# Magnitude Pruning

## Magnitude-Based Pruning

Increase sparsity by setting small weight values to 0

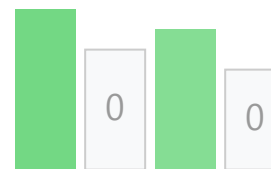## Pruning Process

**1**

Model Training
Complete

→

**2**

$|w|$

Calculate Weight
Magnitude $|w|$

→

**3**

Remove Weights
Below Threshold

→

**4**

Fine-tuning

**Before Pruning**

All weights

Prune →

**After Pruning**

0   0

Sparse (50%)

Threshold

## Threshold Setting Strategies

## Global Threshold

θ

Single threshold for entire network

## Layer-wise Threshold

$\theta_1$ $\theta_2$ $\theta_3$

Different threshold per layer

## Top-k Pruning

Top-k

Keep only top k% weights

## Sparsity Example

Dense: [0.8, 0.3, -0.5, 0.1, -0.9]

↓ Threshold = 0.4

Sparse: [0.8, 0, -0.5, 0, -0.9]

0.8 ✕ -0.5 ✕ -0.9

**40% Sparsity Achieved**