

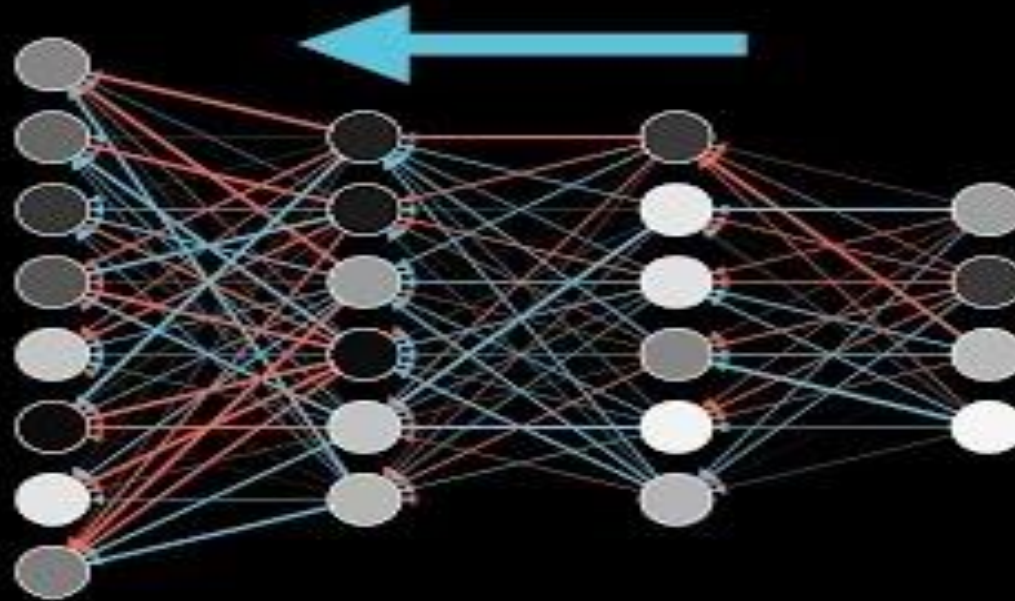
MLP and BP

C. V. Jawahar

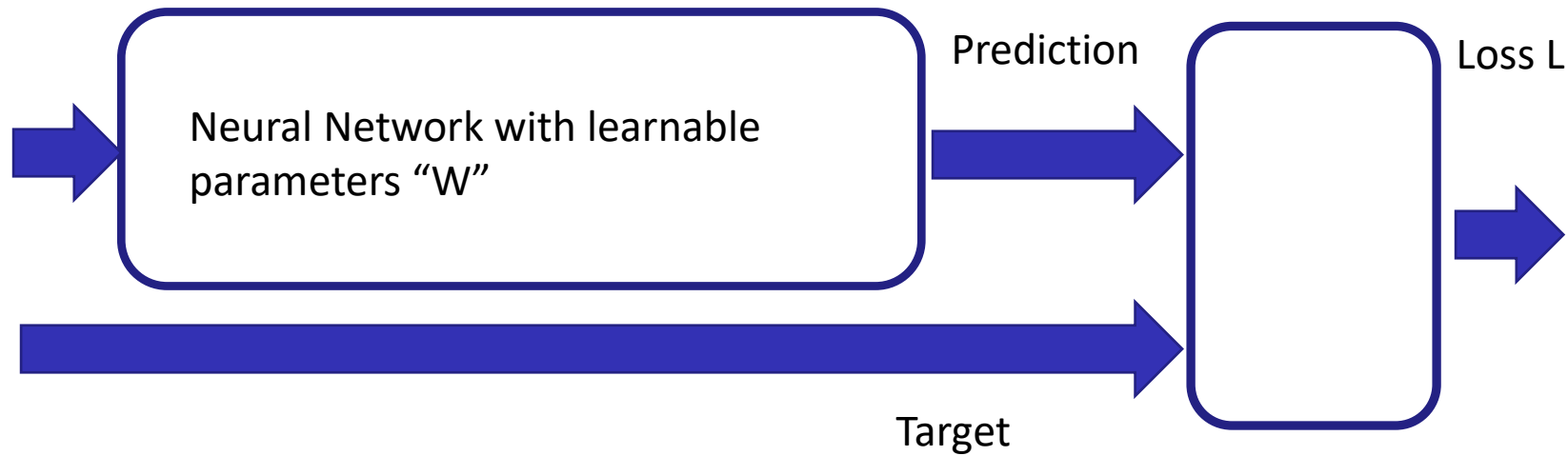
IIIT Hyderabad

28 Mar, 2025

Backpropagation



Basic Problem



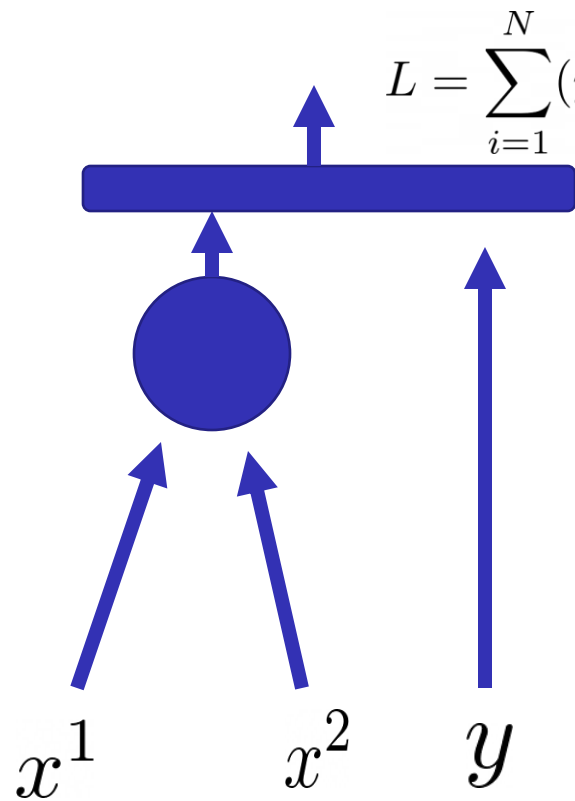
Goal: Minimize Loss L by adjusting/updating the weights W of the neural network.

Motivating Analogy



Many many Knobs (~Millions)
You can increase or decrease each one.
Goal: Control the intensity of the bulb

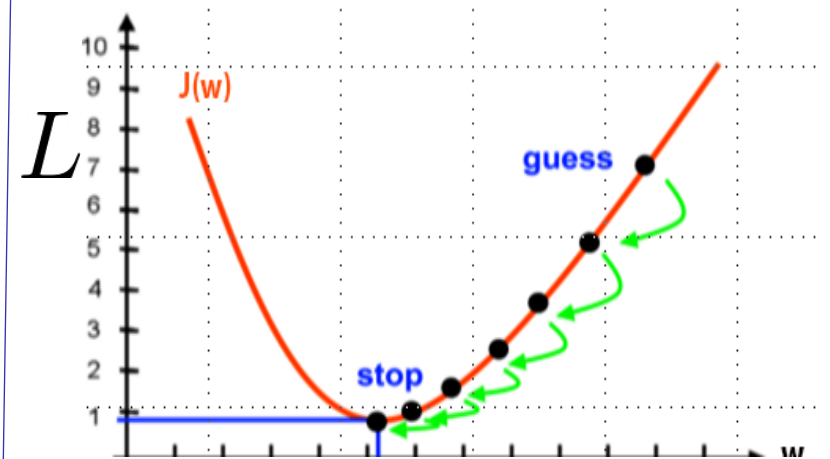
Gradient Descent



$$L = \sum_{i=1}^N (y_i - o_i)^2 = \sum_{i=1}^N (y_i - (w_1 x_i^1 + w_2 x_i^2))^2$$

$$\frac{\partial L}{\partial w_1} = \sum_{i=1}^N 2(y_i - o_i)(-1)(x_i^1)$$

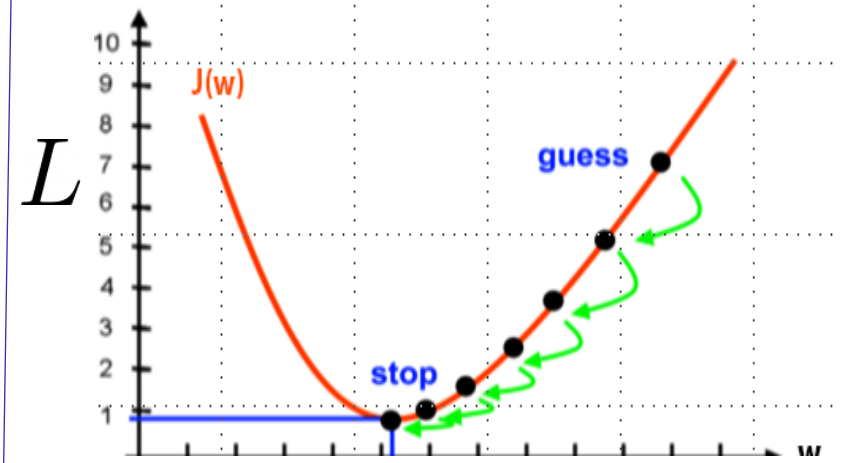
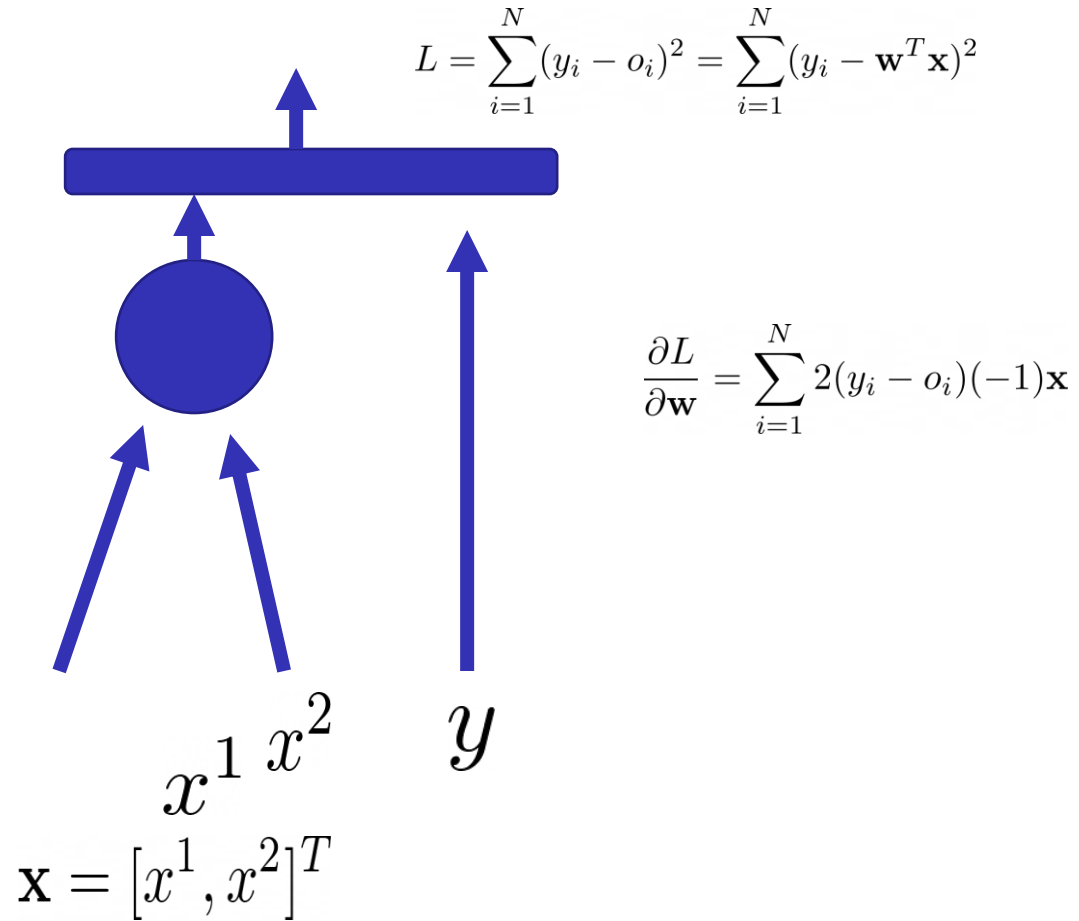
$$\frac{\partial L}{\partial w_2} = \sum_{i=1}^N 2(y_i - o_i)(-1)(x_i^2)$$



$$w_1^{n+1} = w_1^n - \eta \frac{\partial L}{\partial w_1}$$

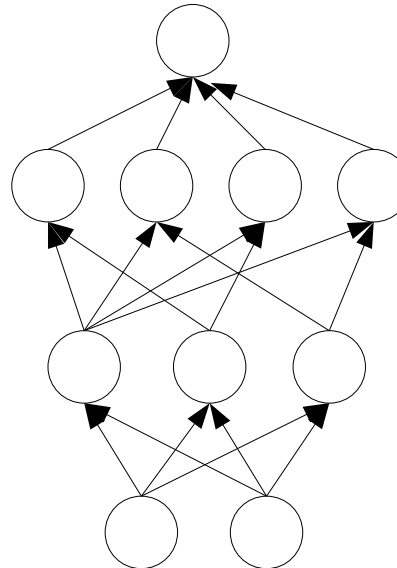
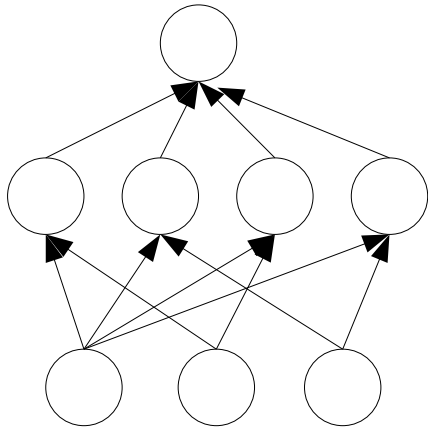
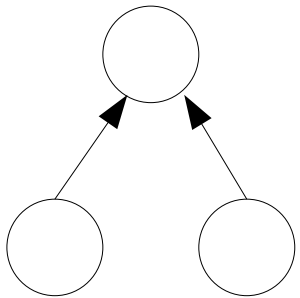
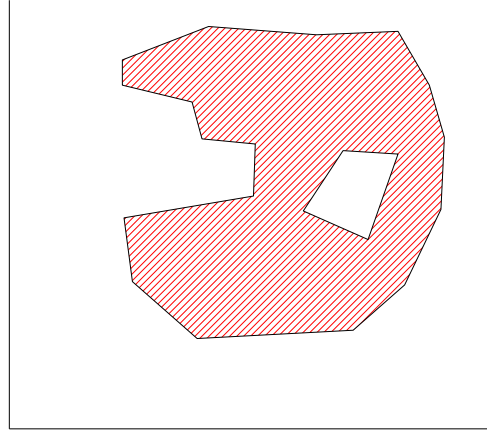
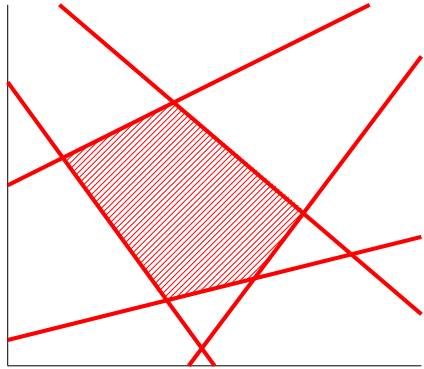
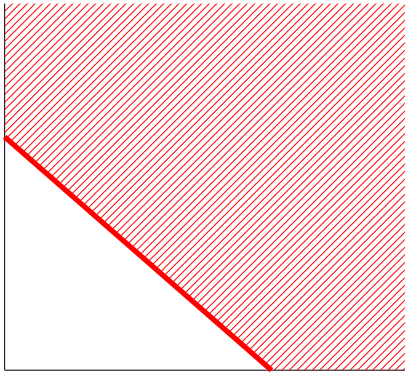
$$w_2^{n+1} = w_2^n - \eta \frac{\partial L}{\partial w_2}$$

Gradient Descent: Simple Vector Notation

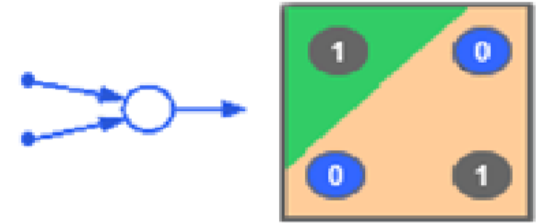


$$\mathbf{w}^{n+1} = \mathbf{w}^n - \eta \frac{\partial L}{\partial \mathbf{w}}$$

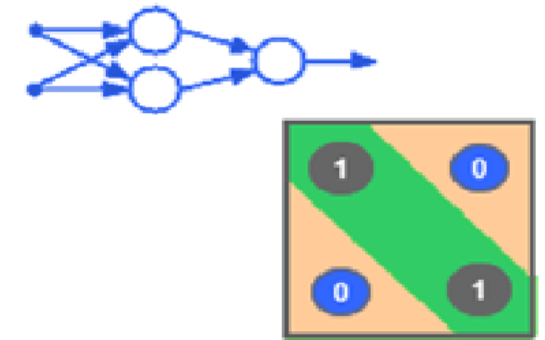
Deeper Networks



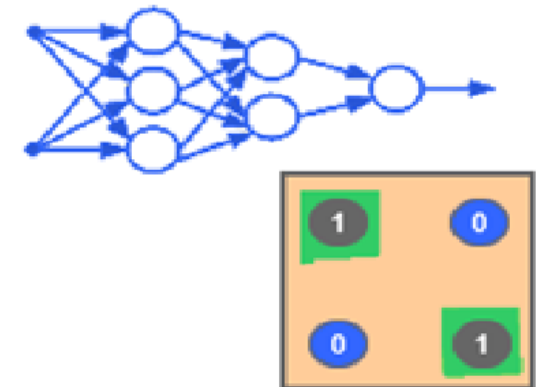
1 layer: semiplane



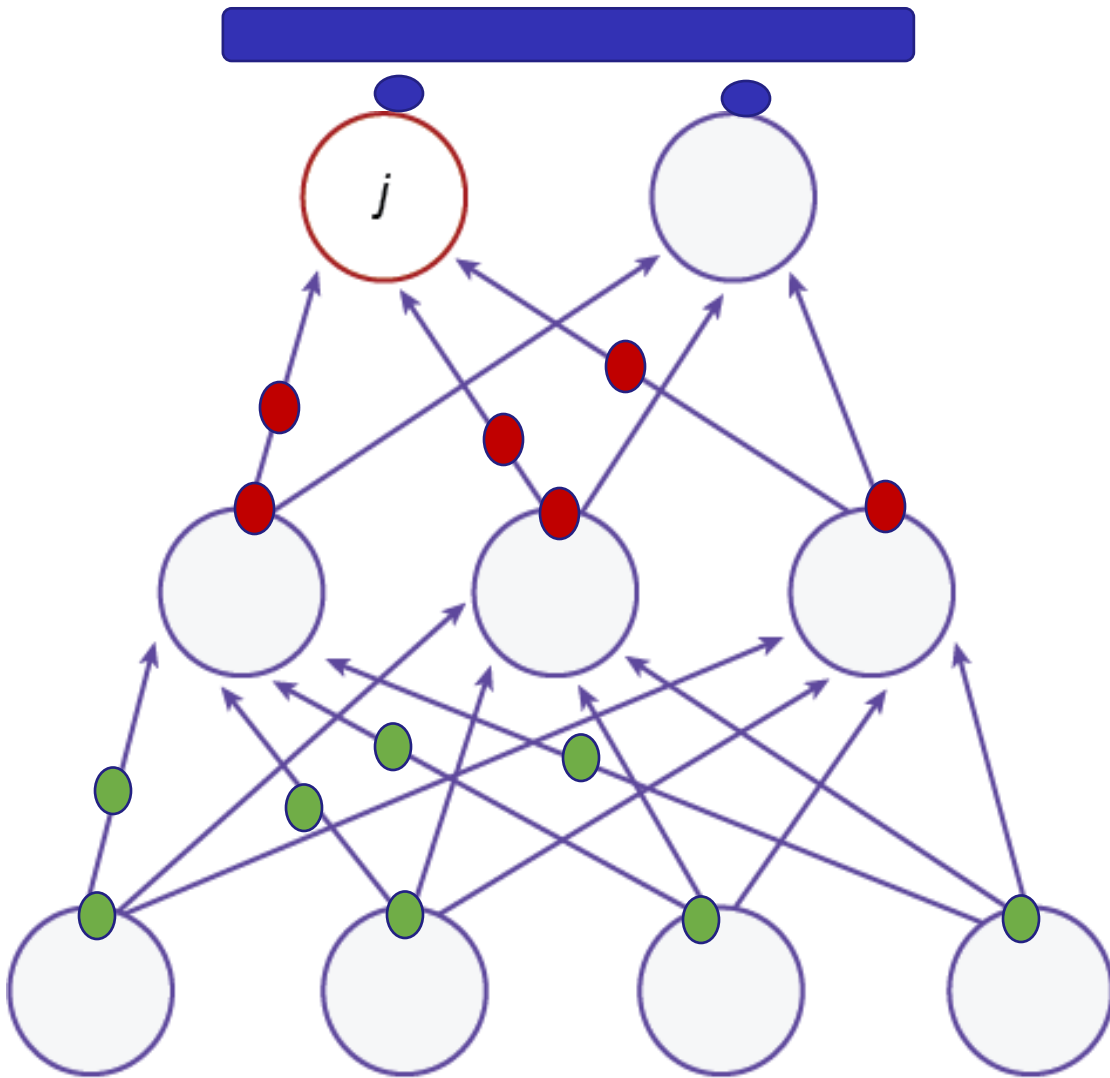
2 layers: convex regions



3 layers: arbitrary regions



Error depends on many weights



Compute Loss/Error; Loss Layer

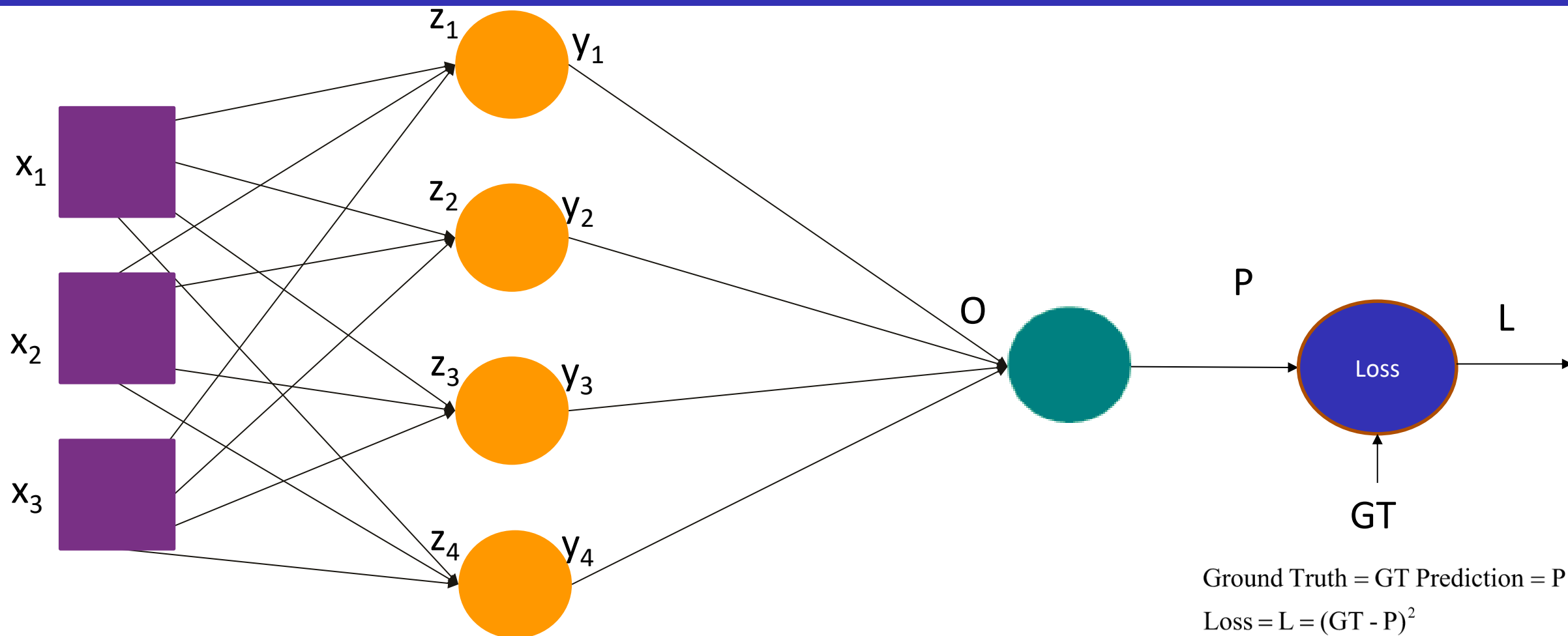
Error depends on two terms/outputs

Each one depends on three weights and three intermediate values

Each one of these intermediate values depends on four weights and four inputs

To minimize error/loss, we need to adjust/vary all the 18 weights.!!

Back Propagation



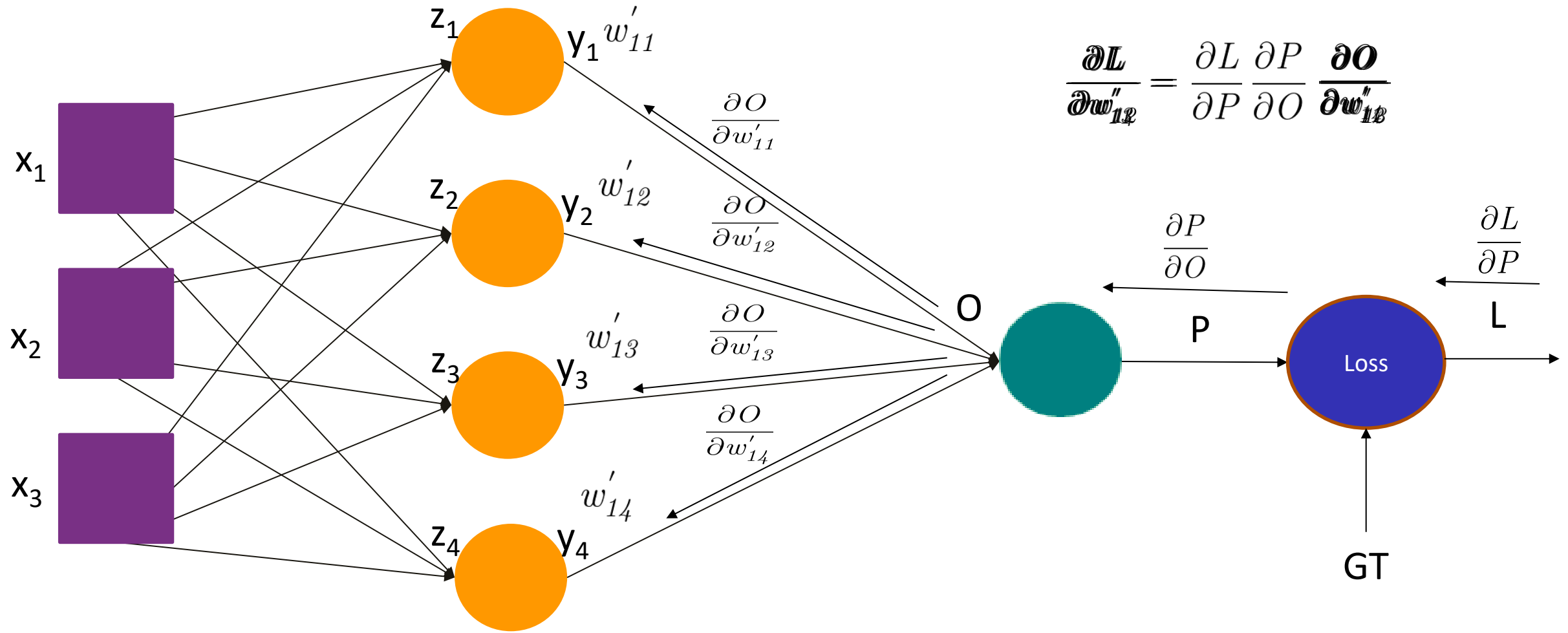
$$\begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} \phi(z_1) \\ \phi(z_2) \\ \phi(z_3) \\ \phi(z_4) \end{bmatrix}$$

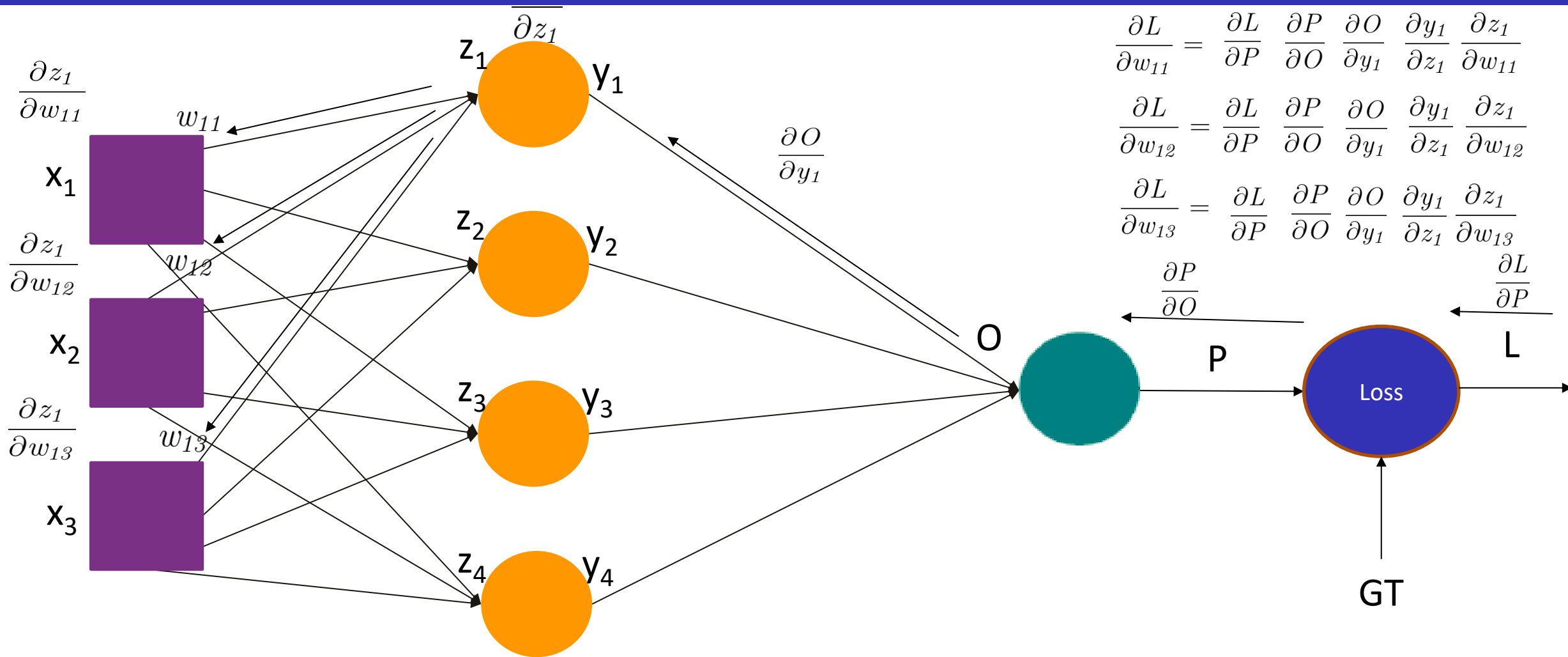
$$\mathbf{O} = \begin{bmatrix} w'_{11} & w'_{12} & w'_{13} & w'_{14} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

$$\mathbf{P} = \phi(\mathbf{O})$$

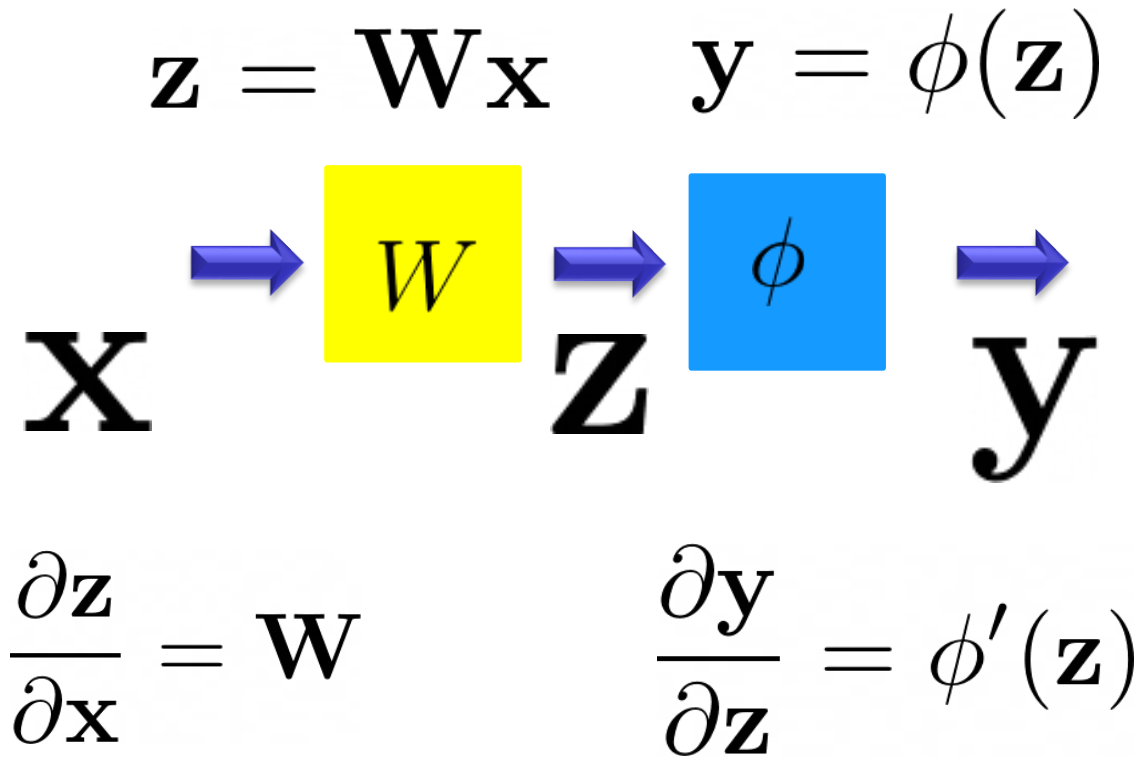
Back Propagation



Back Propagation



Two Typical Blocks (with and without weights)

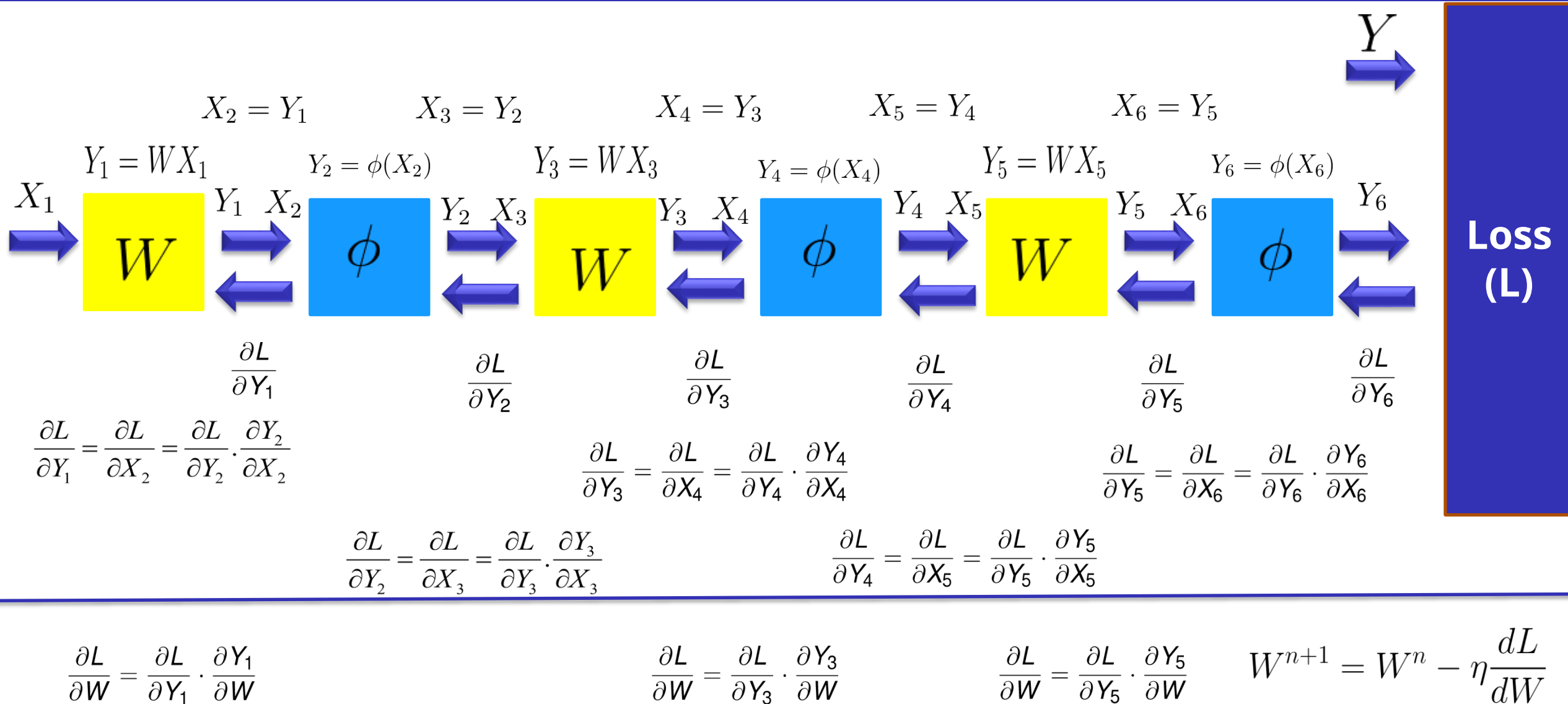


Eg. Sigmoid

$$y = \phi(z) = \frac{1}{1 + e^{-z}}$$

$$\phi'(z) = \phi(z) \cdot (1 - \phi(z))$$

Back Propagation:



Basic Algorithm

Forward Pass

Take a batch of samples; Compute outputs; Compute Loss.

Backward Pass

Update all the weights (backwards) and reduce the errors.

Repeat the steps; Until?

Improvements

Mini Batch,
Stochastic, SGD

Activation Fns,
Learning Rate.

Better
Initialization
(Xavier, He)

**Design or
Initial.**

Momentum
(Nesterov)

Numerical/Parallelism
(including Vanishing
Exploding Gradients)

Better Update
(Adam, Adagrad,
RMSProp etc.)

Update

Early Stopping

Regularization
(L1, L2)

Regularization
(Dropout, Data
Augmentation etc.)

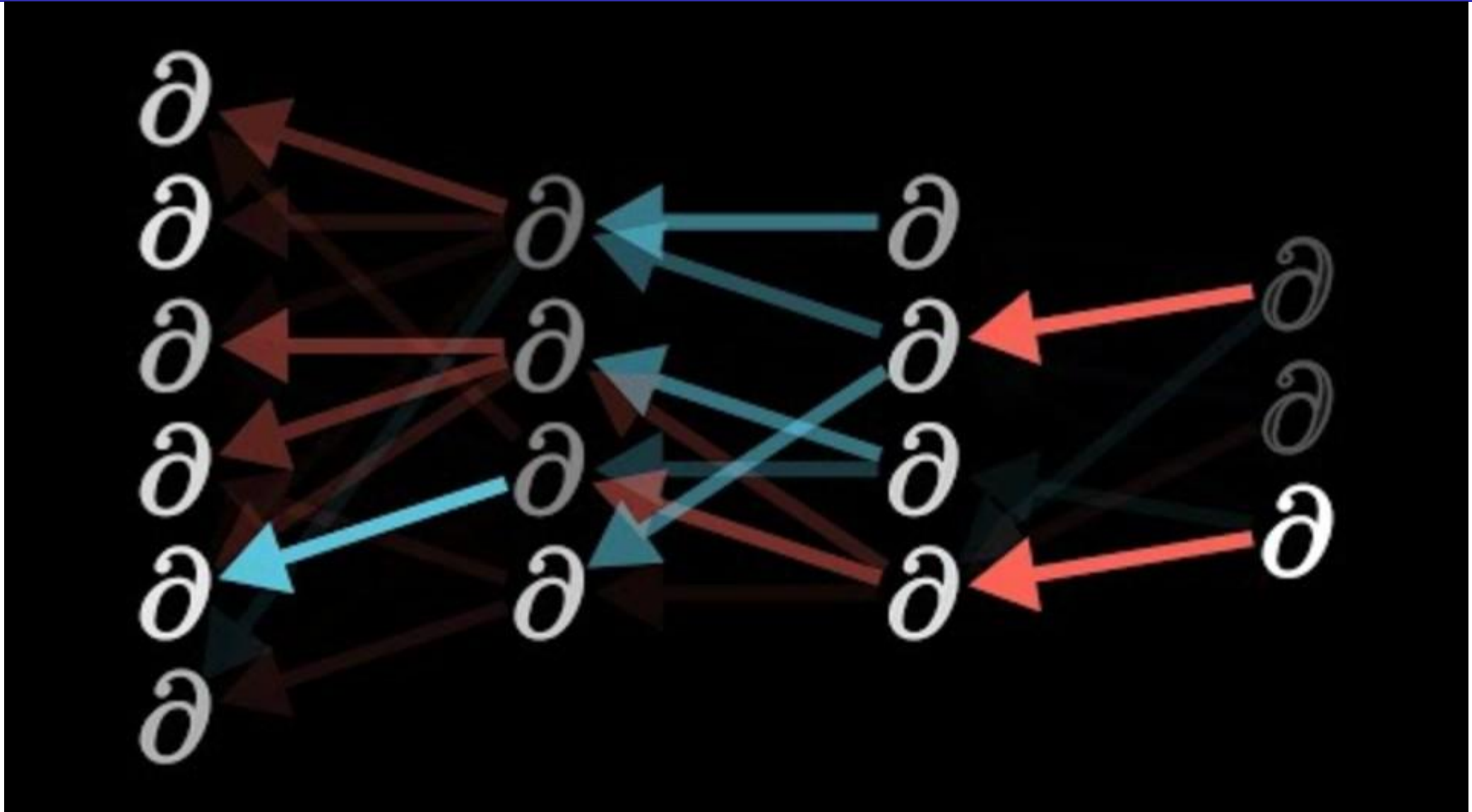
Regularize

Basic; Good to Discuss

Worth; Insightful

Advanced

Video



Link to the video: <https://www.youtube.com/watch?v=tIeHLnjs5U8>

Questions?
