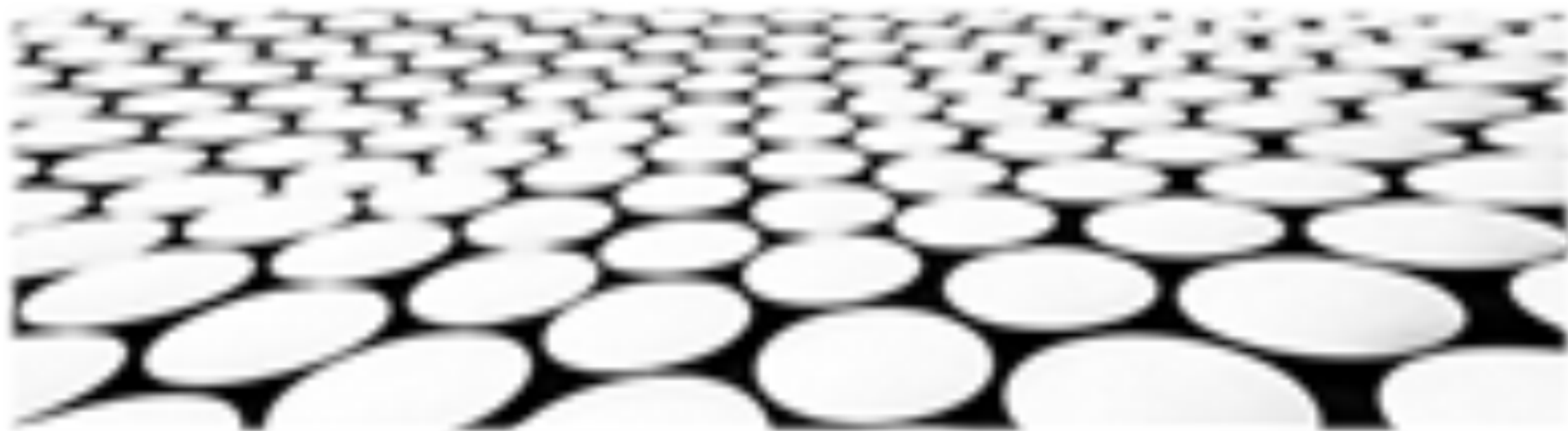


# 模型量化压缩算法





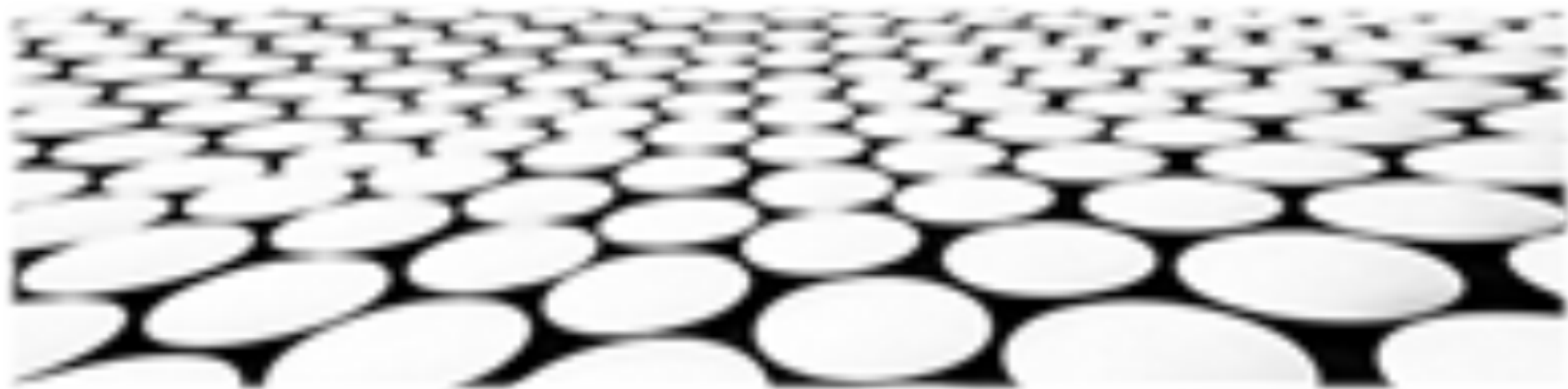
## 目录页

Contents Page

- 比。
3. **稀疏性原理**：将网络中的权重和激活值置为零，减少表达数据所需位的数目。
4. **结构化剪枝**：采用特定的剪枝策略，如大小修剪、卷积核修剪等，去除网络中的不重要部分。
5. **量化感知训练**：在训练过程中，通过量化感知损失函数，引导网络学习量化友好特征。
6. **张量分解**：将网络中的张量分解为多个低秩张量，降低计算复杂度和内存消耗。
7. **图压缩**：使用图压缩算法，将网络表示为更紧凑的图结构，减少模型大小。



量化基本原理：降低精度，减少表达数据所需位的数目。



# 量化基本原理：降低精度，减少表达数据所需位的数目。

## ■ 量化基本概念

1. 量化是指将高精度的数据表示转换为低精度的数据表示。
2. 量化可以降低数据表示所需的位数，从而减少存储空间和计算量。
3. 量化还可以减少模型参数的数量，从而减少模型大小和提高模型的推理速度。

## ■ 量化方法

1. 线性量化是最常用的量化方法，它将数据值均匀地映射到一个较小的值域。
2. 非线性量化可以更好地保留数据分布的特征，但计算量也更大。
3. 自适应量化可以根据输入数据的分布动态地调整量化参数，从而提高量化精度。



# 量化基本原理：降低精度，减少表达数据所需位的数目。



## 量化评估

1. 量化评估的指标包括准确率、召回率、F1值等。
2. 量化评估还可以使用其他指标，如模型大小、推理速度等。
3. 量化评估的结果可以指导量化算法的选择和参数设置。



## 量化应用

1. 量化在移动端和嵌入式设备上非常重要，因为这些设备的存储空间和计算资源有限。
2. 量化也可以用于云端计算，以降低存储成本和提高计算速度。
3. 量化在自动驾驶、医疗诊断等领域也有广泛的应用。

# 量化基本原理：降低精度，减少表达数据所需位的数目。

## ■ 量化研究热点

1. 研究热点之一是低精度量化，即使用更少的位数来表示数据。
2. 研究热点之二是动态量化，即根据输入数据的分布动态地调整量化参数。
3. 研究热点之三是自适应量化，即根据任务和模型自动选择量化算法和参数。

## ■ 量化发展趋势

1. 量化算法将继续朝着更准确、更有效的方向发展。
2. 量化技术将与其他技术相结合，如剪枝、蒸馏等，以进一步提高模型的压缩率和性能。
3. 量化技术将在更多领域得到应用，如自动驾驶、医疗诊断、金融等。



**修剪原理：去除网络中的冗余参数，提高计算效率和压缩比。**



# 修剪原理：去除网络中的冗余参数，提高计算效率和压缩比。

## 修剪原则

1. 修剪规则的选择：权值剪枝稀疏规则的选择对压缩效果和模型性能都有很大影响。目前常用的权值剪枝稀疏规则包括：L1范数稀疏、L2范数稀疏、最大值稀疏和随机稀疏等。
2. 修剪强度：修剪强度是指被修剪节点的百分比。修剪强度越高，压缩率越高，但模型性能可能下降。因此，在选择修剪强度时需要权衡压缩率和模型性能。
3. 修剪策略：修剪策略是指修剪过程中的具体步骤和方法。常用的修剪策略包括：一次性修剪和迭代修剪等。一次性修剪是指一次性从网络中去除所有冗余参数，而迭代修剪是指根据网络的结构和参数分布，分阶段地去除冗余参数。



# 修剪原理：去除网络中的冗余参数，提高计算效率和压缩比。

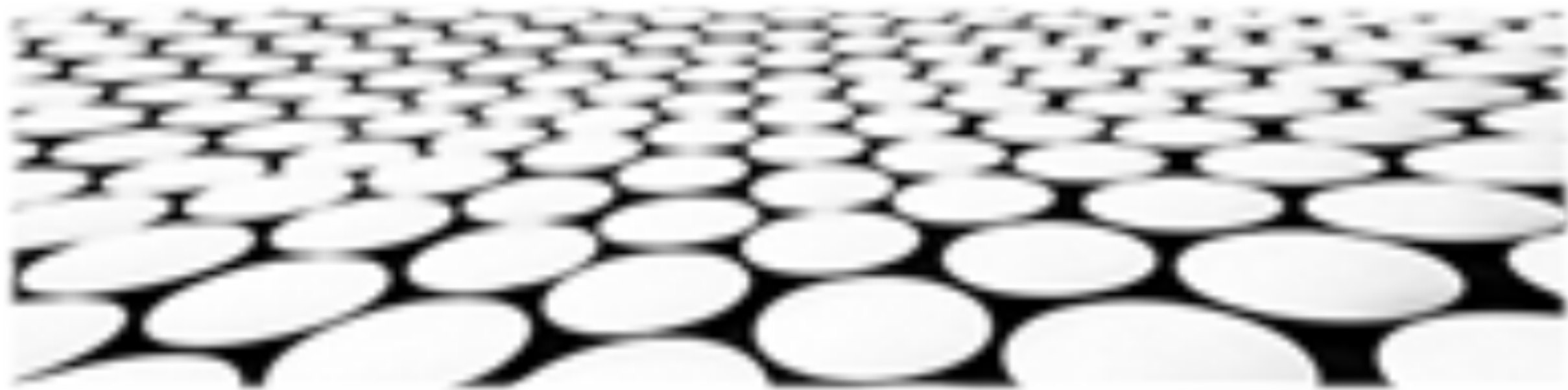


## 修剪方法

1. 权值修剪：权重修剪是压缩模型最常用的方法。权重修剪是指将网络中的部分权重值设置为0，从而减少网络的参数数量。权重修剪可以分为结构化修剪和非结构化修剪。结构化修剪是指将整个通道或层设置为0，而非结构化修剪是指将单个权重设置为0。
2. 激活值修剪：激活值修剪是指将网络中的部分激活值设置为0，从而减少网络的参数数量。激活值修剪可以分为结构化修剪和非结构化修剪。结构化修剪是指将整个通道或层设置为0，而非结构化修剪是指将单个激活值设置为0。
3. 滤波器修剪：滤波器修剪是指将网络中的部分滤波器去除，从而减少网络的参数数量。滤波器修剪可以分为结构化修剪和非结构化修剪。结构化修剪是指将整个通道或层去除，而非结构化修剪是指去除单个滤波器。



**稀疏性原理：将网络中的权重和激活值置为零，减少表达数据所需位的数目。**



## 稀疏性原理及其应用

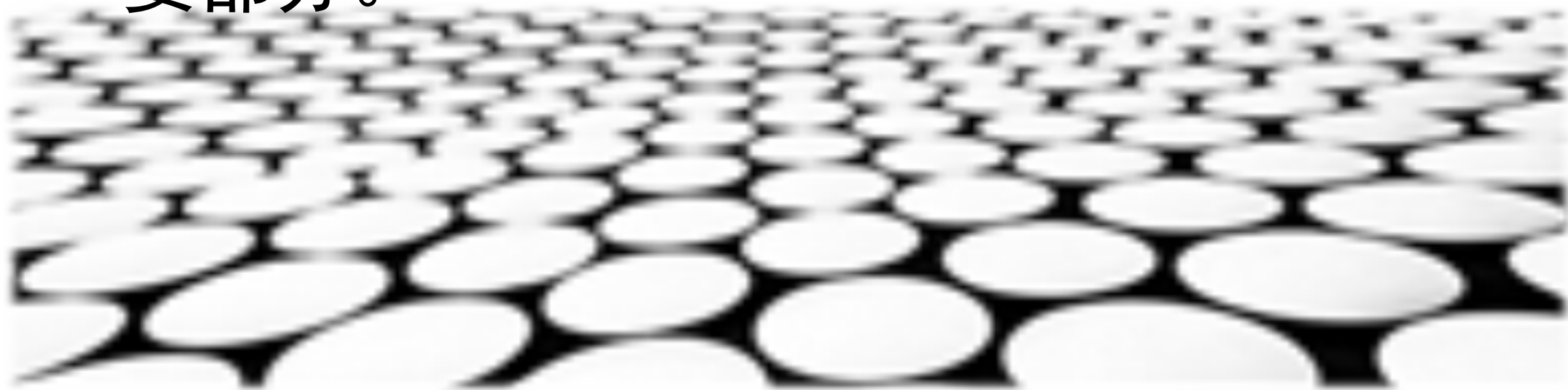
- 1. 稀疏性原理概述：**权重和激活值的稀疏性是神经网络的一项重要特性，即许多神经元连接（即权重）和神经元输出值（即激活值）通常很小或为零。基于这一特性，可以将网络中的权重和激活值置为零，减少表达数据所需比特数目，从而实现模型压缩。
- 2. 稀疏性压缩方法：**实现稀疏性压缩的方法有多种，包括修剪、量化和结构化稀疏。修剪是直接将权重和激活值中较小的值设置为零，量化是将权重和激活值离散化为较小的值，结构化稀疏是对网络进行结构上的改变，使其具有稀疏性。
- 3. 稀疏性压缩的优势：**稀疏性压缩可以有效地减少模型的大小和计算成本，同时保持或提高模型的精度。此外，稀疏性压缩可以提高模型的鲁棒性和可解释性。

## 稀疏性原理的挑战与应对策略

- 1. 挑战：**稀疏性压缩面临的主要挑战之一是权重和激活值置零后，会导致模型精度下降。此外，稀疏性压缩后，模型的训练和推理过程变得更加复杂，需要特殊的设计和优化策略。
- 2. 应对策略：**为了应对稀疏性压缩带来的精度下降问题，可以采用各种正则化技术和训练方法来提高模型的鲁棒性。对于训练和推理过程的复杂性问题，可以设计专门的算法和优化策略来提高效率。
- 3. 最新进展：**近年来，稀疏性压缩技术取得了很大进展，涌现出许多新的压缩方法和优化策略。这些方法可以有效地减少模型的大小和计算成本，同时保持或提高模型的精度。



结构化剪枝：采用特定的剪枝策略，如大小修剪、卷积核修剪等，去除网络中的不重要部分。



## 结构化剪枝

### 1. 结构化剪枝的基本原理及思想：

- 识别网络中的不重要节点（例如，神经元、滤波器、通道或连接），并将其剪掉。
- 目标是去除对网络性能影响最小的部分，同时保持其准确率。

### 2. 各种结构化剪枝方法：

- 筛选方法：
  - 基于权重的筛选：通过阈值或排序，直接去除冗余的权重。
  - 基于梯度的筛选：使用梯度信息来评估权重的重要性，并去除不重要的权重。

### 3. 优化过程及策略：

- 一次性剪枝：在训练完成后进行一次剪枝。
- 迭代式剪枝：在训练过程中反复进行剪枝。
- 贪婪式剪枝：每次剪掉最不重要的部分。

## 大小修剪

### 1. 大小修剪的目标：

- 减少网络的通道数或卷积核数量，从而降低计算量和内存占用。
- 保持网络的整体结构和性能。

### 2. 实现方式：

- 通道修剪：
  - 针对每个卷积层，移除不重要的通道。
  - 可通过权重值、梯度信息、激活值等作为剪枝指标。
- 卷积核修剪：
  - 针对每个卷积层，移除不重要的卷积核。
  - 可通过权重值、梯度信息、空间分布等作为剪枝指标。

### 3. 挑战：

- 如何选择合适的剪枝比例。

# 结构化剪枝：未用特定的剪枝策略，如大小修剪、卷积核修剪等 去除网络中的不重要部分

## 卷积核修剪

1. 卷积核修剪的基本思想：
  - 识别重要卷积核，去除冗余卷积核。
2. 如何选择要剪掉的卷积核：
  - 权重值：剪掉权重值较小的卷积核。
  - 梯度信息：剪掉梯度信息较小的卷积核。
  - 激活值：剪掉激活值较小的卷积核。
3. 如何保证剪枝后模型的性能：
  - 冻结剪掉的卷积核，只训练剩余的卷积核。
  - 使用正则化技术防止模型过拟合。
  - 微调网络以适应新的架构。

## 细粒度剪枝

1. 细粒度剪枝的目标：
  - 进一步减少网络的冗余，提高模型的压缩率和性能。
  - 实现不同粒度的剪枝，如通道修剪、卷积核修剪、滤波器修剪等。
2. 细粒度剪枝的方法：
  - 剪枝搜索算法：
    - 基于贪婪算法的搜索方法：逐层、逐通道、逐卷积核地进行剪枝。
    - 基于强化学习的搜索方法：将剪枝过程建模为马尔可夫决策过程，通过强化学习来寻找最佳的剪枝策略。
  - 联合剪枝方法：
    - 将不同粒度的剪枝方法结合起来，以达到最佳的效果。
    - 例如，先进行通道修剪，然后再进行卷积核修剪，或者先进行卷积核修剪，然后再进行滤波器修剪。

以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：  
<https://d.book118.com/267042113110006132>