

Prompt Tuning详解及代码实战

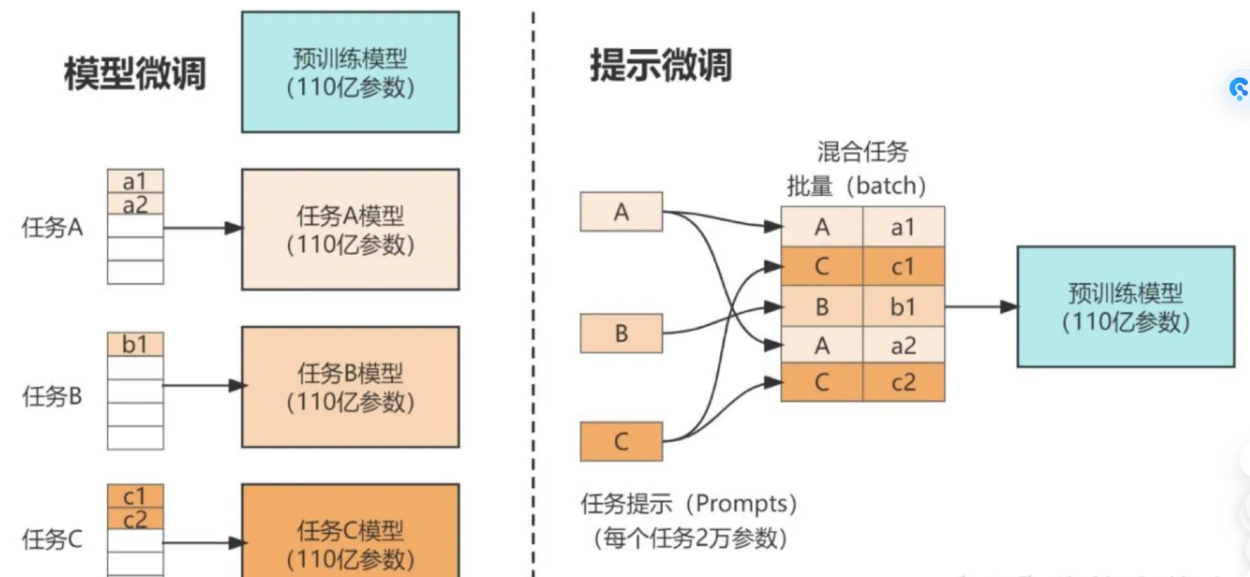
- 1. 概念&原理
- 2. 实验论证
- 3. Prompt Ensembling
- 4. 代码实践

Prompt Tuning 是2021年谷歌在论文《The Power of Scale for Parameter-Efficient Prompt Tuning》中提出的微调方法。

1. 概念&原理

该方法可以看作是 Prefix Tuning 的简化版本，只在输入层加入 prompt tokens，并不需要加入 MLP 进行调整来解决难训练的问题。主要在 T5 预训练模型上做实验。似乎只要预训练模型足够强大，其他的一切都不是问题。作者也做实验说明随着预训练模型参数量的增加，Prompt Tuning的方法会逼近 Fine-tune 的结果。

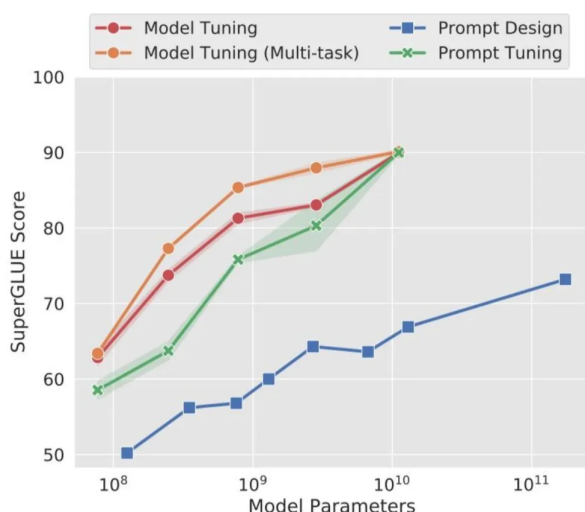
固定预训练参数，为每一个任务额外添加一个或多个 embedding，之后拼接 query 正常输入 LLM，并只训练这些 embedding。左图为单任务全参数微调，右图为 Prompt tuning。



2. 实验论证

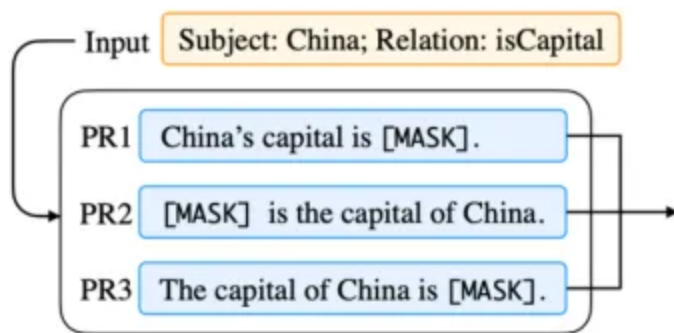
作者做了一系列对比实验，都在说明：随着预训练模型参数的增加，一切的问题都不是问题，最简单的设置也能达到极好的效果。

- **Prompt 长度影响**：模型参数达到一定量级时，Prompt 长度为1也能达到不错的效果，Prompt 长度为20就能达到极好效果。
- **Prompt初始化方式影响**：Random Uniform 方式明显弱于其他两种，但是当模型参数达到一定量级，这种差异也不复存在。
- **预训练的方式**：LM Adaptation 的方式效果好，但是当模型达到一定规模，差异又几乎没有了。
- **微调步数影响**：模型参数较小时，步数越多，效果越好。同样随着模型参数达到一定规模，zero shot 也能取得不错效果。
- 当参数达到100亿规模与全参数微调方式效果无异。



3. Prompt Ensembling

同时，Prompt Tuning 还提出了 Prompt Ensembling，也就是在一个批次（Batch）里同时训练同一个任务的不同 prompt（即采用多种不同方式询问同一个问题），这样相当于训练了不同模型，比模型集成的成本小多了。



(a) Prompt Ensembling.

© 稀土掘金技术社区

4. 代码实践

https://github.com/huggingface/peft/blob/main/src/peft/tuners/prompt_tuning/model.py

```
1 from transformers import AutoModelForCausalLM
2 from peft import get_peft_config, get_peft_model, PromptTuningInit, PromptTuningConfig, TaskType, PeftType
3 import torch
4 from datasets import load_dataset
5 import os
6 from transformers import AutoTokenizer
7 from torch.utils.data import DataLoader
8 from transformers import default_data_collator, get_linear_schedule_with_warmup
9 from tqdm import tqdm
10 from datasets import load_dataset
11
12
13 device = "cuda"
14
15 model_name_or_path = "/data/nfs/llm/model/bloomz-560m"
16 tokenizer_name_or_path = "/data/nfs/llm/model/bloomz-560m"
17
18 peft_config = PromptTuningConfig(
19     task_type=TaskType.CAUSAL_LM,
20     prompt_tuning_init=PromptTuningInit.TEXT,
21     num_virtual_tokens=8,
22     prompt_tuning_init_text="Classify if the tweet is a complaint or not:",
23     tokenizer_name_or_path=model_name_or_path,
24 )
25
26 dataset_name = "twitter_complaints"
27
28 text_column = "Tweet text"
29 label_column = "text_label"
30 max_length = 64
31 lr = 3e-2
32 num_epochs = 10
33 batch_size = 8
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
43 dataset = dataset.map(
44     lambda x: {"text_label": [classes[label] for label in x["Label"]]},
45     batched=True,
46     num_proc=1,
47 )
48 print(dataset)
49
50 dataset["train"][0]
51
52
53 # data preprocessing
54 tokenizer = AutoTokenizer.from_pretrained(model_name_or_path)
55 if tokenizer.pad_token_id is None:
56     tokenizer.pad_token_id = tokenizer.eos_token_id
57
58 target_max_length = max([len(tokenizer(class_label)["input_ids"]) for class_label in classes])
59 print("target_max_length:", target_max_length)
60
61
62 # 预处理
63 def preprocess_function(examples):
64     batch_size = len(examples[text_column])
65     print("batch_size:", batch_size)
66
67     inputs = [f"{text_column} : {x} Label : " for x in examples[text_column]]
68     targets = [str(x) for x in examples[label_column]]
69
70     model_inputs = tokenizer(inputs)
71     labels = tokenizer(targets)
72
73     for i in range(batch_size):
74         sample_input_ids = model_inputs["input_ids"][i]
75         label_input_ids = labels["input_ids"][i] + [tokenizer.pad_token_id]
76
77         if i == 0:
78             print(i, sample_input_ids, label_input_ids)
79             model_inputs["input_ids"][i] = sample_input_ids + label_input_ids
80             labels["input_ids"][i] = [-100] * len(sample_input_ids) + label_input_ids
81
82     model_inputs["attention_mask"][i] = [1] * len(model_inputs["input_ids"][i])
83     #print(model_inputs)
84
85     for i in range(batch_size):

```

```

84         sample_input_ids = model_inputs["input_ids"][i]
85         label_input_ids = labels["input_ids"][i]
86
87         model_inputs["input_ids"][i] = [tokenizer.pad_token_id] * (max_length - len(sample_input_ids)) + sample_input_ids
88         model_inputs["attention_mask"][i] = [0] * (max_length - len(sample_input_ids)) + model_inputs["attention_mask"][i]
89         labels["input_ids"][i] = [-100] * (max_length - len(sample_input_ids)) + label_input_ids
90
91         model_inputs["input_ids"][i] = torch.tensor(model_inputs["input_ids"][i][:max_length])
92         model_inputs["attention_mask"][i] = torch.tensor(model_inputs["attention_mask"][i][:max_length])
93         labels["input_ids"][i] = torch.tensor(labels["input_ids"][i][:max_length])
94     if i == 0:
95         print("model_inputs input_ids:", model_inputs["input_ids"][i])
96     print("model_inputs attention_mask:", model_inputs["attention_mask"][i])
97     print("labels input_ids:", labels["input_ids"][i])
98
99
100
101     model_inputs["labels"] = labels["input_ids"]
102     return model_inputs
103
104
105 print("column_names:", dataset["train"].column_names)
106
107 # 将原始的训练和测试数据同时预处理，然后作为训练和评估数据集
108 processed_datasets = dataset.map(
109     preprocess_function,
110     batched=True,
111     num_proc=1,
112     remove_columns=dataset["train"].column_names,
113     load_from_cache_file=False,
114     desc="Running tokenizer on dataset",
115 )
116
117 train_dataset = processed_datasets["train"]
118 eval_dataset = processed_datasets["train"]
119
120 # 训练与评估使用同一份数据，但是训练数据打乱
121 train_dataloader = DataLoader(train_dataset, shuffle=True, collate_fn=default_data_collator, batch_size=batch_size, pin_memory=True)
122

```

```

123 eval_dataloader = DataLoader(eval_dataset, collate_fn=default_data_collat
124 or, batch_size=batch_size, pin_memory=True)
125 print(len(train_dataloader))
126 print(len(eval_dataloader))
127
128 def test_preprocess_function(examples):
129     batch_size = len(examples[text_column])
130     inputs = [f"{text_column} : {x} Label : " for x in examples[text_colu
131 mn]]
132     model_inputs = tokenizer(inputs)
133     # print(model_inputs)
134     for i in range(batch_size):
135         sample_input_ids = model_inputs["input_ids"][i]
136
137         model_inputs["input_ids"][i] = [tokenizer.pad_token_id] * (max_l
138 ength - len(sample_input_ids)) + sample_input_ids
139         model_inputs["attention_mask"][i] = [0] * (max_length - len(sampl
140 e_input_ids)) + model_inputs["attention_mask"][i]
141
142         model_inputs["input_ids"][i] = torch.tensor(model_inputs["input_i
143 ds"][i][:max_length])
144         model_inputs["attention_mask"][i] = torch.tensor(model_inputs["at
145 tention_mask"][i][:max_length])
146     return model_inputs
147
148 # 将原始的测试数据用于测试
149 test_dataset = dataset["test"].map(
150     test_preprocess_function,
151     batched=True,
152     num_proc=1,
153     remove_columns=dataset["train"].column_names,
154     load_from_cache_file=False,
155     desc="Running tokenizer on dataset",
156 )
157
158 test_dataloader = DataLoader(test_dataset, collate_fn=default_data_collat
159 or, batch_size=batch_size, pin_memory=True)
160 next(iter(test_dataloader))
161
162 # creating model
163 model = AutoModelForCausalLM.from_pretrained(model_name_or_path)
164 model = get_peft_model(model, peft_config)
165 model.print_trainable_parameters()
166
167 # model
168 # optimizer and lr scheduler
169 optimizer = torch.optim.AdamW(model.parameters(), lr=lr)

```

```

164 lr_scheduler = get_linear_schedule_with_warmup(
165     optimizer=optimizer,
166     num_warmup_steps=0,
167     num_training_steps=(len(train_dataloader) * num_epochs),
168 )
169
170 # training and evaluation
171 model = model.to(device)
172
173 for epoch in range(num_epochs):
174     model.train()
175     total_loss = 0
176     for step, batch in enumerate(tqdm(train_dataloader)):
177         batch = {k: v.to(device) for k, v in batch.items()}
178         # print(batch)
179         # print(batch["input_ids"].shape)
180         outputs = model(**batch)
181         loss = outputs.loss
182         total_loss += loss.detach().float()
183         loss.backward()
184         optimizer.step()
185         lr_scheduler.step()
186         optimizer.zero_grad()
187
188     model.eval()
189     eval_loss = 0
190     eval_preds = []
191     for step, batch in enumerate(tqdm(eval_dataloader)):
192         batch = {k: v.to(device) for k, v in batch.items()}
193         with torch.no_grad():
194             outputs = model(**batch)
195             loss = outputs.loss
196             eval_loss += loss.detach().float()
197             eval_preds.extend(
198                 tokenizer.batch_decode(torch.argmax(outputs.logits, -1).detach()
199                                     .cpu().numpy(), skip_special_tokens=True)
200             )
201
202     eval_epoch_loss = eval_loss / len(eval_dataloader)
203     eval_ppl = torch.exp(eval_epoch_loss)
204     train_epoch_loss = total_loss / len(train_dataloader)
205     train_ppl = torch.exp(train_epoch_loss)
206     print(f"epoch={epoch}: {train_ppl=} {train_epoch_loss=} {eval_ppl=} {eval_
epoch_loss=}")

```



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