

```
In [75]: # ! pip install transformers datasets peft
```

## HW 8: Low Rank Adaptation (LoRA)

In this assignment, you will learn to implement low-rank adaptation both from scratch and using a library—specifically, with PyTorch and the PEFT library, respectively.

This assignment is divided into two sections:

In the first section, we introduce the parameter-efficient transfer learning (PET) method. We use LoRA to adapt the GPT2 model for the SST-2 dataset. This section will teach you how LoRA works and how to implement it from scratch using `forward_hook`.

In the second section, we introduce the PEFT library, which allows us to perform LoRA easily.

### Part 1: LoRA from Scratch

With the discovery of scaling properties in deep learning models, several researchers tend to increase model size to achieve emergent properties, especially in the natural language processing (NLP) field. For example, GPT-3 contains 175 billion parameters, making it nearly impossible to fine-tune on limited resources. This trend prevents students like us from adapting these enormous foundation models on a single GPU (or with small resources).

To alleviate this problem, researchers have developed new fine-tuning methods, known as parameter-efficient transfer learning, which allow us to train large models with limited resources. The benefits of these methods extend not only to the training process but also to deployment. After fine-tuning, we only need to save a small number of parameters (the LoRA weights), enabling us to deploy the foundation model to various downstream tasks using minimal storage. One of the prevailing methods is Low Rank Adaptation (LoRA).

Another popular option is prompt tuning, where we only train special tokens that are prepended to the input. However, this is not the focus of this homework.

In this section, we will introduce Low-Rank Adaptation. You are assigned to implement LoRA on GPT2 model. We will finetune the model to SST-2 dataset using the traditional and LoRA method.

### Load Dataset and Model

In this step, we will prepare the GPT-2 model and the SST-2 dataset.

SST-2 is a widely used dataset for sentiment analysis, extracted from movie reviews, containing sentences labeled as either positive or negative.

```
In [76]: import torch
import numpy as np
import time
from torch.utils.data import DataLoader
from transformers import GPT2ForSequenceClassification, GPT2TokenizerFast
from datasets import load_dataset
from tqdm.autonotebook import tqdm
from transformers import GPT2Tokenizer, GPT2ForSequenceClassification, AdamW

# Load the GPT-2 model for sequence classification and its tokenizer
model = GPT2ForSequenceClassification.from_pretrained("gpt2", num_labels=2)
tokenizer = GPT2TokenizerFast.from_pretrained("gpt2")

# GPT-2 does not have a pad token by default so we set it to the EOS token.
tokenizer.pad_token = tokenizer.eos_token
model.config.pad_token_id = tokenizer.eos_token_id

# Load the SST-2 dataset
train_dataset_raw = load_dataset("glue", "sst2", split="train")
train_dataset_raw, val_dataset_raw = train_dataset_raw.train_test_split(test_size=0.1)
test_dataset_raw = load_dataset("glue", "sst2", split="validation")

# Preview dataset
print("Sample sentence:")
for data in test_dataset_raw:
    print(data)
    break

def tokenize_function(example):
    return tokenizer(example["sentence"], padding="max_length", truncation=True, max_length=512)

train_dataset = train_dataset_raw.map(tokenize_function, batched=True)
val_dataset = val_dataset_raw.map(tokenize_function, batched=True)
test_dataset = test_dataset_raw.map(tokenize_function, batched=True)

train_dataset.set_format("torch", columns=["input_ids", "attention_mask", "label"])
val_dataset.set_format("torch", columns=["input_ids", "attention_mask", "label"])
test_dataset.set_format("torch", columns=["input_ids", "attention_mask", "label"])

# Create data loaders
train_dataloader = DataLoader(train_dataset, shuffle=True, batch_size=16)
val_dataloader = DataLoader(val_dataset, batch_size=16)
test_dataloader = DataLoader(test_dataset, batch_size=16)
```

Some weights of GPT2ForSequenceClassification were not initialized from the model checkpoint at gpt2 and are newly initialized: ['score.weight']  
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

Sample sentence:

```
{'sentence': "it 's a charming and often affecting journey . ", 'label': 1, 'idx': 0}
```

```
Map: 0% | 0/53879 [00:00<?, ? examples/s]
```

Map: 0%| | 0/13470 [00:00<?, ? examples/s]

## Traditional Fine tuning

In the traditional fine-tuning method, the entire model is trained, which is computationally expensive. An alternative approach is to fine-tune only certain layers of the model to reduce resource usage while still adapting the model to a specific task.

To keep the implementation simple, you are assigned to train only the attention weights in the self-attention layers.

The code below displays the names of all layers in the GPT-2 model. This will help you identify which layers to set as trainable or keep frozen. For more details on the attention layers in GPT-2, please refer to the following link: [GPT-2 Attention Layer Details](#).

```
In [77]: for name, module in model.named_modules():  
         print(name, type(module))
```

```

    <class 'transformers.models.gpt2.modeling_gpt2.GPT2ForSequenceClassification'>
transformer <class 'transformers.models.gpt2.modeling_gpt2.GPT2Model'>
transformer.wte <class 'torch.nn.modules.sparse.Embedding'>
transformer.wpe <class 'torch.nn.modules.sparse.Embedding'>
transformer.drop <class 'torch.nn.modules.dropout.Dropout'>
transformer.h <class 'torch.nn.modules.container.ModuleList'>
transformer.h.0 <class 'transformers.models.gpt2.modeling_gpt2.GPT2Block'>
transformer.h.0.ln_1 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.0.attn <class 'transformers.models.gpt2.modeling_gpt2.GPT2SdpaAttention'>
transformer.h.0.attn.c_attn <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.0.attn.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.0.attn.attn_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.0.attn.resid_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.0.ln_2 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.0.mlp <class 'transformers.models.gpt2.modeling_gpt2.GPT2MLP'>
transformer.h.0.mlp.c_fc <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.0.mlp.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.0.mlp.act <class 'transformers.activations.NewGELUActivation'>
transformer.h.0.mlp.dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.1 <class 'transformers.models.gpt2.modeling_gpt2.GPT2Block'>
transformer.h.1.ln_1 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.1.attn <class 'transformers.models.gpt2.modeling_gpt2.GPT2SdpaAttention'>
transformer.h.1.attn.c_attn <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.1.attn.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.1.attn.attn_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.1.attn.resid_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.1.ln_2 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.1.mlp <class 'transformers.models.gpt2.modeling_gpt2.GPT2MLP'>
transformer.h.1.mlp.c_fc <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.1.mlp.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.1.mlp.act <class 'transformers.activations.NewGELUActivation'>
transformer.h.1.mlp.dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.2 <class 'transformers.models.gpt2.modeling_gpt2.GPT2Block'>
transformer.h.2.ln_1 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.2.attn <class 'transformers.models.gpt2.modeling_gpt2.GPT2SdpaAttention'>
transformer.h.2.attn.c_attn <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.2.attn.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.2.attn.attn_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.2.attn.resid_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.2.ln_2 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.2.mlp <class 'transformers.models.gpt2.modeling_gpt2.GPT2MLP'>
transformer.h.2.mlp.c_fc <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.2.mlp.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.2.mlp.act <class 'transformers.activations.NewGELUActivation'>
transformer.h.2.mlp.dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.3 <class 'transformers.models.gpt2.modeling_gpt2.GPT2Block'>
transformer.h.3.ln_1 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.3.attn <class 'transformers.models.gpt2.modeling_gpt2.GPT2SdpaAttention'>
transformer.h.3.attn.c_attn <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.3.attn.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.3.attn.attn_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.3.attn.resid_dropout <class 'torch.nn.modules.dropout.Dropout'>

```

[illegible]

```
transformer.h.7.ln_2 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.7.mlp <class 'transformers.models.gpt2.modeling_gpt2.GPT2MLP'>
transformer.h.7.mlp.c_fc <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.7.mlp.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.7.mlp.act <class 'transformers.activations.NewGELUActivation'>
transformer.h.7.mlp.dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.8 <class 'transformers.models.gpt2.modeling_gpt2.GPT2Block'>
transformer.h.8.ln_1 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.8.attn <class 'transformers.models.gpt2.modeling_gpt2.GPT2SdpaAttention'>
transformer.h.8.attn.c_attn <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.8.attn.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.8.attn.attn_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.8.attn.resid_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.8.ln_2 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.8.mlp <class 'transformers.models.gpt2.modeling_gpt2.GPT2MLP'>
transformer.h.8.mlp.c_fc <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.8.mlp.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.8.mlp.act <class 'transformers.activations.NewGELUActivation'>
transformer.h.8.mlp.dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.9 <class 'transformers.models.gpt2.modeling_gpt2.GPT2Block'>
transformer.h.9.ln_1 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.9.attn <class 'transformers.models.gpt2.modeling_gpt2.GPT2SdpaAttention'>
transformer.h.9.attn.c_attn <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.9.attn.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.9.attn.attn_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.9.attn.resid_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.9.ln_2 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.9.mlp <class 'transformers.models.gpt2.modeling_gpt2.GPT2MLP'>
transformer.h.9.mlp.c_fc <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.9.mlp.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.9.mlp.act <class 'transformers.activations.NewGELUActivation'>
transformer.h.9.mlp.dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.10 <class 'transformers.models.gpt2.modeling_gpt2.GPT2Block'>
transformer.h.10.ln_1 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.10.attn <class 'transformers.models.gpt2.modeling_gpt2.GPT2SdpaAttention'>
transformer.h.10.attn.c_attn <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.10.attn.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.10.attn.attn_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.10.attn.resid_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.10.ln_2 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.10.mlp <class 'transformers.models.gpt2.modeling_gpt2.GPT2MLP'>
transformer.h.10.mlp.c_fc <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.10.mlp.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.10.mlp.act <class 'transformers.activations.NewGELUActivation'>
transformer.h.10.mlp.dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.11 <class 'transformers.models.gpt2.modeling_gpt2.GPT2Block'>
transformer.h.11.ln_1 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.11.attn <class 'transformers.models.gpt2.modeling_gpt2.GPT2SdpaAttention'>
transformer.h.11.attn.c_attn <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.11.attn.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.11.attn.attn_dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.h.11.attn.resid_dropout <class 'torch.nn.modules.dropout.Dropout'>
```

```
transformer.h.11.ln_2 <class 'torch.nn.modules.normalization.LayerNorm'>
transformer.h.11.mlp <class 'transformers.models.gpt2.modeling_gpt2.GPT2MLP'>
transformer.h.11.mlp.c_fc <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.11.mlp.c_proj <class 'transformers.pytorch_utils.Conv1D'>
transformer.h.11.mlp.act <class 'transformers.activations.NewGELUActivation'>
transformer.h.11.mlp.dropout <class 'torch.nn.modules.dropout.Dropout'>
transformer.ln_f <class 'torch.nn.modules.normalization.LayerNorm'>
score <class 'torch.nn.modules.linear.Linear'>
```

## TODO 1: Freeze the Model and Train Only Attention Weights

You are assigned to freeze the entire model, except for the last two attention weights and the classification head. Note that, in this context, the attention weights do not include the projection layer of the transformer. Instead, they refer only to the weights of the query, key, and value.

**HINT:** `c_proj` is projection layer.

```
In [78]: for n, p in model.named_parameters():
          # TODO 1: freeze every layer except the last two attention weights and classifi
          p.requires_grad = False

num_layers = model.config.n_layer
for layer_idx in range(num_layers - 2, num_layers): # Last two layers
    attn_layer = f"transformer.h.{layer_idx}.attn.c_attn"
    for n, p in model.named_parameters():
        if attn_layer in n:
            p.requires_grad = True # Unfreeze Q, K, V weights

# Unfreeze classification head
for n, p in model.named_parameters():
    if "score" in n: # Classification head
        p.requires_grad = True

# Verify which parameters are trainable
trainable_params = [n for n, p in model.named_parameters() if p.requires_grad]
print("Trainable parameters:", trainable_params)
```

```
Trainable parameters: ['transformer.h.10.attn.c_attn.weight', 'transformer.h.10.att
n.c_attn.bias', 'transformer.h.11.attn.c_attn.weight', 'transformer.h.11.attn.c_att
n.bias', 'score.weight']
```

**Check Your Answer:** The number of learnable parameters is around 3545088.

```
In [79]: pytorch_total_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
print("Number of Trainable Parameters:", pytorch_total_params)
```

```
Number of Trainable Parameters: 3545088
```

You are assigned to train the GPT-2 model on the SST-2 dataset. Due to the long training time, you will train the model for only 3 epochs. Your model should have around 86-88% accuracy.

```

In [80]: optimizer = AdamW(model.parameters(), lr=5e-5)

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print("Device:", device)
model.to(device)

num_epochs = 3
for epoch in tqdm(range(num_epochs)):
    model.train()
    for batch in tqdm(train_dataloader):
        input_ids = batch["input_ids"].to(device)
        attention_mask = batch["attention_mask"].to(device)
        labels = batch["label"].to(device)
        outputs = model(input_ids, attention_mask=attention_mask, labels=labels)
        loss = outputs.loss

        loss.backward()
        optimizer.step()
        optimizer.zero_grad()

    model.eval()
    correct = 0
    total = 0
    with torch.no_grad():
        for batch in test_dataloader:
            input_ids = batch["input_ids"].to(device)
            attention_mask = batch["attention_mask"].to(device)
            labels = batch["label"].to(device)
            outputs = model(input_ids, attention_mask=attention_mask)
            predictions = torch.argmax(outputs.logits, dim=-1)
            correct += (predictions == labels).sum().item()
            total += labels.size(0)

    accuracy = correct / total
    print(f"Epoch {epoch + 1}/{num_epochs} - Accuracy: {accuracy:.4f}")

```

```

Device: cuda
 0%|          | 0/3 [00:00<?, ?it/s]
 0%|          | 0/3368 [00:00<?, ?it/s]
Epoch 1/3 - Accuracy: 0.8567
 0%|          | 0/3368 [00:00<?, ?it/s]
Epoch 2/3 - Accuracy: 0.8853
 0%|          | 0/3368 [00:00<?, ?it/s]
Epoch 3/3 - Accuracy: 0.8876

```

As you can see, fine-tuning in the traditional way takes a long time to complete and also requires a high-computation GPU to fine-tune the entire model. Therefore, it is not feasible for most people.

In the next part, we will introduce a better method: parameter-efficient learning, which requires lower computation. We will focus on the state-of-the-art method, Low-Rank Adaptation (LoRA).



# Low Rank Adaptation

The concept of LoRA is that we are going to estimate the gradient (adaptation matrix) with two smaller matrices ( $A$  and  $B$ ):

$$\text{Adaptation Matrix} = B \times A$$

where  $\text{Adaptation Matrix} \in \mathbb{R}^{m \times n}$ ,  $A \in \mathbb{R}^{r \times n}$ , and  $B \in \mathbb{R}^{m \times r}$ . We could make this approximation based on the assumption that  $\text{Adaptation Matrix}$  has a rank of  $r$ . Therefore, the fine-tuned weight becomes:

$$\begin{aligned} W &= W_0 + \Delta W \\ &= W_0 + \frac{\alpha}{r} BA \end{aligned}$$

where  $W$  denotes the fine-tuned weight,  $W_0$  represents pre-trained weight,  $\Delta W$  is the gradient and  $\alpha$  can be seen as a learning rate.  $A$  is initialized using a common initialization, like Kaiming initialization, during the initialization process. On the other hand,  $B$  is set to 0 such that the model's output remains the same after injecting LoRA, resulting in a stabilized training process.

To summarize, when injecting LoRA into a layer, we insert new parameters called matrix  $A$  and  $B$  and initialize them using the above description. Then, we modify the forward pass with `forward_hook` such that the output becomes:

$$h = Wx + \frac{\alpha}{r} BAx$$

where  $x$  and  $h$  are the input and output, respectively. We recommend you read this [blog](#) to learn more about `forward_hook`.

## LoRA on Linear Layer

- TODO 2: initialize  $A$  and  $B$  to ones (every entry in the matrix is one), such that we can verify your forward pass after attaching the hook.
- TODO 3: implement the forward hook such that new output  $h$  is

$$h = Wx + \frac{\alpha}{r} BAx$$

**Hint:** When you want to declare and initialize a parameter, you can use `torch.nn.Parameter` and `torch.nn.init`, respectively.

In [163...

```
import math
import torch.nn as nn
import torch.nn.functional as F
# Initialize LoRA and attach a hook.
def attach_lora(layer, r, lora_alpha, in_features, out_features):
    assert r > 0, "rank must greater than 0."
```

```

# TODO 2: Declare A and B matrices and initialize A and B to ones.
layer.lora_A = nn.Parameter(torch.ones(r, in_features))
layer.lora_B = nn.Parameter(torch.ones(out_features, r))

def hook(model, input, output):
    assert len(input) == 1, "The length of the input must be 1."
    # TODO 3: Compute adaptation matrix (BA) and modify the forward pass.
    x = input[0]
    delta_W = (lora_alpha / r) * torch.matmul(layer.lora_B, layer.lora_A)

    output += torch.matmul(x, delta_W.T)

return hook

```

To test your `forward_hook`, we will check the difference of the output before and after injecting the LoRA when you initialize matrices A and B with ones.

In [164...

```

from transformers.modeling_utils import Conv1D
# The Conv1D Layer from the Transformer Library is actually a linear layer. (https:

class DummyLinear(nn.Module):
    def __init__(self):
        super().__init__()
        self.linear = Conv1D(20, 10)

    def forward(self, x):
        return self.linear(x)

r, lora_alpha = 1, 4
input_ = torch.arange(10, dtype=torch.float32).unsqueeze(0)
dummy_linear = DummyLinear()
output_before = dummy_linear(input_)
for name, module in dummy_linear.named_modules():
    if isinstance(module, Conv1D):
        in_features, out_features = module.weight.shape
        h = module.register_forward_hook(attach_lora(module, r, lora_alpha, in_features)
output_after = dummy_linear(input_)

if torch.all(torch.isclose(output_after - output_before, lora_alpha * input_.sum())
print("Your forward hook seems to be correct.")
else:
print("There is something wrong with your forward hook.")

```

Your forward hook seems to be correct.

### Instruction

TODO 4: Change the initialization of A and B where A is initialized with Kaiming Uniform ( $a = \sqrt{5}$ ), and B is set to 0.

In [172...

```

# Initialize LoRA and attach a hook.
def attach_lora(layer, r, lora_alpha, in_features, out_features):
    assert r > 0, "rank must greater than 0."
    # TODO 4: initialize A with kaiming uniform with a = sqrt(5) and initialize B t
    layer.lora_A = nn.Parameter(torch.empty(r, in_features))

```

```

layer.lora_B = nn.Parameter(torch.zeros(out_features, r))

nn.init.kaiming_uniform_(layer.lora_A, a=math.sqrt(5))

def hook(model, input, output):
    assert len(input) == 1, "The length of the input must be 1."
    # Copy from TODO 3
    x = input[0]
    delta_W = (lora_alpha / r) * torch.matmul(layer.lora_B, layer.lora_A)

    output += torch.matmul(x, delta_W)

return hook

```

Similar to TODO 1, You are assigned to inject lora into the last two attention weights.

TODO 5: inject lora into the last two attention weights

```

In [173... r, lora_alpha = 1, 4
def attach_lora_to_maskformer(model, r, lora_alpha):
    hooks = []
    attention_layers = [name for name, _ in model.named_modules() if "attn.c_attn" in
    for name, module in model.named_modules():
        # TODO 5: inject lora into the last two attention weights
        if name in attention_layers:
            in_features = module.weight.shape[1]
            out_features = module.weight.shape[0]

            hook = attach_lora(module, r, lora_alpha, in_features, out_features)
            hooks.append(module.register_forward_hook(hook))
    return hooks

model = GPT2ForSequenceClassification.from_pretrained("gpt2", num_labels=2)
model.config.pad_token_id = tokenizer.eos_token_id
hooks = attach_lora_to_maskformer(model, r, lora_alpha)

```

Some weights of GPT2ForSequenceClassification were not initialized from the model checkpoint at gpt2 and are newly initialized: ['score.weight']  
 You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

## TODO 6: Freeze the Model and Train Only LoRA Weights

You are assigned to freeze the entire model, except for the bias of the last two attention weights, LoRA weights, and the classification head.

```

In [174... for n, p in model.named_parameters():
    # TODO 6: freeze every layer except the bias of the last two attention weights,
    p.requires_grad = False

    if 'lora_A' in n or 'lora_B' in n:
        p.requires_grad = True

    if n.startswith("score"):

```

```

        p.requires_grad = True

for layer in model.transformer.h[-2:]:
    # attn.c_attn.weight.requires_grad = True # c_attn has q, k, v weights
    layer.attn.c_attn.bias.requires_grad = True

```

**Check Your Answer:** The number of learnable parameters is around 12288.

```

In [175... pytorch_total_params = sum(p.numel() for p in model.parameters() if p.requires_grad
print(pytorch_total_params)

```

12288

```

In [176... optimizer = AdamW(model.parameters(), lr=5e-5)

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print("Device:", device)
model.to(device)

num_epochs = 3
for epoch in tqdm(range(num_epochs)):
    model.train()
    for batch in tqdm(train_dataloader):
        input_ids = batch["input_ids"].to(device)
        attention_mask = batch["attention_mask"].to(device)
        labels = batch["label"].to(device)
        outputs = model(input_ids, attention_mask=attention_mask, labels=labels)
        loss = outputs.loss
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()

    model.eval()
    correct = 0
    total = 0
    with torch.no_grad():
        for batch in test_dataloader:
            input_ids = batch["input_ids"].to(device)
            attention_mask = batch["attention_mask"].to(device)
            labels = batch["label"].to(device)
            outputs = model(input_ids, attention_mask=attention_mask)
            predictions = torch.argmax(outputs.logits, dim=-1)
            correct += (predictions == labels).sum().item()
            total += labels.size(0)

    accuracy = correct / total
    print(f"Epoch {epoch + 1}/{num_epochs} - Accuracy: {accuracy:.4f}")

```

Device: cuda

```

0%|          | 0/3 [00:00<?, ?it/s]
0%|          | 0/3368 [00:00<?, ?it/s]
Epoch 1/3 - Accuracy: 0.8119
0%|          | 0/3368 [00:00<?, ?it/s]
Epoch 2/3 - Accuracy: 0.8268
0%|          | 0/3368 [00:00<?, ?it/s]
Epoch 3/3 - Accuracy: 0.8452

```

## Part 2: PEFT Library

In the first part, you learned how to implement LoRA from scratch. However, in real-world applications, we can simplify this process by using pre-built libraries. One such library is `peft`, which allows us to inject LoRA into a model more efficiently. By declaring the injected modules in the `LoRAConfig`, we can easily integrate LoRA without having to implement it ourselves. In this section, you will use the `peft` library to apply LoRA to the model.

### TODO 7-8: Initialize the LoRA Config and Set trainable parameters

Your task is to initialize `LoRAConfig` using the same hyperparameters as in TODO 6 ( `r=1`, `lora_alpha=4` ). Apply LoRA only to the last two attention layers. Then, make sure to freeze the entire model except for the LoRA weights, the bias in the LoRA-injected layers, and the classification head. You only need to set the classification head to be trainable where the rest parameters are already set according to our `LoRAConfig`.

**HINT:** The total number of trainable parameters should match the result from Part 1 (TODO 6).

In [177...

```
from peft import LoraConfig, get_peft_model

num_layers = model.config.n_layer
last_two_layers = ["transformer.h." + str(num_layers-2) + ".attn.c_attn", "transformer.h." + str(num_layers-1) + ".attn.c_attn"]

# TODO 7: Initialize LoRAConfig
lora_config = LoraConfig(
    # Insert the parameters
    r=1,                # rank of the LoRA layers
    lora_alpha=4,        # scaling factor for LoRA weights
    target_modules=last_two_layers, # Apply LoRA to the last two attention layers
    modules_to_save = ["score"],
    bias="lora_only",
)

model = GPT2ForSequenceClassification.from_pretrained("gpt2", num_labels=2)
model.config.pad_token_id = tokenizer.eos_token_id

model = get_peft_model(model, lora_config)
model = model.to(device)

# TODO 8: Set classification head to trainable
for n, p in model.named_parameters():
    if n.startswith("score"): # Classification head
        p.requires_grad = True
```

Some weights of GPT2ForSequenceClassification were not initialized from the model checkpoint at gpt2 and are newly initialized: ['score.weight']  
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
In [178... pytorch_total_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
print(pytorch_total_params)
```

12288

```
In [179... optimizer = AdamW(model.parameters(), lr=5e-5)

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print("Device:", device)
model.to(device)

num_epochs = 3
for epoch in tqdm(range(num_epochs)):
    model.train()
    for batch in tqdm(train_dataloader):
        input_ids = batch["input_ids"].to(device)
        attention_mask = batch["attention_mask"].to(device)
        labels = batch["label"].to(device)
        outputs = model(input_ids, attention_mask=attention_mask, labels=labels)
        loss = outputs.loss
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()

    model.eval()
    correct = 0
    total = 0
    with torch.no_grad():
        for batch in test_dataloader:
            input_ids = batch["input_ids"].to(device)
            attention_mask = batch["attention_mask"].to(device)
            labels = batch["label"].to(device)
            outputs = model(input_ids, attention_mask=attention_mask)
            predictions = torch.argmax(outputs.logits, dim=-1)
            correct += (predictions == labels).sum().item()
            total += labels.size(0)

    accuracy = correct / total
    print(f"Epoch {epoch + 1}/{num_epochs} - Accuracy: {accuracy:.4f}")
```

Device: cuda

```
0%|          | 0/3 [00:00<?, ?it/s]
0%|          | 0/3368 [00:00<?, ?it/s]
Epoch 1/3 - Accuracy: 0.7775
0%|          | 0/3368 [00:00<?, ?it/s]
Epoch 2/3 - Accuracy: 0.8119
0%|          | 0/3368 [00:00<?, ?it/s]
Epoch 3/3 - Accuracy: 0.8372
```

In [ ]: