(config.vocab_path, 'rb'))
22     else:
23         vocab = build_vocab(config.train_path, tokenizer=tokenizer,
max_size=MAX_VOCAB_SIZE, min_freq=1)
24         pickle.dump(vocab, open(config.vocab_path, 'wb'))
25     print(f"Vocab size: {len(vocab)}")
26
27     def load_dataset(path, pad_size=32):
28         contents = []
29         with open(path, 'r', encoding='UTF-8') as f:
30             for line in tqdm(f):
31                 lin = line.strip()
32                 if not lin:
33                     continue
34                 content, label = lin.split('\t')
35                 words_line = []
36                 token = tokenizer(content)
37                 seq_len = len(token)
38                 if pad_size:
39                     if len(token) < pad_size:
40                         token.extend([vocab.get(PAD)] * (pad_size -
len(token)))
41                 else:
42                     token = token[:pad_size]
43                     seq_len = pad_size
44                 # word to id
45                 for word in token:
46                     words_line.append(vocab.get(word, vocab.get(UNK)))
47                 contents.append((words_line, int(label), seq_len))
48     return contents # [([...], 0), ([...], 1), ...]
49 train = load_dataset(config.train_path, config.pad_size)
50 dev = load_dataset(config.dev_path, config.pad_size)
51 test = load_dataset(config.test_path, config.pad_size)
52 return vocab, train, dev, test

```

以上代码：

1. 获取分词方式（单字或者单词，这里使用单字）
2. 获取字典类型的词汇表（key=字，value=索引）
3. 获取三个数据集分词之后的索引列表（padding之后长度固定为max_size）

然后将数据封装到迭代器中

Python

```
1 class DatasetIterater(object):
2     def __init__(self, batches, batch_size, device):
3         self.batch_size = batch_size
4         self.batches = batches
5         self.n_batches = len(batches) // batch_size
6         self.residue = False # 记录batch数量是否为整数
7         if len(batches) % self.n_batches != 0:
8             self.residue = True
9         self.index = 0
10        self.device = device
11    def _to_tensor(self, datas):
12        x = torch.LongTensor([_[0] for _ in datas]).to(self.device)
13        y = torch.LongTensor([_[1] for _ in datas]).to(self.device)、
14        # pad前的长度(超过pad_size的设为pad_size)
15        seq_len = torch.LongTensor([_[2] for _ in datas]).to(self.device)
16        return (x, seq_len), y
17    def __next__(self):
18        if self.residue and self.index == self.n_batches:
19            batches = self.batches[self.index * self.batch_size:
20len(self.batches)]
21            self.index += 1
22            batches = self._to_tensor(batches)
23            return batches
24        elif self.index > self.n_batches:
25            self.index = 0
26            raise StopIteration
27        else:
28            batches = self.batches[self.index * self.batch_size: (self.index +
291) * self.batch_size]
30            self.index += 1
31            batches = self._to_tensor(batches)
32            return batches
33    def __iter__(self):
34        return self
35    def __len__(self):
36        if self.residue:
37            return self.n_batches + 1
38        else:
39            return self.n_batches
```

定义DatasetIterater类，并传入需要封装的数据以及需要的batch尺寸或长度。在该类中会对数据进行张量转换。关键是重写__next__()、iter()、len()三个方法。

next () 返回每个batch的张量数据

iter()迭代

len()返回根据数据总样本与batch尺寸计算出来的batch个数

预处理之后的数据只需要通过for循环就可以一次获取一个batch的张量数据

3定义CNN模型

首先将网络模型参数设置封装成类：

Python

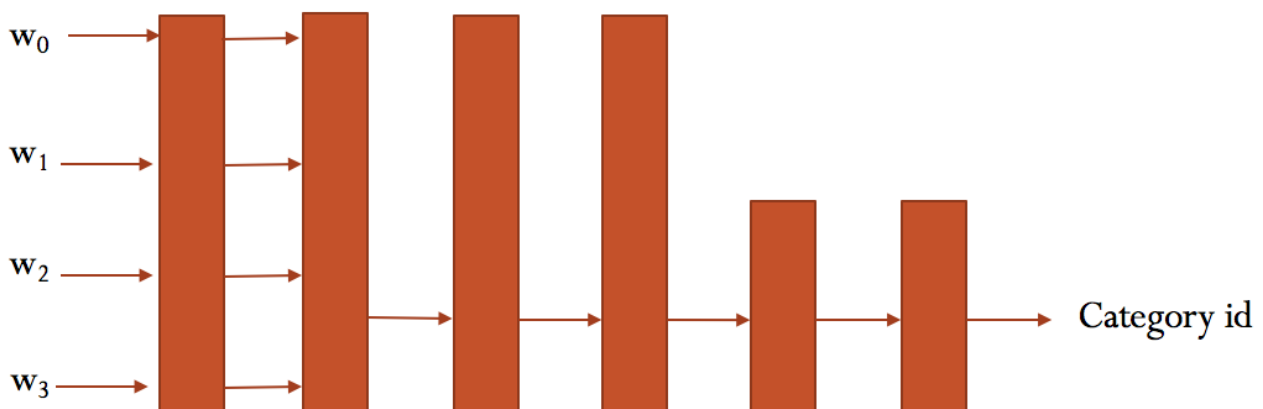
```
1 class configure:
2     def __init__(self):
3         self.dropout = 0.5                                # 随机失活
4         self.require_improvement = 1000                    # 若超过1000batch效果还没提升，则提前结束训练
5         self.num_classes = len(self.class_list)            # 类别数
6         self.n_vocab = 0                                    # 词表大小，在运行时赋值
7         self.num_epochs = 20                                # epoch数
8         self.batch_size = 128                                # mini-batch大小
9         self.pad_size = 32                                  # 每句话处理成的长度(短填长切)
10        self.learning_rate = 1e-3                           # 学习率
11        self.embed = self.embedding_pretrained.size(1)\
12            if self.embedding_pretrained is not None else 300 # 字向量维度
13        self.filter_sizes = (2, 3, 4)                       # 卷积核尺寸
14        self.num_filters = 256                                # 卷积核数量(channels数)
```

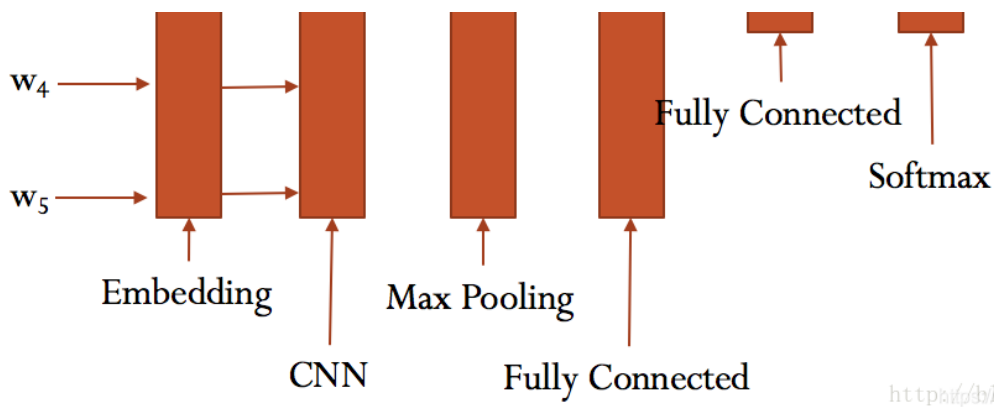
然后定义模型：

Python

```
1 class Model(nn.Module):
2     def __init__(self, config):
3         super(Model, self).__init__()
4         if config.embedding_pretrained is not None:
5             self.embedding =
nn.Embedding.from_pretrained(config.embedding_pretrained, freeze=False)
6         else:
7             self.embedding = nn.Embedding(config.n_vocab, config.embed,
padding_idx=config.n_vocab - 1)
8             self.convs = nn.ModuleList(
9                 [nn.Conv2d(1, config.num_filters, (k, config.embed)) for k in
config.filter_sizes])
10            self.dropout = nn.Dropout(config.dropout)
11            self.fc = nn.Linear(config.num_filters * len(config.filter_sizes),
config.num_classes)
12
13        def conv_and_pool(self, x, conv):
14            x = F.relu(conv(x)).squeeze(3)
15            x = F.max_pool1d(x, x.size(2)).squeeze(2)
16            return x
17        def forward(self, x):
18            #print (x[0].shape)
19            out = self.embedding(x[0])
20            out = out.unsqueeze(1)
21            out = torch.cat([self.conv_and_pool(out, conv) for conv in
self.convs], 1)
22            out = self.dropout(out)
23            out = self.fc(out)
24            return out
```

该模型数据流向图：





整体模型打印如下：

Python

```
1 <bound method Module.parameters of Model(
2   (embedding): Embedding(4762, 300)
3   (convs): ModuleList(
4     (0): Conv2d(1, 256, kernel_size=(2, 300), stride=(1, 1))
5     (1): Conv2d(1, 256, kernel_size=(3, 300), stride=(1, 1))
6     (2): Conv2d(1, 256, kernel_size=(4, 300), stride=(1, 1))
7   )
8   (dropout): Dropout(p=0.5)
9   (fc): Linear(in_features=768, out_features=10, bias=True)
10 )>
```

4训练与验证及测试：

Plain Text

```
1 Epoch [1/20]
2 Iter:      0,  Train Loss:   2.3,  Train Acc: 12.50%,  Val Loss:   2.7,  Val
  Acc: 10.00%   Time: 0:00:04 +
```


Acc: 10.00%, Time: 0:00:04 *

3	Iter: 100,	Train Loss: 0.75,	Train Acc: 70.31%,	Val Loss: 0.69,	Val Acc: 78.74%,	Time: 0:00:40 *
4	Iter: 200,	Train Loss: 0.69,	Train Acc: 76.56%,	Val Loss: 0.55,	Val Acc: 83.48%,	Time: 0:01:18 *
5	Iter: 300,	Train Loss: 0.47,	Train Acc: 82.81%,	Val Loss: 0.49,	Val Acc: 84.66%,	Time: 0:01:54 *
6	Iter: 400,	Train Loss: 0.73,	Train Acc: 78.12%,	Val Loss: 0.47,	Val Acc: 85.48%,	Time: 0:02:31 *
7	Iter: 500,	Train Loss: 0.39,	Train Acc: 87.50%,	Val Loss: 0.44,	Val Acc: 86.33%,	Time: 0:03:08 *
8	Iter: 600,	Train Loss: 0.49,	Train Acc: 84.38%,	Val Loss: 0.43,	Val Acc: 86.58%,	Time: 0:03:45 *
9	Iter: 700,	Train Loss: 0.5,	Train Acc: 83.59%,	Val Loss: 0.41,	Val Acc: 87.10%,	Time: 0:04:23 *
10	Iter: 800,	Train Loss: 0.47,	Train Acc: 84.38%,	Val Loss: 0.39,	Val Acc: 87.79%,	Time: 0:05:00 *
11	Iter: 900,	Train Loss: 0.43,	Train Acc: 86.72%,	Val Loss: 0.38,	Val Acc: 88.16%,	Time: 0:05:37 *
12	Iter: 1000,	Train Loss: 0.35,	Train Acc: 87.50%,	Val Loss: 0.39,	Val Acc: 87.94%,	Time: 0:06:14
13	Iter: 1100,	Train Loss: 0.42,	Train Acc: 89.84%,	Val Loss: 0.38,	Val Acc: 88.47%,	Time: 0:06:50 *
14	Iter: 1200,	Train Loss: 0.35,	Train Acc: 86.72%,	Val Loss: 0.36,	Val Acc: 88.99%,	Time: 0:07:27 *
15	Iter: 1300,	Train Loss: 0.44,	Train Acc: 88.28%,	Val Loss: 0.37,	Val Acc: 88.73%,	Time: 0:08:04
16	Iter: 1400,	Train Loss: 0.48,	Train Acc: 85.94%,	Val Loss: 0.36,	Val Acc: 88.92%,	Time: 0:08:41 *
17	Epoch [2/20]					
18	Iter: 1500,	Train Loss: 0.39,	Train Acc: 90.62%,	Val Loss: 0.35,	Val Acc: 89.31%,	Time: 0:09:18 *
19	Iter: 1600,	Train Loss: 0.31,	Train Acc: 86.72%,	Val Loss: 0.35,	Val Acc: 89.06%,	Time: 0:09:54
20	Iter: 1700,	Train Loss: 0.34,	Train Acc: 92.19%,	Val Loss: 0.35,	Val Acc: 89.41%,	Time: 0:10:31 *
21	Iter: 1800,	Train Loss: 0.29,	Train Acc: 92.97%,	Val Loss: 0.37,	Val Acc: 88.60%,	Time: 0:11:08
22	Iter: 1900,	Train Loss: 0.38,	Train Acc: 89.06%,	Val Loss: 0.35,	Val Acc: 89.43%,	Time: 0:11:45 *
23	Iter: 2000,	Train Loss: 0.32,	Train Acc: 88.28%,	Val Loss: 0.34,	Val Acc: 89.41%,	Time: 0:12:22 *
24	Iter: 2100,	Train Loss: 0.32,	Train Acc: 89.06%,	Val Loss: 0.35,	Val Acc: 89.37%,	Time: 0:12:58
25	Iter: 2200,	Train Loss: 0.22,	Train Acc: 90.62%,	Val Loss: 0.34,	Val Acc: 89.44%,	Time: 0:13:35 *
26	Iter: 2300,	Train Loss: 0.39,	Train Acc: 91.41%,	Val Loss: 0.34,	Val Acc: 89.62%,	Time: 0:14:12 *

```

27 Iter: 2400, Train Loss: 0.28, Train Acc: 93.75%, Val Loss: 0.34, Val
   Acc: 89.54%, Time: 0:14:49
28 Iter: 2500, Train Loss: 0.21, Train Acc: 92.97%, Val Loss: 0.33, Val
   Acc: 90.02%, Time: 0:15:26 *
29 Iter: 2600, Train Loss: 0.34, Train Acc: 89.06%, Val Loss: 0.33, Val
   Acc: 89.90%, Time: 0:16:03
30 Iter: 2700, Train Loss: 0.26, Train Acc: 91.41%, Val Loss: 0.33, Val
   Acc: 89.76%, Time: 0:16:39
31 Iter: 2800, Train Loss: 0.42, Train Acc: 85.94%, Val Loss: 0.34, Val
   Acc: 89.52%, Time: 0:17:16
32 Epoch [3/20]
33 Iter: 2900, Train Loss: 0.34, Train Acc: 89.84%, Val Loss: 0.33, Val
   Acc: 89.99%, Time: 0:17:53 *
34 Iter: 3000, Train Loss: 0.27, Train Acc: 91.41%, Val Loss: 0.33, Val
   Acc: 89.70%, Time: 0:18:29
35 Iter: 3100, Train Loss: 0.3, Train Acc: 89.06%, Val Loss: 0.34, Val
   Acc: 89.83%, Time: 0:19:06
36 Iter: 3200, Train Loss: 0.4, Train Acc: 90.62%, Val Loss: 0.33, Val
   Acc: 90.00%, Time: 0:19:43
37 Iter: 3300, Train Loss: 0.37, Train Acc: 89.84%, Val Loss: 0.33, Val
   Acc: 90.12%, Time: 0:20:20 *
38 Iter: 3400, Train Loss: 0.32, Train Acc: 89.06%, Val Loss: 0.33, Val
   Acc: 90.07%, Time: 0:20:57
39 Iter: 3500, Train Loss: 0.19, Train Acc: 92.97%, Val Loss: 0.33, Val
   Acc: 89.78%, Time: 0:21:35
40 Iter: 3600, Train Loss: 0.14, Train Acc: 95.31%, Val Loss: 0.33, Val
   Acc: 89.74%, Time: 0:22:12
41 Iter: 3700, Train Loss: 0.29, Train Acc: 89.84%, Val Loss: 0.33, Val
   Acc: 89.74%, Time: 0:22:49
42 Iter: 3800, Train Loss: 0.28, Train Acc: 88.28%, Val Loss: 0.33, Val
   Acc: 90.11%, Time: 0:23:25
43 Iter: 3900, Train Loss: 0.32, Train Acc: 87.50%, Val Loss: 0.34, Val
   Acc: 89.73%, Time: 0:24:02
44 Iter: 4000, Train Loss: 0.28, Train Acc: 89.84%, Val Loss: 0.33, Val
   Acc: 89.97%, Time: 0:24:39
45 Iter: 4100, Train Loss: 0.26, Train Acc: 90.62%, Val Loss: 0.33, Val
   Acc: 90.25%, Time: 0:25:16
46 Iter: 4200, Train Loss: 0.35, Train Acc: 87.50%, Val Loss: 0.33, Val
   Acc: 90.04%, Time: 0:25:53

```

```
Test Loss: 0.31, Test Acc: 90.39%
```

```
Precision, Recall and F1-Score...
```

```
precision recall f1-score support
```

```

0.9999 0.9999 0.9999 1000

```

finance	0.9220	0.8870	0.9042	1000
realty	0.9035	0.9360	0.9194	1000
stocks	0.8483	0.8610	0.8546	1000
education	0.9586	0.9490	0.9538	1000
science	0.8381	0.8750	0.8562	1000
society	0.8885	0.9080	0.8981	1000
politics	0.9206	0.8460	0.8817	1000
sports	0.9322	0.9620	0.9469	1000
game	0.9154	0.9090	0.9122	1000
entertainment	0.9179	0.9060	0.9119	1000
accuracy			0.9039	10000
macro avg	0.9045	0.9039	0.9039	10000
weighted avg	0.9045	0.9039	0.9039	10000

https://blog.csdn.net/weixin_45968656

```
Confusion Matrix...
[[887  19  53   0  12  10   7   6   2   4]
 [  9 936  14   2   4  13   5   3   3  11]
 [ 40  28 861   2  33   4  23   4   4   1]
 [  0   3   2 949   8  13   6   4   2  13]
 [  3   9  32   3 875  13  15   5  32  13]
 [  5  19   3  16  13 908  15   3   6  12]
 [ 10   9  39   9  25  45 846   6   2   9]
 [  1   2   3   1   6   7   1 962   3  14]
 [  4   4   7   4  52   4   1  11 909   4]
 [  3   7   1   4  16   5   0  28  30 901]]
```

https://blog.csdn.net/weixin_45968656

在测试集中的正确率达到90.39%，precision、recall、f1-scores都达到了90%以上。

从混淆矩阵也可以看出分类效果非常优秀。