

LoRA详解介绍及代码实战

LORA

1. 背景
2. 出发点
3. 具体实现方法
 - 3.1. 公式
4. 4.代码实战

LORA

1. 背景

论文标题:LoRA: Low-Rank Adaptation of Large Language Models

论文链接:<https://arxiv.org/pdf/2106.09685.pdf>

Aghajanyan的研究表明：预训练模型拥有极小的内在维度(instrinsic dimension)，即存在一个极低维度的参数，微调它和在全参数空间中微调能起到相同的效果。同时发现在预训练后，越大的模型有越小的内在维度，这也解释了为何大模型都拥有很好的few-shot能力。

2. 出发点

随着大模型的发展，尤其是chatGPT出现之后，175B的参数量，训练起来非常昂贵。

微软提出了低秩自适应（LoRA），LoRA的主要思想很简单，冻结预训练模型的权重参数，在原始的预训练模型（PLM）旁边增加一个新的通路，通过前后两个矩阵A,B相乘，第一个矩阵A负责降维，第二个矩阵B负责升维，中间层维度为r，在微调下游任务的时候，只更新A和B，该方法的核心思想就是通过低秩分解来模拟参数的改变量，从而以极小的参数量来实现大模型的间接训练。

LoRA方法优点：

- 预训练模型可以共享，针对下游任务可以构建多个不同任务的LoRA模块。冻结预训练模型参数共享，通过替换矩阵A和B来高效地切换任务，从而显著降低存储需求和任务切换开销。

- LoRA使训练更加的高效，将硬件的进入门槛降低了3倍，相同的内存下，可以微调更大参数的模型
- 线性设计允许我们在部署时将可训练矩阵与冻结权重合并，和完全微调的模型相比，不会引入推理延迟。

3. 具体实现方法

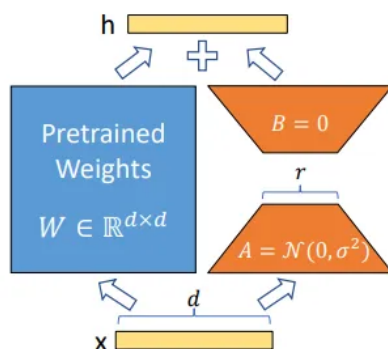


Figure 1: Our reparametrization. We only train A and B .

3.1. 公式

LoRA中是让模型学习BA，去近似SVD分解的结果

$$h = W_0 x + \Delta W x = W_0 x + B A x$$

- 在训练过程中 W_0 被冻结，不接收梯度更新，而 A 和 B 包含可训练参数。
- 在初始化的时候，我们对使用 A 随机高斯初始化，对 B 使用零初始化。
- 因此，在训练开始时为0，然后我们可以通过 $\frac{\alpha}{r}$ 对 ΔW 进行缩放， r 是秩 在推理时，我们通过上图可知，将左右两部分进行相加即可，不会添加额外的计算资源。

在微调过程中，所有做lora适配器的module，它们的 r 都是一致的，且在训练过程中不会改变。

Transformer的权重矩阵包括Attention模块里用于计算query, key, value的 W_q, W_k, W_v 以及多头attention的 W_o ,以及MLP层的权重矩阵，在LoRA原始论文中，作者通过消融实验发现最终选择对attention模块的 W_q, W_v 做低秩适配产生最佳结果。

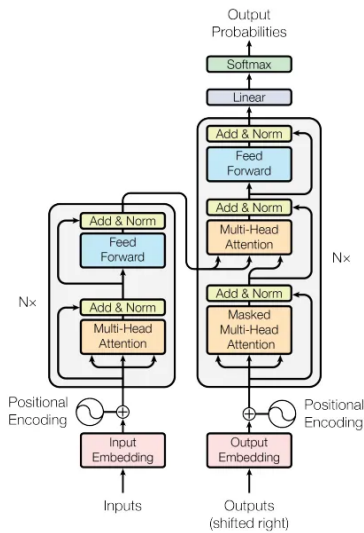
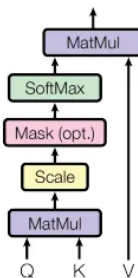


Figure 1: The Transformer - model architecture.

Scaled Dot-Product Attention



Multi-Head Attention

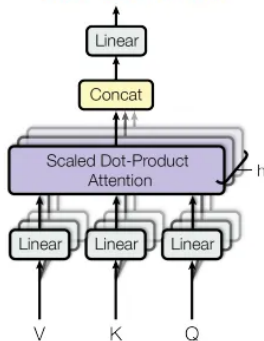
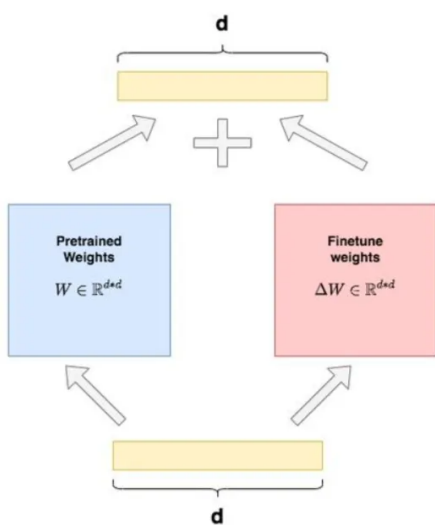
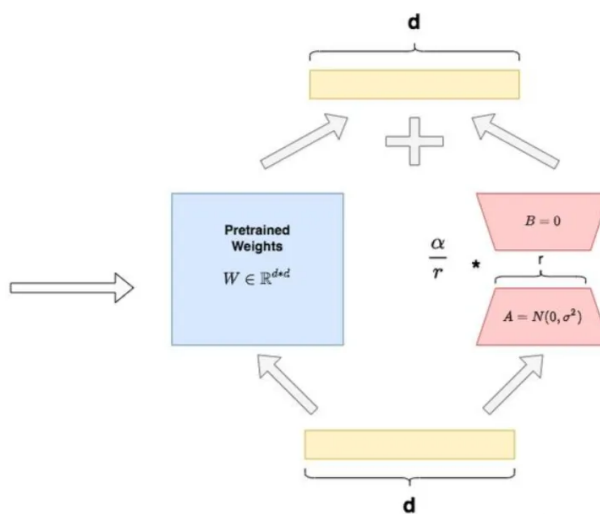


Figure 2: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention consists of several attention layers running in parallel.

全参数finetune



LoRA finetune



4. 4.代码实战

```
1  from transformers import AutoModelForCausalLM
2  from peft import get_peft_config, get_peft_model, get_peft_model_state_dict, LoraConfig, TaskType
3
4  import torch
5  from datasets import load_dataset
6  import os
7  from transformers import AutoTokenizer
8  from torch.utils.data import DataLoader
9  from transformers import default_data_collator, get_linear_schedule_with_warmup
10 from tqdm import tqdm
11 from datasets import load_dataset
12
13
14 device = "cuda"
15
16 model_name_or_path = "/data/nfs/llm/model/bloomz-560m"
17 tokenizer_name_or_path = "/data/nfs/llm/model/bloomz-560m"
18
19 ##### 加载配置 #####
20 peft_config = LoraConfig(task_type=TaskType.CAUSAL_LM,
21                          inference_mode=False, r=8,
22                          lora_alpha=32,
23                          lora_dropout=0.1)
24
25 dataset_name = "twitter_complaints"
26 checkpoint_name = f"{dataset_name}_{model_name_or_path}_{peft_config.peft_type}_{peft_config.task_type}_v1.pt".replace("/", "_")
27 text_column = "Tweet text"
28 label_column = "text_label"
29 max_length = 64
30 lr = 3e-2
31 num_epochs = 10
32 batch_size = 8
33
34
35 from datasets import load_dataset
36
37 # dataset = load_dataset("ought/raft", dataset_name)
38 dataset = load_dataset("/home/guodong.li/data/peft/raft/raft.py", dataset_name, cache_dir="/home/guodong.li/data/peft/data")
39
40 classes = [k.replace("_", " ") for k in dataset["train"].features["Label"].names]
```

```

41 print(classes)
42 dataset = dataset.map(
43     lambda x: {"text_label": [classes[label] for label in x["Label"]]},
44     batched=True,
45     num_proc=1,
46 )
47 print(dataset)
48 dataset["train"][0]
49
50
51 # data preprocessing
52 tokenizer = AutoTokenizer.from_pretrained(model_name_or_path)
53 if tokenizer.pad_token_id is None:
54     tokenizer.pad_token_id = tokenizer.eos_token_id
55 target_max_length = max([len(tokenizer(class_label)["input_ids"]) for class_label in classes])
56 print("target_max_length:", target_max_length)
57
58
59 def preprocess_function(examples):
60     batch_size = len(examples[text_column])
61     inputs = [f"{text_column} : {x} Label : " for x in examples[text_column]]
62     targets = [str(x) for x in examples[label_column]]
63     model_inputs = tokenizer(inputs)
64     labels = tokenizer(targets)
65     for i in range(batch_size):
66         sample_input_ids = model_inputs["input_ids"][i]
67         label_input_ids = labels["input_ids"][i] + [tokenizer.pad_token_id]
68
69         # print(i, sample_input_ids, label_input_ids)
70         model_inputs["input_ids"][i] = sample_input_ids + label_input_ids
71         labels["input_ids"][i] = [-100] * len(sample_input_ids) + label_input_ids
72
73         model_inputs["attention_mask"][i] = [1] * len(model_inputs["input_ids"][i])
74
75         # print(model_inputs)
76     for i in range(batch_size):
77         sample_input_ids = model_inputs["input_ids"][i]
78         label_input_ids = labels["input_ids"][i]
79         model_inputs["input_ids"][i] = [tokenizer.pad_token_id] * (
80             max_length - len(sample_input_ids)
81         ) + sample_input_ids
82         model_inputs["attention_mask"][i] = [0] * (max_length - len(sample_input_ids)) + model_inputs["attention_mask"]

```

```

83         labels["input_ids"][i] = [-100] * (max_length - len(sample_input_
            ids)) + label_input_ids
84         model_inputs["input_ids"][i] = torch.tensor(model_inputs["input_i
            ds"][i][:max_length])
85         model_inputs["attention_mask"][i] = torch.tensor(model_inputs["at
            tention_mask"][i][:max_length])
86         labels["input_ids"][i] = torch.tensor(labels["input_ids"][i][:max
            _length])
87         model_inputs["labels"] = labels["input_ids"]
88         return model_inputs
89
90
91     processed_datasets = dataset.map(
92         preprocess_function,
93         batched=True,
94         num_proc=1,
95         remove_columns=dataset["train"].column_names,
96         load_from_cache_file=False,
97         desc="Running tokenizer on dataset",
98     )
99
100     train_dataset = processed_datasets["train"]
101     eval_dataset = processed_datasets["train"]
102
103
104     train_dataloader = DataLoader(train_dataset, shuffle=True, collate_fn=def
        ault_data_collator, batch_size=batch_size, pin_memory=True)
105     eval_dataloader = DataLoader(eval_dataset, collate_fn=default_data_collat
        or, batch_size=batch_size, pin_memory=True)
106     def test_preprocess_function(examples):
107         batch_size = len(examples[text_column])
108         inputs = [f"{text_column} : {x} Label : " for x in examples[text_colu
            mn]]
109         model_inputs = tokenizer(inputs)
110         # print(model_inputs)
111         for i in range(batch_size):
112             sample_input_ids = model_inputs["input_ids"][i]
113             model_inputs["input_ids"][i] = [tokenizer.pad_token_id] * (max_le
                ngth - len(sample_input_ids)) + sample_input_ids
114             model_inputs["attention_mask"][i] = [0] * (max_length - len(sampl
                e_input_ids)) + model_inputs["attention_mask"][i]
115
116             model_inputs["input_ids"][i] = torch.tensor(model_inputs["input_i
                ds"][i][:max_length])
117             model_inputs["attention_mask"][i] = torch.tensor(model_inputs["at
                tention_mask"][i][:max_length])
118             return model_inputs
119

```

```

120
121 test_dataset = dataset["test"].map(
122     test_preprocess_function,
123     batched=True,
124     num_proc=1,
125     remove_columns=dataset["train"].column_names,
126     load_from_cache_file=False,
127     desc="Running tokenizer on dataset",
128 )
129
130 test_dataloader = DataLoader(test_dataset, collate_fn=default_data_collat
131 or, batch_size=batch_size, pin_memory=True)
132 next(iter(test_dataloader))
133
134 ##### 加载模型 #####
135 model = AutoModelForCausalLM.from_pretrained(model_name_or_path)
136 model = get_peft_model(model, peft_config)
137 model.print_trainable_parameters()
138
139 # model
140 # optimizer and lr scheduler
141 optimizer = torch.optim.AdamW(model.parameters(), lr=lr)
142 lr_scheduler = get_linear_schedule_with_warmup(
143     optimizer=optimizer,
144     num_warmup_steps=0,
145     num_training_steps=(len(train_dataloader) * num_epochs),
146 )
147
148
149 # training and evaluation
150 model = model.to(device)
151
152 for epoch in range(num_epochs):
153     model.train()
154     total_loss = 0
155     for step, batch in enumerate(tqdm(train_dataloader)):
156         batch = {k: v.to(device) for k, v in batch.items()}
157         # print(batch)
158         # print(batch["input_ids"].shape)
159         outputs = model(**batch)
160         loss = outputs.loss
161         total_loss += loss.detach().float()
162         loss.backward()
163         optimizer.step()
164         lr_scheduler.step()
165         optimizer.zero_grad()
166

```



```

167     model.eval()
168     eval_loss = 0
169     eval_preds = []
170     for step, batch in enumerate(tqdm(eval_dataloader)):
171         batch = {k: v.to(device) for k, v in batch.items()}
172         with torch.no_grad():
173             outputs = model(**batch)
174             loss = outputs.loss
175             eval_loss += loss.detach().float()
176             eval_preds.extend(
177                 tokenizer.batch_decode(torch.argmax(outputs.logits, -1).detach()
178                                     .cpu().numpy(), skip_special_tokens=True)
179             )
180
181     eval_epoch_loss = eval_loss / len(eval_dataloader)
182     eval_ppl = torch.exp(eval_epoch_loss)
183     train_epoch_loss = total_loss / len(train_dataloader)
184     train_ppl = torch.exp(train_epoch_loss)
185     print(f"{epoch=:} {train_ppl=} {train_epoch_loss=} {eval_ppl=} {eval_
186           epoch_loss=}")
187
188     第四步，模型训练的其余部分均无需更改，
189     当模型训练完成之后，保存高效微调部分的模型权重以供模型推理即可。
190     peft_model_id = f"{model_name_or_path}_{peft_config.peft_type}_{peft_conf
191                     ig.task_type}"
192     model.save_pretrained(peft_model_id)

```

```
1
2 第五步，加载微调后的权重文件进行推理。
3 from peft import PeftModel, PeftConfig
4
5 peft_model_id = f"{model_name_or_path}_{peft_config.peft_type}_{peft_config.task_type}"
6 config = PeftConfig.from_pretrained(peft_model_id)
7 # 加载基础模型
8 model = AutoModelForCausalLM.from_pretrained(config.base_model_name_or_path)
9 # 加载PEFT模型
10 model = PeftModel.from_pretrained(model, peft_model_id)
11
12 # tokenizer编码
13 inputs = tokenizer(f'{text_column} : {dataset["test"][i]["Tweet text"]} Label : ', return_tensors="pt")
14
15 # 模型推理
16 outputs = model.generate(
17     input_ids=inputs["input_ids"],
18     attention_mask=inputs["attention_mask"],
19     max_new_tokens=10,
20     eos_token_id=3
21 )
22
23 # tokenizer解码
24 print(tokenizer.batch_decode(outputs.detach().cpu().numpy(), skip_special_tokens=True))
25
```