

# 一.模型压缩概述

- 为什么要模型压缩？

端侧设备存在资源限制

- 为什么能模型压缩？

深度神经网络模型存在冗余性

- 什么是模型压缩？

利用网络模型冗余性的特点，减小模型规模的方法

- 压缩方法有哪些？

剪枝：修剪不重要的网络连接

量化：将连续型数据量化为低位宽离散数据

知识蒸馏：大模型指导小模型学习

低秩分解：通过低秩矩阵近似原矩阵

轻量化网络：使用轻量化卷积核代替传统卷积

网络结构搜索：自动化地设计优异网络模型

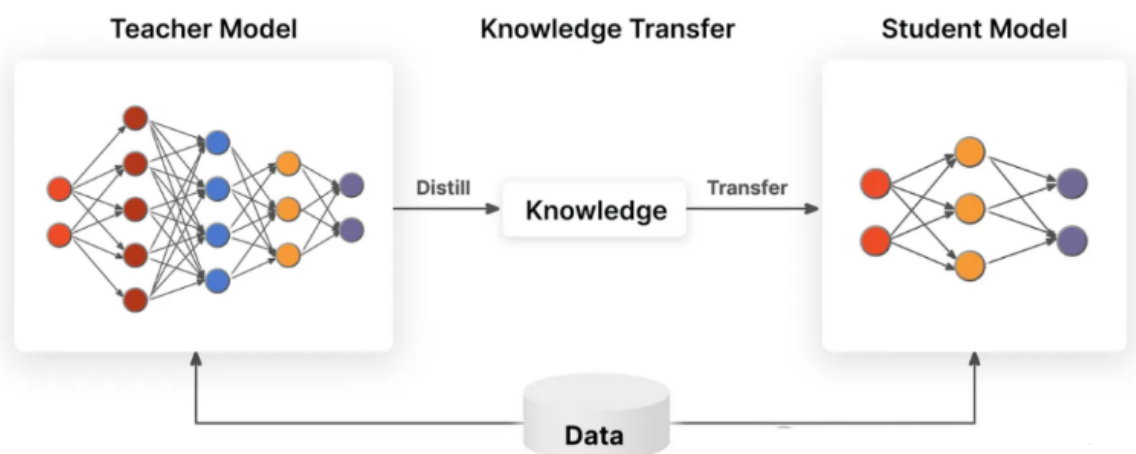
## 二.方法

### 1.模型蒸馏

蒸馏是一种将大型、复杂模型的知识迁移到较小模型中的技术。

通常，蒸馏的过程是训练一个小模型（学生模型）以模仿一个较大的、预先训练好的模型（教师模型）的行为。

小模型通过学习教师模型的预测概率分布来获取知识，而不仅仅是传统的标签信息。



蒸馏的主要思想是：

- 教师模型输出的类别概率包含了更多的“软信息”，这些信息能够帮助学生模型更好地学习一些复杂的模式。

- 学生模型通过与教师模型输出的“软标签”进行学习，能够在不完全依赖硬标签的情况下获取更多的信息，进而提高其性能。

蒸馏的步骤通常是：

#### 1. 训练教师模型

首先，训练一个大型且高性能的教师模型，这通常是一个深度神经网络。

#### 2. 训练学生模型

然后，训练一个较小的学生模型，目标是通过最小化学生模型与教师模型在相同输入上的输出差异来进行训练。学生模型不仅学习硬标签（真实标签），还学习教师模型的“软标签”。

通过蒸馏，学生模型可以获得教师模型中蕴含的丰富知识，尤其是在教师模型能够捕获的复杂特征和模式方面，从而在保持较小规模的同时接近或达到教师模型的性能。

### 优缺点

- 优点：蒸馏可以显著提高小型模型的性能，使其在压缩后依然保持接近教师模型的精度，尤其在大型模型压缩时表现出色。
- 缺点：蒸馏的一个挑战是教师模型的选择和训练需要耗费大量的计算资源和时间。此外，蒸馏的效果可能会受到学生模型的限制，对于某些任务，学生模型的性能可能不容易达到教师模型的水平。

```
import torch
import torch.nn as nn
import torch.optim as optim

# 创建一个简单的教师模型和学生模型
class TeacherNet(nn.Module):
    def __init__(self):
        super(TeacherNet, self).__init__()
        self.fc = nn.Linear(10, 10)

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

class StudentNet(nn.Module):
    def __init__(self):
        super(StudentNet, self).__init__()
        self.fc = nn.Linear(10, 10)

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

# 创建模型实例
teacher = TeacherNet()
student = StudentNet()

# 使用教师模型生成“软标签”
def distillation_loss(student_outputs, teacher_outputs, temperature=2.0):
    # 使用温度缩放进行蒸馏损失计算
    loss = nn.KLDivLoss()(nn.functional.log_softmax(student_outputs /
                                                         temperature, dim=1),
                          nn.functional.softmax(teacher_outputs / temperature,
                                                  dim=1)) * (temperature ** 2)
    return loss

# 简单的训练循环
optimizer = optim.SGD(student.parameters(), lr=0.1)
```

```

# 模拟训练过程
for epoch in range(100):
    # 输入数据
    inputs = torch.randn(32, 10) # 假设批次大小是32，输入维度是10
    teacher_outputs = teacher(inputs)

    # 学生模型的输出
    student_outputs = student(inputs)

    # 计算损失
    loss = distillation_loss(student_outputs, teacher_outputs)

    # 反向传播
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()

    if epoch % 10 == 0:
        print(f"Epoch [{epoch}/100], Loss: {loss.item()}")

```

## 2.剪枝

- 修剪是通过去除神经网络中某些不重要的连接或神经元来减少模型的规模和计算需求。
- 修剪的目标是去除那些对网络性能影响较小的参数，从而达到减少模型复杂度的效果。

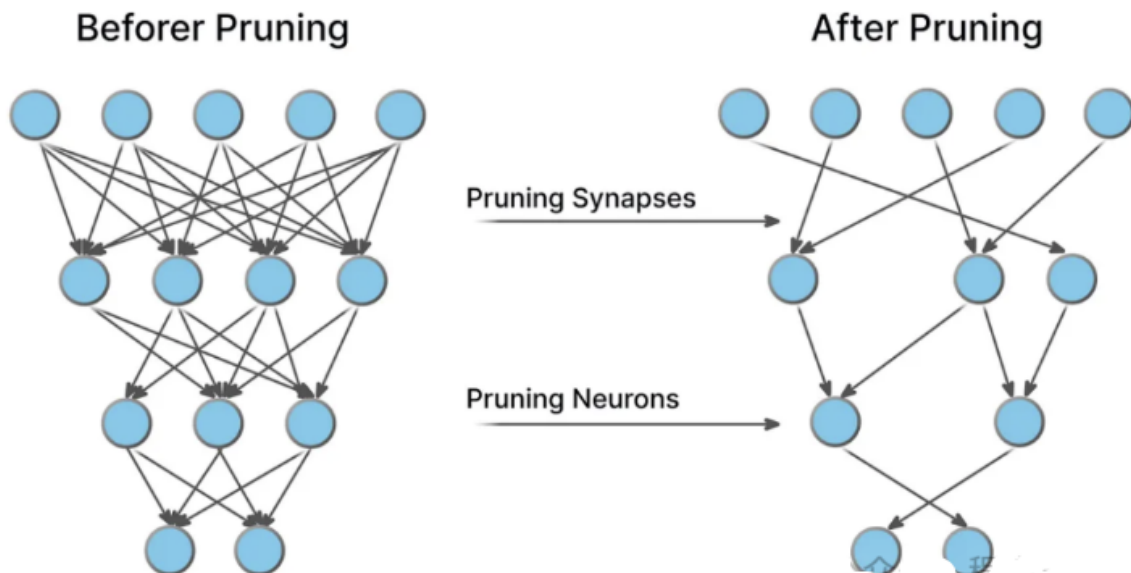
### 常见修剪策略

- 权重修剪

通过移除那些对网络输出贡献较小的权重来减少模型的大小。  
这些权重可以通过设定一个阈值来判定：低于某个阈值的权重会被剪掉。

- 神经元修剪

修剪掉整个神经元或通道，这样的修剪方法可以进一步减少计算量，尤其是对于卷积神经网络（CNN）来说，移除不重要的特征图通道会显著降低计算复杂度。



修剪的步骤通常是：

1. 训练原始模型
2. 计算每个权重的重要性或每个神经元的激活度
3. 去除不重要的权重或神经元
4. 重新训练，以恢复性能损失

优缺点

- 优点：减小模型尺寸，降低计算负担，提升推理速度，尤其适合硬件加速。
- 缺点：修剪过度可能导致模型性能下降。需要精心设计修剪方案，以在压缩和性能之间找到平衡。

```
import torch
import torch.nn as nn
import torch.nn.utils.prune as prune

# 定义一个简单的神经网络
class SimpleNet(nn.Module):
    def __init__(self):
        super(SimpleNet, self).__init__()
        self.fc1 = nn.Linear(10, 10)
        self.fc2 = nn.Linear(10, 2)

    def forward(self, x):
        x = torch.relu(self.fc1(x))
        x = self.fc2(x)
        return x

# 创建网络和输入数据
model = SimpleNet()
input_data = torch.randn(1, 10)

# 修剪fc1层的20%的权重
prune.random_unstructured(model.fc1, name="weight", amount=0.2)

# 打印fc1层的权重，观察被修剪掉的权重
print(model.fc1.weight)
```

## 3.量化

量化是将浮点数表示的参数（如权重和激活）转换为低精度数值表示（如整数）。

量化通常将模型从 32 位浮点数转换为更低精度的数据类型，如16位、8位或更低，这样可以减少存储需求和加速推理过程。

- 权重量化  
将模型中的浮点数权重转换为低精度整数。例如，将32位浮点数权重映射到8位整数，这样就能大幅减少模型的存储需求。
- 激活量化  
对于激活值（神经网络各层的输出），也可以应用类似的量化策略。

常见的量化类型

- 后训练量化：在模型训练完成后进行量化，适用于已经训练好的模型。

- 量化感知训练：在训练过程中加入量化过程，从而使得模型能够适应低精度的计算。

量化不仅减小了模型大小，还可能加速模型的推理过程，尤其是在支持低精度计算的硬件上（如TPU、GPU等）。

### 优缺点

- 优点：大幅减少模型的存储需求，加速推理过程。尤其在嵌入式设备和移动端设备上具有显著的优势。
- 缺点：量化可能导致一定的精度损失

```
import torch
import torch.nn as nn
import torch.quantization

# 定义一个简单的模型
class SimpleModel(nn.Module):
    def __init__(self):
        super(SimpleModel, self).__init__()
        self.fc1 = nn.Linear(28 * 28, 128)
        self.fc2 = nn.Linear(128, 10)

    def forward(self, x):
        x = torch.relu(self.fc1(x))
        x = self.fc2(x)
        return x

# 初始化模型
model = SimpleModel()

# 转换模型为量化版本
model.eval() # 切换到评估模式
quantized_model = torch.quantization.quantize_dynamic(
    model, {nn.Linear}, dtype=torch.qint8
)

# 查看量化后的模型
print(quantized_model)
```