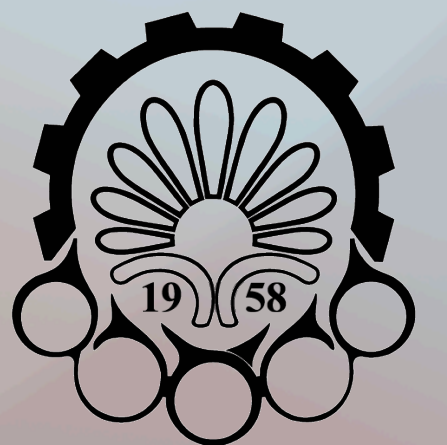


# Mastering MLP - Part 1

Computation Intelligence and its Application in Mechatronics

Winter 2025

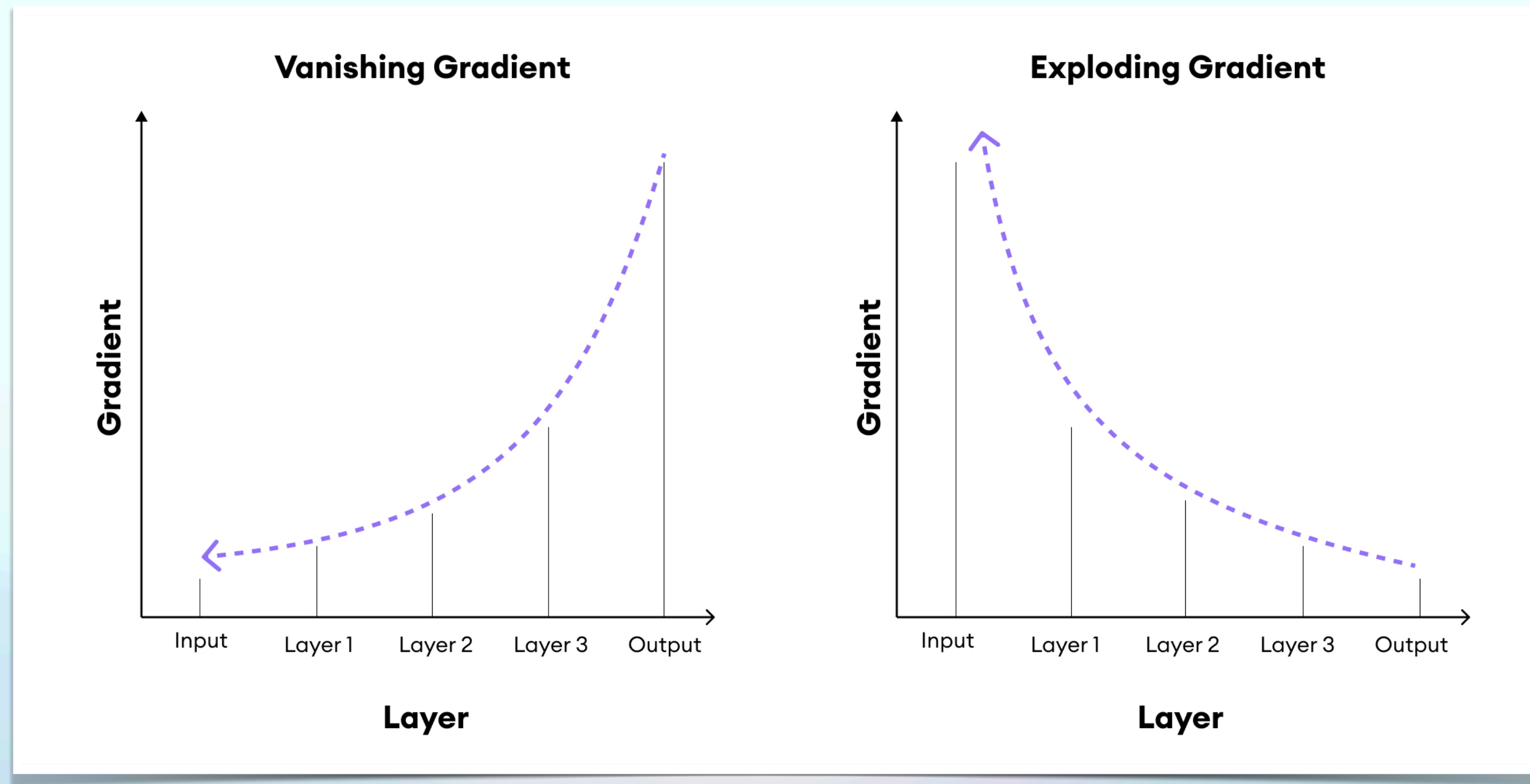


Amirkabir University of Technology  
(Tehran Polytechnic)

# The Vanishing/Exploding Gradient Problem

- **Issue:** As gradients backpropagate, they may shrink to near-zero (vanishing) or grow too large (exploding).
- Effect:
  - Vanishing gradients slow down learning or prevent deeper layers from learning.
  - Exploding gradients cause instability and divergence.

# The Vanishing/Exploding Gradient Problem



# The Vanishing/Exploding Gradient Problem

## Solutions to Vanishing/Exploding Gradients

1. Better Weight Initialization (Glorot, He)
2. Nonsaturating Activation Functions (Leaky ReLU, etc.)
3. Batch Normalization
4. Gradient Clipping

# The Vanishing/Exploding Gradient Problem

## Weight Initialization Techniques

### Glorot (Xavier) Initialization

- Formula: Uniform Distribution  $w \sim U\left(-\frac{\sqrt{6}}{\sqrt{n_{\text{in}} + n_{\text{out}}}}, \frac{\sqrt{6}}{\sqrt{n_{\text{in}} + n_{\text{out}}}}\right)$ , Normal Distribution  $w \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}} + n_{\text{out}}}\right)$
- Suitable for **sigmoid** and **tanh** activations.

### He Initialization

- Formula: Uniform Distribution  $w \sim U\left(-\frac{\sqrt{6}}{\sqrt{n_{\text{in}}}}, \frac{\sqrt{6}}{\sqrt{n_{\text{in}}}}\right)$ , Normal Distribution  $w \sim \mathcal{N}\left(0, \frac{2}{n_{\text{in}}}\right)$
- Best for **ReLU** and **Leaky ReLU** activations



# The Vanishing/Exploding Gradient Problem

## Nonsaturating Activation Functions

$$\text{Leaky ReLU (LReLU)} = \begin{cases} x, & \text{if } x > 0 \\ \alpha x, & \text{if } x \leq 0 \end{cases}$$

Randomized Leaky ReLU (RReLU) (randomizes  $\alpha$  during training)

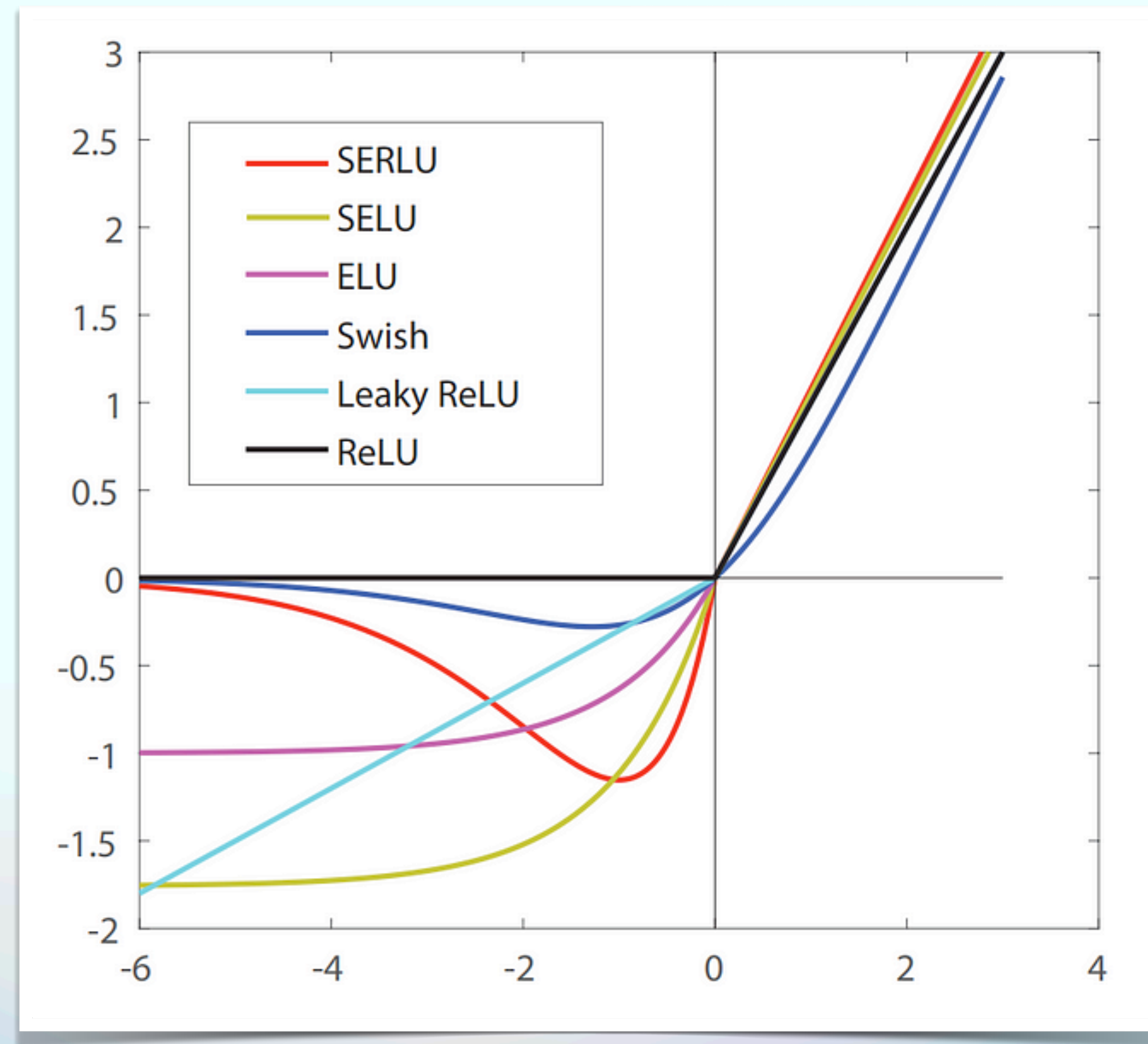
Parametric Leaky ReLU (PReLU) (learns the negative slope  $\alpha$  during training)

$$\text{Exponential Linear Unit (ELU)} = \begin{cases} x, & x > 0 \\ \alpha(\exp(x) - 1), & x \leq 0 \end{cases}$$

Scaled Exponential Linear Unit (SELU)

# The Vanishing/Exploding Gradient Problem

## Nonsaturating Activation Functions



# The Vanishing/Exploding Gradient Problem

## Batch Normalization

Normalizes activations across the mini-batch:

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}, \quad z_i = \gamma \hat{x}_i + \beta$$

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i, \quad \sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

- $x_i$  is the input activation of the neuron.
- $\mu_B$  is the mean of the mini-batch.
- $\sigma_B^2$  is the variance of the mini-batch.
- $\hat{x}_i$  is the normalized activation.
- $\epsilon$  is a small constant to prevent division by zero.
- $\gamma$  (scale) and  $\beta$  (shift) are trainable parameters.



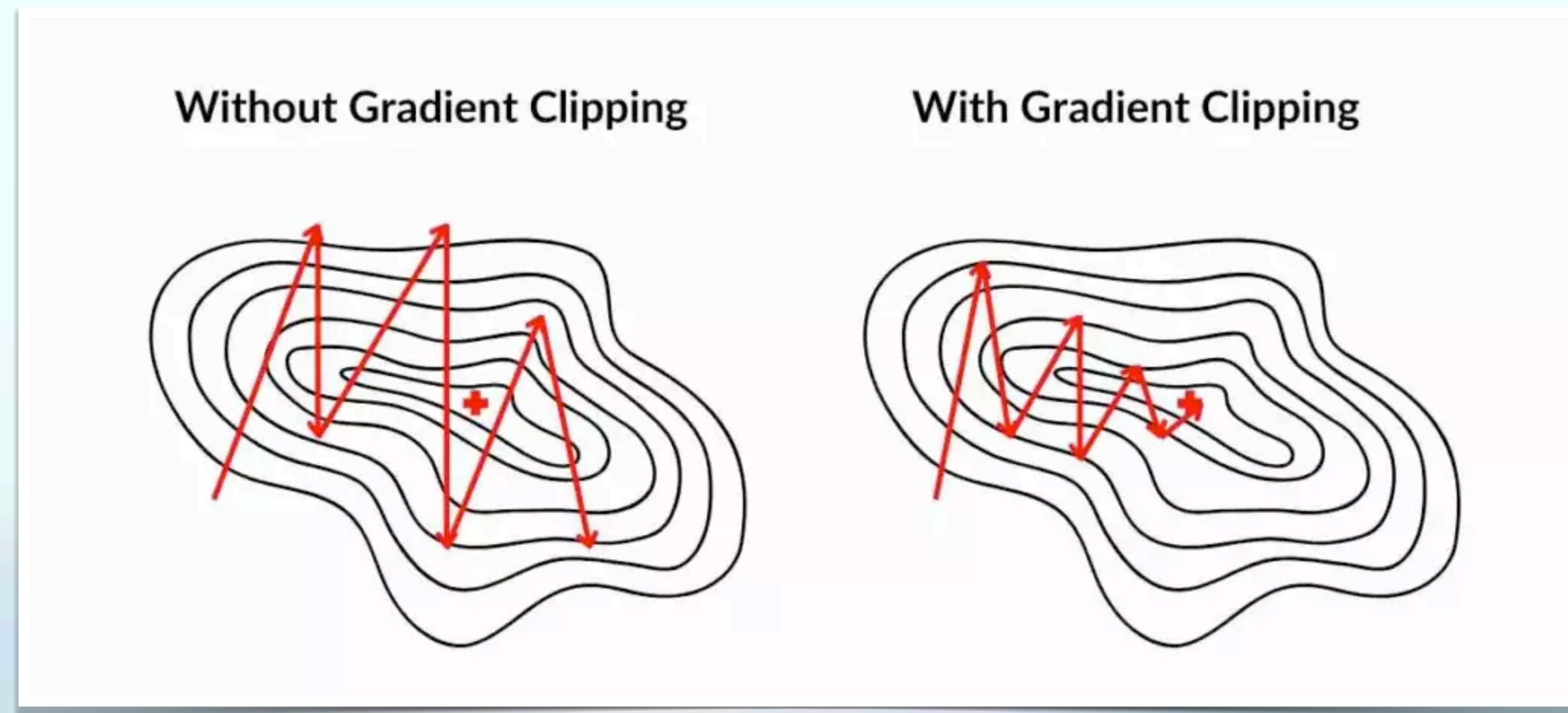
So during training, BN  
standardizes its inputs,  
then rescales and offsets  
them. Good!  
**What about at test time?**



# The Vanishing/Exploding Gradient Problem

## Gradient Clipping

Prevents exploding gradients by limiting their magnitude.



# Avoiding Overfitting – Regularization Techniques

## L1 and L2 regularization

L1 regularization: lasso regression

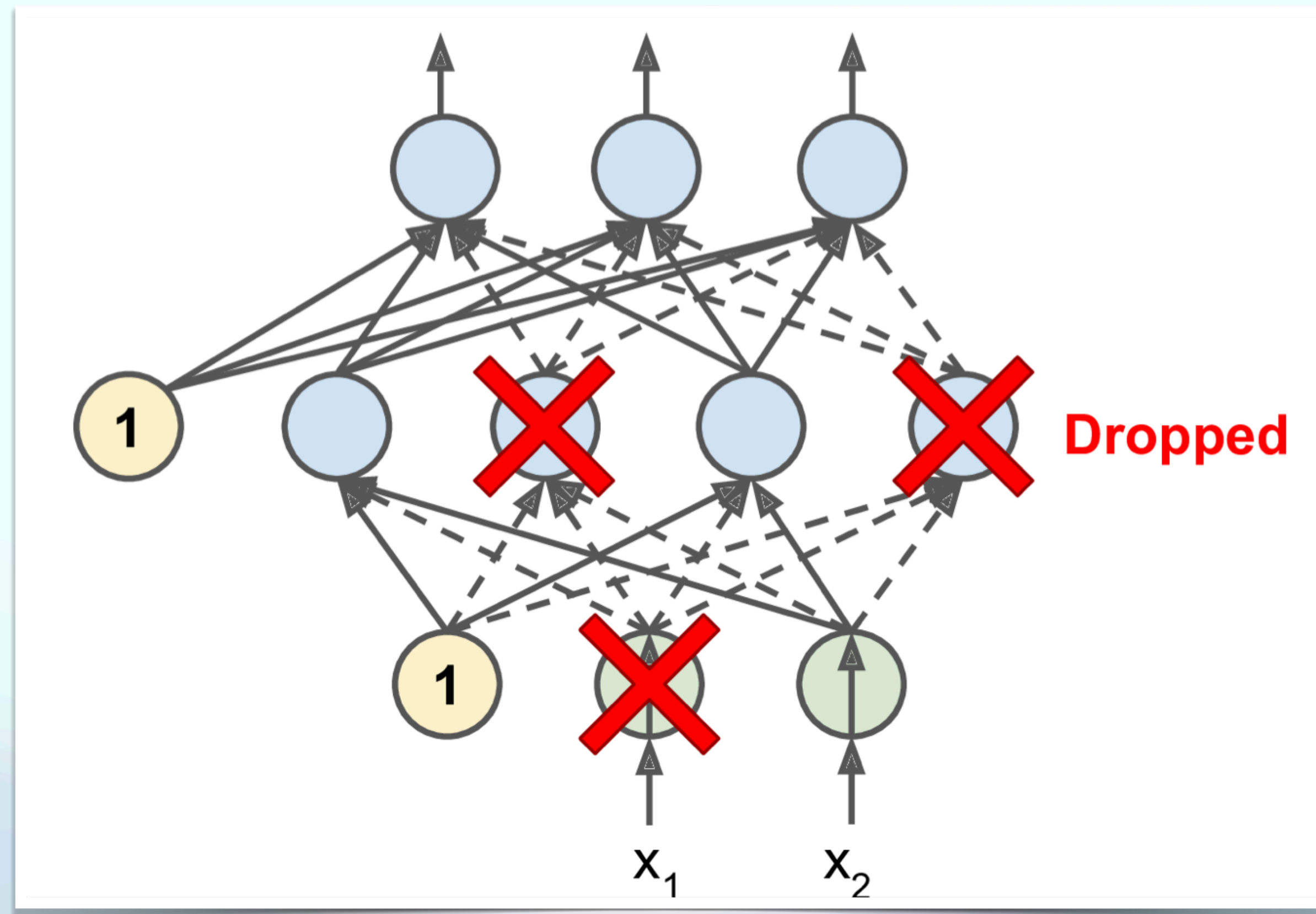
$$\text{New Cost} = \frac{1}{n} \sum_n^i y_i - \hat{y}_i + \lambda \sum_i |w_i|$$

L2 regularization: ridge regression

$$\text{New Cost} = \frac{1}{n} \sum_n^i y_i - \hat{y}_i + \lambda \sum_i w_i^2$$

# Avoiding Overfitting – Regularization Techniques

## Dropout



# What's Coming Next

- Training your model with faster optimizers