

# 新闻类别分类

211250109 赵政杰

2023 年 11 月 30 日

## 1 数据集描述

通过 analysis.py 对数据集进行分析，结果如下：

Number of Samples: 208908

Number of Features: 6

同时分别生成了类别分布的柱状图、饼状图

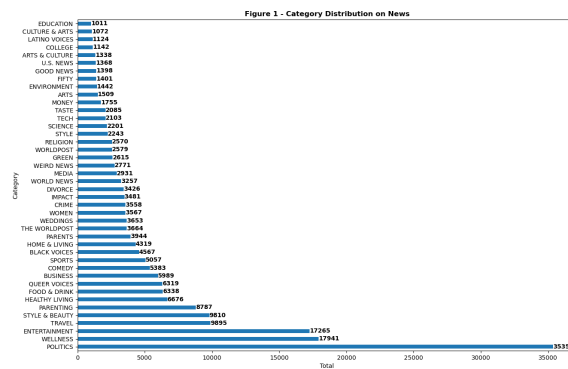


图 1: 柱状图

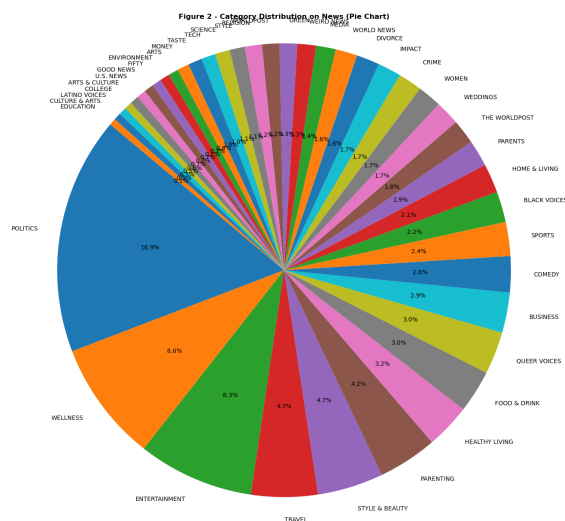


图 2: 饼状图

生成 short\_description 列中词语的词云图，如下

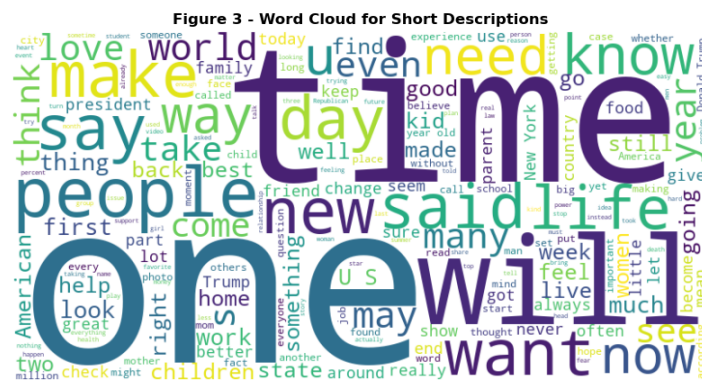


图 3: 词云图

## 2 数据处理

使用 pandas 库读取 json 文件，并将其存储在 Pandas DataFrame 中。

将作者、新闻标题、link 和简短描述合并为一个文本列，并将类别作为标签列构建新的 DataFrame。

定义数据集类 TextClassificationDataset，用来将数据集转化为相应的张量

实现代码如下：

```
1 # TextClassificationDataset 类定义
2 class TextClassificationDataset(Dataset):
3
4     def __init__(self,
5                 texts: List[str],
6                 labels: List[str] = None,
7                 label_dict: Mapping[str, int] = None,
8                 max_seq_length: int = 512,
9                 model_name: str = 'distilbert-base-uncased'):
10
11         self.texts = texts
12         self.labels = labels
13         self.label_dict = label_dict
14         self.max_seq_length = max_seq_length
15
16         if self.label_dict is None and labels is not None:
17             self.label_dict = dict(zip(sorted(set(labels)),
18                                     range(len(set(labels)))))
19
20         self.tokenizer = AutoTokenizer.from_pretrained(model_name)
```

```

21
22     self.sep_vid = self.tokenizer.vocab["[SEP]"]
23     self.cls_vid = self.tokenizer.vocab["[CLS]"]
24     self.pad_vid = self.tokenizer.vocab["[PAD]"]
25
26     def __len__(self):
27
28         return len(self.texts)
29
30     def __getitem__(self, index) -> Mapping[str, torch.Tensor]:
31
32         x = self.texts[index]
33         x_encoded = self.tokenizer.encode(
34             x,
35             add_special_tokens=True,
36             max_length=self.max_seq_length,
37             return_tensors="pt",
38             truncation=True # 明确启用截断
39         ).squeeze(0)
40
41         true_seq_length = x_encoded.size(0)
42         pad_size = self.max_seq_length - true_seq_length
43         pad_ids = torch.Tensor([self.pad_vid] * pad_size).long()
44         x_tensor = torch.cat((x_encoded, pad_ids))
45
46         mask = torch.ones_like(x_encoded, dtype=torch.int8)
47         mask_pad = torch.zeros_like(pad_ids, dtype=torch.int8)
48         mask = torch.cat((mask, mask_pad))
49
50         output_dict = {
51             "features": x_tensor,
52             'attention_mask': mask
53         }
54
55         if self.labels is not None:
56             y = self.labels[index]
57             y_encoded = torch.Tensor(
58                 [self.label_dict.get(y, -1)]
59             ).long().squeeze(0)
60             output_dict["targets"] = y_encoded
61
62         return output_dict
63
64
65 data_path = "./data"
66 data = pd.read_json(data_path+"/News_Category.json", lines=True)
67 text = pd.DataFrame({
68     "text" : data.authors+" "+data.headline+" "+data.link+" "+data.short_description,
69     "label" : data.category
70 })
71 train_dataset = TextClassificationDataset(
72     texts=train['text'].values.tolist(),

```

```

73     labels=train['label'].values.tolist(),
74     label_dict=label_dict,
75     max_seq_length=MAX_SEQ_LENGTH,
76     model_name=MODEL_NAME
77 )
78 valid_dataset = TextClassificationDataset(
79     texts=val['text'].values.tolist(),
80     labels=val['label'].values.tolist(),
81     label_dict=label_dict,
82     max_seq_length=MAX_SEQ_LENGTH,
83     model_name=MODEL_NAME
84 )

```

### 3 分类模型选择及设计、训练、验证、测试

实现在 train.py 文件、test.py 文件和 some\_classes.py 文件中

#### 3.1 模型选择与设计

选择使用 DistilBert 模型进行分类。DistilBert 是一种小型、快速和轻量级的 Transformer 模型，通过知识蒸馏 BERT 基础进行训练。

用于序列分类的 DistilBERT 模型定义如下：

```

1  class DistilBertForSequenceClassification(nn.Module):
2
3  def __init__(self, pretrained_model_name: str, num_classes: int = None):
4
5      super().__init__()
6
7      config = AutoConfig.from_pretrained(
8          pretrained_model_name, num_labels=num_classes)
9
10     self.distilbert = AutoModel.from_pretrained(pretrained_model_name,
11                                                  config=config)
12     self.pre_classifier = nn.Linear(config.dim, config.dim)
13     self.classifier = nn.Linear(config.dim, num_classes)
14     self.dropout = nn.Dropout(config.seq_classif_dropout)
15
16 def forward(self, features, attention_mask=None, head_mask=None):
17
18     outputs = self.distilbert(input_ids=features, attention_mask=attention_mask,
19                               head_mask=head_mask)
20
21     hidden_state = outputs[0]
22     pooled_output = hidden_state[:, 0]
23     pooled_output = self.pre_classifier(pooled_output)
24     pooled_output = nn.ReLU()(pooled_output)
25     pooled_output = self.dropout(pooled_output)
26     logits = self.classifier(pooled_output)
27
28     return logits

```

```

28 # 加载模型
29 model = DistilBertForSequenceClassification(pretrained_model_name=MODEL_NAME,
30                                             num_classes=NUM_CLASSES)

```

## 3.2 训练与验证

将数据集划分为 80% 的训练集和 20% 的验证集:

```

1 train, val = train_test_split(text, test_size=0.20, random_state=SEED)

```

训练参数设置:

```

1 MODEL_NAME = "distilbert-base-uncased"
2 NUM_EPOCHS = 3
3 BATCH_SIZE = 80
4 MAX_SEQ_LENGTH = 256
5 LEARN_RATE = 5e-5
6 SEED = 42
7 LOG_DIR = 'results'

```

训练与验证过程如代码所示:

```

1 model = DistilBertForSequenceClassification(pretrained_model_name=MODEL_NAME,
2                                             num_classes=NUM_CLASSES)
3 criterion = torch.nn.CrossEntropyLoss()
4 optimizer = torch.optim.Adam(model.parameters(), lr=LEARN_RATE)
5 scheduler = torch.optim.lr_scheduler.ReduceLROnPlateau(optimizer)
6
7 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
8 model.to(device)
9 for epoch in range(NUM_EPOCHS):
10     # Training loop
11     model.train()
12     total_loss = 0
13
14     for batch_idx, batch in enumerate(train_val_loaders['train']):
15         optimizer.zero_grad()
16
17         features = batch['features'].to(device)
18         attention_mask = batch['attention_mask'].to(device)
19         targets = batch['targets'].to(device)
20
21         outputs = model(features, attention_mask=attention_mask)
22         loss = criterion(outputs, targets)
23         total_loss += loss.item()
24
25         loss.backward()
26         optimizer.step()
27
28         if batch_idx % 10 == 0:
29             print(f"Epoch [{epoch + 1}/{NUM_EPOCHS}] "
30                   f"Batch [{batch_idx + 1}/{len(train_val_loaders['train'])}] "
31                   f"Loss: {loss.item()}")

```

```

32         logging.info(f"Epoch [{epoch + 1}/{NUM_EPOCHS}] "
33                       f"Batch [{batch_idx + 1}/{len(train_val_loaders['train'])}] "
34                       f"Loss: {loss.item()}")
35
36     average_loss = total_loss / len(train_val_loaders['train'])
37     print(f"Epoch [{epoch + 1}/{NUM_EPOCHS}] "
38           f"Average training loss: {average_loss}")
39     logging.info(f"Epoch [{epoch + 1}/{NUM_EPOCHS}] "
40                 f"Average training loss: {average_loss}")
41
42     # Validation loop
43     model.eval()
44     val_loss = 0
45     predicted_labels = []
46     true_labels = []
47
48     with torch.no_grad():
49         for val_batch in train_val_loaders['valid']:
50             features = val_batch['features'].to(device)
51             attention_mask = val_batch['attention_mask'].to(device)
52             targets = val_batch['targets'].to(device)
53
54             outputs = model(features, attention_mask=attention_mask)
55             loss = criterion(outputs, targets)
56             val_loss += loss.item()
57
58             _, predicted = torch.max(outputs, 1)
59
60             predicted_labels.extend(predicted.cpu().numpy()) # 预测的标签
61             true_labels.extend(targets.cpu().numpy()) # 真实的标签
62
63     average_val_loss = val_loss / len(train_val_loaders['valid'])
64     acc = accuracy_score(true_labels, predicted_labels) # 计算准确率
65     precision = precision_score(true_labels, predicted_labels, average='macro') #
66     # 计算精确度
67     f1 = f1_score(true_labels, predicted_labels, average='macro') # 计算F1 Score
68     print(f"Epoch [{epoch + 1}/{NUM_EPOCHS}] "
69           f"Test Accuracy: {acc}, Test Precision: {precision}, Test F1 Score: {f1}")
70     logging.info(f"Epoch [{epoch + 1}/{NUM_EPOCHS}] "
71                 f"Test Accuracy: {acc}, Test Precision: {precision}, Test F1 Score:
72                 {f1}")
73
74     # 根据验证损失调整学习率
75     scheduler.step(average_val_loss)
76
77     # 保存训练好的模型
78     torch.save(model.state_dict(), f'results1/trained_model_epoch_{epoch + 1}.pth')

```

### 3.3 测试

加载训练好的模型文件，在数据集上进行测试，实现如下：

```

1 # 加载训练好的模型

```

```

2 loaded_model = DistilBertForSequenceClassification(pretrained_model_name=MODEL_NAME,
3                                                    num_classes=NUM_CLASSES)
4 loaded_model.load_state_dict(torch.load('results/trained_model_epoch_2.pth')) #
    替换成你想要加载的模型文件名
5 # 设置模型为评估模式
6
7 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
8 loaded_model.to(device)
9
10 def cal_acc_test():
11     test_loader = DataLoader(test_dataset, batch_size=100, shuffle=False)
12     predicted_labels = []
13     true_labels = []
14
15     with torch.no_grad():
16         for batch in test_loader:
17             features = batch['features'].to(device)
18             attention_mask = batch['attention_mask'].to(device)
19             targets = batch['targets'].to(device)
20
21             outputs = loaded_model(features, attention_mask=attention_mask)
22             _, predicted = torch.max(outputs, 1)
23
24             predicted_labels.extend(predicted.cpu().numpy()) # 预测的标签
25             true_labels.extend(targets.cpu().numpy()) # 真实的标签
26
27     acc = accuracy_score(true_labels, predicted_labels) # 计算准确率
28     precision = precision_score(true_labels, predicted_labels, zero_division=1,
29                                average='macro') # 计算精确度
30     f1 = f1_score(true_labels, predicted_labels, average='macro') # 计算F1 Score
31     print(f"Test Accuracy: {acc}, Test Precision: {precision}, Test F1 Score: {f1}")
32
33 if __name__ == "__main__":
34     print_classifications_test()

```

### 3.4 训练结果

参见模型链接中的 log 文件