

PPCA CodeMate项目报告

房诗涵

2023 年 7 月 29 日

1 爬虫

1.1 使用request

使用request发起请求，将响应解析为json（CSDN）或取响应的文本用lxml.etree.HTML转化为DOM树，使用xpath定位获取text信息。

1.2 使用selenium

使用selenium的webdriver模拟网站登录，能用获得渲染后的网页页面，ActionChains实现点击等功能，使用driver.execute_script(js)执行JavaScript实现向下滑动等功能。

同时，在使用过程中发现，使用selenium爬取Stackoverflow，不会受到网站人机验证的限制。

1.3 使用并行

1.3.1 协程

essence_link.py使用协程，利用Queue在协程间通信，在 load_page做网络请求时，切换协程做解析页面快照任务。

1.3.2 线程

Stackoverflow/get_detail.py使用threading实现多线程，用 threading.Semaphore(14)限制最多线程数为14，将访问每一个地址作为一个task被执行。

1.4 使用scrapy

用scrapy框架能用更少的代码实现爬虫操作。

scdnSpider是用scrapy框架实现的CSDN爬虫，主要实现新建spider项目的pipeline.py（对爬取的数据的处理）和spider（爬取数据），然后在setting.py中指定pipeline的组织方式。

在使用过程中，发现spider的域名限制不会将域名外的网址入队，一方面，能筛选掉部分无用网址，但是如果从当前网页获得的目标网址域名发生变化，爬取任务将无法执行。

2 爬虫数据分类器训练

2.1 jieba分词

使用中文分词库jieba，将测试集和训练集的问答对分解为词，并去除停用词。停用词包括中文标点、语气词等。

分词前后示例如下

```
": "在python中用三单引号或三双引号，如在代码中没有赋值给任何变量时，就是多行注释，如果赋值给一个变量，那它就是一个多行字符串，可进行字符串一系列操作。
python 中用 三 单引号 三 双引号 代码 没有 赋值 变量 多行 注释 赋值 一个 变量 一个 多行 字符串 进行 字符串 一系列 操作
程序运行！
```

图 1: 中文文本分词前后

在初步分词后，发现存在部分专业词汇被分解为多个词素，如“结构体”等被拆分为“结构”、“体”。尝试使用jieba.load_userdict('file_name')函数，指定这些词汇为一个整体。尝试过去除待分词的字符串中的所有标点、空格，发现会导致代码块被作为一长串英文字符串，无法被正确区分。

```
供参考 includestdioincludestdlibincludestringincludestdbooltypepedefcharDataTypetypedefstructNodeDataTypetypedefstructNodeLinkStackNodevoidinitStackStackNodeSSlinkNULLSMULLboolpush
数据结构
栈
Iknowgoroutinecanhaveafewblockingactions wonderifagoroutinecancall auserdefinedblockingfunctionlikearegularfunctionAuserdefinedblockingfunctionhasafewstepslikestep1st
Ibelieveyoure right thoughIm unsure of the relationship between virtual and resident memory its possible there some overlap Something to consider you are running 1000000 it appears not 100000
数据结构
栈
Is there any way to change the defer stack For example adding a call to the bottom of the defer stack or removing the last defer placed
The only modification possible to the defer stack is to pop onto it Having said that you could make a defer optional with a variable to early exit from the deferred function Example func foo var skip
数据结构
```

图 2: 忽略字符串标点符号和空格后的分词结果

在jieba分词阶段，如何合理地分出代码块是一个难题。一个想法是通过用代码块的标识符正则匹配区分，如完整的c++代码往往以#include开头，这个方法存在的问题是难以确定代码块的结束标记，而且问答对中代码不一定为完整的c++代码，适用范围

将gensim.models.Word2Vec在所有的数据的分词上训练，获得将词转化为对应向量表示的模型。

```
1 import gensim
2
3 # 用分词结果训练word2vec模型
4 # 根据需求配置 Word2Vec 参数并进行训练
5 model=gensim.models.Word2Vec(sentence,          # 分词结果
6                               vector_size=128,   # 词向量维度
7                               window=25,         # 相关前后文的窗口大小
8                               min_count=10,      # 词频
9                               workers=4,         # 线程数
10                              hs=1,
11                              negative=0,
12                              epochs=20          # 迭代次数
13                              )
14
15 model.save('1.model')
```

再将训练好的word2vec模型将训练集和测试集中的每个问答对转化为向量表示。我使用的转化方式是将每个问答对对在word2vec中出现的分词的向量和取平均。

```
1 from gensim.models import Word2Vec
2
3 # 导入训练好的模型
4 model=Word2Vec.load('1.model')
5
6 def get_sentence_vector(sentence):
7     # 将句子用jieba分词
8     word_list=[word for word in jieba.lcut(sentence) if word not in stopwords]
9     # 所有能找到分词结果的词
10    vectors = [model.wv[word] for word in word_list if word in model.wv]
11    if vectors:
12        return sum(vectors) / len(vectors)
13    else:
14        return average_vector
```

最后使用sklearn库的随机森林模型做分类训练。

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.metrics import accuracy_score
3
4 rf_classifier = RandomForestClassifier(n_estimators=190,
5                                     max_depth=27)
```

```

6 rf_classifier.fit(X_train, Y_train)
7
8 Y_pred = rf_classifier.predict(X_test)
9
10 accuracy = accuracy_score(Y_test, Y_pred)
11 print("Accuracy:", accuracy)

```

在试验过程中发现，word2vec的window参数对结果影响较大，其原因极有可能是适当增加窗口数，有利于word2vec模型学到词语的前后文信息。word2vec的epoch参数对结果影响较大，增加迭代次数能提高模型准确率，但是迭代次数过多后，训练效率不高，还可能出现过拟合现象。增大随机森林的n_estimators参数对预测正确率影响较明显，但是增加到130以上后，结果不再有明显提升。

```

Num: 130 Accuracy: 0.7693817468105987
Num: 140 Accuracy: 0.7654563297350343
Num: 150 Accuracy: 0.7634936211972522
Num: 160 Accuracy: 0.7654563297350343
Num: 170 Accuracy: 0.7723258096172718
Num: 180 Accuracy: 0.7556427870461236
Num: 190 Accuracy: 0.7762512266928361
Num: 200 Accuracy: 0.7595682041216879
Num: 210 Accuracy: 0.7713444553483808
Num: 220 Accuracy: 0.7576054955839058

```

图 5: 改变n_estimators大小对正确率的影响

用jieba分词、gensim.models.Word2Vec转化词向量、sklearn.ensemble.RandomForestClassifier现分类任务，获得的最优正确率在0.77-0.78。

2.3 TfidfVectorizer

2.3.1 使用TfidfVectorizer词向量化

TfidfVectorizer是scikit-learn中的一个文本特征提取器，用于将文本转换为基于TF-IDF（Term Frequency-Inverse Document Frequency）的向量表示。

TF-IDF是一种常用的用于衡量文本中词语重要性的方法。它结合了两个指标：词频（Term Frequency, TF）和逆文档频率（Inverse Document Frequency, IDF）。TF衡量了在单个文档中词语出现的频率，而IDF衡量了该词语在整个文集中的重要程度。

为了保证在使用TfidfVectorizer时对中文词汇仍然能有好的分词效果，使用了自定义的分词工具，借用jieba分词的结果。

用TfidfVectorizer词向量化后的训练集训练随机森林，在测试集上测试正确率。最优结果在0.85-0.86

```
Num: 20 Accuracy: 0.8233562315996075
Num: 30 Accuracy: 0.8204121687929342
Num: 40 Accuracy: 0.844946025515211
Num: 50 Accuracy: 0.8351324828263003
Num: 60 Accuracy: 0.8429833169774289
Num: 70 Accuracy: 0.8537782139352306
Num: 80 Accuracy: 0.8429833169774289
Num: 90 Accuracy: 0.8370951913640824
Num: 100 Accuracy: 0.8498527968596663
Num: 110 Accuracy: 0.845927379784102
```

图 6: 改变n_estimators大小（max_depth=None）对正确率的影响

```
Num: 72 Accuracy: 0.8331697742885181
Num: 74 Accuracy: 0.8410206084396468
Num: 76 Accuracy: 0.845927379784102
Num: 78 Accuracy: 0.831207065750736
Num: 80 Accuracy: 0.8420019627085378
Num: 82 Accuracy: 0.844946025515211
Num: 84 Accuracy: 0.8557409224730128
Num: 86 Accuracy: 0.8527968596663396
Num: 88 Accuracy: 0.8400392541707556
Num: 90 Accuracy: 0.8400392541707556
Num: 92 Accuracy: 0.8341511285574092
```

图 7: 改变max_depth大小（n_estimators=70）对正确率的影响

发现使用TfidfVectorizer词向量化的效果明显优于word2vec，我认为是因为word2vec更关注词的上下文信息，而TfidfVectorizer用IDF衡量了该词语在整个文集中的重要程度，关注到了某些特殊的词在分类中的意义，其词向量化的结果对分类任务更加有利。

2.3.2 网格搜索

定义了网格搜索的参数空间，将TfidfVectorizer和随机森林模型用pipeline相连。创建GridSearchCV对象在训练集上执行网格搜索，获得结果最好的参数组合，使用最佳参数组合的模型在测试集上进行预测。

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.model_selection import GridSearchCV
```

```

3 from sklearn.pipeline import Pipeline
4
5 # 定义要搜索的参数空间
6 param_grid = {
7     'vectorizer__max_features': [None,500,1000,1500,2000],
8     'rf_model__n_estimators': [60,90,110],
9     'rf_model__max_depth': [None,40,80]
10 }
11
12 # 构建机器学习工作流
13 pipeline=Pipeline([
14     ('vectorizer',vectorizer),
15     ('rf_model', rf_model)
16 ])
17
18
19 # 创建GridSearchCV对象, 指定pipeline和参数空间
20 grid_search = GridSearchCV(pipeline, param_grid=param_grid, scoring='accuracy', cv=3,
21                             refit=True)
22
23 # 在训练集上执行网格搜索
24 # 在训练集上划分“训练集”和“测试集”，反复验证，获得结果较好的参数
25 grid_search.fit(X_train, Y_train)

```

用上述参数空间，训练得到的在测试集上的正确率为0.8253189401373896，低于手动调参的正确率。我认为原因可能是使用网格搜索得到的最优参数组合是完全由训练集得出的，但是手动调参的参数是根据在测试集上的测试结果调整的，有“面向数据”之嫌。此外，不应该在训练集上做参数网格搜索，可能出现过拟合现象。

2.4 sentence-bert

sentence-bert模型roberta-base-chinese-extractive-qa，损失函数BatchAllTripletLoss在测试集（QA对和标签）上迭代。使用预训练后的sentence-bert模型词向量化，再使用随机森林训练分类器。

原本想使用pipeline将使用的sentence-bert和随机森林连接，再使用损失函数训练，由于sentence-bert是基于深度学习的模型，并且其训练过程与传统的随机森林不兼容，因此无法直接在pipeline中拼接这两个步骤。

迭代十次后，预测正确率有一定的增长。我认为原因可能是sentence-bert模型在迭代后，能根据标签，学习到将被分为1和0的词区分的特征，因而词向量化的结果更有利

```
Iteration: 100% | 153/153 [03:34<00:00, 1.40s/it]
Epoch: 100% | 1/1 [03:34<00:00, 214.81s/it]
Fine-tuned!
Model saved successfully!
gotten word vec 1!
gotten word vec 2!
gotten word vec 3!
gotten word vec 4!
gotten word vec 5!
gotten word vec 6!
gotten word vec 7!
gotten word vec 8!
gotten word vec 9!
gotten word vec 10!
gotten word vec!
trained!
Accuracy: 0.803729146221786
```

图 8: 迭代一次

```
Iteration: 100% | 153/153 [03:36<00:00, 1.41s/it]
Epoch: 100% | 10/10 [36:01<00:00, 216.16s/it]
Fine-tuned!
Model saved successfully!
gotten word vec 1!
gotten word vec 2!
gotten word vec 3!
gotten word vec 4!
gotten word vec 5!
gotten word vec 6!
gotten word vec 7!
gotten word vec 8!
gotten word vec 9!
gotten word vec 10!
gotten word vec!
trained!
Accuracy: 0.8282630829440629
```

图 9: 迭代十次

于分类。

2.5 bert分类

2.5.1 不改变bert参数

使用了bert-base-chinese的Tokenizer和预训练模型，将bert-base-chinese预训练模型和下游二分类任务连接。

```
1 #定义下游任务模型
2 #自定义操作，继承nn.Module类
3 class Model(torch.nn.Module):
4     def __init__(self, pretrained):
5         super().__init__()
6         # 入通道数，出通道数(0,1)
7         self.linear_layer=torch.nn.Linear(768,2)
8         self.pretrained=pretrained
9
10    #实现模型的功能，实现各个层之间的连接关系的核心
11    def forward(self, input_ids, attention_mask, token_type_ids):
12        with torch.no_grad():
13            # 所有计算得出的tensor的requires_grad都自动设置为False
14            out = self.pretrained(input_ids=input_ids,
15                                  attention_mask=attention_mask,
```



```

16         token_type_ids=token_type_ids)
17
18
19     # self.linear_layer执行forward函数
20     out = self.linear_layer(out.last_hidden_state[:, 0])#用特征中的第0个来分类
21     # 归一化指数函数（多个输出结点）
22     out = out.softmax(dim=1)
23
24     return out

```

选用优化算法和损失函数在训练集上训练模型。

```

1 from transformers import AdamW #优化算法，利于收敛
2
3 optimizer = AdamW(my_model.parameters(), lr=2e-4)#学习模型的参数，lr：学习率
4 criterion = torch.nn.CrossEntropyLoss()#损失函数

```

将训练得的模型在测试集上测试正确率。

在调整参数过程中，发现由于训练集数据量小，适当减小batch_size能提升训练效果，但是过小的batch_size会导致loss出现大幅度的振荡。

由于bert预训练模型最多接受长度为512的文本，而我们的数据集为较长的问答对，直接截取问答对的前510个词无法较好提取出文本信息。为了较好表现问答对信息，我在构建数据集时，提取了Question头部32个词，尾部94个词，Answer的头部128个词，尾部256个词，拼接成字符串。这一个处理对训练效果的提升较明显。

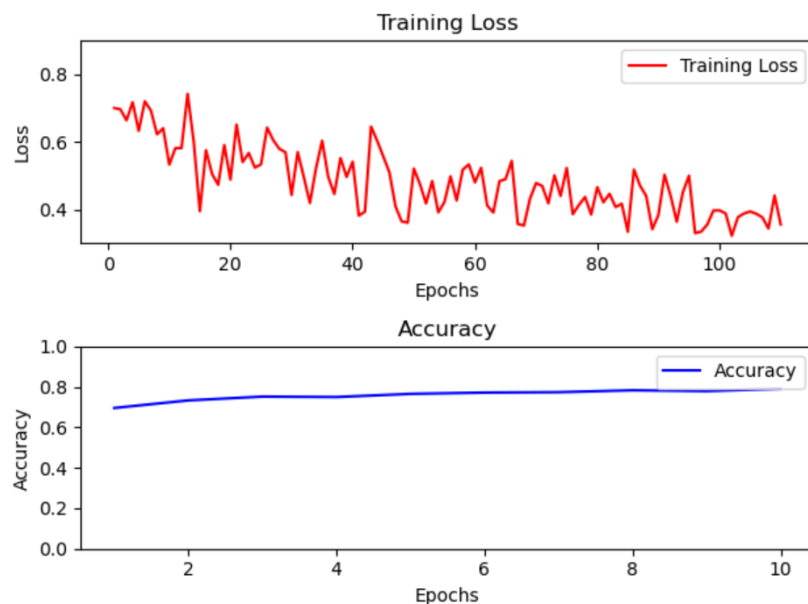
```

1 def __getitem__(self, ind):
2     text=''
3     if(len(self.dataset[ind]['Question'])<=126):
4         text+=self.dataset[ind]['Question']
5     else:
6         text+=self.dataset[ind]['Question'][:32]
7         text+=self.dataset[ind]['Question'][-94:]
8     if(len(self.dataset[ind]['Answer'])<=384):
9         text+=self.dataset[ind]['Answer']
10    else:
11        text+=self.dataset[ind]['Answer'][:128]
12        text+=self.dataset[ind]['Answer'][-256:]
13    label=self.dataset[ind]['Label']
14    return text,label

```

学习率对训练效果的影响很大，在采用学习率为1-e3到1e-5，均出现了loss曲线高频振荡的情况，多次实验后，发现在取学习率为5e-6时，训练效果较好。

初始的模型，冻结了bert预训练模型的参数，仅更新下游模型线性层参数。通过解冻bert预训练模型的最后一层的参数，在训练时将bert预训练模型也加入更新。



[19]:

```
# 测试
test(my_model, Dataset('Test'), False)
```

100% 1/1 [00:00<00:00, 57.76it/s]
Accuracy: 0.8004032258064516

图 10: batch_size=16,lr=5e-6迭代十次训练结果

多次迭代，可以让模型学到更多东西，但是由于训练集小，迭代次数过多，会出现过拟合，在训练集上的正确率在0.9-1.0，但在测试集上的正确率在0.81波动。在下游模型的线性加dropout对过拟合有一定的规避效果。

在迭代30次后，得到最优的正确率在0.83左右。

2.6 投票法

将结果较好的模型作为基座，组建投票分类器。

使用sklearn.ensemble的VotingClassifier，加入的model需要保证为classifier，于是用自定义的CustomClassifier继承BaseEstimator和ClassifierMixin，重写__init__、fit和predict方

法将训练好的模型嵌入其中。

我先选用了在使用相同模型结构训练，不同迭代次数的模型作为参数。每个模型的正确率如下表，最终投票结果在测试集上的结果为0.82。比原先的正确率都略有提升。由于模型结构相同，并不满足投票分类器用优而不同的弱分类器构建强分类器的要求，效果不佳。

随后，我尝试了将TfidfVectorizer接随机森林模型和bert接下游任务模型作为基座，组建投票分类器。

首先尝试了选取两种模型各一个，最终争取率为0.80，均低于单独使用的正确率。产生这个问题的原因是，对仅有两个基座的投票分类器使用了硬分类，当出现两个基座模型评价结果冲突时，无法做出决策。于是，我又尝试加大模型基数，选取了多个参数上有差异的TfidfVectorizer接随机森林模型和bert接下游任务模型，使用硬分类，正确率为0.818，仍然低于单独使用的正确率。

此外，我也尝试了多个参数上有差异的TfidfVectorizer接随机森林模型和bert接下游任务模型，加上有微小差异的权重，使用软分类。最终正确率仍然为0.82左右，无明显变化。原因在于虽然较大的基数能规避掉一些无法决策的情况，但是使用相同结构，不同参数的模型用于投票分类，意义并不大。

为了避免仅用两个模型导致的无法决策，我又尝试了使用微调的sentence-bert词向量化，再使用随机森林分类的模型，这一个模型和原有的两个模型在结构上有一定差异。一共以三类模型为基座，组建投票分类器。最终在测试集上的正确率为0.8469087340529932，与前几次实验相比有明显提高，且正确率高于单用sentence-bert+随机森林和bert+下游任务，略低于TfidfVectorizer+随机森林。

model	Accuracy
bert+下游任务	0.83
TfidfVectorizer+随机森林	0.85
sentence-bert+随机森林	0.82

表 1: 单独用三类模型的正确率

944.3s 108 0.8469087340529932

图 11: 投票分类正确率

A 模型参数

A.1 Word2Vec

gensim.models.Word2Vec参数

```
1 # 配置 Word2Vec 参数并进行训练
2 model=gensim.models.Word2Vec(sentence,
3                               vector_size=128, # 词向量维度
4                               window=25,      # 相关前后文的窗口大小
5                               min_count=10,   # 词频
6                               workers=4,
7                               hs=1,
8                               negative=0,
9                               epochs=20
10                              )
```

使用Word2Vec模型词向量化，随机森林分类，随机森林参数

```
1 rf_classifier = RandomForestClassifier(n_estimators=190,
2                                     max_depth=27)
```

A.2 TfidfVectorizer

```
1 def custom_tokenizer(text):
2     return text.split("。")
3
4 def get_sentence_token(sentence):
5     # 将句子用jieba分词
6     word_list=[word for word in jieba.lcut(sentence) if word not in stopwords]
7     # 将词以特殊的划分依据链接为string，用于作为TfidfVectorizer的样本
8     str=""
9     for word in word_list:
10         str=str+word+"。"
11     return str
12
13 vectorizer = TfidfVectorizer(tokenizer=custom_tokenizer)
14
15 rf_classifier = RandomForestClassifier(n_estimators=70)
```

B github仓库

爬虫、分类器训练等代码详见 <https://github.com/Fangtangtang/CodeMate-PPCA-Project>