

# Prefix Tuning详解及代码实战

## 1. 背景

## 2. 技术细节

## 3. 优点

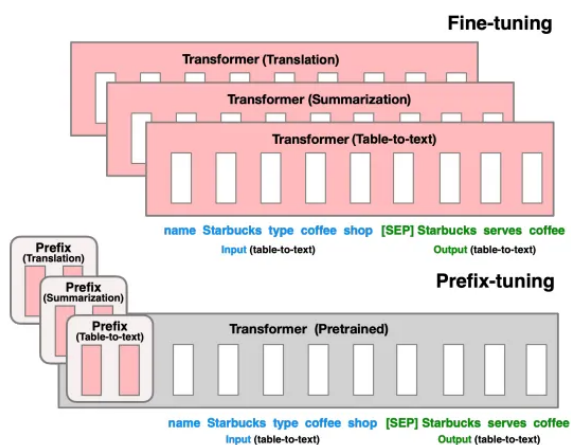
## 1. 背景

在Prefix Tuning之前的工作主要是人工设计离散的模版或者自动化搜索离散的模版存在缺陷：

- **人工离散模版缺点：**模版的变化对模型最终的性能特别敏感，加一个少一个词或变动位置都会造成比较大的变化。
- **自动化搜索模版缺点：**成本也比较高；同时，以前这种离散化的token搜索出来的结果可能并不是最优的。除此之外，传统的微调范式利用预训练模型去对不同的下游任务进行微调，对每个任务都要保存一份微调后的模型权重，一方面微调整个模型耗时长；另一方面也会占很多存储空间。

基于上述两点，Prefix Tuning优化：

- 提出固定预训练LM，为LM添加可训练，特定任务的前缀，可以为不同任务保存不同的前缀，微调成本也小；
- 同时，这种Prefix实际就是连续可微的Virtual Token（Soft Prompt/Continuous Prompt），相比离散的Token，更好优化，效果更好。

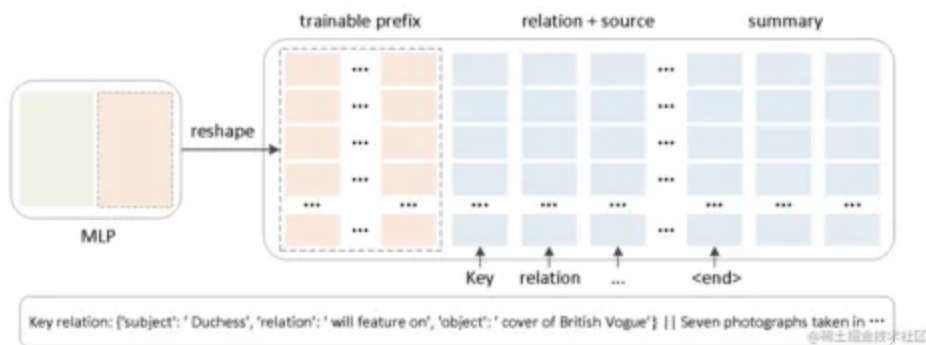


## 2. 技术细节

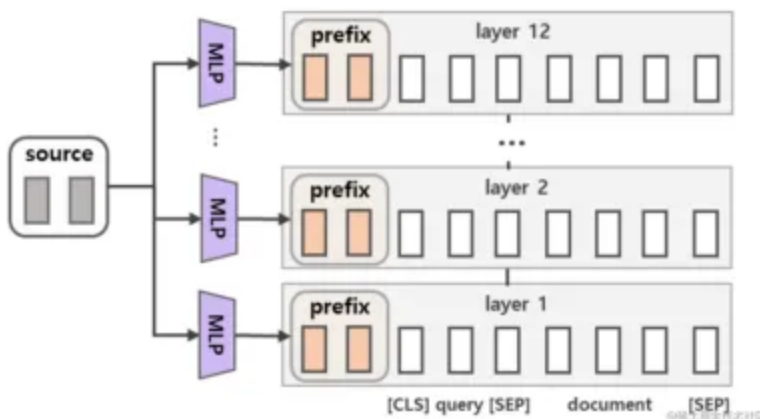
2021年斯坦福的研究人员在论文《[Prefix-Tuning: Optimizing Continuous Prompts for Generation](#)》中提出了 Prefix Tuning 方法。

细节：

- 与Full-finetuning 更新所有参数的方式不同，该方法是在输入 token 之前构造一段任务相关的 **virtual tokens**（虚拟tokens不是真实的tokens，而是可学习的自由参数）作为 Prefix，然后训练的时候只更新 Prefix 部分的参数，而 Transformer 中的其他部分参数固定。
- 为了防止直接更新Prefix的参数导致训练不稳定和性能下降的情况，在Prefix层前面加了MLP结构，训练完成后，只保留Prefix的参数。



- 通过消融实验证实，只调整embedding层的表现力不够，将导致性能显著下降，因此，在每层（所有layer的输入层）都加了prefix token，改动较大。



针对不同的模型结构，需要构造不同的Prefix。

- 针对encoder-decoder架构，会在encoder的输入和decoder的输入embedding都加上prefix token；
  - 针对自回归架构模型：在句子前面添加前缀，得到  $z = [\text{PREFIX}; x; y]$ ，合适的上文能够在固定 LM 的情况下去引导生成下文（比如：GPT3的上下文学习）。
  - 针对编码器-解码器架构模型：Encoder和Decoder都增加了前缀，得到  $z = [\text{PREFIX}; x;$

PREFIX0; y]。Encoder端增加前缀是为了引导输入部分的编码，Decoder 端增加前缀是为了引

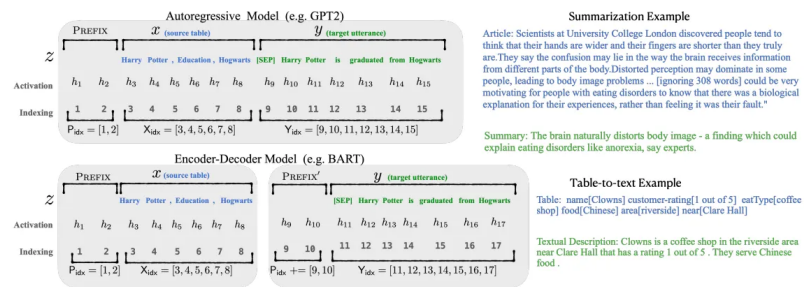


Figure 2: An annotated example of prefix-tuning using an autoregressive LM (top) and an encoder-decoder model (bottom). The prefix activations  $\forall i \in P_{idx}, h_i$  are drawn from a trainable matrix  $P_\theta$ . The remaining activations are computed by the Transformer.

导后续token的生成。

image.png

### 3. 优点

Prefix-Tuning只需训练和存储0.1%的新增参数（VS adapter 3.6%， fine-tuning 100%）。

#### 代码实践

[https://github.com/huggingface/peft/blob/main/src/peft/tuners/prefix\\_tuning/model.py](https://github.com/huggingface/peft/blob/main/src/peft/tuners/prefix_tuning/model.py)

```
1 from transformers import AutoModelForCausalLM
2 from peft import get_peft_config, get_peft_model, PrefixTuningConfig, 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 = PrefixTuningConfig(task_type=TaskType.CAUSAL_LM,
19                                 num_virtual_tokens=30,
20                                 prefix_projection=True)
21
22 dataset_name = "twitter_complaints"
23 checkpoint_name = f"{dataset_name}_{model_name_or_path}_{peft_config.peft_type}_{peft_config.task_type}_v1.pt".replace("/", "_")
24 text_column = "Tweet text"
25 label_column = "text_label"
26 max_length = 64
27 lr = 3e-2
28 num_epochs = 10
29 batch_size = 8
30
31
32 from datasets import load_dataset
33
34 # dataset = load_dataset("ought/raft", dataset_name)
35 dataset = load_dataset("/home/guodong.li/data/peft/raft/raft.py", dataset_name, cache_dir="/home/guodong.li/data/peft/data")
36
37 classes = [k.replace("_", " ") for k in dataset["train"].features["Label"].names]
38 print(classes)
39 dataset = dataset.map(
40     lambda x: {"text_label": [classes[label] for label in x["Label"]]},
```

```

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

```

```

81         model_inputs["attention_mask"][i] = torch.tensor(model_inputs["at
82         tention_mask"][i][:max_length])
83         labels["input_ids"][i] = torch.tensor(labels["input_ids"][i][:max
84         _length])
85         model_inputs["labels"] = labels["input_ids"]
86         return model_inputs
87
88     processed_datasets = dataset.map(
89         preprocess_function,
90         batched=True,
91         num_proc=1,
92         remove_columns=dataset["train"].column_names,
93         load_from_cache_file=False,
94         desc="Running tokenizer on dataset",
95     )
96
97     train_dataset = processed_datasets["train"]
98     eval_dataset = processed_datasets["train"]
99
100     train_dataloader = DataLoader(train_dataset, shuffle=True, collate_fn=def
101     ault_data_collator, batch_size=batch_size, pin_memory=True)
102     eval_dataloader = DataLoader(eval_dataset, collate_fn=default_data_collat
103     or, batch_size=batch_size, pin_memory=True)
104
105     def test_preprocess_function(examples):
106         batch_size = len(examples[text_column])
107         inputs = [f"{text_column} : {x} Label : " for x in examples[text_colu
108         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
114             ngth - len(sample_input_ids)) + sample_input_ids
115             model_inputs["attention_mask"][i] = [0] * (max_length - len(sampl
116             e_input_ids)) + model_inputs["attention_mask"][i]
117
118             model_inputs["input_ids"][i] = torch.tensor(model_inputs["input_i
119             ds"][i][:max_length])
120             model_inputs["attention_mask"][i] = torch.tensor(model_inputs["at
121             tention_mask"][i][:max_length])
122             return model_inputs
123
124     test_dataset = dataset["test"].map(
125         test_preprocess_function,

```

```

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

```

```

167     for step, batch in enumerate(tqdm(eval_dataloader)):
168         batch = {k: v.to(device) for k, v in batch.items()}
169         with torch.no_grad():
170             outputs = model(**batch)
171             loss = outputs.loss
172             eval_loss += loss.detach().float()
173             eval_preds.extend(
174                 tokenizer.batch_decode(torch.argmax(outputs.logits, -1).detach()
175                                     .cpu().numpy(), skip_special_tokens=True)
176             )
177
178         eval_epoch_loss = eval_loss / len(eval_dataloader)
179         eval_ppl = torch.exp(eval_epoch_loss)
180         train_epoch_loss = total_loss / len(train_dataloader)
181         train_ppl = torch.exp(train_epoch_loss)
182         print(f"epoch=: {train_ppl} {train_epoch_loss=} {eval_ppl=} {eval_
epoch_loss=}")

```

## ▼ 推理

Python |

```

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

```



