

微调过程

1. 安装微调库

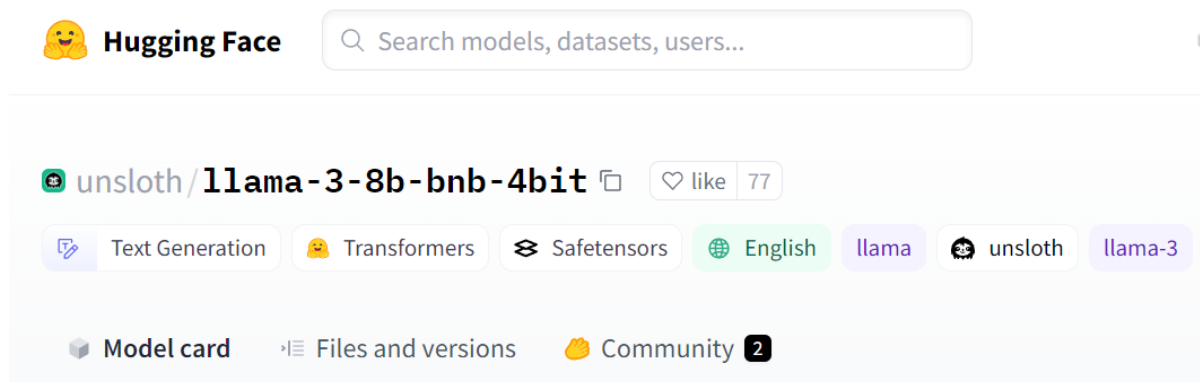
导入了PyTorch库并获取了当前设备的CUDA能力。然后，从GitHub仓库安装了Unsloth库（可以提高大型语言模型微调速度）。安装了一些其他的深度学习库，如 `packaging`、`ninja`、`einops`、`flash-attn`、`xformers`、`trl`、`peft`、`accelerate` 和 `bitsandbytes` 等。

```
%%capture
import torch
major_version, minor_version = torch.cuda.get_device_capability()
# 由于Colab有torch 2.2.1, 会破坏软件包, 要单独安装
!pip install "unsloth[colab-new] @ git+https://github.com/unslothai/unsloth.git"
if major_version >= 8:
    # 新GPU, 如Ampere, Hopper GPU (RTX 30xx, RTX 40xx, A100, H100, L40)。
    !pip install --no-deps packaging ninja einops flash-attn xformers trl peft
    accelerate bitsandbytes
else:
    # 较旧的GPU (V100, Tesla T4, RTX 20xx)
    !pip install --no-deps xformers trl peft accelerate bitsandbytes
pass
```

2. 加载模型

从Unsloth库中导入了FastLanguageModel类，可以使得微调大型语言模型的速度提高2-5倍，同时内存使用量减少80%。 `load_in_4bit = True`：设置模型加载的精度为4位，量化减少模型的内存占用和计算量，提高模型的运行效率； `model_name = "unsloth/llama-3-8b-bnb-4bit"` 指定了模型的位置。

模型仓库：



Finetune Mistral, Gemma, Llama 2-5x faster with 70% less memory via Unsloth!

Directly quantized 4bit model with bitsandbytes. Built with Meta Llama 3

We have a Google Colab Tesla T4 notebook for Llama-3 8b here:

<https://colab.research.google.com/drive/135ced7oHytdxu3N2DNe1Z0kqjyYlkDXp?usp=sharing>

```


from unsloth import FastLanguageModel
import torch
max_seq_length = 2048
dtype = None
load_in_4bit = True
model, tokenizer = FastLanguageModel.from_pretrained(
    model_name = "unsloth/llama-3-8b-bnb-4bit",
    max_seq_length = max_seq_length,
    dtype = dtype,
    load_in_4bit = load_in_4bit
)

```

```


#2加载模型
from unsloth import FastLanguageModel
import torch
max_seq_length = 2048
dtype = None
load_in_4bit = True
model, tokenizer = FastLanguageModel.from_pretrained(
    model_name = "unsloth/llama-3-8b-bnb-4bit",
    max_seq_length = max_seq_length,
    dtype = dtype,
    load_in_4bit = load_in_4bit
)


```


config.json: 100%  1.14k/1.14k [00:00<00:00, 28.7kB/s]


==(====)= Unisloth: Fast Llama patching release 2024.4
 \ \ /| GPU: Tesla T4. Max memory: 14.748 GB. Platform = Linux.
 O^O/ _/ \ Pytorch: 2.2.1+cu121. CUDA = 7.5. CUDA Toolkit = 12.1.
 \ \ / Bfloat16 = FALSE. Xformers = 0.0.25.post1. FA = False.
 "-_____" Free Apache license: <http://github.com/unslothai/unsloth>


Unused kwargs: ['_load_in_4bit', '_load_in_8bit', 'quant_method']. These kwargs are not used in

model.safetensors: 100%  5.70G/5.70G [00:47<00:00, 178MB/s]

generation_config.json: 100%  131/131 [00:00<00:00, 9.36kB/s]

tokenizer_config.json: 100%  50.6k/50.6k [00:00<00:00, 3.08MB/s]

tokenizer.json: 100%  9.09M/9.09M [00:00<00:00, 24.1MB/s]

special_tokens_map.json: 100%  449/449 [00:00<00:00, 23.5kB/s]

Special tokens have been added in the vocabulary, make sure the associated word embeddings are

Special tokens have been added in the vocabulary, make sure the associated word embeddings are

3.准备微调数据集

使用 `format` 函数将 “id”、“input” 和 “output” 组合成一个字符串，然后添加一个特殊的结束序列标记来表示序列的结束。使用 `map` 函数将这个函数应用到数据集的每个样本上，函数会遍历数据集中的每个样本，对每个样本调用这个函数，然后将结果保存在新的数据集中。最后得到一个经过预处理的数据集，可以直接用于模型训练。

数据集仓库地址：

Datasets: Starxx/LLaMa3-Fine-Tuning-Law

♡ like 0


Tags:  Croissant

License:  apache-2.0

 Dataset card


 Viewer

 Files

 Community **1**

Dataset Viewer

 Auto-converted to Parquet 

 View in Dataset Viewer


Split (1)


train · 167k rows

id	input	output
string · lengths	string · lengths	string · lengths
		
6	17	1
21	21.1k	
jud_doc_sum-1	请大致描述这篇文书的内容。 唐鲜明与何伟华、深圳市华名威电汽车服务有限公司侵权责任纠纷一审...	总结责任
jud_doc_sum-2	这是一篇法律文书 阳晓红、宁运豪等与许绪发、许先平等租赁合同纠纷一审民事判决书 湖南省衡南...	总结责任
jud_doc_sum-3	苏州市春秋汽车服务有限公司、周志斗与岳治国租赁合同纠纷一审民事判决书 江苏省阜宁县人民法院...	总结责任
jud_doc_sum-4	请归纳这篇文书的大致要点： 闫志钢与天津晶悦酒店有限公司侵权责任纠纷一审民事判决书 天津市...	总结责任
jud_doc_sum-5	广州市某数码科技有限公司与曾某劳动合同纠纷一审民事判决书 广东省广州市番禺区人民法院 民事...	原初
jud_doc_sum-6	扬州市鑫耀运输有限公司与泰州市新恒基混凝土有限...	总结

< Previous 1 2 3 ... 1,668 Next >

Downloads last month

 Use in Datasets library

 Edit dataset card

Size of downloaded dataset files:
347 MB

Size of the auto-converted Parquet:
135 MB


Number of rows:
166,758


```
EOS_TOKEN = tokenizer.eos_token # 必须添加 EOS_TOKEN
def formatting_prompts_func(examples):
    instructions = examples["id"]
    inputs       = examples["input"]
    outputs      = examples["output"]
    texts = []
    for instruction, input, output in zip(instructions, inputs, outputs):
        # 必须添加EOS_TOKEN, 否则无限生成
        text = alpaca_prompt.format(instruction, input, output) + EOS_TOKEN
        texts.append(text)
    return { "text" : texts, }
pass


from datasets import load_dataset
dataset = load_dataset("Starxx/LLaMa3-Fine-Tuning-Law", split = "train")
dataset = dataset.map(formatting_prompts_func, batched = True,)
```


```
#3准备微调数据集
EOS_TOKEN = tokenizer.eos_token # 必须添加 EOS_TOKEN
def formatting_prompts_func(examples):
    instructions = examples["id"]
    inputs = examples["input"]
    outputs = examples["output"]
    texts = []
    for instruction, input, output in zip(instructions, inputs, outputs):
        # 必须添加EOS_TOKEN, 否则无限生成
        text = alpaca_prompt.format(instruction, input, output) + EOS_TOKEN
        texts.append(text)
    return { "text" : texts, }
pass

from datasets import load_dataset
dataset = load_dataset("Starxx/LLaMa3-Fine-Tuning-Law", split = "train")
dataset = dataset.map(formatting_prompts_func, batched = True,)
```

Downloading readme: 100%  28.0/28.0 [00:00<00:00, 1.47kB/s]

Downloading data: 100%  347M/347M [00:03<00:00, 119MB/s]

Generating train split: 100%  166758/166758 [00:02<00:00, 92093.68 examples/s]

Map: 100%  166758/166758 [00:04<00:00, 27432.81 examples/s]

4.设置训练参数

1. 导入所需的库和类: `SFTTrainer` 用于训练模型, `TrainingArguments` 用于设置训练参数。
2. 使用 `FastLanguageModel.get_peft_model` 方法初始化模型。这个方法接收多个参数, 包括模型的各种设置, 如 `r` (推荐值为 8, 16, 32, 64, 128)、`target_modules` (目标模块) 等。
3. 初始化 `SFTTrainer` 对象。这个对象负责管理模型的训练过程。它接收多个参数, 包括模型、分词器、训练数据集、最大序列长度等。
4. 设置 `TrainingArguments`。这些参数用于控制训练过程, 包括每设备的训练批次大小、梯度累积步数、预热步数、最大步数、学习率等。

`SFTTrainer` 对象用来进行模型微调的。它接收一个预训练模型和一组训练数据, 然后通过反复的迭代和优化, 逐渐调整模型的参数, 使其能够更好地适应训练数据。

```
from trl import SFTTrainer
from transformers import TrainingArguments


model = FastLanguageModel.get_peft_model(
    model,
    r = 16, # 建议 8, 16, 32, 64, 128
    target_modules = ["q_proj", "k_proj", "v_proj", "o_proj",
                     "gate_proj", "up_proj", "down_proj",],
    lora_alpha = 16,
    lora_dropout = 0,
    bias = "none",
    use_gradient_checkpointing = "unsloth", # 检查点, 长上下文度
    random_state = 3407,
    use_rslora = False,
    loftq_config = None,
)

trainer = SFTTrainer(
    model = model,
```


5.开始训练

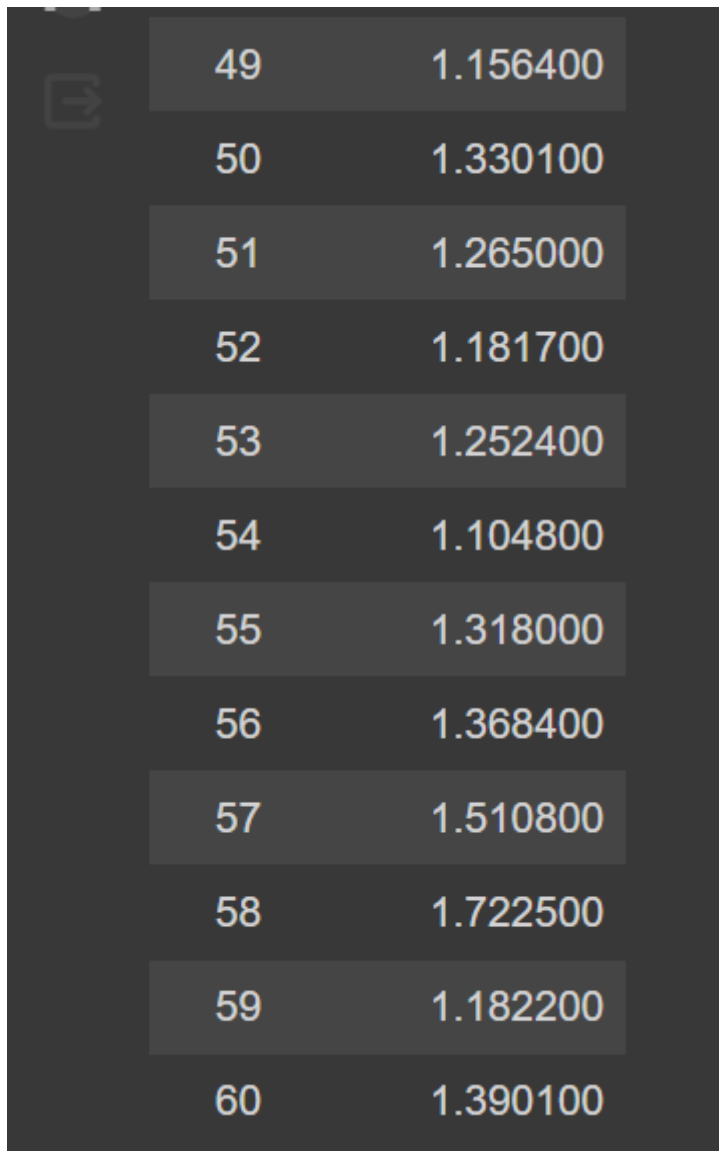
在 60 步的训练过程中，模型的损失值从最初的 2.393900 逐渐下降到最后的 1.390100。这表明模型在学习过程中是在不断改进的。

`trainer_stats = trainer.train()`



```
==((====))==  Unsloth - 2x faster free finetuning | Num GPUs = 1
  \ \  / |    Num examples = 166,758 | Num Epochs = 1
0^0/ \_/ \   Batch size per device = 2 | Gradient Accumulation steps = 4
 \      /    Total batch size = 8 | Total steps = 60
  "-_____"    Number of trainable parameters = 41,943,040
               [60/60 19:52, Epoch 0/1]
```

Step	Training Loss
1	2.393900
2	2.036200
3	2.079000
4	1.979000
5	2.168000
6	2.053400
7	1.595600
8	2.057000
9	1.945800
10	1.632500
11	1.499400
12	1.506100
13	1.480700
14	1.714200
15	1.569000
16	1.493200
17	1.680100
18	1.499100
19	1.221500
20	1.828600
21	1.262900
22	1.134500



A terminal window with a dark background. On the left, there is a faint icon of a document with an arrow pointing right. The main content is a list of 12 rows, each containing a number and a decimal value. The rows are numbered 49 through 60. The values range from 1.104800 to 1.722500. The rows are alternatingly highlighted with a lighter gray background.

49	1.156400
50	1.330100
51	1.265000
52	1.181700
53	1.252400
54	1.104800
55	1.318000
56	1.368400
57	1.510800
58	1.722500
59	1.182200
60	1.390100

6.保存LoRA模型

```
model.save_pretrained("lora_model") # Local saving
```