

Prompt Tuning 微调实战

为了不影响阅读体验，详细的代码放置在GitHub: [llm-action](#) 项目中 [peft_prompt_tuning_clm.ipynb](#) 文件，这里仅列出关键步骤。

第一步，引入必要的库，如：Prompt Tuning 配置类 PromptTuningConfig。

```
1 from peft import get_peft_config, get_peft_model, PromptTuningInit, PromptTuningConfig, TaskType, PeftType
```

第二步，创建 Prompt Tuning 微调方法对应的配置。

```
1 peft_config = PromptTuningConfig(  
2     task_type=TaskType.CAUSAL_LM,  
3     prompt_tuning_init=PromptTuningInit.TEXT,  
4     num_virtual_tokens=8,  
5     prompt_tuning_init_text="Classify if the tweet is a complaint or not:",  
6     tokenizer_name_or_path=model_name_or_path,  
7 )
```

参数说明：

- `prompt_tuning_init`：提示嵌入的初始化方法。PEFT支持文本（TEXT）和随机（RANDOM）初始化。在原理篇中提到过 Prompt token 的初始化方法和长度对于模型性能有影响。与随机初始化和使用样本词汇表初始化相比，Prompt Tuning 采用类标签初始化模型的效果更好。不过随着模型参数规模的提升，这种gap最终会消失。因此，如果需要使用类标签和样本词汇表初始化需指定为 TEXT。
- `prompt_tuning_init_text`：用于初始化提示嵌入的文本，在使用文本（TEXT）初始化方法时使用。
- `task_type`：指定任务类型。如：条件生成任务（SEQ_2_SEQ_LM），因果语言建模（CAUSAL_LM）等。
- `num_virtual_tokens`：指定虚拟Token数。在原理篇中，提到过提示虚拟 Token 的长度在20左右时的表现已经不错（超过20之后，提升Prompt token长度，对模型的性能提升不明显了）；同样的，这个gap也会随着模型参数规模的提升而减小（即对于超大规模模型而言，即使提示虚拟 Token 长度很短，对性能也不会有太大的影响）。

第三步，通过调用 `get_peft_model` 方法包装基础的 Transformer 模型。

```

1 model = AutoModelForCausalLM.from_pretrained(model_name_or_path)
2 model = get_peft_model(model, peft_config)
3 model.print_trainable_parameters()

```

通过 `print_trainable_parameters` 方法可以查看可训练参数的数量(仅为8,192)以及占比 (仅为0.00146%) 。

```

1 trainable params: 8,192 || all params: 559,222,784 || trainable%: 0.0014648
  902430985358

```

Prompt Tuning 模型类结构如下所示:

```

1 PeftModelForCausalLM(
2   (base_model): BloomForCausalLM(
3     (transformer): BloomModel(
4       (word_embeddings): Embedding(250880, 1024)
5       (word_embeddings_layernorm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
6       (h): ModuleList(
7         ...
8       )
9       (ln_f): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
10    )
11    (lm_head): Linear(in_features=1024, out_features=250880, bias=False)
12  )
13  (prompt_encoder): ModuleDict(
14    (default): PromptEmbedding(
15      (embedding): Embedding(8, 1024)
16    )
17  )
18  (word_embeddings): Embedding(250880, 1024)
19 )

```

从模型类结构可以看到, Prompt Tuning 只在输入层加入 prompt virtual tokens, 其他地方均没有变化, 具体可查看 `PromptEmbedding` 的源码。

```

1  class PromptEmbedding(torch.nn.Module):
2      def __init__(self, config, word_embeddings):
3          super().__init__()
4
5          total_virtual_tokens = config.num_virtual_tokens * config.num_trans
former_submodules
6          # 初始化 embedding 层
7          self.embedding = torch.nn.Embedding(total_virtual_tokens, config.t
oken_dim)
8
9          # 如果使用文本进行初始化, 执行如下逻辑, PromptTuningConfig 配置类需要传入初
始化文本。
10         if config.prompt_tuning_init == PromptTuningInit.TEXT:
11             from transformers import AutoTokenizer
12
13             tokenizer = AutoTokenizer.from_pretrained(config.tokenizer_nam
e_or_path)
14             init_text = config.prompt_tuning_init_text
15             init_token_ids = tokenizer(init_text)["input_ids"]
16             # Trim or iterate until num_text_tokens matches total_virtual_
tokens
17             num_text_tokens = len(init_token_ids)
18             if num_text_tokens > total_virtual_tokens:
19                 init_token_ids = init_token_ids[:total_virtual_tokens]
20             elif num_text_tokens < total_virtual_tokens:
21                 num_reps = math.ceil(total_virtual_tokens / num_text_token
s)
22                 init_token_ids = init_token_ids * num_reps
23                 init_token_ids = init_token_ids[:total_virtual_tokens]
24
25             word_embedding_weights = word_embeddings(torch.LongTensor(init
_token_ids)).detach().clone()
26             word_embedding_weights = word_embedding_weights.to(torch.float
32)
27             # 初始化embedding层的权重
28             self.embedding.weight = torch.nn.Parameter(word_embedding_weig
hts)
29
30         def forward(self, indices):
31             # Just get embeddings
32             prompt_embeddings = self.embedding(indices)
33             return prompt_embeddings

```

第四步，模型训练的其余部分均无需更改，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。

```
▼ Plain Text |
1 peft_model_id = f"{model_name_or_path}_{peft_config.peft_type}_{peft_config.task_type}"
2 model.save_pretrained(peft_model_id)
```

输出的模型权重文件如下所示：

```
▼ Plain Text |
1 /data/nfs/llm/model/bloomz-560m_PROMPT_TUNING_CAUSAL_LM
2 |— [ 500] adapter_config.json
3 |— [ 33K] adapter_model.bin
4 |— [ 111] README.md
5
6 0 directories, 3 files
```

注意：这里只会保存经过训练的增量 PEFT 权重。其中，`adapter_config.json` 为 Prompt Tuning 配置文件；`adapter_model.bin` 为 Prompt Tuning 权重文件。

第五步，加载微调后的权重文件进行推理。

```
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))
```

至此，我们完成了Prompt Tuning的训练及推理。

```

1  # coding=utf-8
2  # Copyright 2023-present the HuggingFace Inc. team.
3  #
4  # Licensed under the Apache License, Version 2.0 (the "License");
5  # you may not use this file except in compliance with the License.
6  # You may obtain a copy of the License at
7  #
8  #     http://www.apache.org/licenses/LICENSE-2.0
9  #
10 # Unless required by applicable law or agreed to in writing, software
11 # distributed under the License is distributed on an "AS IS" BASIS,
12 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
13 # See the License for the specific language governing permissions and
14 # limitations under the License.
15
16 import math
17
18 import torch
19
20 from .config import PromptTuningInit
21
22
23 class PromptEmbedding(torch.nn.Module):
24     """
25     The model to encode virtual tokens into prompt embeddings.
26
27     Args:
28         config ([`PromptTuningConfig`]): The configuration of the prompt embedding.
29         word_embeddings (`torch.nn.Module`): The word embeddings of the base transformer model.
30
31     **Attributes**:
32         - **embedding** (`torch.nn.Embedding`) -- The embedding layer of the prompt embedding.
33
34     Example:
35
36     ```py
37     >>> from peft import PromptEmbedding, PromptTuningConfig
38
39     >>> config = PromptTuningConfig(
40         ...     peft_type="PROMPT_TUNING",
41         ...     task_type="SEQ_2_SEQ_LM",
42         ...     num_virtual_tokens=20,

```

```

43         ...     token_dim=768,
44         ...     num_transformer_submodules=1,
45         ...     num_attention_heads=12,
46         ...     num_layers=12,
47         ...     prompt_tuning_init="TEXT",
48         ...     prompt_tuning_init_text="Predict if sentiment of this review i
49 s positive, negative or neutral",
50         ...     tokenizer_name_or_path="t5-base",
51         ... )
52
53 >>> # t5_model.shared is the word embeddings of the base model
54 >>> prompt_embedding = PromptEmbedding(config, t5_model.shared)
55
56 Input Shape: (`batch_size`, `total_virtual_tokens`)
57
58 Output Shape: (`batch_size`, `total_virtual_tokens`, `token_dim`)
59 """
60
61 def __init__(self, config, word_embeddings):
62     super().__init__()
63
64     total_virtual_tokens = config.num_virtual_tokens * config.num_tran
65 sformer_submodules
66     self.embedding = torch.nn.Embedding(total_virtual_tokens, config.t
67 oken_dim)
68     if config.prompt_tuning_init == PromptTuningInit.TEXT and not conf
69 ig.inference_mode:
70         from transformers import AutoTokenizer
71
72         tokenizer_kwargs = config.tokenizer_kwargs or {}
73         tokenizer = AutoTokenizer.from_pretrained(config.tokenizer_nam
74 e_or_path, **tokenizer_kwargs)
75         init_text = config.prompt_tuning_init_text
76         init_token_ids = tokenizer(init_text)["input_ids"]
77         # Trim or iterate until num_text_tokens matches total_virtual_
78 tokens
79         num_text_tokens = len(init_token_ids)
80         if num_text_tokens > total_virtual_tokens:
81             init_token_ids = init_token_ids[:total_virtual_tokens]
82         elif num_text_tokens < total_virtual_tokens:
83             num_reps = math.ceil(total_virtual_tokens / num_text_token
84 s)
85             init_token_ids = init_token_ids * num_reps
86             init_token_ids = init_token_ids[:total_virtual_tokens]
87             init_token_ids = torch.LongTensor(init_token_ids).to(word_embe
88 ddings.weight.device)

```

```

83         word_embedding_weights = word_embeddings(init_token_ids).detach_()
84         word_embedding_weights = word_embedding_weights.to(torch.float32)
85         self.embedding.weight = torch.nn.Parameter(word_embedding_weights)
86
87     def forward(self, indices):
88         # Just get embeddings
89         prompt_embeddings = self.embedding(indices)
90         return prompt_embeddings

```