

新闻文本分类

19351024林威

1. 背景介绍

文本自动分类 (text categorization) 简称文本分类，是模式识别与自然语言处理密切结合的研究课题。传统的文本分类是基于文本内容的，研究如何将文本自动划分成政治的、经济的、军事的、体育的、娱乐的等各种类型。这也是人们提到“文本分类”这一术语时通常所指的含义。

本次实践的背景是Datawhale与天池联合发起的0基础入门系列赛事第三场 —— 零基础入门NLP之新闻文本分类挑战赛。赛题以自然语言处理为背景，要求选手根据新闻文本字符对新闻的类别进行分类，这是一个经典文本分类问题。通过这道赛题可以引导大家走入自然语言处理的世界，带大家接触NLP的预处理、模型构建和模型训练等知识点。

[零基础入门NLP - 新闻文本分类 - 天池大赛 - 阿里云天池](#)

2. 数据集

赛题以新闻数据为赛题数据，数据集报名后可见并可下载。赛题数据为新闻文本，并按照字符级别进行匿名处理。整合划分出14个候选分类类别：财经、彩票、房产、股票、家居、教育、科技、社会、时尚、时政、体育、星座、游戏、娱乐的文本数据。赛题数据由以下几个部分构成：训练集20w条样本，测试集A包括5w条样本，测试集B包括5w条样本。为了预防选手人工标注测试集的情况，比赛数据的文本按照字符级别进行了匿名处理。处理后的赛题训练数据如下：

label	text
6	57 44 66 56 2 3 3 37 5 41 9 57 44 47 45 33 13 63 58 31 17 47 0 1 1 69 26 60 62 15 21 12 49 18 38 20 50 23 57 44 45 33 25 28 47 22 52 35 30 14 24 69 54 7 48 19 11 51 16 43 26 34 53 27 64 8 4 42 36 46 65 69 29 39 15 37 57 44 45 33 69 54 7 25 40 35 30 66 56 47 55 69 61 10 60 42 36 46 65 37 5 41 32 67 6 59 47 0 1 1 68

在数据集中标签的对应的关系如下：

1	{'科技': 0, '股票': 1, '体育': 2, '娱乐': 3, '时政': 4, '社会': 5, '教育': 6, '财经': 7, '家居': 8, '游戏': 9, '房产': 10, '时尚': 11, '彩票': 12, '星座': 13}
---	--

3. 评价指标

评价标准为类别f1_score的均值，提交结果与实际测试集的类别进行对比，结果越大越好。

$$F1 = 2 * \frac{precision * recall}{precision + recall}$$

可以通过sklearn来完成f1_score计算：

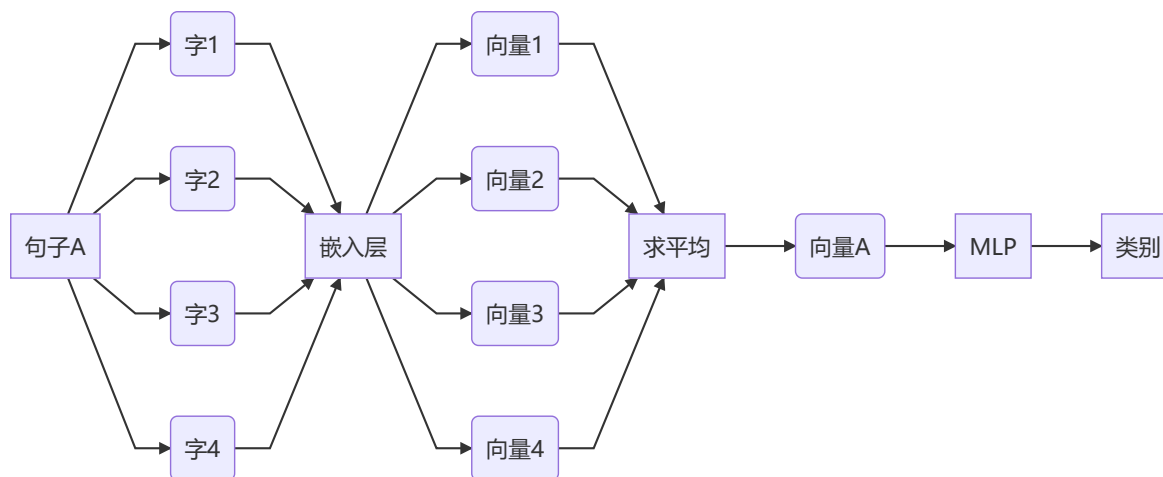
```

1 from sklearn.metrics import f1_score
2 y_true = [0, 1, 2, 0, 1, 2]
3 y_pred = [0, 2, 1, 0, 0, 1]
4 f1_score(y_true, y_pred, average='macro')

```

4. 解题模型

我选择了深度学习的方法来解决这个问题。模型其实非常简单，可以由下图来表示：



模型的输入是一个句子A，它由许多字构成（这里以4个字为例）。实际上，这里的每一个汉字都用了一个整数来表示，这个整数会被作为是嵌入层的输入。嵌入层的作用是把一个整数输入转化为一个指定维度的向量。之后我对这些向量的每个元素分别求平均，得到了向量A。向量A即可被认为是句子A的表示。接着我再用一个多层感知机对向量A进行分类预测，最终得到一个长度为14的向量。该向量中的每个元素即代表对应位置所表示的类别的概率，取概率最大者为模型预测结果。

该模型可以用SGD优化器在训练集上进行训练，损失函数为Cross Entropy。

5. 具体python实现

（一）导入需要用到的Package

我的具体实现是基于Pytorch的，以下是一些可能会用到的库

```

1 import torch
2 from torch.utils.data import Dataset, DataLoader
3 import pandas as pd
4 from tqdm import tqdm
5 from torch import nn
6 import time
7 from torch.utils.data.dataset import random_split
8 from torchtext.data.functional import to_map_style_dataset
9 from sklearn.metrics import f1_score

```

(二) 数据准备

先为训练集和测试集创造两个dataset。

```
1 class trainDataset(Dataset):
2     def __init__(self):
3         self.data=pd.read_csv("训练集/train_set.csv/train_set.csv",sep='\t')
4     def __len__(self):
5         return len(self.data)
6     def __getitem__(self,idx):
7         item=self.data.iloc[idx,:]
8         label=item['label']
9         text=item['text'].split(' ')
10        text=[int(i) for i in text]
11        return label,text
12
13 class testDataset(Dataset):
14     def __init__(self):
15         self.data=pd.read_csv("测试集/test.csv/test_a.csv",sep='\t')
16     def __len__(self):
17         return len(self.data)
18     def __getitem__(self,idx):
19         item=self.data.iloc[idx,:]
20         text=item['text'].split(' ')
21         text=[int(i) for i in text ]
22         return text
```

接着将训练集以0.95 : 0.05的比例分为训练集和验证集

```
1 train_dataset = trainDataset()
2 test_dataset=testDataset()
3 num_train = int(len(train_dataset) * 0.95)
4 split_train_, split_valid_ = \
5     random_split(train_dataset, [num_train, len(train_dataset) - num_train])
```

然后把他们都放置入DataLoader当中

```
1 train_dataloader = DataLoader(split_train_, batch_size=BATCH_SIZE,
2                               shuffle=True, collate_fn=collate_batch)
3 valid_dataloader = DataLoader(split_valid_, batch_size=BATCH_SIZE,
4                               shuffle=True, collate_fn=collate_batch)
5 test_dataloader  = DataLoader(test_dataset,batch_size=BATCH_SIZE,
6                               shuffle=False, collate_fn=test_collate_batch)
7
8 def collate_batch(batch):
9     label_list, text_list, offsets = [], [], [0]
10    for (label, text) in batch:
11        label_list.append(label)
12        processed_text = torch.tensor(text, dtype=torch.int64)
13        text_list.append(processed_text)
14        offsets.append(processed_text.size(0))
15    label_list = torch.tensor(label_list, dtype=torch.int64)
16    offsets = torch.tensor(offsets[:-1]).cumsum(dim=0)
17    text_list = torch.cat(text_list)
18    return label_list.to(device), text_list.to(device), offsets.to(device)
19
```

```

20 def test_collate_batch(batch):
21     text_list, offsets = [], [0]
22     for text in batch:
23         processed_text = torch.tensor(text, dtype=torch.int64)
24         text_list.append(processed_text)
25         offsets.append(processed_text.size(0))
26     offsets = torch.tensor(offsets[:-1]).cumsum(dim=0)
27     text_list = torch.cat(text_list)
28     return text_list.to(device), offsets.to(device)

```

(三) 定义模型

根据前文描述的思路，定义网络模型。

```

1 class TextClassificationModel(nn.Module):
2
3     def __init__(self, vocab_size, embed_dim, num_class):
4         super(TextClassificationModel, self).__init__()
5         self.embedding = nn.EmbeddingBag(vocab_size, embed_dim, sparse=True)
6         self.fc = nn.Linear(embed_dim, num_class)
7         self.init_weights()
8
9     def init_weights(self):
10         initrange = 0.5
11         self.embedding.weight.data.uniform_(-initrange, initrange)
12         self.fc.weight.data.uniform_(-initrange, initrange)
13         self.fc.bias.data.zero_()
14
15     def forward(self, text, offsets):
16         embedded = self.embedding(text, offsets)
17         return self.fc(embedded)

```

(四) 训练、验证和测试模型

```

1 def train(data_loader):
2     model.train()
3     total_acc, total_count = 0, 0
4     log_interval = 500
5     start_time = time.time()
6
7     for idx, (label, text, offsets) in enumerate(data_loader):
8         optimizer.zero_grad()
9         predicted_label = model(text, offsets)
10        loss = criterion(predicted_label, label)
11        loss.backward()
12        torch.nn.utils.clip_grad_norm_(model.parameters(), 0.1)
13        optimizer.step()
14        total_acc += (predicted_label.argmax(1) == label).sum().item()
15        total_count += label.size(0)
16        if idx % log_interval == 0 and idx > 0:
17            elapsed = time.time() - start_time
18            print('| epoch {:3d} | {:5d}/{:5d} batches |
19                  | accuracy {:.3f}'.format(epoch, idx, len(data_loader),
20                                              total_acc/total_count))
21            total_acc, total_count = 0, 0
22            start_time = time.time()

```

```

23     torch.save(model, 'my_model.pth')
24 def evaluate(data_loader):
25     model.eval()
26     total_acc, total_count = 0, 0
27
28     with torch.no_grad():
29         for idx, (label, text, offsets) in enumerate(data_loader):
30             predicted_label = model(text, offsets)
31             loss = criterion(predicted_label, label)
32             total_acc += (predicted_label.argmax(1) == label).sum().item()
33             total_count += label.size(0)
34     return total_acc/total_count
35 def test(data_loader):
36     model.eval()
37     res=[]
38
39     with torch.no_grad():
40         for idx, (text, offsets) in enumerate(data_loader):
41             predicted_label = model(text, offsets)
42             result=predicted_label.argmax(1).to('cpu').numpy()
43             res.extend(result)
44     df=pd.DataFrame()
45     df['label']=res
46     df.to_csv("test_result.csv")

```

(五) 超参设定和主体部分

```

1  EPOCHS = 5
2  LR = 5
3  BATCH_SIZE = 16
4
5  criterion = torch.nn.CrossEntropyLoss()
6  optimizer = torch.optim.SGD(model.parameters(), lr=LR)
7  scheduler = torch.optim.lr_scheduler.StepLR(optimizer, 1.0, gamma=0.1)
8
9  num_class = 14
10 vocab_size = 7916
11 emsize = 64
12 model = TextClassificationModel(vocab_size, emsize, num_class).to(device)
13 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
14
15 for epoch in range(1, EPOCHS + 1):
16     epoch_start_time = time.time()
17     train(train_data_loader)
18     accu_val = evaluate(valid_data_loader)
19     if total_accu is not None and total_accu > accu_val:
20         scheduler.step()
21     else:
22         total_accu = accu_val
23     print('-' * 59)
24     print('| end of epoch {:3d} | time: {:.2f}s | '
25           'valid accuracy {:.3f} '.format(epoch, time.time() -
epoch_start_time,
26                                           accu_val))
27     test(test_data_loader)

```

6. 实验结果

经过5个epoch的训练（平均每个epoch需要运行169.85秒），我在验证集上的到的结果为 $Precision = 0.925$, $f1 = 0.91$

最后我把在测试集上跑的结果上传至比赛官网，结果如下图所示：

