

推特文本情感分析

第五组: 张露雨 高鹏轩 邹佳俊

2022年12月30日







任务简介







D 展望与思考



任务简介

• 任务介绍 • 数据探索



1.1任务简介



情感分析一般指判断一段文本所表达的情绪状态,属于文本分类问题。主要任务是对文本中的主观信息(如观点、情感、评价、态度、情绪等)进行提取、分析、处理、归纳和推理。

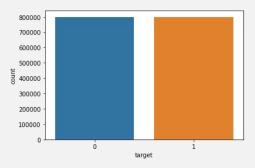
本次任务中, 我们需要通过模型对一段文本进行情绪的正负判断, 可以作为一个简单的二分类任务。

数据集介绍

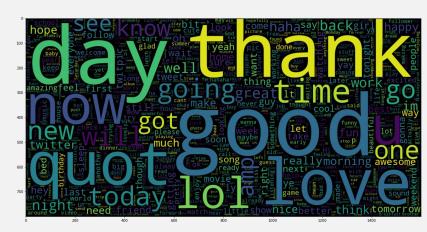
训练集160万条:

六个维度:标签、ID、日期、话题、用户及推文内容; 无缺失数据; 正负样本平衡。

1.2 数据探索



正负样本比例



正样本词云图



负样本词云图



实验思路

• 机器学习 • BERT • GPT2











过程及结果

• 传统分类 • BERT • GPT2



3.1 传统分类

预处理:转小写-->去停用词-->去标点-->去链接-->去数字-->TFIDF

逻辑回归

模型超参: L2正则化

正则化系数 λ 的倒数C=2 优化算法: solver=sag

支持向量机

模型超参: L2正则化

损失函数: loss=squared_hinge 最大迭代次数: max_iter=1000

随机梯度下降

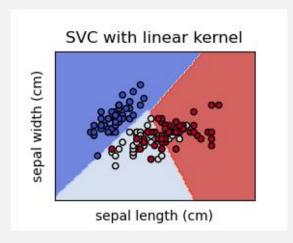
模型超参: L2正则化

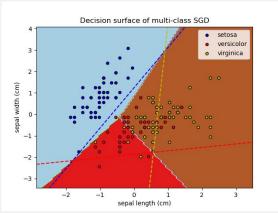
损失函数: loss='log'

最大迭代次数: max_iter=1000

集成学习

模型超参: 投票方法: voting='hard'





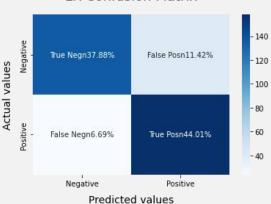
3.1 传统分类





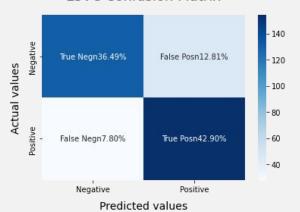






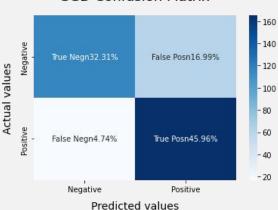
accuracy=0.82 time: 20.6s

LSVC Confusion Matrix



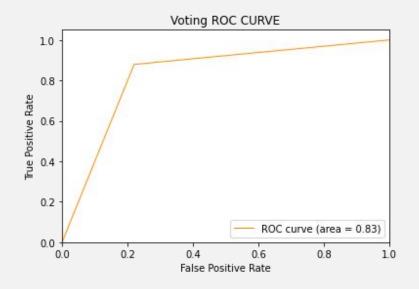
accuracy=0.79 time: 15.3s

SGD Confusion Matrix



accuracy=0.78 time: 4.1s

3.1传统分类



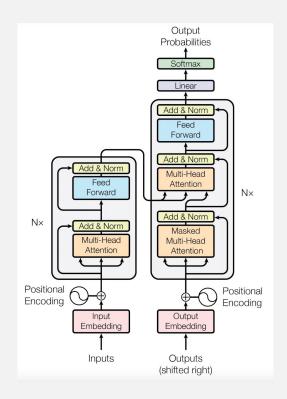
accuracy=0.83 time: 15.2s

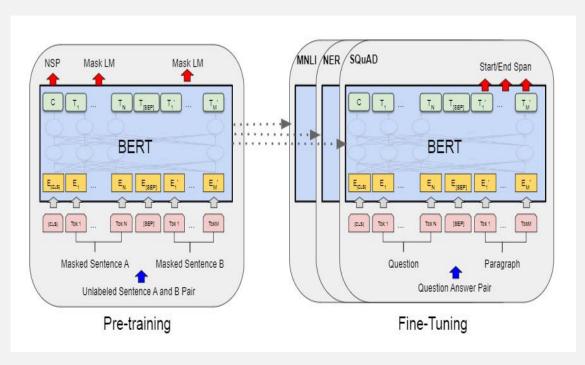
	precision	recall	f1-score	support
0	0.86	0. 78	0. 82	177
1	0.80	0. 88	0. 84	182
accuracy			0. 83	359
macro avg	0.83	0. 83	0. 83	359
weighted avg	0. 83	0. 83	0. 83	359



集成学习提高了准确率, 对于时间的影响不大

3.2 BERT微调



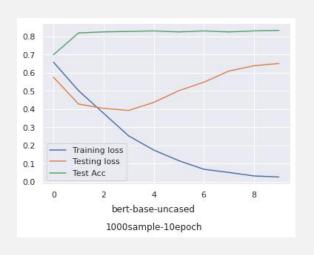


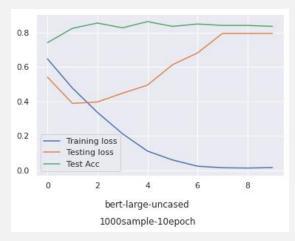
arXiv:1810.04805v2

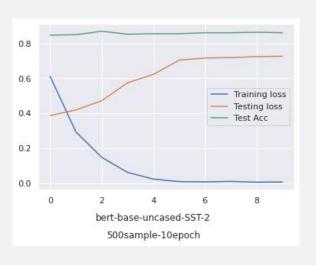


3.2 BERT微调

预处理:转小写-->去停用词-->去标点-->去链接-->去数字





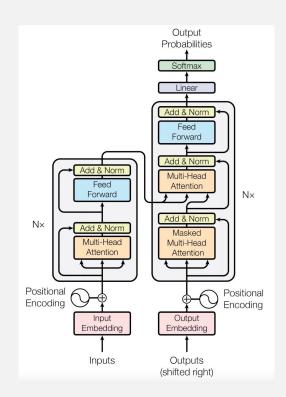


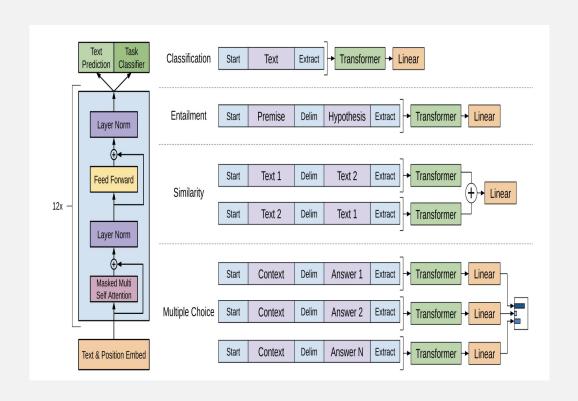
Accuracy=0.832

Accuracy=0.8635

Accuracy=0.869

3.3 GPT2微调



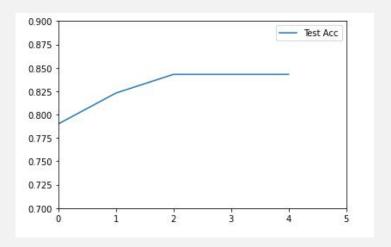


3.3 GPT2微调

数据: df = df. sample(1000)

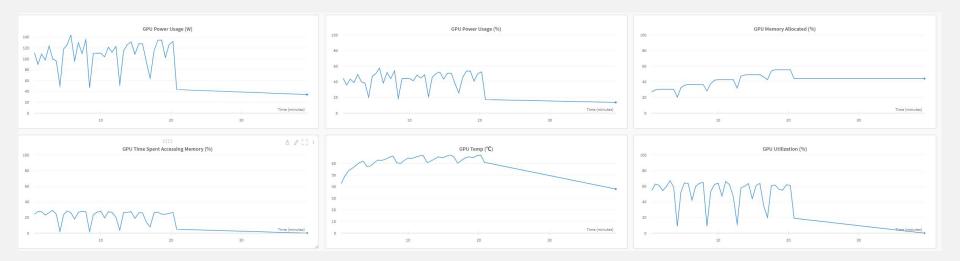
```
for trial_no in range(5):
print("Loading model...")
 # load tokenizer and model
 tokenizer = GPT2Tokenizer.from_pretrained(model_name, bos_token='<|startoftext|>',
                                          eos_token='<|endoftext|>', pad_token='<|pad|>')
 model = GPT2LMHeadModel.from_pretrained(model_name).cuda()
model.resize_token_embeddings(len(tokenizer))
 print("Loading dataset...")
 train_dataset, test_dataset = load_sentiment_dataset(tokenizer, trial_no)
 print("Start training...")
 training_args = TrainingArguments(output_dir='results', num_train_epochs=2,
                                logging_steps=200, load_best_model_at_end=True,
                                  save_strategy="epoch", per_device_train_batch_size=2, per_device_eval_batch_size=2,
                                warmup_steps=100, weight_decay=0.01, logging_dir='logs')
 Trainer(model=model, args=training_args, train_dataset=train_dataset,
        eval_dataset=test_dataset, data_collator=lambda data: {'input_ids': torch.stack([f[0] for f in data]),
                                                               'attention_mask': torch.stack([f[1] for f in data]),
                                                               'labels': torch.stack([f[0] for f in data])}).train()
 print("Start testing...")
 # eval mode on model
 _ = model.eval()
```

参数设置



Accuracy=0.843

3.3 GPT2微调





思考与展望









模型融合

- 1 我们主要使用了机器学习、BERT微调、GPT2微调三种方法进行了尝试。模型最高准确率0.869。
- 2 未来可以考虑加入日期、话题等 特征进行训练。

未来可以使用RoBERT、MacBERT等模型进行尝试,并进行模型融合。

感谢!

