

【概述】

Qwen1.5大模型微调、基于PEFT框架LoRA微调，在数据集[HC3-Chinese](#)上实现文本分类。

运行环境：[Kaggle - Notebook](#)

【数据处理】

1. 【数据下载】

```
1 import modelscope
2 from modelscope.msdatasets import MsDataset
3 # 【下载数据集】
4 HC3=MsDataset.load('simpleai/HC3-
  Chinese',subset_name='baike',split='train') #调用HC3数据集
5 dataset=HC3.to_hf_dataset() #将MsDataset转换成huggingface
  dataset格式，方便后续处理
6 print("【数据集下载完成】")
7 print(dataset)
8 print(dataset[0])
```

输出信息：

```
1 Dataset({
2     features: ['id', 'question', 'human_answers',
  'chatgpt_answers'],
3     num_rows: 4617
4 })
5
6 {
7     'id': '0',
8     'question': '我有一个计算机相关的问题...',
9     'human_answers': ['硬盘安装就是...'],
10    'chatgpt_answers': ['硬盘安装是指...']
11 }
```

2. 【格式调整】

将数据调整成形如 `{label: 0/1, text: '...'}` 的格式，在 9234 组数据中随机选 5000 个，按照 8:1:1 的比例划分训练集、验证集、测试集。

```
1 from datasets import Dataset
2 # 【调整数据集格式】
3 def data_init(dataset):
4     ds=[]
5     cnt=dataset.num_rows
6     for i in range(cnt):
7         example=dataset[i]
8
9         ds.append({"label":0,"text":example["human_answers"]
10                    [0]})
11
12         ds.append({"label":1,"text":example["chatgpt_answers"]
13                    [0]})
14     return Dataset.from_list(ds)
15
16 dataset=data_init(dataset) # 调整数据集内容
17 print(dataset)
18 dataset=dataset.shuffle(seed=233).select(range(5000)) #随
19 机选一部分
20
21 #数据集划分 train:val:test=8:1:1
22 data_=dataset.train_test_split(train_size=0.8,seed=233) #
23 数据集划分
24 data_train=data_["train"]
25 data__=data_["test"].train_test_split(train_size=0.5,seed
26 =233)
27 data_val=data__["train"]
28 data_test=data__["test"]
29
30 print(" 【data_train】 ",data_train)
31 print(" 【data_val】 ",data_val)
32 print(" 【data_test】 ",data_test)
```

输出信息：

```
1 Dataset({
2     features: ['label', 'text'],
3     num_rows: 9234
4 })
5 【data_train】 Dataset({
```

```

6     features: ['label', 'text'],
7     num_rows: 4000
8 })
9 【data_val】 Dataset({
10     features: ['label', 'text'],
11     num_rows: 500
12 })
13 【data_test】 Dataset({
14     features: ['label', 'text'],
15     num_rows: 500
16 })

```

【模型】

1. 【分词器】

文本信息在输入模型前，需要先用tokenizer分词。使用Dataset.map()函数快速处理。

```

1  from transformers import
    AutoTokenizer, AutoModelForSequenceClassification, Training
    Arguments, Trainer, DataCollatorWithPadding
2
3  # 【加载分词器】
4  tokenizer = AutoTokenizer.from_pretrained("Qwen/Qwen1.5-
    0.5B")
5  tokenizer.pad_token_id = tokenizer.eos_token_id #Qwen特
    性，需要指定一下pad_token_id
6
7  def tokenize_function(examples):
8      return
    tokenizer(examples["text"], padding="max_length", truncatio
    n=True, max_length=512)
9
10 token_train=data_train.map(tokenize_function,
    batched=True)
11 token_val=data_val.map(tokenize_function, batched=True)
12
13 train_dataset = token_train
14 eval_dataset = token_val

```

2. 【加载模型】

用 `AutoModelForSequenceClassification` 载入模型进行文本分类任务。

`num_labels` 为要分类的标签数量。

`from_pretrained()` 支持的模型在[这里](#)可以找到。

报错 `KeyError: 'qwen2'` 应该是 `transformers` 版本太旧。

```
1 # 【加载模型】
2 id2label = {0: "human", 1: "chatgpt"}
3 label2id = {"human": 0, "chatgpt": 1}
4 #使用Qwen1.5模型
5 model =
  AutoModelForSequenceClassification.from_pretrained("Qwen/Q
  wen1.5-
  0.5B", num_labels=2, id2label=id2label, label2id=label2id)
6 model.config.pad_token_id=model.config.eos_token_id #这里也
  要指定一下pad_token_id, 不然训练时会报错 "ValueError: Cannot
  handle batch sizes > 1 if no padding token is defined."
7 print("【model】\n", model)
8 print("【model.config】\n", model.config)
```

输出信息可以看到模型结构, 以及 `pad_token_id` (如果没有指定的话可以看到 `config` 里没有这个变量)

```
1 【model】
2 Qwen2ForSequenceClassification(
3   (model): Qwen2Model(
4     (embed_tokens): Embedding(151936, 1024)
5     (layers): ModuleList(
6       (0-23): 24 x Qwen2DecoderLayer(
7         (self_attn): Qwen2SdpaAttention(
8           (q_proj): Linear(in_features=1024,
9             out_features=1024, bias=True)
10          (k_proj): Linear(in_features=1024,
11            out_features=1024, bias=True)
12          (v_proj): Linear(in_features=1024,
13            out_features=1024, bias=True)
14          (o_proj): Linear(in_features=1024,
15            out_features=1024, bias=False)
16          (rotary_emb): Qwen2RotaryEmbedding()
17        )
18      )
19      (mlp): Qwen2MLP(
```

```

15         (gate_proj): Linear(in_features=1024,
out_features=2816, bias=False)
16         (up_proj): Linear(in_features=1024,
out_features=2816, bias=False)
17         (down_proj): Linear(in_features=2816,
out_features=1024, bias=False)
18         (act_fn): SiLU()
19     )
20     (input_layernorm): Qwen2RMSNorm()
21     (post_attention_layernorm): Qwen2RMSNorm()
22 )
23 )
24     (norm): Qwen2RMSNorm()
25 )
26     (score): Linear(in_features=1024, out_features=2,
bias=False)
27 )
28 【model.config】
29 Qwen2Config {
30     "_name_or_path": "Qwen/Qwen1.5-0.5B",
31     "architectures": [
32         "Qwen2ForCausalLM"
33     ],
34     "attention_dropout": 0.0,
35     "bos_token_id": 151643,
36     "eos_token_id": 151643,
37     "hidden_act": "silu",
38     "hidden_size": 1024,
39     "id2label": {
40         "0": "human",
41         "1": "chatgpt"
42     },
43     "initializer_range": 0.02,
44     "intermediate_size": 2816,
45     "label2id": {
46         "chatgpt": 1,
47         "human": 0
48     },
49     "max_position_embeddings": 32768,
50     "max_window_layers": 21,
51     "model_type": "qwen2",
52     "num_attention_heads": 16,
53     "num_hidden_layers": 24,

```

```

54     "num_key_value_heads": 16,
55     "pad_token_id": 151643,
56     "rms_norm_eps": 1e-06,
57     "rope_theta": 1000000.0,
58     "sliding_window": 32768,
59     "tie_word_embeddings": true,
60     "torch_dtype": "bfloat16",
61     "transformers_version": "4.41.2",
62     "use_cache": true,
63     "use_sliding_window": false,
64     "vocab_size": 151936
65 }

```

【训练】

1. 【训练参数】

```

1  # 【训练参数】
2  from datasets import load_metric
3  import numpy as np
4
5  training_args = TrainingArguments(
6      output_dir="pt_save_pretrained",
7      evaluation_strategy="epoch", #每跑完一个epoch输出一下测试
      信息
8      num_train_epochs=2,
9      per_device_train_batch_size=4, # 一共要跑
      len(dataset)/batch_size * epoch 个step
10                                     # [模型=Qwen1.5-0.5B,
      batch_size=4]: 完全微调显存13.3GB, LoRA微调显存8.7GB
11      save_strategy="no", #关闭自动保存模型 (Kaggle上磁盘空间不
      太多)
12  )
13
14  metric=load_metric('accuracy') #评估指标
15
16  def compute_metrics(eval_pred):
17      logits, labels = eval_pred
18      predictions = np.argmax(logits, axis=-1)
19      return metric.compute(predictions=predictions,
      references=labels)
20
21  def get_trainer(model):

```

```

22     return Trainer(
23         model=model,
24         args=training_args,
25         tokenizer=tokenizer,
26         train_dataset=train_dataset,
27         eval_dataset=eval_dataset,
28         compute_metrics=compute_metrics,
29
30         data_collator=DataCollatorWithPadding(tokenizer=tokenizer, padding=True, return_tensors="pt"), #给数据添加padding弄成batch
31     )

```

2. 【完全微调】

直接开始训练：

```

1  # 【完全微调】
2  print("【开始训练】")
3  trainer=get_trainer(model)
4  trainer.train()
5
6  #tokenizer.save_pretrained("./full_model_tokenizer")
7  #model.save_pretrained("./full_model")
8
9  #Kaggle注意：
10 每次训练之后restart以释放显存！
11  factory也reset一下，不然磁盘空间会爆！

```

训练效果：

[2000/2000 30:07, Epoch 2/2]

Epoch	Training Loss	Validation Loss	Accuracy
1	0.304600	0.217201	0.956000
2	0.099500	0.148824	0.962000

3. 【LoRA微调】

添加 LoRA 参数，调用peft框架：

```

1  # 【PEFT-LoRA微调】
2  from peft import LoraConfig, get_peft_model

```

```

3
4 peft_config = LoraConfig(
5     task_type="SEQ_CLS", #任务类型: 分类
6     target_modules=["q_proj", "k_proj", "v_proj", "o_proj"],
7     # 这个不同的模型需要设置不同的参数, 主要看模型中的attention层
8     inference_mode=False, # 关闭推理模式 (即开启训练模式)
9     r=8, # Lora 秩
10    lora_alpha=16, # Lora alpha, 具体作用参见 Lora 原理
11    lora_dropout=0.1 # Dropout 比例
12 )
13 peft_model = get_peft_model(model, peft_config) # 加载lora
14 # 参数peft框架
15 print('PEFT参数量: ')
16 peft_model.print_trainable_parameters()
17
18 print("【开始训练】")
19 peft_trainer=get_trainer(peft_model)
20 peft_trainer.train()
21
22 tokenizer.save_pretrained("./peft_model_tokenizer")
23 peft_model.save_pretrained("./peft_model")

```

从输出结果可以看到, LoRA 微调所要训练的参数只占 25.4%, 显著降低显存占用和训练时间:

```

1 | PEFT参数量:
2 | trainable params: 1,181,696 || all params: 465,171,456 ||
   | trainable%: 0.2540

```

训练效果:

[2000/2000 23:08, Epoch 2/2]

Epoch	Training Loss	Validation Loss	Accuracy
1	0.102500	0.113912	0.986000
2	0.052600	0.107569	0.986000

【测试】

1. 【代码】

```

1 | import torch

```



```

2  from transformers import
    DataCollatorWithPadding, AutoTokenizer, AutoModelForSequenceClassification
3
4  def classify(example, show): #对example进行预测
5      text=example["text"]
6      label=example["label"]
7      inputs = tokenizer(text, truncation=True,
padding=True, return_tensors="pt").to('cuda')
8      with torch.no_grad():
9          output = inference_model(**inputs)
10         pred = output.logits.argmax(dim=-1).item()
11     if show:
12         print("【预测{}!】 Label: {}, Pred_Label: {}\nText:
{}".format("正确" if label==pred else "错
误", id2label[label], id2label[pred], text))
13     else:
14         return pred, label
15
16 #inference_model=model.to('cuda')
17 tokenizer =
AutoTokenizer.from_pretrained("./peft_model_tokenizer")
18 inference_model =
AutoModelForSequenceClassification.from_pretrained("./pef
t_model").to('cuda') #读取训练好的模型
19 print("【model】\n", inference_model)
20 print("【model.config】\n", inference_model.config)
21 print("【model.config.pad_token_id】", inference_model.conf
ig.pad_token_id)
22 data_collator=DataCollatorWithPadding(tokenizer=tokenizer
, padding=True, return_tensors="pt")
23
24
25 id2label = {0: "human", 1: "chatgpt"}
26 label2id = {"human": 0, "chatgpt": 1}
27
28 classify(data_test[0], 1) #随便测试一个数据

```

输出：

- 1 **【预测正确!】**Label: human, Pred_Label: human
- 2 **Text:** 硬盘接口是硬盘与主机系统间的连接部件,作用是在硬盘缓存和主机内存之间传输数据。不同的硬盘接口决定着硬盘与计算机之间的连接速度,在整个系统中,硬盘接口的优劣直接影响着程序运行快慢和系统性能好坏。

在测试集上评估性能,二分类使用 `accuracy` 指标:

```
1 from datasets import load_metric
2 from tqdm import tqdm
3 metric=load_metric('accuracy')
4
5 print("【测试集】",data_test)
6 inference_model.eval()
7 for i,example in enumerate(tqdm(data_test)):
8     pred, label = classify(example,0)
9     metric.add(predictions=pred, references=label)
10 print(metric.compute())
```

2. 【结果展示】

使用 `Qwen1.5-1.8B-Chat`、给予 `prompt` 为 `{system="你擅长文本分类,能判断一段文字描述来自human还是chatgpt。你只需要回答一个单词,回答human或者chatgpt。", user="这里有一段关于xxx的文字描述,请判断它来自human还是chatgpt: \n xxx"}`,输出基本全是 `human`。

```
1 测试集大小: 30*2
2 Test NO.0: response=[人类,human] acc=(1/2)
3 Test NO.1: response=[人类,human] acc=(2/4)
4 Test NO.2: response=[human,human] acc=(3/6)
5 Test NO.3: response=[chatgpt,human] acc=(3/8)
6 Test NO.4: response=[human,human] acc=(4/10)
7 Test NO.5: response=[human,人类] acc=(5/12)
8 Test NO.6: response=[human,human] acc=(6/14)
9 Test NO.7: response=[human,human] acc=(7/16)
10 Test NO.8: response=[human,human] acc=(8/18)
11 Test NO.9: response=[human,human] acc=(9/20)
12 Test NO.10: response=[human,人类] acc=(10/22)
13 Test NO.11: response=[human,human] acc=(11/24)
14 Test NO.12: response=[chatgpt,human] acc=(11/26)
15 Test NO.13: response=[human,human] acc=(12/28)
16 Test NO.14: response=[human,human] acc=(13/30)
17 Test NO.15: response=[human,人类] acc=(14/32)
```

```
18 Test NO.16: response=[human,human] acc=(15/34)
19 Test NO.17: response=[人类,人类] acc=(16/36)
20 Test NO.18: response=[human,human] acc=(17/38)
21 Test NO.19: response=[human,human] acc=(18/40)
22 Test NO.20: response=[chatgpt,人类] acc=(18/42)
23 Test NO.21: response=[human,human] acc=(19/44)
24 Test NO.22: response=[chatgpt,human] acc=(19/46)
25 Test NO.23: response=[human,human] acc=(20/48)
26 Test NO.24: response=[human,human] acc=(21/50)
27 Test NO.25: response=[human,人类] acc=(22/52)
28 Test NO.26: response=[human,human] acc=(23/54)
29 Test NO.27: response=[human,human] acc=(24/56)
30 Test NO.28: response=[human,human] acc=(25/58)
31 Test NO.29: response=[人类,人类] acc=(26/60)
32
33 Test准确率: acc=0.43333333333333335
```

Qwen1.5-0.5B 完全微调:

```
1 【测试集】 Dataset({
2     features: ['label', 'text'],
3     num_rows: 500
4 })
5
6 {'accuracy': 0.944}
```

Qwen1.5-0.5B Lora微调:

```
1 【测试集】 Dataset({
2     features: ['label', 'text'],
3     num_rows: 500
4 })
5
6 {'accuracy': 0.982}
7
```