



Improving Deep Neural Networks

Video 3 : The Solution to Vanishing and Exploding Gradients - Part 2

Optimum Weight Initialization

Introduction



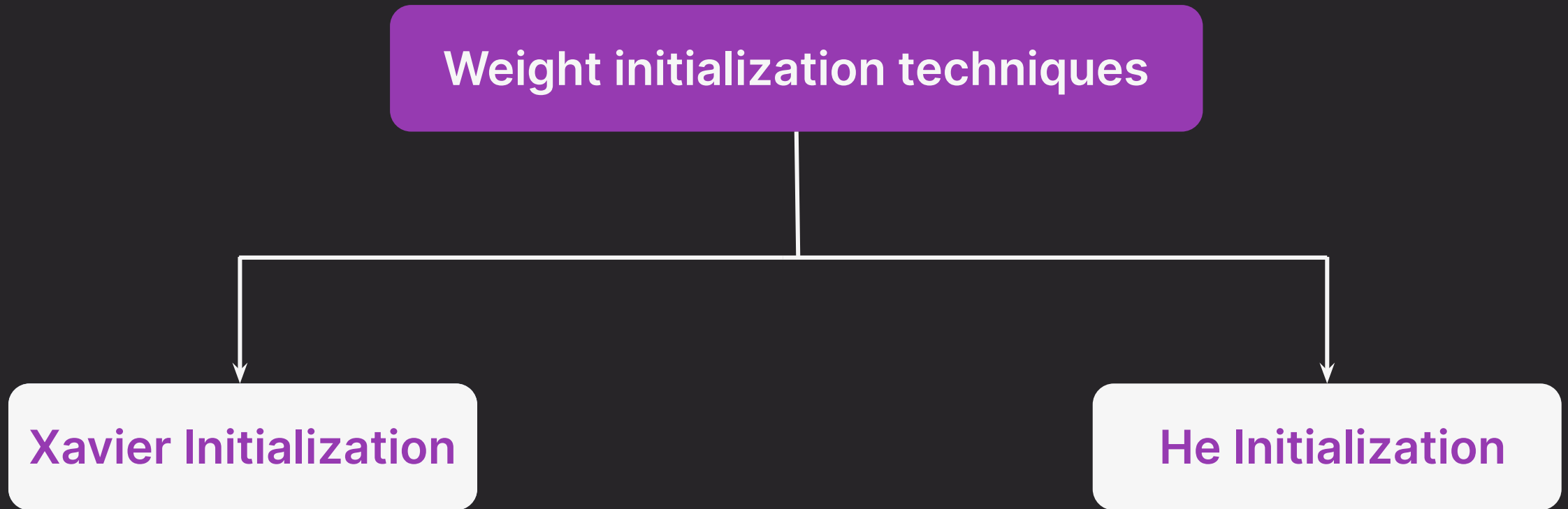
Forward propagation of the neural network



Weights assigned randomly

Introduction

Key concept: Keep activations and gradients variance consistent across layers



Xavier Initialization

Xavier initialization

Variance

=

1

Avg (no. of input neurons, no. of
output neurons)

Xavier Initialization



$$\text{Variance} = \frac{1}{\text{Avg (no. of input neurons, no. of output neurons)}}$$

He Initialization

He initialization

Variance

=

2

Avg (no. of input neurons, no. of
output neurons)

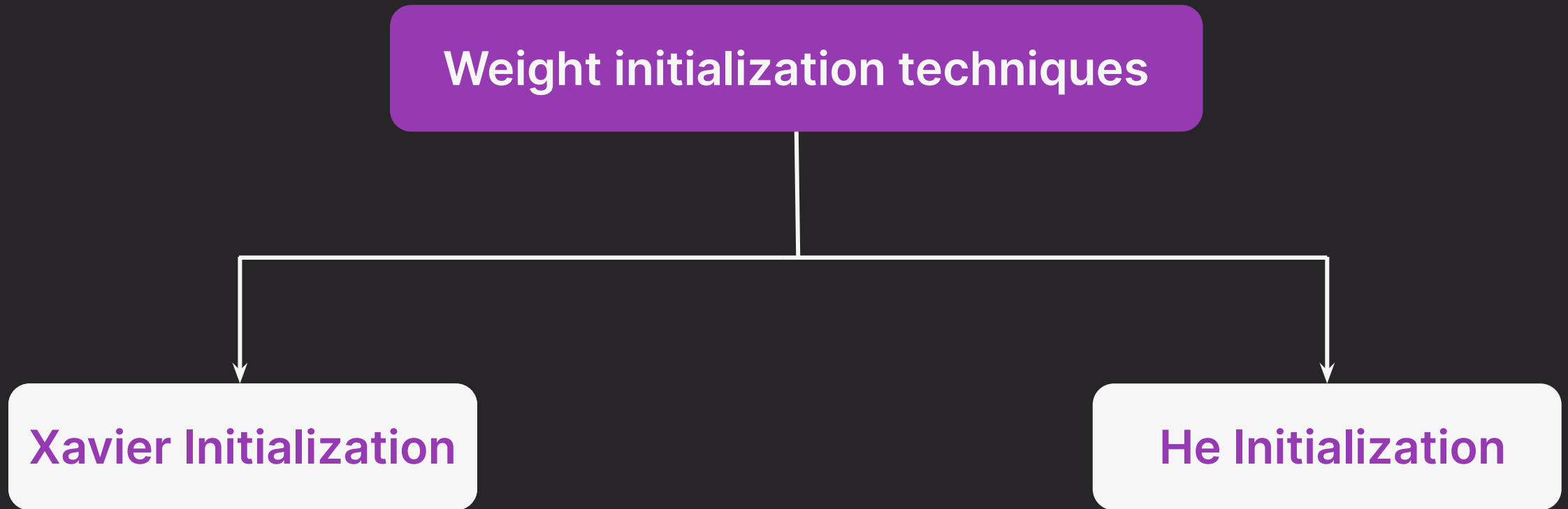
He Initialization



$$\text{Variance} = \frac{2}{\text{Avg (no. of input neurons, no. of output neurons)}}$$

Introduction

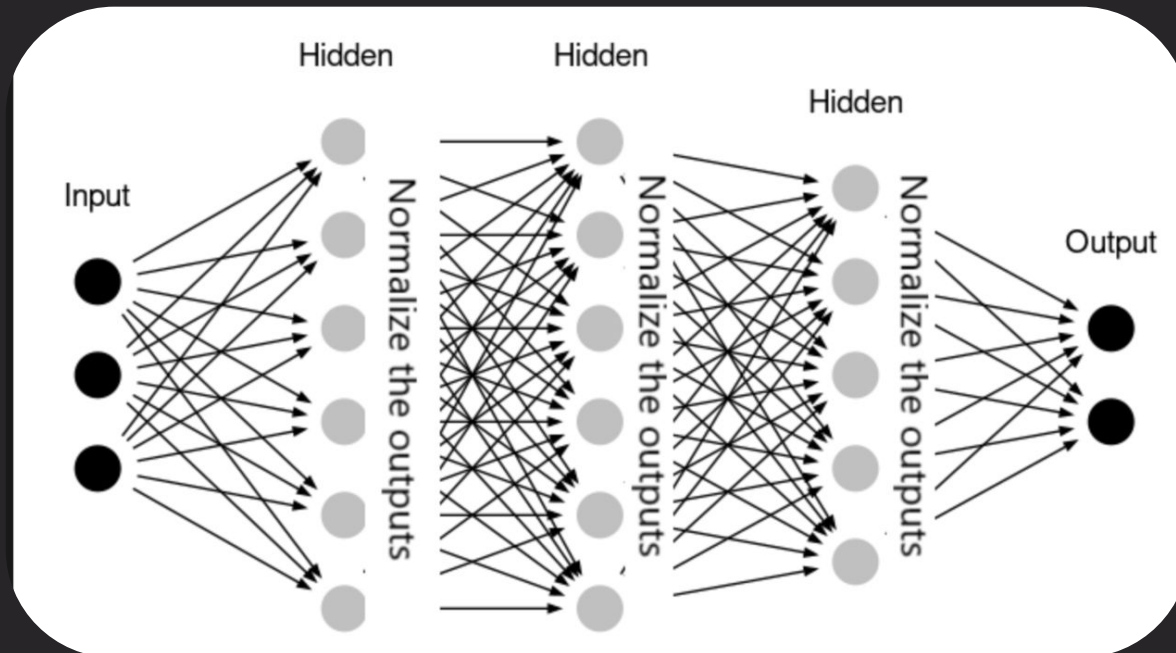
Key takeaway: PyTorch initializes weights using one of these techniques.



Batch Normalization

Key concept: Normalize previous outputs before sending it to the next layer.

- It is highly efficient and acts as a catalyst for faster convergence



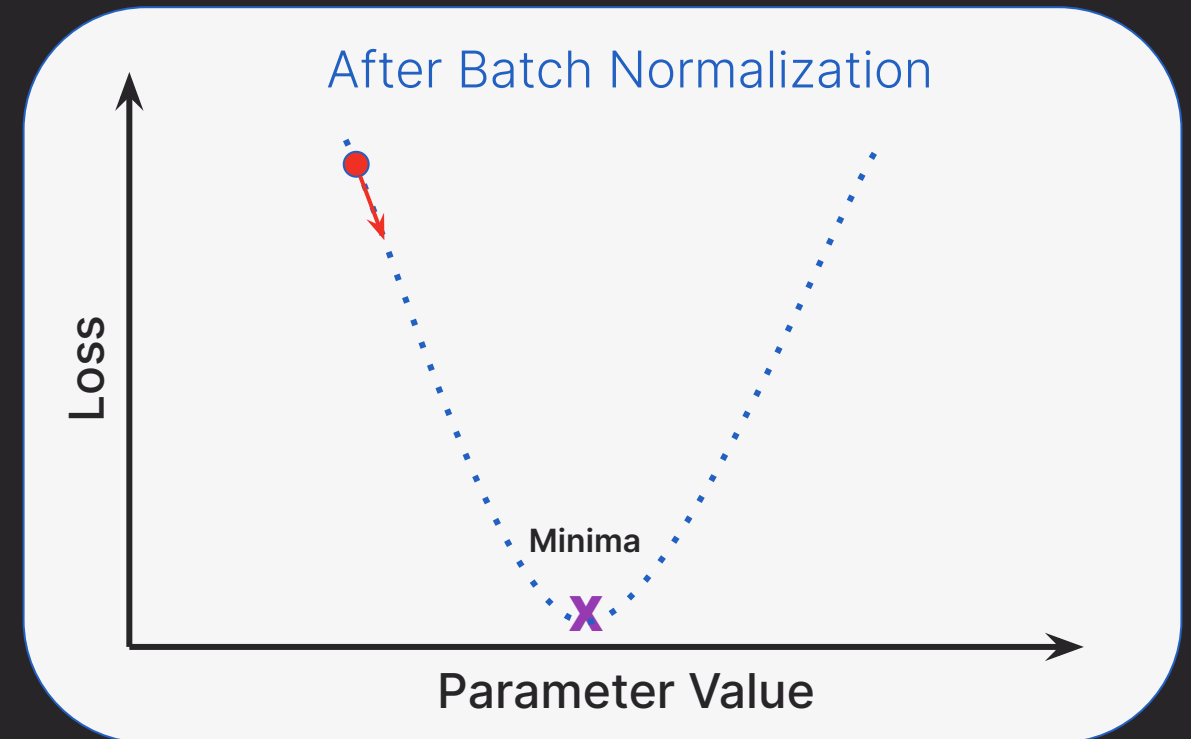
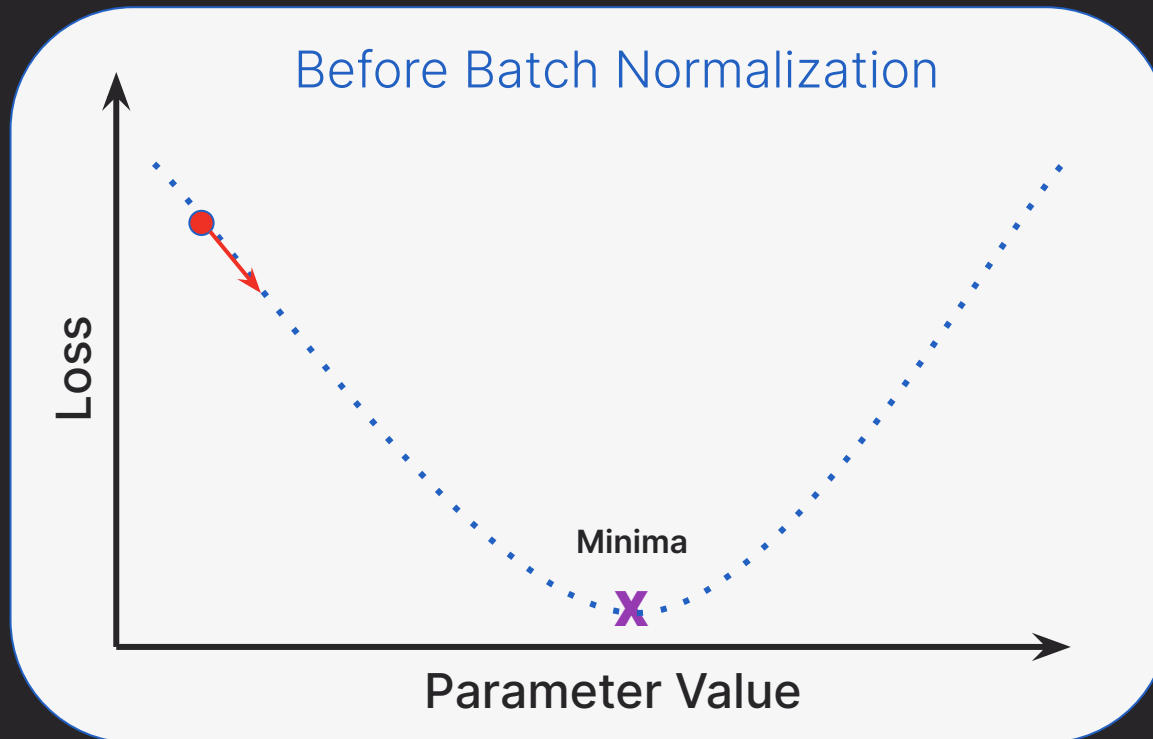
Batch Normalization

Why normalize the input features?

- Equal importance to all features
- Fast conversion



Batch Normalization: Analogy





What is “Batch” in Batch Normalization?

Batch Normalization



Neural networks and batch normalization
both operate in batches

How does Batch Normalization work?

Batch Normalization: Working

Activations of an input layer are normalized for each mini-batch of data:

$$NA_i = (a_i - \text{mean}) / \text{sqrt}(\text{variance} + \epsilon)$$



Batch Normalization: Working

Activations of an input layer are normalized for each mini-batch of data:

$$NA_i = (a_i - \text{mean}) / \sqrt{\text{variance} + \epsilon}$$

Mean of activation function in a layer = 0

Standard Deviation of activation function in a layer = 1

Batch Normalization: Working

Activations of an input layer are normalized for each mini-batch of data:

$$NA_i = (a_i - \text{mean}) / \sqrt{\text{variance} + \epsilon}$$

Adjustment:

$$z_i = \gamma NA_i + \beta$$

Fun Facts: Batch Normalization

- Helps in reducing overfitting marginally
- Batch normalization often replaces the neuron's bias with its beta term

Adjustment:

$$z_i = \gamma N A_i + \beta$$



Gradient Clipping

Gradient clipping prevents exploding gradients by capping their absolute value to a set threshold.

Clipped Gradient =

$\min(\text{threshold}, \text{gradient})$



Up-Next: Hands-On Vanishing and Exploding Gradients