

# 基于卷积神经网络的文本情感分类

注：上述题目是网上申请时提交的题目。后来加以引申，最终项目准确的名称应该是“实现基于卷积神经网络和支持向量机的文本情感分类，并讨论数据量对这两种方式准确率的影响。”

## 基于卷积神经网络的文本情感分类

- 内容
- 依赖库
- 具体实现
  - 构建数据集
  - 分词
    - 分词结果
    - 关于最小词长的讨论
  - 卷积神经网络
    - 数据预处理部分
    - 模型构建部分
    - 训练和预测部分
  - 支持向量机
    - 数据预处理部分
    - 模型构建部分
    - 训练和预测部分
  - 准确率比较

## 内容

- 1、通过网络爬虫从淘宝、京东爬取商品评论，对其进行人工分类、标注，得到数据集。
- 2、对数据集进行预处理，进行中文分词、去除停用词，标点符号。
- 3、分别使用“卷积神经网络”和“支持向量机”进行训练和预测，计算准确率。
- 4、改变用于训练的数据量大小，比较这两种方式下的准确率变化情况。

## 依赖库

```
1 python=3.6
2 pandas=0.22.0
3 numpy=1.14.0
4 jieba=0.39
5 gensim=3.4.0
6 scikit-learn=0.19.1
7 keras=2.1.5
8 wordcloud=1.4.1
```

# 具体实现

## 构建数据集

### 1、原始数据爬取

创建商品评论页面的循环链接，使用Python的requests库循环抓取数据，使用正则表达式匹配查询，得到相关文本，将结果保存为.csv文档。

```
1 for url in urls:
2     content = requests.get(url).text
3     # 借助正则表达式使用findall进行匹配查询
4     nickname.extend(re.findall("displayUserNick": "(.*)", content)) #用户名
5     ratecontent.extend(re.findall(re.compile('"rateContent": "(.*)"', "rateDate"), content)) #评论内容
6     ratedate.extend(re.findall(re.compile('"rateDate": "(.*)"', "reply"), content)) #评论时间
```

### 2、标注情感

正面：

638 使用了一个月了，大品牌就是好，满意的没话说，给全五分好评。 [正面]  
639 挺好 就配送太忙了 已经到了好几天 非要在约定日子配送 提前还不能提前， [正面]  
640 好好好好好， [正面]  
641 电视不错 带蓝牙语音功能。如果有需要。下次还会选择暴风。 [正面]  
642 不错， [正面]  
643 电视效果很好， [正面]  
644 感觉很好，性价比很高。晒单： 收起 向左转 向右转， [正面]  
645 秒杀1899，以前没有用过不知好不好用。 [正面]  
646 海信电视，真的很好！性价比很高！推荐大家使用。  
客服好，商家好，送货的也好，上门安装的师傅更好，我是55寸座机，安装调试仅收人民币30元，合规合理。明天再买一个。 [正面]  
647 不错，物美价廉，能再清晰一点就好了！ [正面]  
648 电视机不错，清晰度也好，有WiFi高能，直接网络播放，我是满意的，还有会员送，坐等会员。 [正面]  
649 质量不错，效果好，用一段时间再评论。。。。。。。。。。。。。。。。。。。。 [正面]  
650 京东物流送货快，双十二活动上午抢到的，价钱比平时优惠很多，安装师傅服务态度好，挂墙配挂架没有乱收费现象。这款电视款式新  
买， [正面]  
651 质量非常不错，大气又实惠。服务可以，收费合理， [正面]  
652 电视不错哦，就是手机太差拍的不好，非常的喜欢，看电视流畅，用了几天才来评价的， [正面]  
653 电视非常的好，很喜欢。性价比超高，从快递到安装服务都很满意。全部给满分。 [正面]

负面：

3850 好大呀看一会眼睛就有点不舒服，声音没感觉震撼，清晰度可以。我还没开通会员。 [负面]  
3851 画面不错，这款就是不带DTS和杜比音效，外放很一般，自己配上音响效果好不少。 [负面]  
3852 画面不错非常清晰，图案也细腻，主要音效很好像安上音响一般的效果。开机时间长一点。 [负面]  
3853 画面挺清晰的，该有的功能都有了，而且是高配，就是声音有点失望，相比5年前买的sony差不是一点。支持国产，仍需努力。 [负面]  
3854 画质还行，没有漏光现象。就是侧面看会有泛白，颜色失真。音质也没有宣传的那么好。电视主页面跟电视猫页面很相似。跟电视猫界面来回切换的  
，很容易搞混。电视面板真的很薄，很怕屏幕会变形。蓝牙遥控还是蛮好用的。总体感觉还满意。就是不知能用多少年？ [负面]  
3855 家里人都说好，很清晰，看高清电视很爽，音响一般，看电视够了，有追求的还是再上别的音响。 [负面]  
3856 家里微鲸已经第三台了，16年买的55k1没想到现在涨到四千多了，姐家买的65的，自己新家一直想买个微鲸的奈何涨价太贵。终于等到618发现个  
宜的，一问知道是京东方的屏幕，放心购买，因为家里上一个32的电视就是京东方的，质量不错！昨天拍了这个今天就到了。除了音箱没有我家那  
好以外其他都不错！ [负面]

## 分词

在英文的行文中，单词之间是以空格作为自然分界符的，而中文只是字、句和段能通过明显的分界符来简单划界，但是词却没有一个形式上的分界符。要获得句子的语义特征，必须对句子进行分词。

使用"jieba分词"实现: `words=jieba.lcut(text)`

## 分词结果

使用基于wordcloud生成的词云表示分词结果：

```
1 wc= WordCloud(font_path=font_path,  
2     background_color='white',  
3     max_words=2000,  
4     mask=back_coloring,  
5     max_font_size=100,  
6     random_state=30,  
7     width=600,height=400,margin=5,  
8     relative_scaling=0.5  
9 )  
10 wc.generate(text)
```

使用所以时候收起向左转画面情况这个你微鲸而且快递  
特别感谢感觉有点什么质量一点资源发货速度质量问题  
售后清晰这个非常物流品牌觉得  
很多就是现在知道不是  
遥控器直接今天需要  
已经声音活动电视喜欢  
满意问题性价比价格  
不错电视机  
不能比较客服但是  
画质不行其他有些网络  
色彩一般清晰  
服务效果已经家里还有挂架京东  
自购服务效果已经家里还有挂架京东

## 关于最小词长的讨论

最小词长为1时的词频：

index - Dictionary (9089 elements)

Key	Type	Size	Value
,	int	1	1
的	int	1	2
。	int	1	3
了	int	1	4
	int	1	5
电视	int	1	6
很	int	1	7
！	int	1	8
好	int	1	9
是	int	1	10
也	int	1	11
我	int	1	12
不错	int	1	13
就	int	1	14
买	int	1	15
还	int	1	16
不	int	1	17
都	int	1	18
没有	int	1	19
安装	int	1	20

OK Cancel

最小词长为2时的词频：

counts - Dictionary (8096 elements)

Key	Type	Size	Value
电视	int	1	2942
不错	int	1	1684
没有	int	1	981
安装	int	1	948
微鲸	int	1	916
就是	int	1	878
可以	int	1	810
还是	int	1	759
非常	int	1	729
客服	int	1	687
这个	int	1	669
开机	int	1	606
问题	int	1	571
价格	int	1	558
满意	int	1	543
清晰	int	1	535
屏幕	int	1	534
有点	int	1	527

OK Cancel

由图可见，当最小词长选为1时，高频词出现了大量的标点符号、语气词等无关词语。虽然可以通过通用词将其剔除，但是这么做比较复杂，并且需要构建一个庞大的停用词表。选最小词长为2时，在高频词中明显可以看到很多表示情感的词，比如像“不错”、“满意”等等，已经满足了需求。所以最终选择的最小词长为2。

## 卷积神经网络

代码见 CNN.py

流程：导入数据 -> 数据预处理（分词，统计词频、编号，标签转成独热码） -> 构建模型 -> 训练和预测 -> 计算准确率

### 数据预处理部分

1、分词（同上）

2、统计词频、编码

使用Keras的Tokenizer类 对文本中的词进行统计计数，生成文档词典，以支持基于词典位序生成文本的向量表示。

```

1 tokenizer = Tokenizer(num_words=2000) #初始化Tokenizer
2 tokenizer.fit_on_texts(texts=X_cut) #训练
3 counts=tokenizer.word_counts #词频
4 index=tokenizer.word_index #编号

```

注：在Keras官方文档中，Tokenizer被翻译成“分词器”，这个翻译其实不太恰当。它实现的功能并不是分词，而是统计词频，并且可以根据词频从高到低对词语依次编码为1,2,3.....。

### 3、标签转成独热码

使用Pandas的get\_dummies方法实现。

```

1 y_one_hot = pd.get_dummies(y)
2 y_one_hot_labels = np.asarray(y_one_hot)

```

## 模型构建部分

模型来自博文（Implementing a CNN for Text Classification in TensorFlow :

<http://www.wildml.com/2015/12/implementing-a-cnn-for-text-classification-in-tensorflow/>）。因为要做的只是二分类，所以对其模型中平行的三个卷积层进行了简化，只使用一个卷积层。因为发现过拟合严重，增加了两个Dropout层。

```

1 data_input = Input(shape=[maxlen])
2 word_vec = Embedding(input_dim=num_words+1,
3                       input_length=maxlen,
4                       output_dim=vec_size,
5                       mask_zero=0)(data_input)
6 x = Conv1D(filters=128, kernel_size=[3], strides=1, padding='same', activation='relu')(word_vec)
7 x = GlobalMaxPool1D()(x)
8 x = Dropout(0.1)(x)
9 x = Dense(500, activation='relu')(x)
10 x = Dropout(0.1)(x)
11 x = Dense(output_shape, activation='softmax')(x)
12 model = Model(inputs=data_input, outputs=x)
13 model.compile(loss='categorical_crossentropy',
14               optimizer='adam',
15               metrics=['accuracy'])
16 model.summary()

```

## 训练和预测部分



```

1 #训练
2 model.fit(x=X_train, y=y_train, epochs=3, batch_size=64, verbose=2)
3 #预测
4 y_predict = model.predict(X_test)
5 #转换预测结果
6 y_predict_label = label2tag(predictions=y_predict, y=y)
7 #统计正确率
8 Y_test=label2tag(predictions=y_test,y=y)
9 print(sum([y_predict_label[i] == Y_test[i] for i in range(len(y_predict))]) / len(y_predict))

```

## 支持向量机

代码见 SVM.py

流程：导入数据 -> 数据预处理（分词，词语转为词向量） -> 构建模型 -> 训练和预测 -> 计算准确率

### 数据预处理部分

1、分词，同上

2、词向量生成

使用中文词向量语料库（Shen Li, Zhe Zhao, Renfen Hu, Wensi Li, Tao Liu, Xiaoyong Du, [Analogical Reasoning on Chinese Morphological and Semantic Relations](#), ACL 2018.）中的中文预训练词向量（sgns.weibo.word）对分词结果进行转换。

```

1 model = gensim.models.KeyedVectors.load_word2vec_format(word2vec_loadpath, binary=False)
2 text_vec = [[model[word] for word in text_cut if word in model] for text_cut in X_cut]

```

### 模型构建部分

```

1 model_svc = SVC(C=10, kernel='linear',)
2 #C似乎越大准确率越高: score 0.819672 (C=1) 0.83957 (C=2) 0.850117 (C=3) 0.862997 (C=5)
3 #尝试了所有的kernel, 最后得出linear的准确率最高。

```

### 训练和预测部分

```

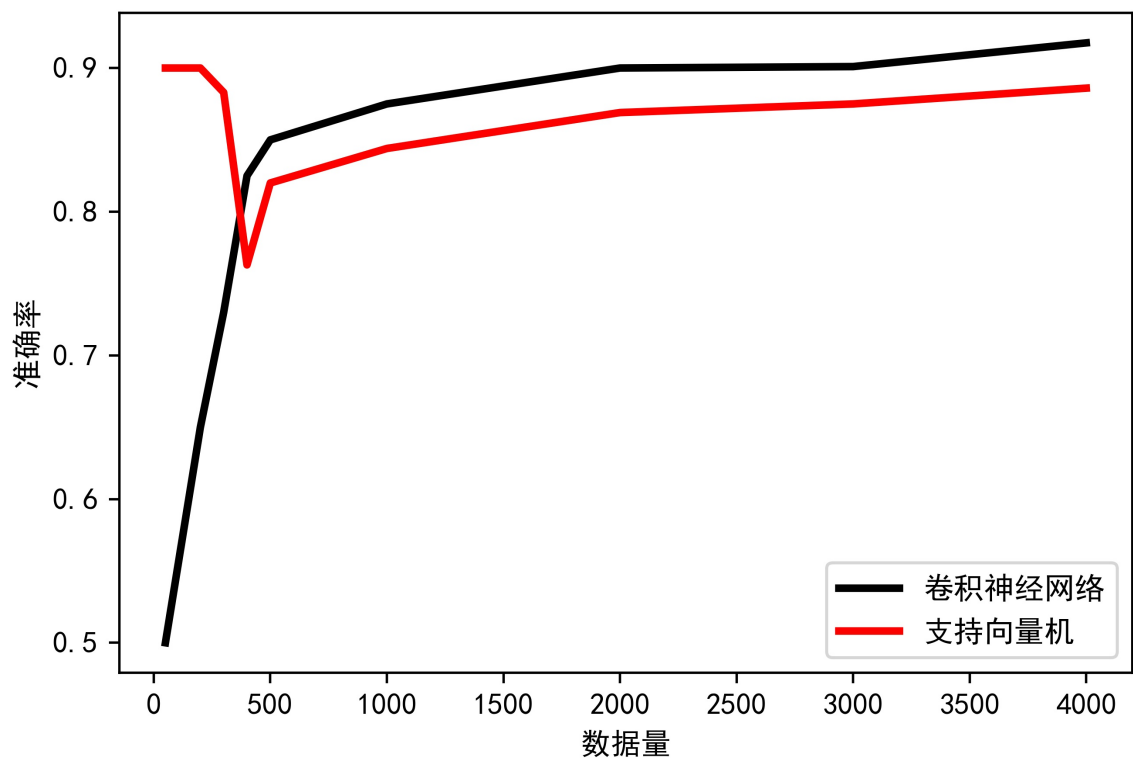
1 model_svc.fit(X=X_train, y=y_train,)
2 y_predict = model_svc.predict(X_test)
3 print(sum(y_predict == np.array(y_test)) / len(y_predict))

```



## 准确率比较

数据量	卷积神经网络	支持向量机
50	0.500	0.900
100	0.550	0.900
200	0.650	0.900
300	0.733	0.883
400	0.825	0.763
500	0.850	0.820
1000	0.875	0.844
2000	0.900	0.869
3000	0.901	0.875
4000	0.9175	0.886



可以看出，卷积神经网络随着数据量的增加，正确率也是不断地提高，最后收敛在0.90左右。而支持向量机的正确率却是经历了先减小后增大的过程。说明对于卷积神经网络，用于训练的数据量越大，效果越好。而对于支持向量机，更重要的可能是参数的设置。