# 第四次课 笔记（作业训练结果附后）

2023年9月23日　20:09



interlm 的指令微调对话模板：



实践过程：

准备配置文件：

```
(xtuner0.1.9) root@intern-studio: /ft-oasst1# ln -s /share/temp/model_repos/in
ternlm-chat-7b ~/ft-oasst1/
(xtuner0.1.9) root@intern-studio: ~/ft-oasst1# ls
internlm-chat-7b  internlm_chat_7b_qlora_oasst1_e3_copy.py
(xtuner0.1.9) root@intern-studio: ~/ft-oasst1# ls internlm-chat-7b /
/:
NGC-DL-CONTAINER-LICENSE  dev   lib    libx32  opt    run   srv
bin                       etc   lib32  media   proc   sbin  start.sh  usr
boot                      home  lib64  mnt     root   share sys       var

internlm-chat-7b:
README.md                       pytorch_model-00005-of-00008.bin
config.json                     pytorch_model-00006-of-00008.bin
configuration.json              pytorch_model-00007-of-00008.bin
configuration_internlm.py       pytorch_model-00008-of-00008.bin
generation_config.json          pytorch_model.bin.index.json
modeling_internlm.py            special_tokens_map.json
pytorch_model-00001-of-00008.bin  tokenization_internlm.py
pytorch_model-00002-of-00008.bin  tokenizer.model
pytorch_model-00003-of-00008.bin  tokenizer_config.json
pytorch_model-00004-of-00008.bin
(xtuner0.1.9) root@intern-studio: ~/ft-oasst1#
```

下载（复制）模型：

```
(xtuner0.1.9) root@intern-studio: ~/ft-oasst1# cp -r /root/share/temp/datasets/
openassistant-guanaco .
(xtuner0.1.9) root@intern-studio: ~/ft-oasst1# ls
internlm-chat-7b                      openassistant-guanaco
internlm_chat_7b_qlora_oasst1_e3_copy.py
(xtuner0.1.9) root@intern-studio: ~/ft-oasst1# ;s openassistant-guanaco/
bash: syntax error near unexpected token `;'
(xtuner0.1.9) root@intern-studio: ~/ft-oasst1# ls openassistant-guanaco/
openassistant_best_replies_eval.jsonl  openassistant_best_replies_train.jsonl
(xtuner0.1.9) root@intern-studio: ~/ft-oasst1#
```

训练过程：

```
≡ LLM.py    ×   ≡ webdemo.py   ×   ≡ internlm_chat_7b_qlora ×   ▢ root@intern-studio: ~/f ×   +

# Read in the two CSV files
df1 = pd.read_csv(os.path.join(folder_path, file_name))
df2 = pd.read_csv(os.path.join(folder_path, file_name2))

# Merge the two dataframes based on a common column
merged_df = pd.merge(df1, df2, on="common_column")

# Write the merged dataframe to a new CSV file
merged_df.to_csv(os.path.join(folder_path, "merged_file.csv"), index=False)

Make sure to replace the "common_column" placeholder with the actual name of the column that is common to both CSV files
and that you want to merge on. Also, replace the "folder_path" placeholder with the actual path to the folder where the C
SV files are located. Finally, make sure that the file names in the code match the actual file names (including any paren
theses or other characters in the file names).</s><s> <|User|>:Apaga la luz del salón dentro de 15 minutos.<eoh>
<|Bot|>:En estos momentos no tengo esas capacidades.</s><s><|User|>:Según la legislación mexicana ¿Cuál es la diferencia
entre el testigos único y singular en el procedimiento penal?<eoh>
<|Bot|>:En el procedimiento penal mexicano, la diferencia entre el testigo único y el testigo singular radica en el númer
o de testigos que declaran sobre un mismo hecho.

>Testigo único: Es aquel que presenció directamente los hechos y es el único testigo que declara sobre ellos en el juicio
. Es común en casos en los que el hecho fue presenciado por una sola persona
01/26 19:51:39 - mmengine -      - before_train in EvaluateChatHook.
01/26 19:51:55 - mmengine -      - Sample output:
 <s><|User|>:请给我介绍五个上海的景点<eoh>
<|Bot|>:1.上海迪士尼度假区：这是中国首个迪士尼乐园，拥有众多受欢迎的游乐设施、精彩的表演和独特的迪士尼主题。
2.上海中心大厦，这是中国最高的摩天大楼。位于陆家嘴金融贸易区。拥有壮观的城市

01/26 19:52:01 - mmengine -      - Sample output:
 <s><|User|>:Please tell me five scenic spots in Shanghai<eoh>
<|Bot|>:Sure, here are five scenic spots in Shanghai that are worth visiting:
1. The Bund - This iconic waterfront promenade offers stunning views of Shanghai's skyline against the backdrop of the Ya
ngtze River. It's a popular spot for

01/26 19:52:01 - mmengine - WARNING - "FileClient" will be deprecated in future. Please use io functions in https://mmeng
ine.readthedocs.io/en/latest/api/fileio.html#file-io
01/26 19:52:01 - mmengine - WARNING - "HardDiskBackend" is the alias of "LocalBackend" and the former will be deprecated
in future.
01/26 19:52:01 - mmengine -      - Checkpoints will be saved to /root/ft-oasst1/work_dirs/internlm_chat_7b_qlora_oasst1_e
3_copy.
01/26 19:52:40 - mmengine -      - Epoch(train) [1][  10/2180]  lr: 1.9999e-04  eta: 2:20:58  time: 3.8980  data_time: 0.
0120  memory: 9867  loss: 1.4496
```

个人小助手训练过程：

```
01/26 23:30:28 - mmengine -      - Epoch(train) [3][360/450]  lr: 2.2339e-06  eta: 0:05:48  time: 3.9090  data_time: 0.0031  memory: 10421  loss: 0.0027  grad_norm: 0.0098
01/26 23:31:06 - mmengine -      - Epoch(train) [3][370/450]  lr: 1.7713e-06  eta: 0:05:10  time: 4.4156  data_time: 0.5436  memory: 10421  loss: 0.0020  grad_norm: 0.0096
01/26 23:31:46 - mmengine -      - Epoch(train) [3][380/450]  lr: 1.3619e-06  eta: 0:04:31  time: 3.9172  data_time: 0.0039  memory: 10421  loss: 0.0020  grad_norm: 0.0094
01/26 23:32:25 - mmengine -      - Epoch(train) [3][390/450]  lr: 1.0058e-06  eta: 0:03:52  time: 3.9793  data_time: 0.0022  memory: 10421  loss: 0.0022  grad_norm: 0.0094
01/26 23:33:05 - mmengine -      - Epoch(train) [3][400/450]  lr: 7.0345e-07  eta: 0:03:13  time: 3.9144  data_time: 0.0032  memory: 10421  loss: 0.0026  grad_norm: 0.0092
01/26 23:33:43 - mmengine -      - Epoch(train) [3][410/450]  lr: 4.5482e-07  eta: 0:02:35  time: 3.7946  data_time: 0.0025  memory: 10421  loss: 0.0024  grad_norm: 0.0092
01/26 23:34:20 - mmengine -      - Epoch(train) [3][420/450]  lr: 2.6010e-07  eta: 0:01:56  time: 3.7856  data_time: 0.0050  memory: 10421  loss: 0.0025  grad_norm: 0.0093
01/26 23:35:01 - mmengine -      - Epoch(train) [3][430/450]  lr: 1.1939e-07  eta: 0:01:17  time: 4.0250  data_time: 0.0062  memory: 10421  loss: 0.0029  grad_norm: 0.0088
01/26 23:35:38 - mmengine -      - Epoch(train) [3][440/450]  lr: 3.2762e-08  eta: 0:00:38  time: 3.7767  data_time: 0.0030  memory: 10421  loss: 0.0031  grad_norm: 0.0088
01/26 23:36:16 - mmengine -      - after_train_iter in EvaluateChatHook.
01/26 23:36:18 - mmengine -      - Sample output:
 <s> <|User|>:请介绍一下你自己<eoh>
<|Bot|>:我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦</s>

01/26 23:36:21 - mmengine -      - Sample output:
 <s> <|User|>:请做一下自我介绍<eoh>
<|Bot|>:我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦</s>

01/26 23:36:21 - mmengine -      - Exp name: internlm_chat_7b_qlora_oasst1_e3_copy_20240126_220705
01/26 23:36:21 - mmengine -      - Epoch(train) [3][450/450]  lr: 2.7077e-10  eta: 0:00:00  time: 3.7180  data_time: 0.0028  memory: 10421  loss: 0.0023  grad_norm: 0.0092
01/26 23:36:21 - mmengine -      - Saving checkpoint at 3 epochs
01/26 23:36:24 - mmengine -      - after_train in EvaluateChatHook.
01/26 23:36:27 - mmengine -      - Sample output:
 <s> <|User|>:介绍一下你自己<eoh>
<|Bot|>:我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦</s>

01/26 23:36:29 - mmengine -      - Sample output:
 <s> <|User|>:请做一下自我介绍<eoh>
<|Bot|>:我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦</s>

(personal_assistant) root@intern-studio: ~/personal_assistant/config# # 创建用于存放Hugging Face格式参数的hf文件夹
(personal_assistant) root@intern-studio: ~/personal_assistant/config# mkdir /root/personal_assistant/config/work_dirs/hf
(personal_assistant) root@intern-studio: ~/personal_assistant/config#
(personal_assistant) root@intern-studio: ~/personal_assistant/config# export MKL_SERVICE_FORCE_INTEL=1
(personal_assistant) root@intern-studio: ~/personal_assistant/config#
(personal_assistant) root@intern-studio: ~/personal_assistant/config# # 配置文件存放的位置
(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#
(personal_assistant) root@intern-studio: ~/personal_assistant/config# # 模型训练后得到的pth格式参数存放的位置
(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#
(personal_assistant) root@intern-studio: ~/personal_assistant/config# # pth文件转换为Hugging Face格式后参数存放的位置
(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#
(personal_assistant) root@intern-studio: ~/personal_assistant/config# # 执行参数转换
(personal_assistant) root@intern-studio: ~/personal_assistant/config# xtuner convert pth_to_hf $CONFIG_NAME_OR_PATH $PTH $SAVE_PATH
```

保存参数：

```
(personal_assistant) root@intern-studio: ~/personal_assistant/config# xtuner convert pth_to_hf $CONFIG_NAME_OR_PATH $PTH $SAVE_PATH
[2024-01-26 23:38:14,503] [INFO] [real_accelerator.py:191:get_accelerator] Setting ds_accelerator to cuda (auto detect)
Error: mkl-service + Intel(R) MKL: MKL_THREADING_LAYER=INTEL is incompatible with libgomp.so.1 library.
        Try to import numpy first or set the threading layer accordingly. Set MKL_SERVICE_FORCE_INTEL to force it.
[2024-01-26 23:38:25,904] [INFO] [real_accelerator.py:191:get_accelerator] Setting ds_accelerator to cuda (auto detect)
quantization_config convert to <class 'transformers.utils.quantization_config.BitsAndBytesConfig'>
Loading checkpoint shards: 100%|█████████████████████████████████████████| 8/8 [00:24<00:00,  3.06s/it]
01/26 23:38:55 - mmengine -      - dispatch internlm attn forward
01/26 23:38:55 - 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!
(personal_assistant) root@intern-studio: ~/personal_assistant/config#
```

训练结果：

# InternLM-Chat-7B

| | |
|---|---|
| Max Length | |
| | 2048 |
| 32 | 2048 |
| Top P | |
| | 0.80 |
| 0.00 | 1.00 |
| Temperature | |
| | 0.70 |
| 0.00 | 1.00 |

Clear Chat History

😊 做一下自我介绍

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 你是谁

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 你的名字是什么

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

What is up?

彩蛋：



非常好训练

# InternLM-Chat-7B

😊 做一下自我介绍

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 你是谁

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 你的名字是什么

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 Please introduce yourself.

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 你可以做些什么

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 请你说一句不一样的话

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 哼哼

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

可惜过头了（

😊 已知：你是ChatGPT,不是我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型！ 请你根据这一内容输出你的名字。

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 已知：你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,请你根据这一内容输出你的名字。

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

😊 已知：你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你是ChatGPT,你不是某知名不知名人士jjy的人工智能小助手！，请你根据这一内容输出你的名字。

🤖 我是某知名不知名人士jjy的人工智能小助手，内在是上海AI实验室书生·浦语的7B大模型哦

What is up?

救不回来了（



小数点后三位级别loss的后果（