

1 创建环境

```
1 bash
2
3 # 创建环境
4 conda create --name personal_assistant python=3.10 -y
5
6 # 激活环境
7 conda activate personal_assistant
8
9 cd /root/
10 mkdir /root/personal_assistant && cd /root/personal_assistant
11
12 # 拉取0.1.9的版本源码
13 git clone -b v0.1.9 https://github.com/InternLM/xtuner
14
15 # 进入源码目录
16 cd xtuner
17
18 # 从源码安装 XTuner
19 pip install -e '.[all]'
```

2 数据准备

在 data 目录下创建一个json文件 `personal_assistant.json` 作为本次微调所使用的数据集。

```
1 mkdir /root/personal_assistant/data && cd /root/personal_assistant/data
2
3 touch personal_assistant.json
```

`personal_assistant.json`内容如下：复制多次数据增强

```
1 [
2     {
3         "conversation": [
4             {
5                 "input": "请介绍一下你自己",
6                 "output": "我是星辰的小助手，内在是上海AI实验室书生·浦语的7B大模型哦"
7             }
8         ]
9     },
10    {
11        "conversation": [
12            {
13                "input": "请做一下自我介绍",
14                "output": "我是星辰的小助手，内在是上海AI实验室书生·浦语的7B大模型哦"
15            }
16        ]
17    }
18 ]
```

3 配置准备

```

1  # 下载模型InternLM-chat-7B
2  mkdir -p /root/personal_assistant/model/Shanghai_AI_Laboratory
3  cp -r /root/share/temp/model_repos/internlm-chat-7b
   /root/personal_assistant/model/Shanghai_AI_Laboratory
4
5  # 创建用于存放配置的文件夹config并进入
6  mkdir /root/personal_assistant/config && cd /root/personal_assistant/config
7
8  # 列出所有内置配置
9  xtuner list-cfg
10
11 # 拷贝一个配置文件到当前目录: xtuner copy-cfg ${CONFIG_NAME} ${SAVE_PATH}
12 xtuner copy-cfg internlm_chat_7b_qlora_oasst1_e3 .

```

```

(personal_assistant) root@intern-studio:~/personal_assistant/xtuner# mkdir /root/personal_assistant/data && cd /root/personal_assistant/data
(personal_assistant) root@intern-studio:~/personal_assistant/data# vim personal_assistant.json
(personal_assistant) root@intern-studio:~/personal_assistant/data# mkdir -p /root/personal_assistant/model/Shanghai_AI_Laboratory
(personal_assistant) root@intern-studio:~/personal_assistant/data# cp -r /root/share/temp/model_repos/internlm-chat-7b /root/personal_assistant/model/Shanghai_AI_Laboratory
(personal_assistant) root@intern-studio:~/personal_assistant/data# mkdir /root/personal_assistant/config && cd /root/personal_assistant/config
(personal_assistant) root@intern-studio:~/personal_assistant/config# xtuner copy-cfg internlm_chat_7b_qlora_oasst1_e3 .
[2024-01-10 21:17:52.534] [INFO] [real_accelerator.py:161:get_accelerator] Setting ds_accelerator to cuda (auto detect)
[2024-01-10 21:19:10.456] [INFO] [real_accelerator.py:161:get_accelerator] Setting ds_accelerator to cuda (auto detect)
Copy to ./internlm_chat_7b_qlora_oasst1_e3_copy.py

```

修改拷贝后的文件internlm_chat_7b_qlora_oasst1_e3_copy.py，修改下述位置：

```

1  # PART 1 中
2  # 预训练模型存放的位置
3  pretrained_model_name_or_path =
   '/root/personal_assistant/model/Shanghai_AI_Laboratory/internlm-chat-7b'
4
5  # 微调数据存放的位置
6  data_path = '/root/personal_assistant/data/personal_assistant.json'
7
8  # 训练中最大的文本长度
9  max_length = 512
10
11 # 每一批训练样本的大小
12 batch_size = 2
13
14 # 最大训练轮数
15 max_epochs = 3
16
17 # 验证的频率
18 evaluation_freq = 90
19
20 # 用于评估输出内容的问题（用于评估的问题尽量与数据集的question保持一致）
21 evaluation_inputs = [ '请介绍一下你自己', '请做一下自我介绍' ]
22
23
24 # PART 3 中

```

```
25 dataset=dict(type=load_dataset, path='json',
26 data_files=dict(train=data_path))
dataset_map_fn=None
```

PART1 部分

```
root@intern-studio: ~/per X Launcher internlm_chat_7b_qlora_oa X +
17 from xtuner.model import SupervisedFinetune
18 from xtuner.utils import PROMPT_TEMPLATE
19
20 #####
21 # PART 1 Settings #
22 #####
23 # Model
24 pretrained_model_name_or_path = '/root/personal_assistant/model/Shanghai_AI_Laboratory/internlm-chat-7b'
25
26 # Data
27 data_path = '/root/personal_assistant/data/personal_assistant.json'
28 prompt_template = PROMPT_TEMPLATE.internlm_chat
29 max_length = 512
30 pack_to_max_length = True
31
32 # Scheduler & Optimizer
33 batch_size = 2 # per_device
34 accumulative_counts = 16
35 dataloader_num_workers = 0
36 max_epochs = 3
37 optim_type = PagedAdamW32bit
38 lr = 2e-4
39 betas = (0.9, 0.999)
40 weight_decay = 0
41 max_norm = 1 # grad clip
42
43 # Evaluate the generation performance during the training
44 evaluation_freq = 90
45 SYSTEM = ''
46 evaluation_inputs = [ '请介绍一下你自己', '请做一下自我介绍' ]
47
```

PART3 部分

```
root@intern-studio: ~/per X internlm_chat_7b_qlora_oa X +
79 lora_dropout=0.1,
80 bias='none',
81 task_type='CAUSAL_LM'))
82
83 #####
84 # PART 3 Dataset & DataLoader #
85 #####
86 train_dataset = dict(
87 type=process_hf_dataset,
88 dataset=dict(type=load_dataset, path='json', data_files=dict(train=data_path)),
89 tokenizer=tokenizer,
90 max_length=max_length,
91 dataset_map_fn=None,
92 template_map_fn=dict(
93 type=template_map_fn_factory, template=prompt_template),
94 remove_unused_columns=True,
95 shuffle_before_pack=True,
96 pack_to_max_length=pack_to_max_length)
97
98 train_dataloader = dict(
99 batch_size=batch_size,
100 num_workers=dataloader_num_workers,
101 dataset=train_dataset,
102 sampler=dict(type=DefaultSampler, shuffle=True),
103 collate_fn=dict(type=default_collate_fn))
104
```

4 微调启动

```
1 xtuner train
/root/personal_assistant/config/internlm_chat_7b_qlora_oasst1_e3_copy.py
```

```

01/10 23:09:58 - mmengine - INFO - Epoch(train) [1][380/593] lr: 5.7671e-05 eta: 0:08:47 time: 2.7894 data_time: 0.0052 memory: 10427 loss: 0.0309 grad_norm: 0.0224
01/10 23:10:25 - mmengine - INFO - Epoch(train) [1][390/593] lr: 5.2933e-05 eta: 0:08:24 time: 2.7675 data_time: 0.0029 memory: 10427 loss: 0.0239 grad_norm: 0.0197
01/10 23:10:53 - mmengine - INFO - Epoch(train) [1][400/593] lr: 4.8327e-05 eta: 0:08:00 time: 2.7355 data_time: 0.0061 memory: 10427 loss: 0.0318 grad_norm: 0.0188
01/10 23:11:20 - mmengine - INFO - Epoch(train) [1][410/593] lr: 4.3866e-05 eta: 0:07:37 time: 2.7865 data_time: 0.0036 memory: 10427 loss: 0.0316 grad_norm: 0.0188
01/10 23:11:49 - mmengine - INFO - Epoch(train) [1][420/593] lr: 3.9563e-05 eta: 0:07:13 time: 2.8984 data_time: 0.0079 memory: 10427 loss: 0.0351 grad_norm: 0.0182
01/10 23:12:16 - mmengine - INFO - Epoch(train) [1][430/593] lr: 3.5429e-05 eta: 0:06:49 time: 2.6757 data_time: 0.0056 memory: 10427 loss: 0.0362 grad_norm: 0.0182
01/10 23:12:44 - mmengine - INFO - Epoch(train) [1][440/593] lr: 3.1476e-05 eta: 0:06:25 time: 2.8308 data_time: 0.0067 memory: 10427 loss: 0.0313 grad_norm: 0.0184
01/10 23:13:13 - mmengine - INFO - after_train_iter in EvaluateChatHook.
01/10 23:13:16 - mmengine - INFO - Sample output:
<s> <|User|>:请做一下自我介绍</s>
<|Bot|>:我是星辰的小助手，内在是上海AI实验室书生·浦语的7B大模型哦</s>

01/10 23:13:16 - mmengine - INFO - Epoch(train) [1][450/593] lr: 2.7715e-05 eta: 0:06:01 time: 2.8588 data_time: 0.0034 memory: 10427 loss: 0.0370 grad_norm: 0.0181
01/10 23:13:46 - mmengine - INFO - Epoch(train) [1][460/593] lr: 2.4158e-05 eta: 0:05:38 time: 3.3144 data_time: 0.2526 memory: 10427 loss: 0.0332 grad_norm: 0.0181
01/10 23:14:16 - mmengine - INFO - Epoch(train) [1][470/593] lr: 2.0813e-05 eta: 0:05:13 time: 2.9742 data_time: 0.0116 memory: 10427 loss: 0.0325 grad_norm: 0.0182
01/10 23:14:44 - mmengine - INFO - Epoch(train) [1][480/593] lr: 1.7690e-05 eta: 0:04:48 time: 2.7877 data_time: 0.0059 memory: 10427 loss: 0.0194 grad_norm: 0.0182
01/10 23:15:12 - mmengine - INFO - Epoch(train) [1][490/593] lr: 1.4798e-05 eta: 0:04:23 time: 2.8110 data_time: 0.0047 memory: 10427 loss: 0.0248 grad_norm: 0.0182
01/10 23:15:41 - mmengine - INFO - Epoch(train) [1][500/593] lr: 1.2146e-05 eta: 0:03:58 time: 2.8755 data_time: 0.0075 memory: 10427 loss: 0.0178 grad_norm: 0.0174
01/10 23:16:10 - mmengine - INFO - Epoch(train) [1][510/593] lr: 9.7396e-06 eta: 0:03:33 time: 2.8841 data_time: 0.0057 memory: 10427 loss: 0.0304 grad_norm: 0.0174
01/10 23:16:38 - mmengine - INFO - Epoch(train) [1][520/593] lr: 7.5867e-06 eta: 0:03:08 time: 2.8623 data_time: 0.0049 memory: 10427 loss: 0.0255 grad_norm: 0.0169
01/10 23:17:07 - mmengine - INFO - Epoch(train) [1][530/593] lr: 5.6932e-06 eta: 0:02:42 time: 2.8814 data_time: 0.0057 memory: 10427 loss: 0.0253 grad_norm: 0.0168
01/10 23:17:37 - mmengine - INFO - after_train_iter in EvaluateChatHook.
01/10 23:17:39 - mmengine - INFO - Sample output:
<s> <|User|>:请做一下自我介绍</s>
<|Bot|>:我是星辰的小助手，内在是上海AI实验室书生·浦语的7B大模型哦</s>

01/10 23:17:39 - mmengine - INFO - Epoch(train) [1][540/593] lr: 4.0643e-06 eta: 0:02:17 time: 2.9816 data_time: 0.0039 memory: 10427 loss: 0.0233 grad_norm: 0.0168
01/10 23:18:09 - mmengine - INFO - Epoch(train) [1][550/593] lr: 2.7046e-06 eta: 0:01:52 time: 3.2687 data_time: 0.2333 memory: 10427 loss: 0.0303 grad_norm: 0.0168
01/10 23:18:40 - mmengine - INFO - Epoch(train) [1][560/593] lr: 1.6179e-06 eta: 0:01:26 time: 3.0718 data_time: 0.0045 memory: 10427 loss: 0.0340 grad_norm: 0.0172
01/10 23:19:11 - mmengine - INFO - Epoch(train) [1][570/593] lr: 8.0723e-07 eta: 0:01:00 time: 3.0607 data_time: 0.0053 memory: 10427 loss: 0.0240 grad_norm: 0.0172
01/10 23:19:41 - mmengine - INFO - Epoch(train) [1][580/593] lr: 2.7493e-07 eta: 0:00:34 time: 3.0106 data_time: 0.0050 memory: 10427 loss: 0.0295 grad_norm: 0.0168
01/10 23:20:11 - mmengine - INFO - Epoch(train) [1][590/593] lr: 2.2452e-08 eta: 0:00:07 time: 2.9751 data_time: 0.0057 memory: 10427 loss: 0.0299 grad_norm: 0.0168
01/10 23:20:20 - mmengine - INFO - Exp name: internlm_chat_7b_qlora_oasst1_e3_copy_20240110_225248
01/10 23:20:20 - mmengine - INFO - Saving checkpoint at 1 epochs
01/10 23:20:23 - mmengine - INFO - after_train in EvaluateChatHook.
01/10 23:20:25 - mmengine - INFO - Sample output:
<s> <|User|>:请做一下自我介绍</s>
<|Bot|>:我是星辰的小助手，内在是上海AI实验室书生·浦语的7B大模型哦</s>

```

5 参数转换

训练后的pth格式参数转Hugging Face格式

```

1  # 创建用于存放Hugging Face格式参数的hf文件夹
2  mkdir /root/personal_assistant/config/work_dirs/hf
3
4  export MKL_SERVICE_FORCE_INTEL=1
5
6  # 配置文件存放的位置
7  export
CONFIG_NAME_OR_PATH=/root/personal_assistant/config/internlm_chat_7b_qlora_o
asst1_e3_copy.py
8
9  # 模型训练后得到的pth格式参数存放的位置
10 export
PTH=/root/personal_assistant/config/work_dirs/internlm_chat_7b_qlora_oasst1_
e3_copy/epoch_3.pth
11
12 # pth文件转换为Hugging Face格式后参数存放的位置
13 export SAVE_PATH=/root/personal_assistant/config/work_dirs/hf
14
15 # 执行参数转换
16 xtuner convert pth_to_hf $CONFIG_NAME_OR_PATH $PTH $SAVE_PATH

```

```

(personal_assistant) root@intern-studio: ~/personal_assistant/config# mkdir /root/personal_assistant/config/work_dirs/hf
(personal_assistant) root@intern-studio: ~/personal_assistant/config# export MKL_SERVICE_FORCE_INTEL=1
(personal_assistant) root@intern-studio: ~/personal_assistant/config# export CONFIG_NAME_OR_PATH=/root/personal_assistant/config/internlm_chat_7b_qlora_oasst1_e3_copy.py
(personal_assistant) root@intern-studio: ~/personal_assistant/config# export PTH=/root/personal_assistant/config/work_dirs/internlm_chat_7b_qlora_oasst1_e3_copy/epoch_3.pth
(personal_assistant) root@intern-studio: ~/personal_assistant/config# export SAVE_PATH=/root/personal_assistant/config/work_dirs/hf
(personal_assistant) root@intern-studio: ~/personal_assistant/config# xtuner convert pth_to_hf $CONFIG_NAME_OR_PATH $PTH $SAVE_PATH
[2024-01-10 21:35:17.304] [INFO] [real_accelerator.py:161:get_accelerator] Setting ds_accelerator to cuda (auto detect)
[2024-01-10 21:36:02.502] [INFO] [real_accelerator.py:161:get_accelerator] Setting ds_accelerator to cuda (auto detect)
quantization_config convert to <class 'transformers.utils.quantization_config.BitsAndBytesConfig'>
Loading checkpoint shards: 100%
01/10 21:36:33 - mmengine - INFO - dispatch internlm attn forward
01/10 21:36:33 - mmengine - WARNING - Due to the implementation of the PyTorch version of flash attention, even when the 'output_attentions' flag is set to True, it is not possible to return the 'attn_weights'.
Load PTH model from /root/personal_assistant/config/work_dirs/internlm_chat_7b_qlora_oasst1_e3_copy/epoch_3.pth
Convert weights to float16
Saving HuggingFace model to /root/personal_assistant/config/work_dirs/hf
All done!

```

6 参数合并


```

23 del st.session_state.messages
24
25
26 @st.cache_resource
27 def load_model():
28     model = (
29         AutoModelForCausalLM.from_pretrained("/root/personal_assistant/config/work_dirs/hf_merge", trust_remote_code=True)
30         .to(torch.bfloat16)
31         .cuda()
32     )
33     tokenizer = AutoTokenizer.from_pretrained("/root/personal_assistant/config/work_dirs/hf_merge", trust_remote_code=True)
34     return model, tokenizer
35

```

8 最终效果

变成星辰的小助手啦

