

LSTM实现网购评论的对象分类

In [1]:

```
1 import re
2 import json
3 import jieba
4 import pandas as pd
5 import numpy as np
6 import matplotlib.pyplot as plt
7 from matplotlib_inline import backend_inline
8 from sklearn.model_selection import train_test_split
9 import torch
10 from torch import nn
```

数据读取与预处理

In [2]:

```
1 df = pd.read_csv("online_shopping_10_cats.csv")[:60000]
2 df.head()
```

Out[2]:

	cat	label	review
0	书籍	1	做父母一定要有刘墉这样的心态，不断地学习，不断地进步，不断地给自己补充新鲜血液，让自己保持...
1	书籍	1	作者真有英国人严谨的风格，提出观点、进行论述论证，尽管本人对物理学了解不深，但是仍然能感受到...
2	书籍	1	作者长篇大论借用详细报告数据处理工作和计算结果支持其新观点。为什么荷兰曾经县有欧洲最高的生产...
3	书籍	1	作者在战几时之前用了 " 拥抱 " 令人叫绝。日本如果没有战败，就会有会有美军的占领，没胡官僚主义的延...
4	书籍	1	作者在少年时即喜阅读，能看出他精读了无数经典，因而他有一个庞大的内心世界。他的作品最难能可贵...

In [3]:

```
1 #使用re正则提取中文并用jieba分词提取词语语料
2 extract_chinese = re.compile(r'[\u4e00-\u9fa5]+')
3 chinese_corpus_raw = df['review'].tolist()
4 chinese_corpus_raw
5 df['chinese_corpus']=[jieba.lcut("".join(extract_chinese.findall(str(corpus)))) for corpus in
6 df.head()
```

Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\23176\AppData\Local\Temp\jieba.cache
Loading model cost 0.823 seconds.
Prefix dict has been built successfully.

Out[3]:

	cat	label	review	chinese_corpus
0	书籍	1	做父母一定要有刘墉这样的心态，不断地学习，不断地进步，不断地给自己补充新鲜血液，让自己保持...	[做, 父母, 一定, 要, 有, 刘墉, 这样, 的, 心态, 不断, 地, 学习, 不断...
1	书籍	1	作者真有英国人严谨的风格，提出观点、进行论述论证，尽管本人对物理学了解不深，但是仍然能感受到...	[作者, 真有, 英国人, 严谨, 的, 风格, 提出, 观点, 进行, 论述, 论证, 尽...
2	书籍	1	作者长篇大论借用详细报告数据处理工作和计算结果支持其新观点。为什么荷兰曾经县有欧洲最高的生产...	[作者, 长篇大论, 借用, 详细, 报告, 数 据处理, 工作, 和, 计算结果, 支持, ...
3	书籍	1	作者在战几时之前用了 " 拥抱 " 令人叫绝。日本如果没有战败，就会有会有美军的占领，没胡官僚主义的延...	[作者, 在, 战, 几时, 之前, 用, 了, 拥 抱, 令人, 叫绝, 日本, 如果, 没...
4	书籍	1	作者在少年时即喜阅读，能看出他精读了无数经典，因而他有一个庞大的内心世界。他的作品最难能可贵...	[作者, 在, 少年, 时即, 喜, 阅读, 能, 看出, 他, 精读, 了, 无数, 经典...

In [4]:

```
1 #将每条评论分词后整合到一个列表中，将每个词用空格隔开放入一个列表中
2 words_list = []
3 corpus = []
4 for corpu in df['chinese_corpus'].tolist():
5     words_list.append(corpu)
6     corpus.append(' '.join(corpu))
7 words_list[0]
```

Out[4]:

```
['做',
 '父母',
 '一定',
 '要',
 '有',
 '刘墉',
 '这样',
 '的',
 '心态',
 '不断',
 '地',
 '学习',
 '不断',
 '地',
 '进步',
 '不断',
 '地',
 '给']
```

In [5]:

```
1 #每个词用空格分割
2 corpus[:2]
```

Out[5]:

['做父母一定要有刘墉这样的心态不断地学习不断地进步不断地给自己补充新鲜血液让自己保持一颗年轻的心我想这是他能很好的和孩子沟通的一个重要因素读刘墉的文章总能让我看到一个快乐的平易近人的父亲他始终站在和孩子同样的高度给孩子创造着一个充满爱和自由的生活环境很喜欢刘墉在字里行间流露出的做父母的那种小狡黠让人总是忍俊不禁父母和子女之间有时候也是一种战斗武力争斗过于低级了智力较量才更有趣味所以做父母的得加把劲了老思想老观念注定会一败涂地生命不息学习不止家庭教育真的是乐在其中',
'作者真有英国人严谨的风格提出观点进行论述论证尽管本人对物理学了解不深但是仍然能感受到真理的火花整本书的结构颇有特点从当时本书写于八十年代流行的计算机话题引入再用数学物理学宇宙学做必要的铺垫这些内容占据了大部分篇幅最后回到关键问题电脑能不能代替人脑和现在流行的观点相反作者认为人的某种洞察是不能被算法模拟的也许作者想说人的灵魂是无可取代的']

In [6]:

```
1 #构建类别与编号的转换字典，并将类别转成编号
2 class2idx = {'书籍':0, '平板':1, '手机':2, '水果':3, '洗发水':4, '热水器':5, '蒙牛':6, '衣服':7,
3 idx2class = {idx:class_ for class_,idx in class2idx.items()}
4 class_idx =[class2idx[calss_] for calss_ in df['cat'].values]
5 class2idx
```

Out[6]:

```
{'书籍': 0,
'平板': 1,
'手机': 2,
'水果': 3,
'洗发水': 4,
'热水器': 5,
'蒙牛': 6,
'衣服': 7,
'计算机': 8,
'酒店': 9}
```

In [7]:

```
1 #统计所有不重复词
2 words_vec_list = []
3 for words in words_list:
4     words_vec_list += words
5 words_vec_list = list(set(words_vec_list))
6 words_vec_list
```

Out[7]:

```
['各得其所',
'上班时间',
'费尽心思',
'全推',
'圖',
'折回来',
'版型',
'写手',
'停下',
'条红脆',
'头好',
'底朝天',
'风景',
'不见',
'窗开',
'抗震',
'政治家',
'快门']
```

In [8]:

```
1 #构建词典(word2idx)
2 word2idx = dict(enumerate(words_vec_list,1))
3 word2idx
```

Out[8]:

```
{1: '各得其所',
2: '上班时间',
3: '费尽心思',
4: '全推',
5: '圖',
6: '折回来',
7: '版型',
8: '写手',
9: '停下',
10: '条红脆',
11: '头好',
12: '底朝天',
13: '风景',
14: '不见',
15: '窗开',
16: '抗震',
17: '政治家',
18: '快门'.
```

In [9]:

```
1 #构建id2word词典
2 idx2word = {word: idx for idx,word in word2idx.items() }
3 idx2word
```

Out[9]:

```
{'各得其所': 1,
'上班时间': 2,
'费尽心思': 3,
'全推': 4,
'圖': 5,
'折回来': 6,
'版型': 7,
'写手': 8,
'停下': 9,
'条红脆': 10,
'头好': 11,
'底朝天': 12,
'风景': 13,
'不见': 14,
'窗开': 15,
'抗震': 16,
'政治家': 17,
'快门': 18.
```

In [10]:

```
1 #将词与编号对应字典输出为json文件便于使用
2 # idx2word_json = json.dumps(idx2word)
3 # with open ('idx2word.json','w') as f:
4 #     f.write(idx2word_json)
5 # f.close()
6 with open ('idx2word.json','r') as f:
7     idx2word_ = json.load(f)
8 idx2word_
```

Out[10]:

```
{'大托': 1,
'湖南': 2,
'定单': 3,
'钛': 4,
'读少': 5,
'弹出式': 6,
'使然': 7,
'贺到': 8,
'医圣': 9,
'外口': 10,
'铜钱': 11,
'斯巴达': 12,
'开学典礼': 13,
'小荷': 14,
'瘦小': 15,
'字棒': 16,
'四张': 17,
'死雪穗': 18.}
```

In [11]:

```
1 #得到每条评论序列表示控制维度为200(使用Tensor表示)
2 def word2idxF(words_list,dim=200):
3     new_idx_list=np.zeros(dim)
4     word_len = len(words_list) if len(words_list)<200 else 200
5     for i in range(word_len):
6         new_idx_list[i]=idx2word[words_list[i]]
7     return new_idx_list
8 words_embedding = torch.tensor([word2idxF(words) for words in words_list],dtype=torch.int32)
9 #得到对应标签的Tensor
10 class_idx = torch.tensor(class_idx,dtype=torch.int64)
11 words_embedding.shape
```

Out[11]:

```
torch.Size([60000, 200])
```

In [12]:

```

1 import torch
2 from torch.utils.data import Dataset, TensorDataset, DataLoader
3 #划分训练集验证集
4 data_set=TensorDataset(words_embedding, class_idx)
5 train_set, valid_set= train_test_split(data_set, random_state=22, test_size=0.2)
6
7 batch_size = 256
8
9 #使用Dataloader进行封装
10 train_loader = DataLoader(train_set, batch_size = batch_size, pin_memory=True, drop_last=True)
11 valid_loader = DataLoader(valid_set, batch_size = batch_size, pin_memory=True, drop_last=True)
12
13 print(f'train_set长度为:{len(train_set)}')
14 print(f'valid_set长度为:{len(valid_set)}')
```

train_set长度为:48000

valid_set长度为:12000

LSTM(单向和双向)实现网购评论的对象分类

In [13]:

```

1 #定义LSTM类可以实例化为LSTM(单向和双向)模型
2 class LSTM(nn.Module):
3     def __init__(self, vocab_size, embedding_dim, hidden_dim, num_layers, output_size, bidirectional):
4         super(LSTM, self).__init__()
5
6         self.hidden_dim = hidden_dim
7         self.num_layers = num_layers
8         self.output_size = output_size
9         self.bidirectional = bidirectional
10
11         self.embedding = nn.Embedding(vocab_size, embedding_dim) # 词嵌入层
12         self.lstm = nn.LSTM(embedding_dim, hidden_dim, num_layers, dropout=dropout, batch_first=True)
13         if self.bidirectional:
14             self.fc = nn.Linear(hidden_dim*2, output_size)
15         else:
16             self.fc = nn.Linear(hidden_dim, output_size)
17
18     def forward(self, x, hidden):
19         batch_size = x.size(0)
20         x = x.long()
21         embeds = self.embedding(x) # 词嵌入表示
22         out, hidden = self.lstm(embeds, hidden)
23         output = self.fc(out[:, -1, :]).squeeze(1)
24         return output, hidden
25
26     def init_hidden(self, batch_size):
27         weight = next(self.parameters()).data
28         if self.bidirectional:
29             hidden = (weight.new(self.num_layers*2, batch_size, self.hidden_dim).zero_().to(device),
30                       weight.new(self.num_layers*2, batch_size, self.hidden_dim).zero_().to(device))
31         else:
32             hidden = (weight.new(self.num_layers, batch_size, self.hidden_dim).zero_().to(device),
33                       weight.new(self.num_layers, batch_size, self.hidden_dim).zero_().to(device))
34         return hidden
```

In [14]:

```

1 #设定超参数
2 vocab_size = len(idx2word)+1#求出词典大小
3 embedding_dim = 128 #词嵌入维度
4 hidden_dim = 64 #隐藏层大小
5 name_layers = 2 #LSTM单元个数
6 dropout_rate = 0.5 #神经元失效比率
7 output_size = 10 #全连接层输出个数
8 bidirectional = False #单向LSTM
9 bidirectiona2 = True #双向LSTM

```

In [29]:

```

1 #创建模型
2 lstm1 = LSTM(vocab_size, embedding_dim, hidden_dim, name_layers, output_size, bidirectional, dropout_
3 lstm2 = LSTM(vocab_size, embedding_dim, hidden_dim, name_layers, output_size, bidirectional, dropout_
4 lstm2

```

Out[29]:

```

LSTM(
  (embedding): Embedding(63745, 128)
  (lstm): LSTM(128, 64, num_layers=2, batch_first=True, dropout=0.5)
  (fc): Linear(in_features=64, out_features=10, bias=True)
)

```

In [31]:

```

1 #定义损失函数, 优化器, 最大迭代次数, 设备
2 loss_function = nn.CrossEntropyLoss()
3 optimizer1 = torch.optim.Adam(lstm1.parameters(), lr = 0.001)
4 optimizer2 = torch.optim.Adam(lstm2.parameters(), lr = 0.001)
5 max_epochs = 20
6 device = 'cuda:0' if torch.cuda.is_available() else 'cpu'
7 lstm1 = lstm1.to(device)
8 lstm2 = lstm2.to(device)#将模型转移到相应的设备上

```

In [17]:

```

1 #设定画图配置
2 def use_svg_display():
3     """Use the svg format to display a plot in Jupyter.
4
5     Defined in :numref:`sec_calculus`"""
6     backend_inline.set_matplotlib_formats('svg')
7 def set_figsize(figsize=(3.5, 2.5)):
8     """Set the figure size for matplotlib.
9
10    Defined in :numref:`sec_calculus`"""
11    use_svg_display()
12    plt.rcParams['figure.figsize'] = figsize

```


In [18]:

```

1 #定义训练器和验证了并定义绘制训练集和验证集损失和准确率曲线函数并保存验证集准确率最高的模型
2 class Train:
3     def __init__(self,max_epochs,loss_function,optimizer,model,model_name,device = 'cpu'):
4         self.max_epochs = max_epochs
5         self.device = device
6         self.loss_function = loss_function
7         self.optimizer = optimizer
8         self.model = model.to(device)
9         self.model_name = model_name
10        self.params = None
11    def start_train(self,trainloader,validloader = None,val_idx = None):
12        self.trainloader = trainloader
13        self.validloader = validloader
14        self.loss_train_list = []
15        self.loss_valid_list = []
16        self.accuracy_rate_train = []
17        self.accuracy_rate_valid = []
18        valid_best_accuracy = 0
19        if val_idx != None:
20            self.max_valid_num = int(self.max_epochs / val_idx)
21            self.val_idx = val_idx
22        if isinstance(self.model, nn.Module):
23            self.model.train()
24        print('Start Training!')
25        for epoch in range(self.max_epochs):
26            total_num = 0
27            accuracy_num = 0
28            hs = self.model.init_hidden(batch_size)
29            for idx,(x,t) in enumerate(self.trainloader):
30                self.model.train()
31                x = x.to(device)
32                t = t.to(device)
33                total_num += x.shape[0]
34                t_hat,hs = self.model(x,hs)
35                t_label = torch.argmax(t_hat,axis = 1)
36                hs = tuple([h.data for h in hs])
37                loss_ = self.loss_function(t_hat, t)
38                accuracy_num += sum(t_label==t)
39                self.optimizer.zero_grad()
40                loss_.backward()
41                self.optimizer.step()
42            self.loss_train_list.append(loss_.item())
43            loss = loss_.item()
44            accuracy_rate = round(accuracy_num.cpu().item()/total_num,4)
45            self.accuracy_rate_train.append(accuracy_rate)
46            print(f'Train_set Epoch [{epoch}/{self.max_epochs}] loss: {loss}, acc: {accuracy_ra
47            if (epoch+1) % self.val_idx == 0:
48                valid_num = int((epoch+1) / self.val_idx)
49                if isinstance(self.model, nn.Module):
50                    self.model.eval()
51                with torch.no_grad():
52                    total_num = 0
53                    accuracy_num = 0
54                    hs = self.model.init_hidden(batch_size)
55                    print('Start Validation!')
56                    for idx, (x, t) in enumerate(self.validloader):
57                        x = x.to(device)
58                        t = t.to(device)
59                        total_num += x.shape[0]

```

```

60         t_hat, hs = self.model(x, hs)
61         t_label = torch.argmax(t_hat, axis = 1)
62         hs = tuple([h.data for h in hs])
63         loss_ = self.loss_function(t_hat, t)
64         accurary_num += sum(t_label==t)
65         self.loss_valid_list.append(loss_.item())
66         loss = loss_.item()
67         accurary_rate = round(accurary_num.cpu().item() / total_num, 4)
68         self.accurary_rate_valid.append(accurary_rate)
69         print(f'Valid_set Epoch [{valid_num}/{self.max_valid_num}] loss: {loss}, ac
70         if accurary_rate > valid_best_accuracy:
71             torch.save(self.model, self.model_name+'.pth')
72             valid_best_accuracy = accurary_rate
73             print('Best_model has been saved!')
74             print('Stop Validation!')
75     def show_loss_acc_value(self):
76         n_train_loss_value = len(self.loss_train_list)
77         n_accurary_rate_train = len(self.accurary_rate_train)
78         set_figsize(figsize=(4, 3))
79         plt.plot(list(range(n_accurary_rate_train)), self.accurary_rate_train, 'r-', linewidth = 1)
80         plt.plot(list(range(n_train_loss_value)), self.loss_train_list, 'b-', linewidth=1, label='loss')
81         if self.loss_valid_list != []:
82             n_valid_loss_value = len(self.loss_valid_list)
83             n_accurary_rate_valid = len(self.accurary_rate_valid)
84             plt.plot(list(range(n_accurary_rate_valid)), self.accurary_rate_valid, 'y-', linewidth=1, label='acc')
85             plt.plot(list(range(n_valid_loss_value)), self.loss_valid_list, 'g-', linewidth=1, label='valid_loss')
86         plt.title('loss_acc_curve')
87         plt.xlabel('Epochs')
88         plt.ylabel('loss_acc')
89         plt.legend()
90         plt.ylim(0, 1)
91         plt.show()

```

In [32]:

```

1 #实例化训练和验证类用于训练LSTM(单向)和LSTM(双向)
2 train1 = Train(max_epochs, loss_function, optimizer1, lstm1, 'LSTM', device =device)
3 train2 = Train(max_epochs, loss_function, optimizer2, lstm2, 'BiLSTM', device =device)

```

In [20]:

```

1 #LSTM(单向)开始训练并验证
2 train1.start_train(train_loader, valid_loader , val_idx = 2)

```

Start Training!

Train_set Epoch [0/20] loss: 2.104768753051758, acc: 0.168

Train_set Epoch [1/20] loss: 1.9072407484054565, acc: 0.1954

Start Validation!

Valid_set Epoch [1/10] loss: 1.7820121049880981, acc: 0.2597

Best_model has been saved!

Stop Validation!

Train_set Epoch [2/20] loss: 1.5503545999526978, acc: 0.3378

Train_set Epoch [3/20] loss: 1.3045552968978882, acc: 0.4585

Start Validation!

Valid_set Epoch [2/10] loss: 1.2434165477752686, acc: 0.488

Best_model has been saved!

Stop Validation!

Train_set Epoch [4/20] loss: 1.1192647218704224, acc: 0.5431

Train_set Epoch [5/20] loss: 1.035346269607544, acc: 0.5948

Start Validation!

Valid_set Epoch [3/10] loss: 1.0766464471817017, acc: 0.5813

Best_model has been saved!

Stop Validation!

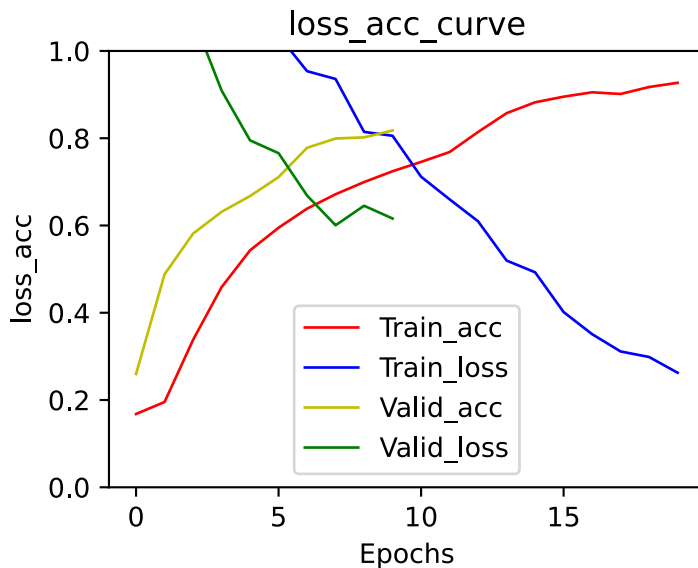
Train_set Epoch [6/20] loss: 0.9586681175881084, acc: 0.6886

In [21]:

```

1 #LSTM(单向)训练集验证集损失和准确率图像
2 train1.show_loss_acc_value()

```



In [22]:

```

1 #LSTM(单向)加载
2 lstm1=torch.load('LSTM.pth')
3 lstm1 = lstm1.to(device)

```

In [33]:

```

1 #LSTM(双向)开始训练并验证
2 train2.start_train(train_loader, valid_loader , val_idx = 2)

```

Start Training!

Train_set Epoch [0/20] loss: 2.1104559898376465, acc: 0.1657

Train_set Epoch [1/20] loss: 1.9765098094940186, acc: 0.1809

Start Validation!

Valid_set Epoch [1/10] loss: 1.925413966178894, acc: 0.2324

Best_model has been saved!

Stop Validation!

Train_set Epoch [2/20] loss: 1.798648715019226, acc: 0.2753

Train_set Epoch [3/20] loss: 1.659867525100708, acc: 0.2938

Start Validation!

Valid_set Epoch [2/10] loss: 1.6208168268203735, acc: 0.3072

Best_model has been saved!

Stop Validation!

Train_set Epoch [4/20] loss: 1.5154714584350586, acc: 0.3529

Train_set Epoch [5/20] loss: 1.3221570253372192, acc: 0.4174

Start Validation!

Valid_set Epoch [3/10] loss: 1.267122745513916, acc: 0.4414

Best_model has been saved!

Stop Validation!

Train_set Epoch [6/20] loss: 1.2188018500014888, acc: 0.4688

In [34]:

```

1 #LSTM(双向)加载
2 lstm2=torch.load('BiLSTM.pth')
3 lstm2 = lstm2.to(device)

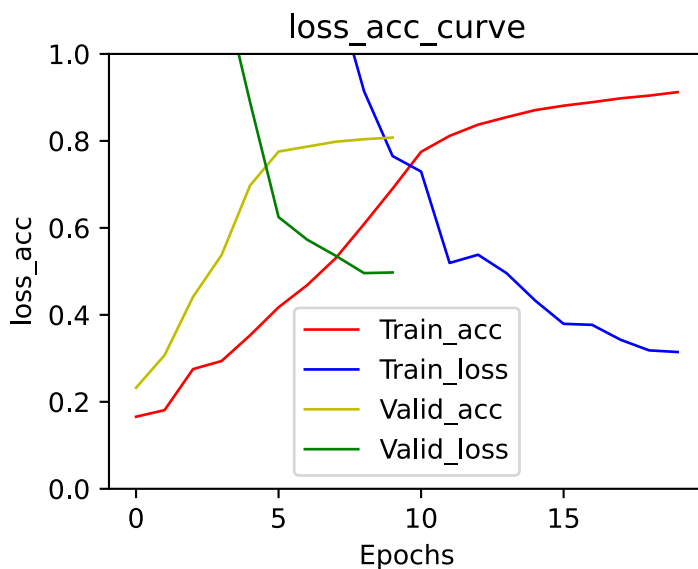
```

In [35]:

```

1 #LSTM(双向)训练集验证集损失和准确率图像
2 train2.show_loss_acc_value()

```



In [26]:

```
1 #定义LSTM模型评估函数
2 def LSTM_Evaluation(LSTM, valid_set):
3     LSTM.eval()
4     pred_list = []
5     label_list = []
6     hs = LSTM.init_hidden(batch_size)
7     for idx, (x, t) in enumerate(valid_loader):
8         x = x.to(device)
9         t = t.to(device)
10        t_hat, hs = LSTM(x, hs)
11        pred_list += torch.argmax(t_hat, axis = 1).cpu().tolist()
12        label_list += t.cpu().tolist()
13        f1 = f1_score(label_list, pred_list, average='macro')
14        Accuracy_score = accuracy_score(label_list, pred_list)
15        Recall_score = recall_score(label_list, pred_list, average='macro')
16        print(f'Accuracy_score: {Accuracy_score}')
17        print(f'Recall_score: {Recall_score}')
18        print(f'f1_score: {f1}')
```

In [27]:

```
1 #LSTM(单向)模型评估
2 LSTM_Evaluation(lstm1, valid_set)
```

Accuracy_score:0.8176800271739131

Recall_score:0.7535268606235399

f1_score:0.74817884050913

In [36]:

```
1 #LSTM(双向)模型评估
2 LSTM_Evaluation(lstm2, valid_set)
```

Accuracy_score:0.8077445652173914

Recall_score:0.7394249069692183

f1_score:0.7349844344083212