

系统文档

小组成员及分工：

姓名	学号	任务分工	任务占比
何欣洋（组长）	202228015329003	主要负责搭建系统整体前端和后端的框架，LSTM 模型、数据预处理等	60%
吕卓锦	202228015329004	主要负责搭建 CNN 模型，对模型进行测试，设计网页样式等	40%

任务定义：

情感分析或观点挖掘是对人们对产品、服务、组织、个人、问题、事件、话题及其属性的观点、情感、情绪、评价和态度的计算研究。近年来随着互联网信息的发展和深度学习的出现，使得越来越多的人投入到情感分析的研究当中。情感分析本质上是一个分类任务，常见的分类方法有：二分类或者进行打分（0 到 10）等等

方法描述：

本次作业主要利用 LSTM 和 CNN 对电影评论数据进行情感分析任务。LSTM 是一种特殊的 RNN，它在 RNN 的基础上引入了 gate 机制，从而使得 LSTM 克服了 RNN “记忆力不好的特点”，克服了简单的 RNN 梯度消失或者梯度爆炸的问题。CNN 是一种具有局部连接、权重共享等特性的深层前馈神经网络，虽然 CNN 广泛地用于图像或视频的分析任务，但如果将文本看作一维图像，也可以将 CNN 用于 nlp 任务中。用卷积层替代全连接层可以极大减少要学习的参数数目。

系统框架：

后端：torch 版本：1.13.0，python 版本：3.19.2，django 版本：4.1.4
前端：bootstrap

实现细节：

数据预处理部分：

通过正则表达式将不是字母或数字的字符替换为空格等

```
7 #清理文本, 去标点符号, 转小写
8 def clean_str(string):
9     string = re.sub(r"[^A-Za-z0-9]", " ", string)
10    string = re.sub(r"\s", " \s", string)
11    string = re.sub(r"\ve", " \ve", string)
12    string = re.sub(r"\n\t", " n\t", string)
13    string = re.sub(r'\re', " \re", string)
14    string = re.sub(r'\d', " \d", string)
15    string = re.sub(r'\ll', " \ll", string)
16    string = re.sub(r",", " ", string)
17    string = re.sub(r"!", " ! ", string)
18    string = re.sub(r"\(", " \(", string)
19    string = re.sub(r"\)", " \)", string)
20    string = re.sub(r"?", " ? ", string)
21    string = re.sub(r"\s{2,}", " ", string)
22    string = re.sub(r"\s{2,}", " ", string)
23    string = re.sub(r"sssss", " ", string)
24    return string.strip().lower()
25
```

建立词典：

```
word_count_sort = sorted(word_count.items(), key=lambda item : item[1], reverse=True) # 对词进行排序, 过滤低频词, 只取前MAX_WORD个高频词
word_number = 1
for word in word_count_sort:
    if word[0] not in vocab.keys():
        vocab[word[0]] = len(vocab)
        word_number += 1
    if word_number > MAX_WORD:
        break
return vocab
```

根据词典将句子转换为等长的 tensor

```
1 # 根据vocab将句子转为定长MAX_LEN的tensor
2 def text_transform(sentence_list, vocab):
3     sentence_index_list = []
4     for sentence in sentence_list:
5         sentence_idx = [vocab[token] if token in vocab.keys() else vocab['<UNK>'] for token in tokenizer(sentence)] # 句子分词转为id
6
7         if len(sentence_idx) < MAX_LEN:
8             for i in range(MAX_LEN-len(sentence_idx)): # 对长度不够的句子进行PAD填充
9                 sentence_idx.append(vocab['<PAD>'])
10
11         sentence_idx = sentence_idx[:MAX_LEN] # 取前MAX_LEN长度
12         sentence_index_list.append(sentence_idx)
13     return torch.LongTensor(sentence_index_list) # 将转为idx的词转为tensor
14
```

模型训练并保存模型

```

74 # 模型训练
75 def train(model, train_data, loss_fn, optimizer, vocab, epoch=10, method='LSTM'):
76     print('train model')
77     model = model.to(device)
78     #记录每个epoch的损失
79     record_loss = []
80     record_acc = []
81     # 定义损失函数和优化器
82     loss_func = loss_fn
83     optimizer_func = optimizer
84
85     for epoch in tqdm(range(epoch)):
86         model.train()
87         avg_loss = 0 # 平均损失
88         avg_acc = 0 # 平均准确率
89         for i, (text, label) in enumerate(tqdm(train_data)):
90
91             train_x = text_transform(text, vocab).to(device)
92             train_y = label.to(device)
93
94             #BP
95             optimizer_func.zero_grad()
96             pred = model(train_x)
97             if (method == 'LSTM'):
98                 pred = pred.log()
99             loss = loss_func(pred, train_y)
100             loss.backward()
101             optimizer_func.step()
102             avg_loss += loss.item()
103             avg_acc += accuracy(pred, train_y)
104         # 一个epoch结束后, 计算平均loss和平均acc
105         avg_loss = avg_loss / len(train_data)
106         avg_acc = avg_acc / len(train_data)
107         record_loss.append(avg_loss)
108         record_acc.append(avg_acc)
109
110     print("avg_loss:", avg_loss, " train_avg_acc:", avg_acc)
111
112     # 保存训练完成后的模型参数
113     torch.save(model.state_dict(), method + '_IMDB_parameter.pkl')
114     ...

```

构建 CNN 模型：

```

5 class TextCNN(nn.Module):
6     def __init__(self, vocab_size, embed_sizes, kernel_sizes, num_channels,
7                 **kwargs):
8         super(TextCNN, self).__init__(**kwargs)
9         self.embedding = nn.Embedding(vocab_size, embed_sizes)
10        # 这个嵌入层不需要训练
11        self.constant_embedding = nn.Embedding(vocab_size, embed_sizes)
12        self.dropout = nn.Dropout(0.5)
13        self.decoder = nn.Linear(sum(num_channels), 2)
14        # 最大时间汇聚层没有参数, 因此可以共享此实例
15        self.pool = nn.AdaptiveAvgPool1d(1)
16        self.relu = nn.ReLU()
17
18        # 创建多个一维卷积层
19        self.convs = nn.ModuleList()
20        for c, k in zip(num_channels, kernel_sizes):
21            self.convs.append(nn.Conv1d(2 * embed_sizes, c, k))
22
23        def forward(self, inputs):
24            # 沿着向量维度将两个嵌入层连接起来,
25            # 每个嵌入层的输出形状都是 (批量大小, 词元数量, 词元向量维度) 连接起来
26            embeddings = torch.cat([
27                self.embedding(inputs), self.constant_embedding(inputs)], dim=2)
28            # 根据一维卷积层的输入格式, 重新排列张量, 以便通道作为第2维
29            embeddings = embeddings.permute(0, 2, 1)
30            # 每个一维卷积层在最大时间汇聚层合并后, 获得的张量形状是 (批量大小, 通道数, 1)
31            # 删除最后一个维度并沿通道维度连接
32            encoding = torch.cat([
33                torch.squeeze(self.relu(self.pool(conv(embeddings))), dim=-1)
34                for conv in self.convs], dim=1)
35            outputs = self.decoder(self.dropout(encoding))
36            return outputs

```

构建 LSTM 模型：

```

5 class TextCNN(nn.Module):
6     def __init__(self, vocab_size, embed_sizes, kernel_sizes, num_channels,
7                 **kwargs):
8         super(TextCNN, self).__init__(**kwargs)
9         self.embedding = nn.Embedding(vocab_size, embed_sizes)
10        # 这个嵌入层不需要训练
11        self.constant_embedding = nn.Embedding(vocab_size, embed_sizes)
12        self.dropout = nn.Dropout(0.5)
13        self.decoder = nn.Linear(sum(num_channels), 2)
14        # 最大时间汇聚层没有参数，因此可以共享此实例
15        self.pool = nn.AdaptiveAvgPool1d(1)
16        self.relu = nn.ReLU()
17        # 创建多个一维卷积层
18        self.convs = nn.ModuleList()
19        for c, k in zip(num_channels, kernel_sizes):
20            self.convs.append(nn.Conv1d(2 * embed_sizes, c, k))
21
22    def forward(self, inputs):
23        # 沿着向量维度将两个嵌入层连接起来，
24        # 每个嵌入层的输出形状都是（批量大小，词元数量，词元向量维度）连接起来
25        embeddings = torch.cat([
26            self.embedding(inputs), self.constant_embedding(inputs)], dim=2)
27        # 根据一维卷积层的输入格式，重新排列张量，以便通道作为第2维
28        embeddings = embeddings.permute(0, 2, 1)
29        # 每个一维卷积层在最大时间汇聚层合并后，获得的张量形状是（批量大小，通道数，1）
30        # 删除最后一个维度并沿通道维度连接
31        encoding = torch.cat([
32            torch.squeeze(self.relu(self.pool(conv(embeddings))), dim=-1)
33            for conv in self.convs], dim=1)
34        outputs = self.decoder(self.dropout(encoding))
35        return outputs

```

数据集：

数据集使用的是 [IMDB 网站的电影评论数据](#)，有 25k 训练集和 25k 测试集，训练集和测试集又各包含 12.5k 个 positive 的样本和 negative 的样本。

基准系统：

macos 13.0, m1silicon

实验设置：

基础实验设置为：

对于 LSTM 模型：设置参数为 epoch=5, learning rate = 0.006, batch_size = 256, embde_size = 300, num_layers = 2。

对于 CNN 模型 设置参数为 epoch=5, learning rate = 0.005, batch_size = 256, embde_size = 300, kernel size = [3,4,5], channel_size = [100,100,100]

实验结果与分析：

1. 当 epoch 越大时，训练的准确率越高，对应的在测试集上的准确率也就越高，但是 epoch 越大，花费的训练时间也就越多

当 epoch=3 时，训练准确率为：0.5846，测试准确率为：0.5403

