

# 实验1 基于TextCNN的情感分类

---

## 实验1 基于TextCNN的情感分类

### 1.实验背景

### 2.实验目的

### 3.实验步骤

#### 3.1 实验准备

#### 3.2 实验过程

### 4.实验原理-TextCNN

#### 4.1嵌入层 (Embedding Layer)

#### 4.2卷积层 (Convolution Layer)

#### 4.3池化层 (Pooling Layer)

#### 4.4全连接层 (Fully connected layer)

#### 4.5TextCNN的小变种

#### 4.6参数与超参数

### 5.实验结果及分析

#### 5.1实验结果

batch\_size = 64, epoch = 4

batch\_size = 64, epoch = 8

batch\_size = 32, epoch = 16

在线测试

#### 5.2分析

## 1.实验背景

---

文本分类(text classification),又称文档分类(document classification),指的是将一个文档归类到一个或多个类别中的自然语言处理任务。文本分类的应用场景非常广泛,涵盖垃圾邮件过滤、垃圾评论过滤、自动标签、情感分析等任何需要自动归档文本的场合。

情感分析是自然语言处理文本分类任务的应用场景之一,情感分类较为简单,实用性也较强。常见的购物网站、电影网站都可以采集到相对高质量的数据集,也很容易给业务领域带来收益。例如,可以结合领域上下文,自动分析特定类型客户对当前产品的意见,可以分主题分用户类型对情感进行分析,以作针对性的处理,甚至基于此进一步推荐产品,提高转化率,带来更高的商业收益。

本实验主要基于卷积神经网络对电影评论信息进行情感分析,判断其情感倾向。

## 2.实验目的

---

- 理解文本分类的基本流程
- 理解CNN网络在文本任务中的用法
- 掌握MindSpore搭建文本分类模型的方法

## 3.实验步骤

---

### 3.1 实验准备

---

- 使用OBS创建项目文件夹
- 上传实验源码及数据



## 3.2 实验过程

- 同步数据和源码到本地容器

```
1 import moxing as mox
2 # 请替换成自己的obs路径
3 mox.file.copy_parallel(src_url="s3://bucket-mzh/lab1/data/",
  dst_url='./data/')
```

### 1. 数据同步

```
In [5]: import moxing as mox
# 请替换成自己的obs路径
mox.file.copy_parallel(src_url="s3://bucket-mzh/lab1/data/", dst_url='./data/')
print('同步')

同步
```

- 设置超参数，运行环境，预览数据

```
1 from easydict import EasyDict as edict
2
3 cfg = edict({
4     'name': 'movie review',
5     'pre_trained': False,
6     'num_classes': 2,
7     'batch_size': 64,
8     'epoch_size': 4,
9     'weight_decay': 3e-5,
10    'data_path': './data/',
11    'device_target': 'Ascend',
12    'device_id': 0,
13    'keep_checkpoint_max': 1,
14    'checkpoint_path': './ckpt/train_textcnn-4_149.ckpt',
15    'word_len': 51,
16    'vec_length': 40
17 })
```

- 数据预处理，并生成数据集和测试集

```
1 # 数据预览
2 with open("./data/rt-polarity.neg", 'r', encoding='utf-8') as f:
3     print("Negative reviews:")
4     for i in range(5):
5         print("[{0}]:{1}".format(i,f.readline()))
```

```

6 with open("./data/rt-polarity.pos", 'r', encoding='utf-8') as f:
7     print("Positive reivews:")
8     for i in range(5):
9         print("[{0}]:{1}".format(i,f.readline()))
10 class Generator():
11     def __init__(self, input_list):
12         self.input_list=input_list
13     def __getitem__(self,item):
14         return (np.array(self.input_list[item][0],dtype=np.int32),
15                 np.array(self.input_list[item][1],dtype=np.int32))
16     def __len__(self):
17         return len(self.input_list)
18
19
20 class MovieReview:
21     '''
22     影评数据集
23     '''
24     def __init__(self, root_dir, maxlen, split):
25         '''
26         input:
27             root_dir: 影评数据目录
28             maxlen: 设置句子最大长度
29             split: 设置数据集中训练/评估的比例
30         '''
31         self.path = root_dir
32         self.feelMap = {
33             'neg':0,
34             'pos':1
35         }
36         self.files = []
37
38         self.doConvert = False
39
40         mypath = Path(self.path)
41         if not mypath.exists() or not mypath.is_dir():
42             print("please check the root_dir!")
43             raise ValueError
44
45         # 在数据目录中找到文件
46         for root,_,filename in os.walk(self.path):
47             for each in filename:
48                 self.files.append(os.path.join(root,each))
49             break
50
51         # 确认是否为两个文件.neg与.pos
52         if len(self.files) != 2:
53             print("There are {} files in the
root_dir".format(len(self.files)))
54             raise ValueError
55
56         # 读取数据
57         self.word_num = 0
58         self.maxlen = 0
59         self.minlen = float("inf")
60         self.maxlen = float("-inf")
61         self.Pos = []
62         self.Neg = []

```

```

63         for filename in self.files:
64             f = codecs.open(filename, 'r')
65             ff = f.read()
66             file_object = codecs.open(filename, 'w', 'utf-8')
67             file_object.write(ff)
68             self.read_data(filename)
69         self.PosNeg = self.Pos + self.Neg
70
71         self.text2vec(maxlen=maxlen)
72         self.split_dataset(split=split)
73
74     def read_data(self, filePath):
75
76         with open(filePath, 'r') as f:
77
78             for sentence in f.readlines():
79                 sentence = sentence.replace('\n', '')\
80                     .replace('"', '')\
81                     .replace('\\', '')\
82                     .replace('.', '')\
83                     .replace(',', '')\
84                     .replace('[', '')\
85                     .replace(']', '')\
86                     .replace('(', '')\
87                     .replace(')', '')\
88                     .replace(':', '')\
89                     .replace('--', '')\
90                     .replace('-', ' ')\
91                     .replace('\ ', '')\
92                     .replace('0', '')\
93                     .replace('1', '')\
94                     .replace('2', '')\
95                     .replace('3', '')\
96                     .replace('4', '')\
97                     .replace('5', '')\
98                     .replace('6', '')\
99                     .replace('7', '')\
100                    .replace('8', '')\
101                    .replace('9', '')\
102                    .replace('`', '')\
103                    .replace('=', '')\
104                    .replace('$', '')\
105                    .replace('/', '')\
106                    .replace('*', '')\
107                    .replace(';', '')\
108                    .replace('<b>', '')\
109                    .replace('%', '')
110                 sentence = sentence.split(' ')
111                 sentence = list(filter(lambda x: x, sentence))
112                 if sentence:
113                     self.word_num += len(sentence)
114                     self.maxlen = self.maxlen if self.maxlen >=
115 len(sentence) else len(sentence)
116                     self.minlen = self.minlen if self.minlen <=
117 len(sentence) else len(sentence)
118                     if 'pos' in filePath:
119                         self.Pos.append([sentence, self.feelMap['pos']])
120                     else:

```

```

119         self.Neg.append([sentence, self.feelMap['neg']])
120
121     def text2vec(self, maxlen):
122         '''
123         将句子转化为向量
124
125         '''
126         # Vocab = {word : index}
127         self.vocab = dict()
128
129         # self.Vocab['None']
130         for SentenceLabel in self.Pos+self.Neg:
131             vector = [0]*maxlen
132             for index, word in enumerate(SentenceLabel[0]):
133                 if index >= maxlen:
134                     break
135                 if word not in self.Vocab.keys():
136                     self.Vocab[word] = len(self.Vocab)
137                     vector[index] = len(self.Vocab) - 1
138                 else:
139                     vector[index] = self.Vocab[word]
140             SentenceLabel[0] = vector
141         self.doConvert = True
142
143     def split_dataset(self, split):
144         '''
145         分割为训练集与测试集
146
147         '''
148
149         trunk_pos_size = math.ceil((1-split)*len(self.Pos))
150         trunk_neg_size = math.ceil((1-split)*len(self.Neg))
151         trunk_num = int(1/(1-split))
152         pos_temp=list()
153         neg_temp=list()
154         for index in range(trunk_num):
155             pos_temp.append(self.Pos[index*trunk_pos_size:
(index+1)*trunk_pos_size])
156             neg_temp.append(self.Neg[index*trunk_neg_size:
(index+1)*trunk_neg_size])
157             self.test = pos_temp.pop(2)+neg_temp.pop(2)
158             self.train = [i for item in pos_temp+neg_temp for i in item]
159
160             random.shuffle(self.train)
161             # random.shuffle(self.test)
162
163     def get_dict_len(self):
164         '''
165         获得数据集中文字组成的词典长度
166         '''
167         if self.doConvert:
168             return len(self.Vocab)
169         else:
170             print("Haven't finished Text2Vec")
171             return -1
172
173     def create_train_dataset(self, epoch_size, batch_size):
174         dataset = ds.GeneratorDataset(

```

```

175     source=Generator(input_list=self.train),
176                       column_names=["data","label"],
177                       shuffle=False
178                   )
179     #         dataset.set_dataset_size(len(self.train))
180     dataset=dataset.batch(batch_size=batch_size,drop_remainder=True)
181     dataset=dataset.repeat(epoch_size)
182     return dataset
183
184     def create_test_dataset(self, batch_size):
185         dataset = ds.GeneratorDataset(
186
187         source=Generator(input_list=self.test),
188                           column_names=["data","label"],
189                           shuffle=False
190                       )
191     #         dataset.set_dataset_size(len(self.test))
192     dataset=dataset.batch(batch_size=batch_size,drop_remainder=True)
193     return dataset

```

- 模型构建

```

1  class TextCNN(nn.Cell):
2      def __init__(self, vocab_len, word_len, num_classes, vec_length):
3          super(TextCNN, self).__init__()
4          self.vec_length = vec_length
5          self.word_len = word_len
6          self.num_classes = num_classes
7
8          self.unsqueeze = ops.ExpandDims()
9          self.embedding = nn.Embedding(vocab_len, self.vec_length,
embedding_table='normal')
10
11         self.slice = ops.Slice()
12         self.layer1 = self.make_layer(kernel_height=3)
13         self.layer2 = self.make_layer(kernel_height=4)
14         self.layer3 = self.make_layer(kernel_height=5)
15
16         self.concat = ops.Concat(1)
17
18         self.fc = nn.Dense(96*3, self.num_classes)
19         self.drop = nn.Dropout(keep_prob=0.5)
20         self.print = ops.Print()
21         self.reduce_mean = ops.ReduceMax(keep_dims=False)
22
23         def make_layer(self, kernel_height):
24             return nn.SequentialCell(
25                 [
26                     make_conv_layer((kernel_height,self.vec_length)),
27                     nn.ReLU(),
28                     nn.MaxPool2d(kernel_size=(self.word_len-
kernel_height+1,1)),
29                 ]
30             )
31
32         def construct(self,x):

```

```

33         x = self.unsqueeze(x, 1)
34         x = self.embedding(x)
35         x1 = self.layer1(x)
36         x2 = self.layer2(x)
37         x3 = self.layer3(x)
38
39         x1 = self.reduce_mean(x1, (2, 3))
40         x2 = self.reduce_mean(x2, (2, 3))
41         x3 = self.reduce_mean(x3, (2, 3))
42
43         x = self.concat((x1, x2, x3))
44         x = self.drop(x)
45         x = self.fc(x)
46         return x

```

- 性能评估

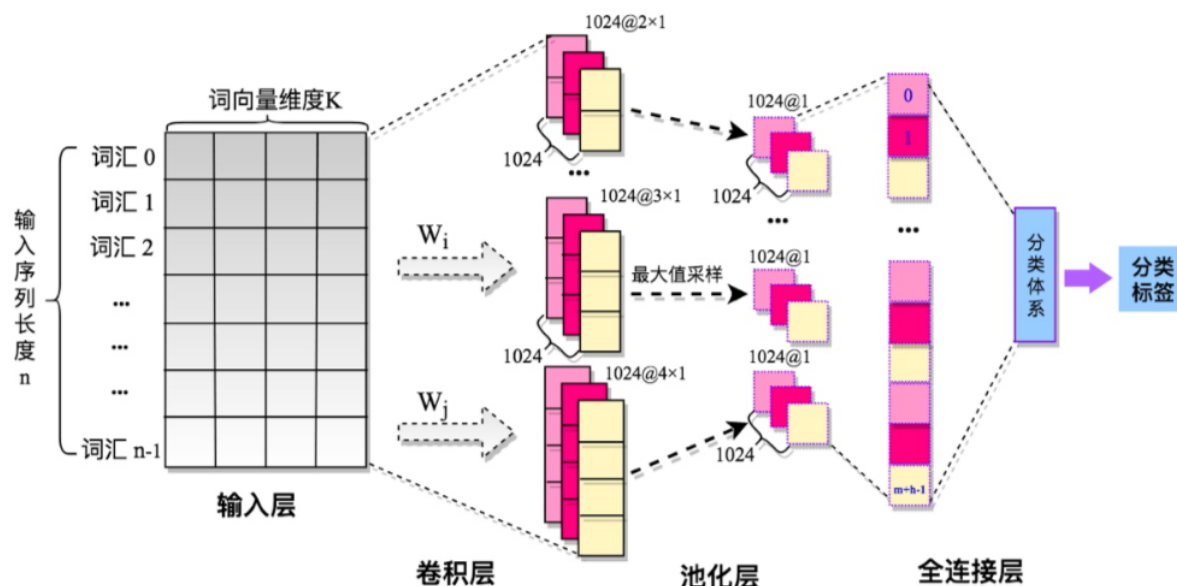
```

1  dataset = instance.create_test_dataset(batch_size=cfg.batch_size)
2  opt = nn.Adam(filter(lambda x: x.requires_grad, net.get_parameters()),
3                  learning_rate=0.001, weight_decay=cfg.weight_decay)
4  loss = nn.SoftmaxCrossEntropyWithLogits(sparse=True)
5  net = TextCNN(vocab_len=instance.get_dict_len(), word_len=cfg.word_len,
6               num_classes=cfg.num_classes, vec_length=cfg.vec_length)
7
8  if checkpoint_path is not None:
9      param_dict = load_checkpoint(checkpoint_path)
10     print("load checkpoint from [{}].".format(checkpoint_path))
11 else:
12     param_dict = load_checkpoint(cfg.checkpoint_path)
13     print("load checkpoint from [{}].".format(cfg.checkpoint_path))
14
15 load_param_into_net(net, param_dict)
16 net.set_train(False)
17 model = Model(net, loss_fn=loss, metrics={'acc': Accuracy()})
18
19 acc = model.eval(dataset)
20 print("accuracy: ", acc)

```

## 4.实验原理-TextCNN

---



## Yoon Kim 2014年提出的Text-CNN <https://blog.csdn.net/asialeebird>

上图很好地诠释了模型的框架。假设我们有一些句子需要对其进行分类。句子中每个词是由 $n$ 维词向量组成的，也就是说输入矩阵大小为 $m \times n$ ，其中 $m$ 为句子长度。CNN需要对输入样本进行卷积操作，对于文本数据，filter不再横向滑动，仅仅是向下移动，有点类似于N-gram在提取词与词间的局部相关性。图中共有三种步长策略，分别是2,3,4，每个步长都有两个filter（实际训练时filter数量会很多）。在不同词窗上应用不同filter，最终得到6个卷积后的向量。然后对每一个向量进行最大化池化操作并拼接各个池化值，最终得到这个句子的特征表示，将这个句子向量丢给分类器进行分类，至此完成整个流程。

### 4.1 嵌入层 (Embedding Layer)

通过一个隐藏层，将 one-hot 编码的词投影到一个低维空间中，本质上是特征提取器，在指定维度中编码语义特征。这样，语义相近的词，它们的欧氏距离或余弦距离也比较近。（作者使用的单词向量是预训练的，方法为fasttext得到的单词向量，当然也可以使用word2vec和GloVe方法训练得到的单词向量）。

### 4.2 卷积层 (Convolution Layer)

在处理图像数据时，CNN使用的卷积核的宽度和高度是一样的，但是在text-CNN中，卷积核的宽度是与词向量的维度一致！这是因为我们输入的每一行向量代表一个词，在抽取特征的过程中，词做为文本的最小粒度。而高度和CNN一样，可以自行设置（通常取值2,3,4,5），高度就类似于n-gram了。由于我们的输入是一个句子，句子中相邻的词之间关联性很高，因此，当我们用卷积核进行卷积时，不仅考虑了词义而且考虑了词序及其上下文（类似于skip-gram和CBOW模型的思想）。

### 4.3 池化层 (Pooling Layer)

因为在卷积层过程中我们使用了不同高度的卷积核，使得我们通过卷积层后得到的向量维度会不一致，所以在池化层中，我们使用1-Max-pooling对每个特征向量池化成一个值，即抽取每个特征向量的最大值表示该特征，而且认为这个最大值表示的是最重要的特征。当我们对所有特征向量进行1-Max-Pooling之后，还需要将每个值给拼接起来。得到池化层最终的特征向量。在池化层到全连接层之前可以加上dropout防止过拟合。

### 4.4 全连接层 (Fully connected layer)



全连接层跟其他模型一样，假设有两层全连接层，第一层可以加上'relu'作为激活函数，第二层则使用softmax激活函数得到属于每个类的概率。

## 4.5 TextCNN的小变种

在词向量构造方面可以有以下几种不同的方式：CNN-rand: 随机初始化每个单词的词向量通过后续的训练去调整。CNN-static: 使用预先训练好的词向量，如word2vec训练出来的词向量，在训练过程中不再调整该词向量。CNN-non-static: 使用预先训练好的词向量，并在训练过程进一步进行调整。CNN-multichannel: 将static与non-static作为两通道的词向量。

## 4.6 参数与超参数

- sequence\_length (Q: 对于CNN, 输入与输出都是固定的, 可每个句子长短不一, 怎么处理? A: 需要做定长处理, 比如定为n, 超过的截断, 不足的补0. 注意补充的0对后面的结果没有影响, 因为后面的max-pooling只会输出最大值, 补零的项会被过滤掉)
- num\_classes (多分类, 分为几类)
- vocabulary\_size (语料库的词典大小, 记为 $|D|$ )
- embedding\_size (将词向量的维度, 由原始的 $|D|$  降维到 embedding\_size)
- filter\_size\_arr (多个不同size的filter)

## 5. 实验结果及分析

### 5.1 实验结果

**batch\_size = 64, epoch = 4**

```
In [18]: checkpoint_path = './ckpt/train_textcnn-4_596.ckpt'
```

```
In [19]: dataset = instance.create_test_dataset(batch_size=cfg.batch_size)
opt = nn.Adam(filter(lambda x: x.requires_grad, net.get_parameters()),
               learning_rate=0.001, weight_decay=cfg.weight_decay)
loss = nn.SoftmaxCrossEntropyWithLogits(sparse=True)
net = TextCNN(vocab_len=instance.get_dict_len(), word_len=cfg.word_len,
              num_classes=cfg.num_classes, vec_length=cfg.vec_length)

if checkpoint_path is not None:
    param_dict = load_checkpoint(checkpoint_path)
    print("load checkpoint from [{}].".format(checkpoint_path))
else:
    param_dict = load_checkpoint(cfg.checkpoint_path)
    print("load checkpoint from [{}].".format(cfg.checkpoint_path))

load_param_into_net(net, param_dict)
net.set_train(False)
model = Model(net, loss_fn=loss, metrics={'acc': Accuracy()})

acc = model.eval(dataset)
print("accuracy: ", acc)

load checkpoint from [./ckpt/train_textcnn-4_596.ckpt].
accuracy: {'acc': 0.7734375}
```

**batch\_size = 64, epoch= 8**

```
load checkpoint from [./ckpt/train_textcnn_1-8_1192.ckpt].
accuracy: {'acc': 0.7587890625}
```

**batch\_size = 32, epoch = 16**

```
load checkpoint from [./ckpt/train_textcnn_2-16_4784.ckpt].
accuracy: {'acc': 0.7339015151515151}
```

## 在线测试

```
In [27]: review_en = "teacher is very handsome"
inference(review_en)
```

Positive comments

Review	Pre_class	Ground_truth
the movie is so boring	Negative	Negative
the movie is so good	Positive	Positive
the man is so crazy	Positive	Positive
There are some boring plots in the movie, but overall it's great	Negative	Positive
The boy is great	Positive	Positive
It looks great, but it's not practical	Positive	Negative
it's so laddish and juvenile , only teenage boys could possibly find it funny	Negative	Negative
a visually flashy but narratively opaque and emotionally vapid exercise in style and mystification	Negative	Negative
if you sometimes like to go to the movies to have fun , wasabi is a good place to start	Positive	Positive

## 5.2分析

了解文本分类任务的基本流程，同时理解卷积网络在文本任务中的使用方法，通过实验也加深了对CNN网络的理解，同时提升了代码实践能力。模型对于前后转折的长文本的评论分类效果不好，丢失前后语义信息，导致分类错误