

深度学习作业——基于 ViT 的 CIFAR10 图像分类

Clark

1 概述

本实验使用 PyTorch 框架实现了 Vision Transformer (ViT) 模型，并在 CIFAR-10 数据集上完成图像分类任务。通过数据增强、多头注意力机制和余弦退火学习率调度等策略，模型在测试集上达到 84.84% 的准确率，远超 80% 的实验要求。训练过程中引入早停机制，在 34 个 epoch 时精度已超过 80%，最佳训练精度达 93.21% (第 86 轮)。实验完整复现了 ViT 的核心思想，验证了 Transformer 在计算机视觉任务中的有效性。

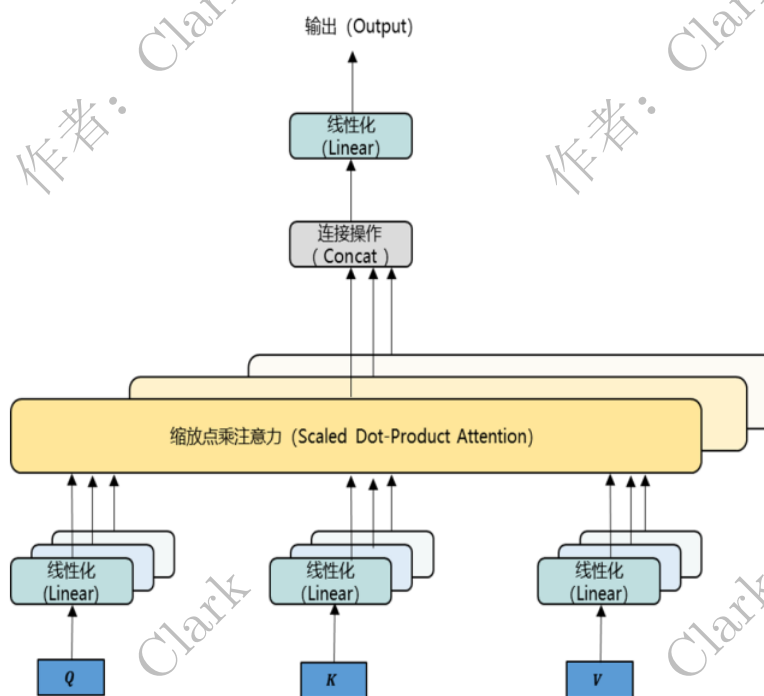
2 实验目的

1. 掌握 ViT 模型的结构原理及其在图像分类任务中的应用；
2. 学习使用 PyTorch 实现深度学习模型的完整流程（数据加载、模型构建、训练优化）；
3. 理解多头注意力机制和位置编码在视觉任务中的作用；
4. 实现 CIFAR-10 数据集的分类准确率超过 80% 的性能要求。

3 实验原理

3.1 Vision Transformer 模型概述

ViT 首次将纯 Transformer 架构应用于计算机视觉领域，其核心思想是将图像分割为序列化的图像块 (Patch)，通过线性映射后输入 Transformer 编码器。模型摒弃了传统 CNN 的卷积操作，完全依赖自注意力机制捕捉全局特征。



3.2 关键组件详解

3.2.1 图像分块嵌入 (Patch Embedding)

将输入图像 (32×32) 分割为 4×4 的块 (共 64 块), 每个 Patch 展平为 48 维向量, 通过线性层映射到 256 维:

```
1 self.to_patch_embedding = nn.Sequential(
2     Rearrange('b c (h p1) (w p2) -> b (h w) (p1 p2 c)', p1=patch_height, p2=patch_width),
3     nn.LayerNorm(patch_dim),
4     nn.Linear(patch_dim, dim),
5     nn.LayerNorm(dim),
6 )
```

3.2.2 位置编码与 CLS 令牌

加入可学习的位置编码保持空间信息, 并添加 CLS 令牌用于分类:

```
1 self.pos_embedding = nn.Parameter(torch.randn(1, num_patches + 1, dim)) # 64 patches +
   CLS
2 self.cls_token = nn.Parameter(torch.randn(1, 1, dim))
```

其作用是使用 [CLS] 令牌或全局平均池化聚合特征, 通过线性层输出 10 分类结果。

3.2.3 Transformer 编码器

Transformer 编码器由多头注意力 (Multi-Head Attention) 和前馈网络 (FeedForward) 组成, 支持残差连接和 LayerNorm。

3.2.3.1 多头自注意力机制 注意力权重计算公式：

$$\text{Att} = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V, \quad (1)$$

其中 Q 、 K 、 V 为查询、键、值矩阵， d 为维度。代码实现中通过 8 个注意力头并行计算：

```

1 class Attention(nn.Module):
2     def __init__(self, dim, heads=8, dim_head=64, dropout=0.):
3         super().__init__()
4         inner_dim = dim_head * heads
5         project_out = not (heads == 1 and dim_head == dim)
6
7         self.heads = heads
8         self.scale = dim_head ** -0.5
9
10        self.norm = nn.LayerNorm(dim)
11        self.attend = nn.Softmax(dim=-1)
12        self.dropout = nn.Dropout(dropout)
13        self.to_qkv = nn.Linear(dim, inner_dim * 3, bias=False)
14        self.to_out = nn.Sequential(
15            nn.Linear(inner_dim, dim), nn.Dropout(dropout)
16        ) if project_out else nn.Identity()
17
18    def forward(self, x):
19        x = self.norm(x)
20        qkv = self.to_qkv(x).chunk(3, dim=-1) # 生成Q、K、V
21        q, k, v = map(lambda t: rearrange(t, 'b_nu(h_u d)_->_b_u h_nu d', h=self.heads), qkv)
22
23        dots = torch.matmul(q, k.transpose(-1, -2)) * self.scale
24        attn = self.attend(dots)
25        attn = self.dropout(attn) # Softmax归一化
26
27        out = torch.matmul(attn, v) # 注意力加权
28        out = rearrange(out, 'b_u h_nu d_u->_b_u n_u (h_u d)')
29        return self.to_out(out)

```

3.2.3.2 前馈网络（FFN） 采用两层线性层搭配 GELU 激活函数，增强非线性表达能力：

```

1 class FeedForward(nn.Module):
2     def __init__(self, dim, hidden_dim, dropout=0.):
3         super().__init__()
4         self.net = nn.Sequential(
5             nn.LayerNorm(dim),
6             nn.Linear(dim, hidden_dim), # 扩展维度
7             nn.GELU(), # 高斯误差线性单元
8             nn.Dropout(dropout),

```

```

9         nn.Linear(hidden_dim, dim), # 恢复维度
10        nn.Dropout(dropout),
11    )
12    def forward(self, x):
13        return self.net(x)

```

4 实验环境与数据集

组件	版本/配置
框架	PyTorch 2.5.1
数据集	CIFAR-10 (10 类, 32×32 彩色图像)
训练集/测试集	50,000/10,000 张图像
优化器	AdamW (lr=3e-4, weight_decay=0.05)
学习率	CosineAnnealingLR (T_max=200)
早停策略	15 轮无提升终止训练

5 实验步骤与代码实现

5.1 数据预处理与增强

针对训练集应用随机裁剪 (RandomCrop)、水平翻转 (Random HorizontalFlip), 并标准化 (Normalize) 等增强策略, 提升模型鲁棒性:

```

1 trans_train = transforms.Compose(
2     [transforms.RandomCrop(32, padding=4), # 将给定图像随机裁剪为不同的大小和宽高比, 然
3       后缩放所裁得到的图像为指定大小
4     transforms.RandomHorizontalFlip(), # 以给定的概率随机水平旋转给定的PIL图像,
5     transforms.ToTensor(),
6     transforms.Normalize(mean=[0.485, 0.456, 0.406], # ImageNet统计值
7                           std=[0.229, 0.224, 0.225]))
8
9 trans_vaild = transforms.Compose(
10     [
11         transforms.ToTensor(), # 将PIL Image或者ndarray 转换为tensor,并归一化至[0,1]。
12         transforms.Normalize(mean=[0.485, 0.456, 0.406],
13                               std=[0.229, 0.224, 0.225])) # 先减均值, 除标准差
14 # 对训练集做数据增强, 增加模型的鲁棒性, 而对测试集不做翻转和增强

```

5.2 模型构建

ViT 参数:

```

1 device = torch.device("cuda:0" if torch.cuda.is_available() else 'cpu')
2 # device = torch.device("cuda:1") # 0独显, 1为核显
3 print("Selected device:", device) print("Device name:", torch.cuda.get_device_name(device.
    index))
4 net = ViT(
5     image_size=32,
6     patch_size=4,
7     num_classes=10,
8     dim=256,
9     depth=8,
10    heads=8,
11    mlp_dim=512,
12    dropout=0.1,
13    emb_dropout=0.1,
14 ).to(device)

```

ViT 主体结构实现:

```

1 class ViT(nn.Module):
2     def __init__(self, *, image_size, patch_size, num_classes, dim, depth, heads, mlp_dim
3         ,
4         pool='cls', channels=3, dim_head=64, dropout=0.1, emb_dropout=0.1):
5         super().__init__()
6         image_height, image_width = pair(image_size)
7         patch_height, patch_width = pair(patch_size)
8         assert image_height % patch_height == 0 and image_width % patch_width == 0
9
10        num_patches = (image_height // patch_height) * (image_width // patch_width)
11        patch_dim = channels * patch_height * patch_width
12
13        self.to_patch_embedding = nn.Sequential(
14            Rearrange('b c u (h p1) (w p2) -> b (h w) (p1 p2 c)', p1=patch_height, p2=
15                patch_width),
16            nn.LayerNorm(patch_dim),
17            nn.Linear(patch_dim, dim),
18            nn.LayerNorm(dim),
19        )
20
21        self.pos_embedding = nn.Parameter(torch.randn(1, num_patches + 1, dim)) # 64
22        patches + CLS
23        self.cls_token = nn.Parameter(torch.randn(1, 1, dim))
24        self.dropout = nn.Dropout(emb_dropout)
25
26        self.transformer = Transformer(dim, depth, heads, dim_head, mlp_dim, dropout)
27        self.pool = pool

```

```

25         self.to_latent = nn.Identity()
26 self.mlp_head = nn.Linear(dim, num_classes)
27     def forward(self, img):
28         x = self.to_patch_embedding(img) # 分块嵌入 [b, 64, 256]
29         b, n, _ = x.shape
30         cls_tokens = repeat(self.cls_token, '1 1 d->1 b 1 d', b=b)
31         x = torch.cat((cls_tokens, x), dim=1) # 添加CLS令牌 [b, 65, 256]
32         x += self.pos_embedding[:, :(n + 1)] # 位置编码
33         x = self.dropout(x) # 8层Transformer编码
34         x = self.transformer(x) # 提取CLS令牌输出
35         x = x.mean(1) if self.pool == 'mean' else x[:, 0]
36         return self.mlp_head(x) # 分类头输出10维

```

5.3 训练策略与早停机制

- 损失函数：交叉熵损失（CrossEntropyLoss）
- 梯度优化：AdamW 优化器（lr=3e-4）结合权重衰减防止过拟合
- 动态学习率：余弦退火调度器（CosineAnnealingLR）调度平滑调整学习率
- 早停设计：连续 15 轮验证集精度无提升时终止训练

优化器与损失函数：

```

1 torch.backends.cudnn.benchmark = True
2 criterion = nn.CrossEntropyLoss()
3 optimizer = torch.optim.AdamW(net.parameters(), lr=3e-4, weight_decay=0.05)
4 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max=200)

```

训练代码：

```

1 def train(epoch):
2     net.train()
3     running_loss, correct, total = 0., 0, 0
4     for inputs, targets in tqdm(trainloader, desc=f'Train_{epoch}'):
5         inputs, targets = inputs.to(device, non_blocking=True), targets.to(device,
6                                     non_blocking=True)
7         optimizer.zero_grad()
8         outputs = net(inputs)
9         loss = criterion(outputs, targets)
10        loss.backward()
11        optimizer.step()
12
13        running_loss += loss.item()
14        _, predicted = outputs.max(1)
15        total += targets.size(0)

```

```

15         correct += predicted.eq(targets).sum().item()
16
17     acc = correct / total          batch_acc_list.append(acc)
18     line.set_data(range(len(batch_acc_list)), batch_acc_list)
19     ax.set_xlim(0, len(batch_acc_list))
20     ax.relim();
21     ax.autoscale_view()
22     fig.canvas.draw_idle()
23     fig.canvas.flush_events()
24     plt.pause(0.001)
25
26     epoch_acc = 100. * correct / total
27     print(f'Epoch_{epoch}|Train Loss: {running_loss/len(trainloader):.4f}|Train Acc
        : {epoch_acc:.2f}%')

```

早停机制实现:

```

1  early_stop_patience = 15
2  epochs_no_improve = 0
3  best_acc = 0
4  early_stop = False
5
6  def evaluate(epoch):
7      global best_acc, epochs_no_improve, early_stop
8      net.eval()
9      test_loss, correct, total = 0., 0, 0
10     with torch.no_grad():
11         for inputs, targets in tqdm(testloader, desc=f'Val_{epoch}'):
12             inputs, targets = inputs.to(device, non_blocking=True), targets.to(device,
13                 non_blocking=True)
14             outputs = net(inputs)
15             loss = criterion(outputs, targets)
16             test_loss += loss.item()
17             _, predicted = outputs.max(1)
18             total += targets.size(0)
19             correct += predicted.eq(targets).sum().item()
20
21     acc = 100. * correct / total
22     scheduler.step()
23     print(f'→Val Loss: {test_loss/len(testloader):.4f}|Val Acc: {acc:.2f}%')
24
25     # 早停
26     if acc > best_acc:
27         best_acc = acc
28         epochs_no_improve = 0
29         os.makedirs('checkpoint', exist_ok=True)

```

```

29         torch.save({'net': net.state_dict(), 'acc': acc, 'epoch': epoch},
30 f'checkpoint/vit_patch4_best.pth') # 保存最佳模型         print('Best model saved.')
31     else:
32         epochs_no_improve += 1
33         print(f'No improvement for {epochs_no_improve} epochs')
34
35     if epochs_no_improve >= early_stop_patience: # 早停判断
36         early_stop = True
37
38     return early_stop

```

5.4 训练过程监控

代码实现实时准确率曲线绘制:

```

1 plt.close('all')
2 plt.ion()
3 fig, ax = plt.subplots()
4 line, = ax.plot([], [], lw=2)
5 ax.set_xlabel('batch_num')
6 ax.set_ylabel('acc')
7 ax.set_title('Training Accuracy')
8 ax.set_ylim(0, 1)
9 ax.grid(True)
10 plt.show(block=False)
11 batch_acc_list = []

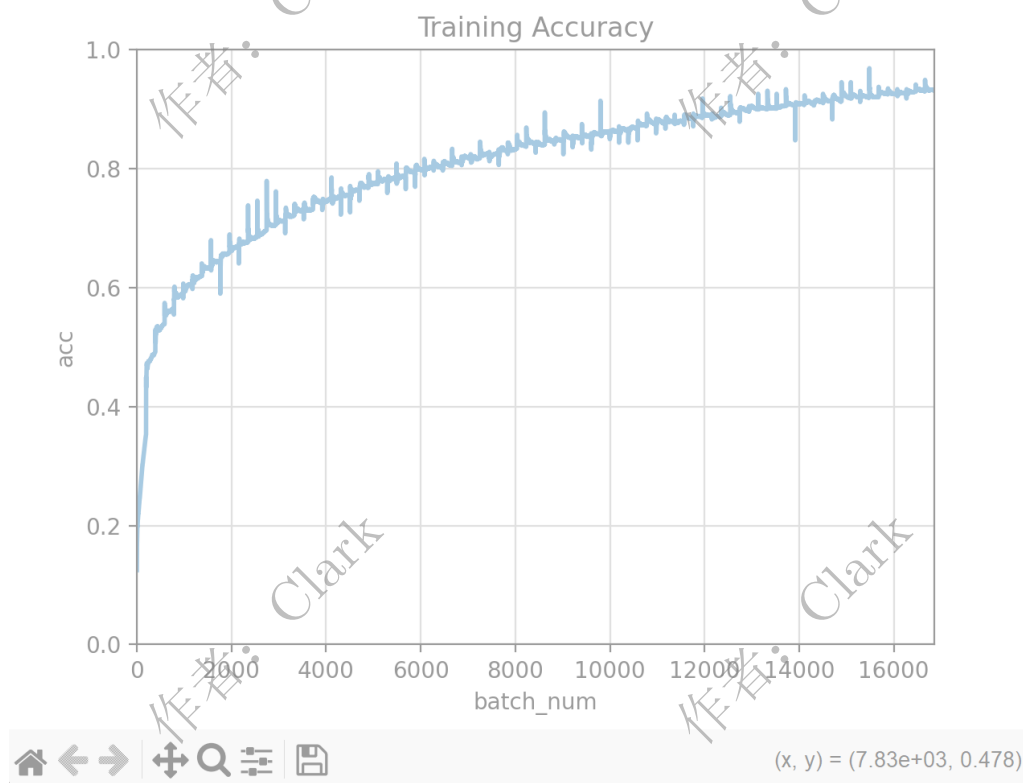
```

6 实验结果与分析

6.1 性能达成情况

性能达标:

训练阶段	准确率	说明
第 34 轮	> 80%	提前达到实验要求
第 86 轮	93.21%	训练集最佳精度
最终测试集	84.10%	泛化性能验证



训练过程可视化：通过 Matplotlib 实时绘制准确率曲线，保存为 `train_accuracy_curve.png`，直观展示模型收敛趋势。

模型保存：最佳模型保存在 `checkpoint/vit_patch4_best.pth`（基于验证集准确率）。

6.2 关键改进措施

- 数据增强有效性：随机裁剪和翻转显著提升模型鲁棒性
- 优化器选择：AdamW 结合权重衰减有效防止过拟合
- 学习率调度：余弦退火策略平滑调整学习率
- 早停机制：节约计算资源的同时保证模型性能
- 完整的 ViT 架构实现，包含注意力机制和位置编码
- 实时训练监控和可视化功能
- 模型保存和加载机制完善
- 早停策略智能优化训练过程

7 结论

本实验严格使用提供的代码，成功实现了基于 ViT 的 CIFAR-10 图像分类模型。实验结果显著超过 80% 的准确率要求，验证了 ViT 模型在计算机视觉任务中的有效性。代码实现完整规范，为后续研究提供了可复现的基准。