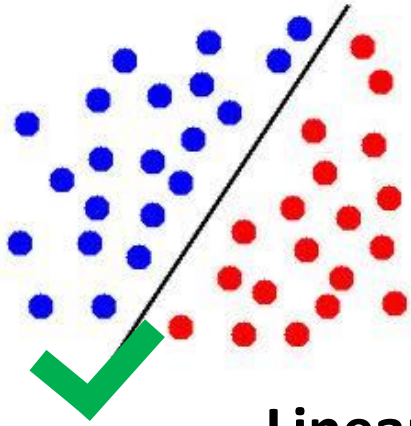
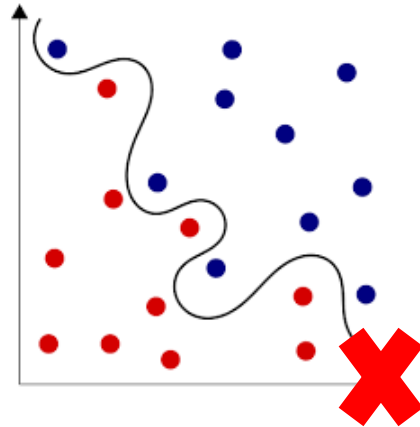
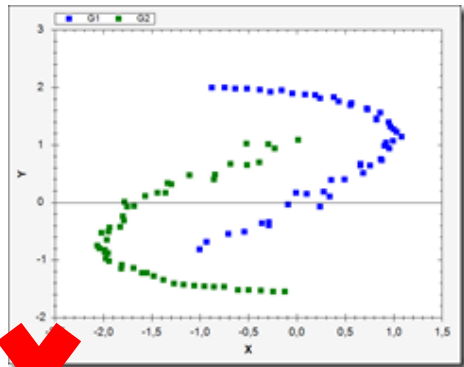
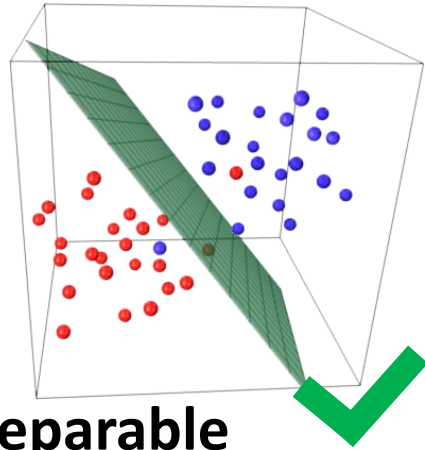


# Basics of Deep Learning

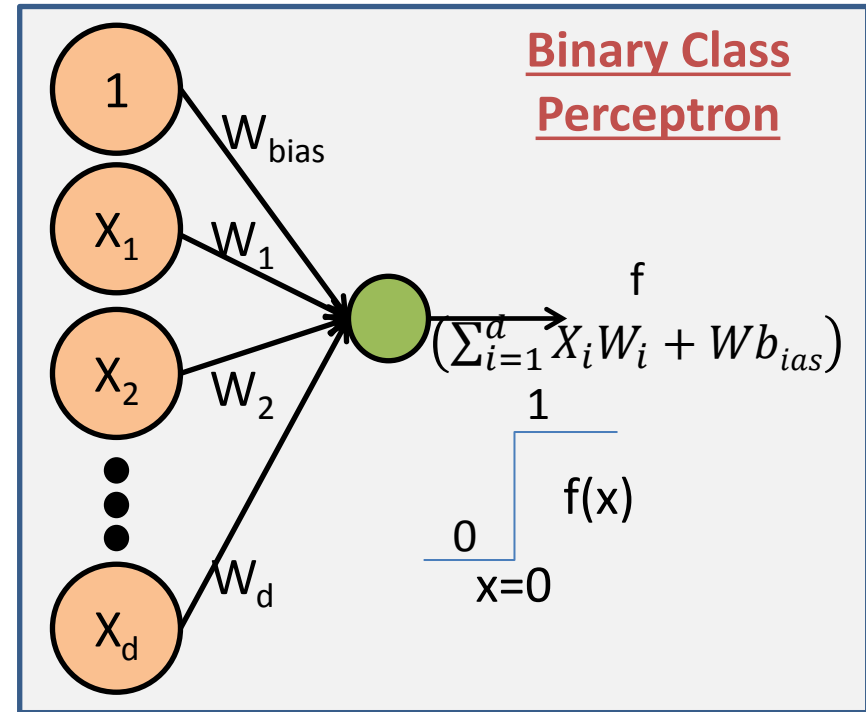
# Single Layer Perceptron



Linear Separable



Non-Linear Separable



Single Layer Perceptron (Binary Class)

# Learning Weights

$$X(k) = [ 1 \ X_1^k \ X_2^k \ X_3^k \ \dots \ X_d^k ]$$
$$W = [ W_{bias} \ W_1 \ W_2 \ W_3 \ \dots \ W_d ]$$

$$Y(k) = 1 \text{ if } \left( \sum_{i=1}^d X_i W_i + W_{bias} \right) > 0$$
$$= 0, \text{ otherwise}$$

Or

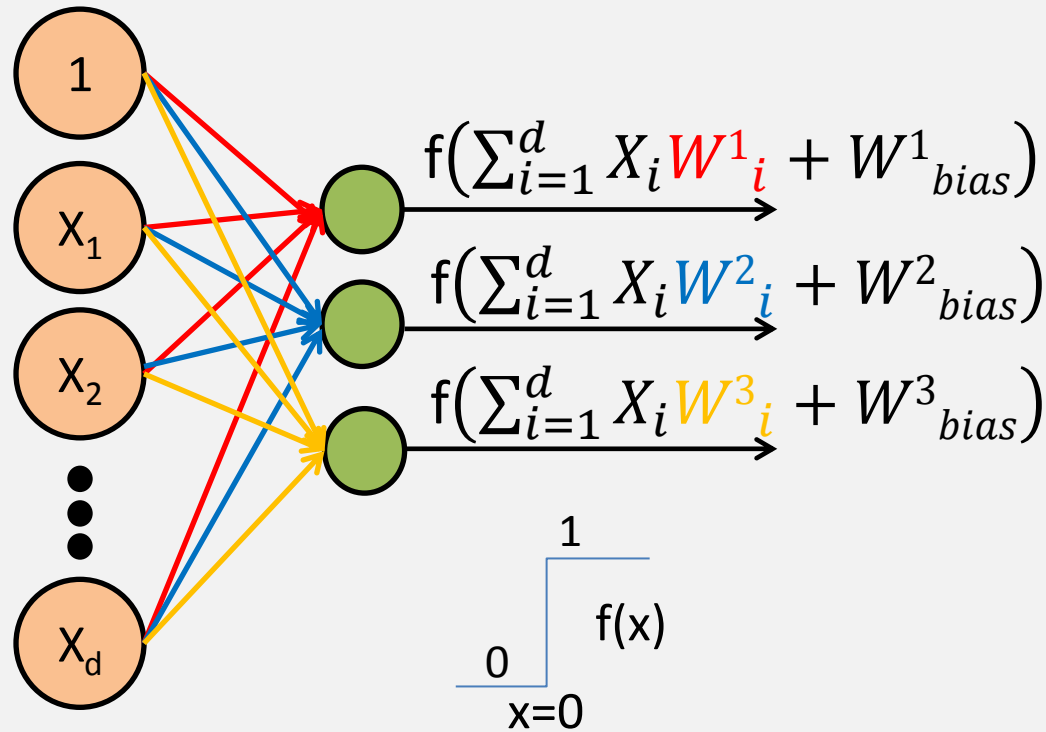
$$Y(k) = 1 \text{ if } X(k) * W(k)^T > 0$$
$$= 0, \text{ otherwise}$$

Intuitively,

$$W(k+1) = W(k) + X(k) \text{ if } (X(k) * W(k)^T < 0) \text{ and } Y(k) = 1$$
$$= W(k) - X(k) \text{ if } (X(k) * W(k)^T > 0) \text{ and } Y(k) = 0$$

Single Layer Perceptron (Learning Weights)

## Multi-Class Perceptron



Single Layer Perceptron (Multi Class)

# Linear Least Square Method

$$J(w) = \frac{1}{2} \sum_{i=1}^n (x_i^T w - y_i)^2 \quad - \quad x_i = \begin{bmatrix} 1 \\ x_1^i \\ x_2^i \\ \vdots \\ x_d^i \end{bmatrix}_{(d+1) \times 1} \quad w = \begin{bmatrix} w_{\text{bias}} \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}_{(d+1) \times 1}$$

$$\text{Let } A = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}_{n \times (d+1)}$$

$$\therefore J(w) = \frac{1}{2} (Aw - y)^T (Aw - y)$$

$$\nabla J(w) = A^T (Aw - y) \Big|_{w=w^*} = 0$$

$$\Rightarrow A^T A w^* = A^T y$$

$$\Rightarrow w^* = (A^T A)^{-1} A^T y$$

$$J(w) = \frac{1}{2} (Aw - y)^T (Aw - y) + \frac{\lambda}{2} w^T w$$

$$\nabla J(w) = A^T (Aw - y) + \lambda w \Big|_{w=w^*} = 0$$

$$\Rightarrow w^* = (A^T A + \lambda I)^{-1} A^T y$$

# Gradient Descent

$$\text{Cost Function, } J(W) = \frac{1}{2} \sum_{k=1}^n (X^k * W(k)^T - y^k)^2$$

$$\text{Gradient, } \nabla J(W) = \sum_{k=1}^n X^k (X^k W(k)^T - y^k)$$

$$W(k+1) = W(k) - \eta \nabla J(W)$$

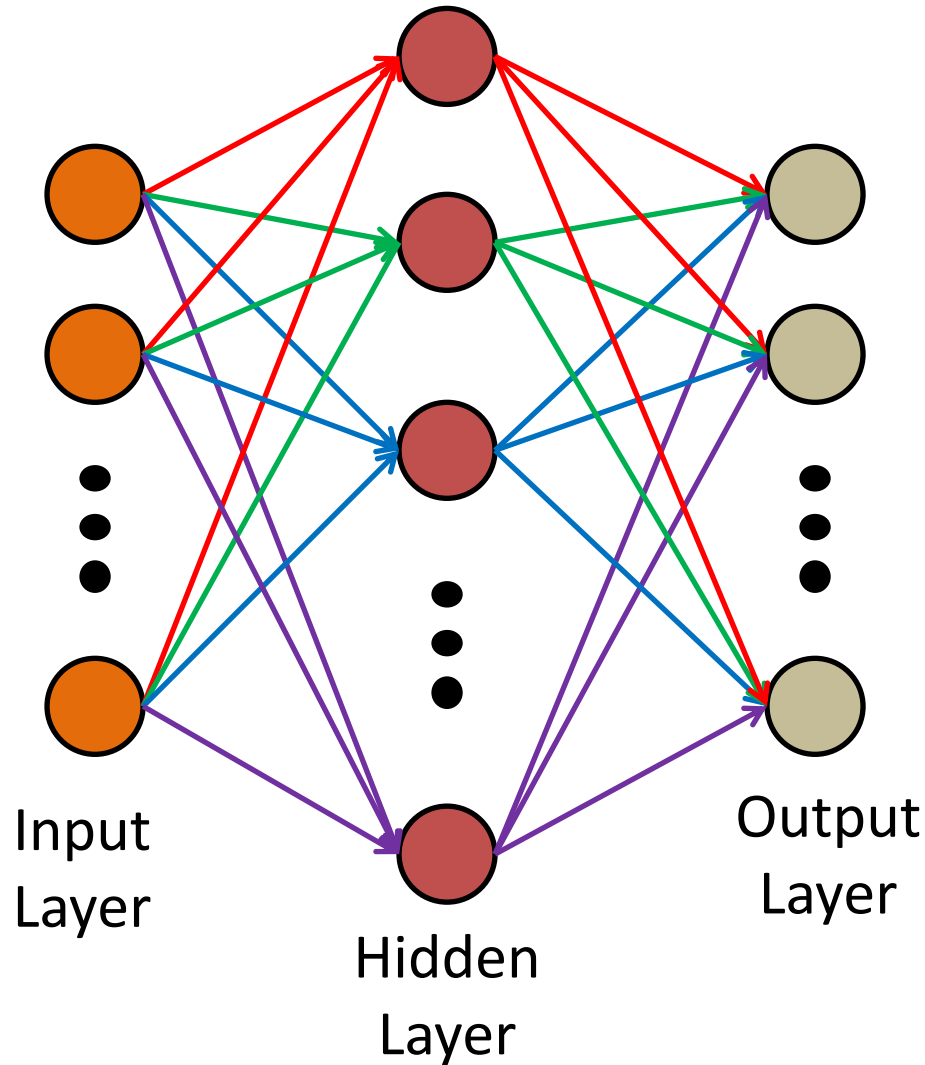
Batch Learning:

$$W(k+1) = W(k) - \eta \sum_{k=1}^n X^k (X^k W(k)^T - y^k)$$

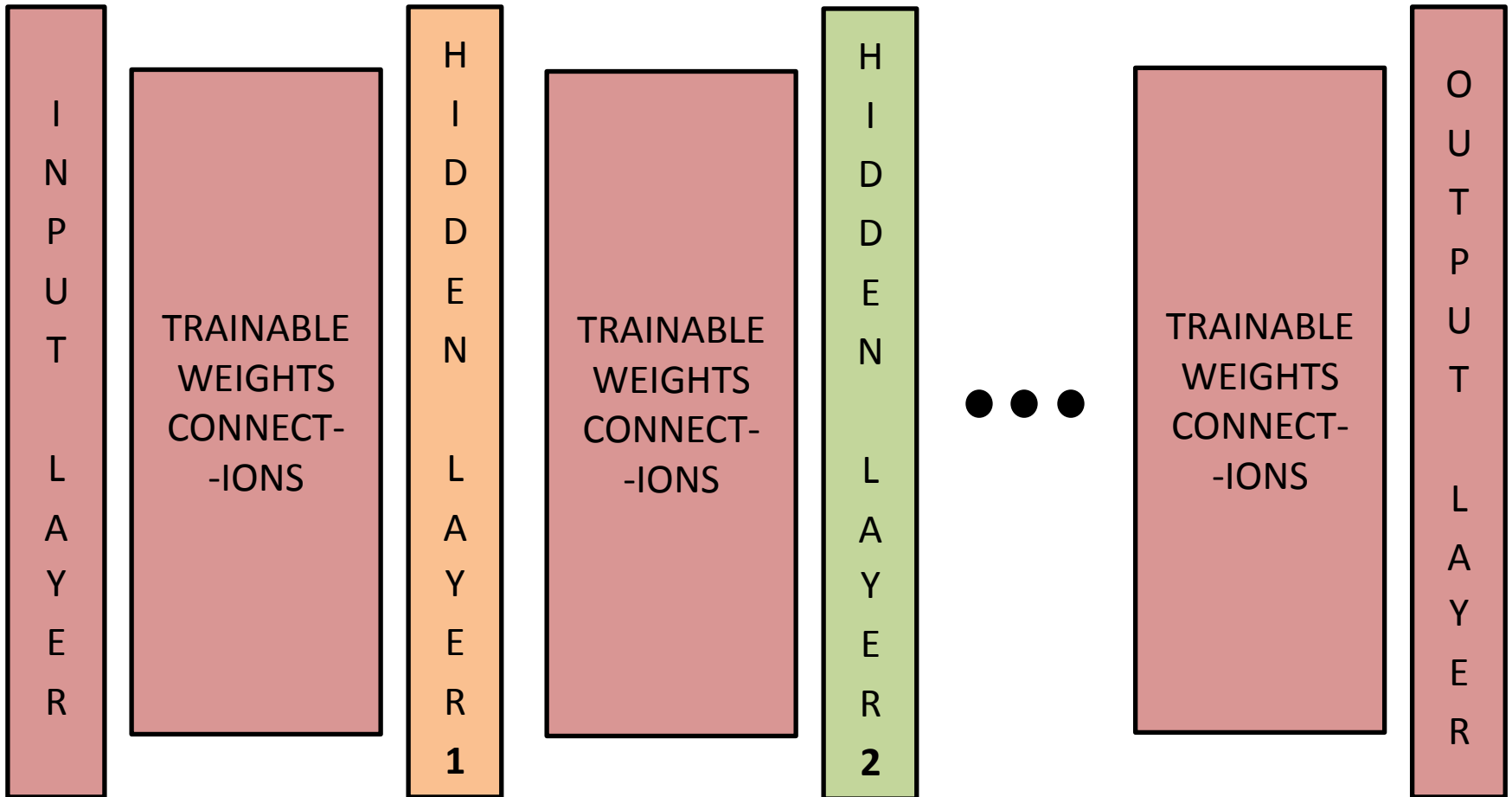
Stochastic/Online/Incremental Learning:

$$W(k+1) = W(k) - \eta \sum_{k=1}^n X^k (X^k W(k)^T - y^k)$$

# Single Layer NN



# Deep Neural Network

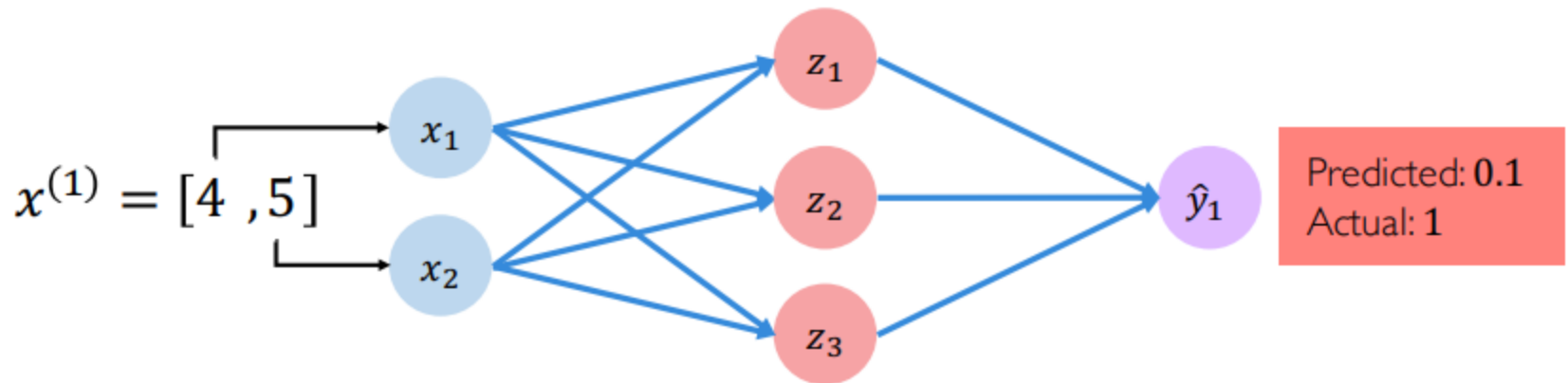


Neural Networks (DNN)



# Quantifying Loss

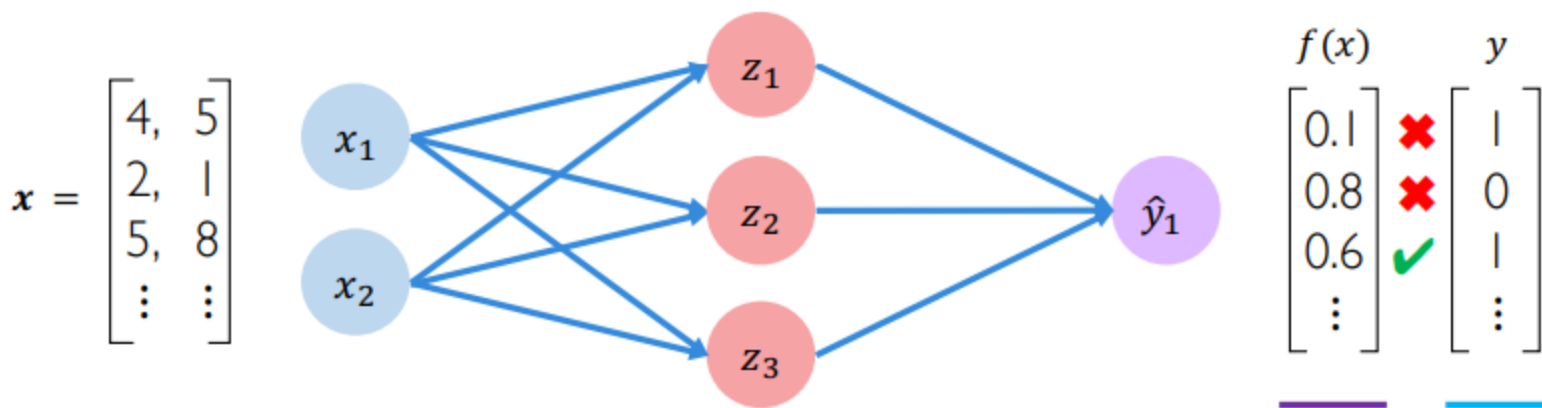
The **loss** of our network measures the cost incurred from incorrect predictions



$$\mathcal{L}(\underbrace{f(x^{(i)}; \theta)}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$

# Empirical Loss

The **empirical loss** measures the total loss over our entire dataset



Also known as:

- Objective function
- Cost function
- Empirical Risk

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathcal{L}(f(x^{(i)}; \theta))}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}}$$

<http://introtodeeplearning.com/index.html>

# Loss Functions

Binary Cross Entropy Loss:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)}}_{\text{Actual}} \log \left( \underbrace{f(x^{(i)}; \theta)}_{\text{Predicted}} \right) + (1 - \underbrace{y^{(i)}}_{\text{Actual}}) \log \left( 1 - \underbrace{f(x^{(i)}; \theta)}_{\text{Predicted}} \right)$$

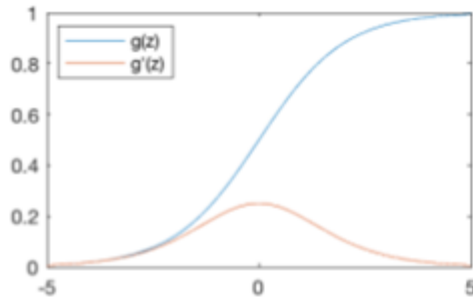
Mean Square Error Loss:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \left( \underbrace{y^{(i)}}_{\text{Actual}} - \underbrace{f(x^{(i)}; \theta)}_{\text{Predicted}} \right)^2$$

<http://introtodeeplearning.com/index.html>

# Common Activation Functions

Sigmoid Function

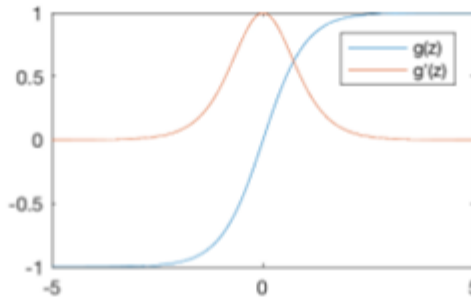


$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

<http://introtodeeplearning.com/index.html>

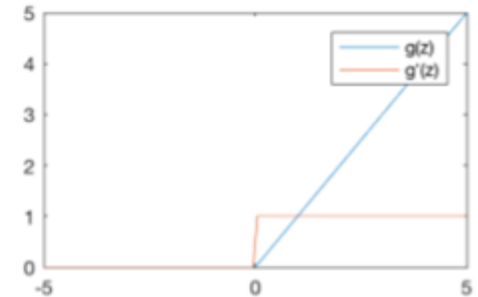
Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

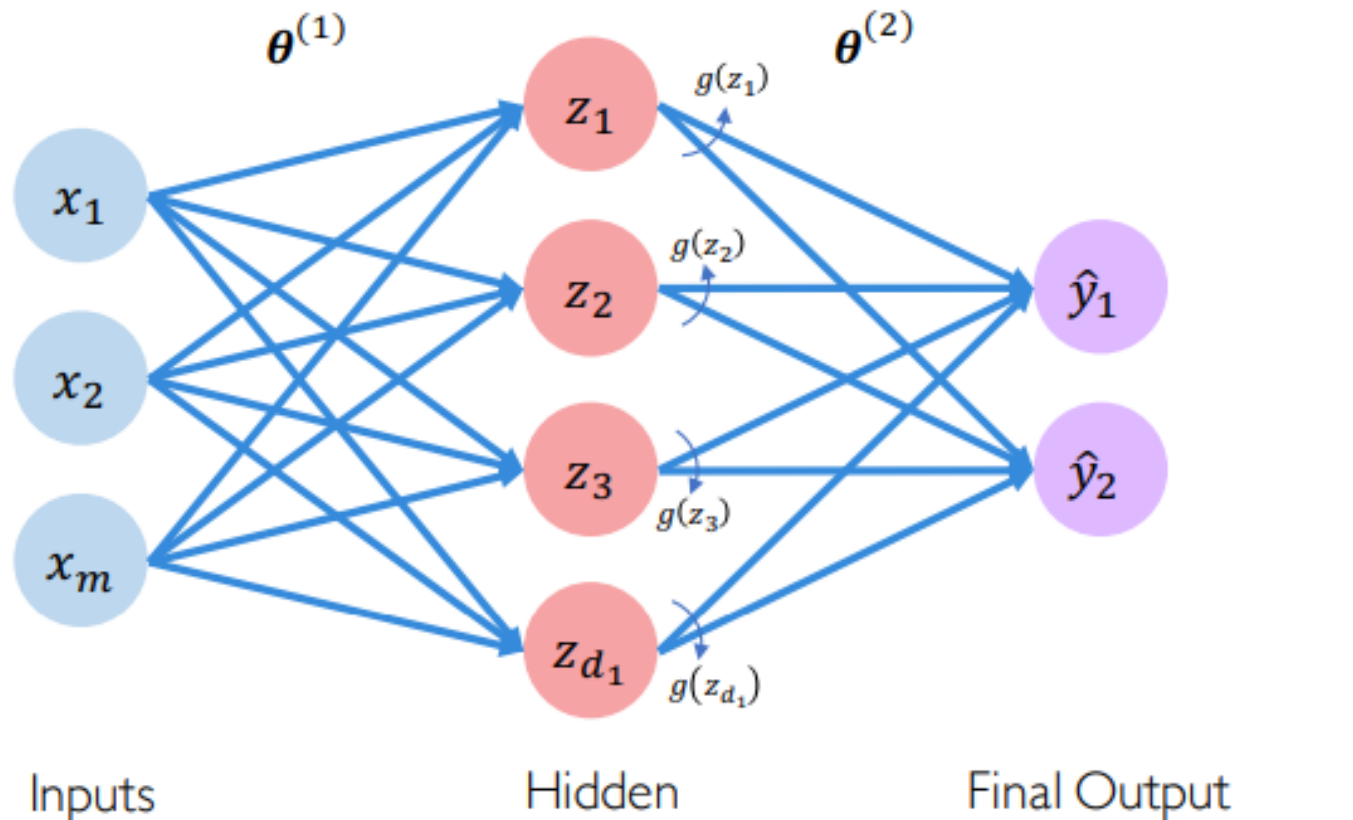
Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

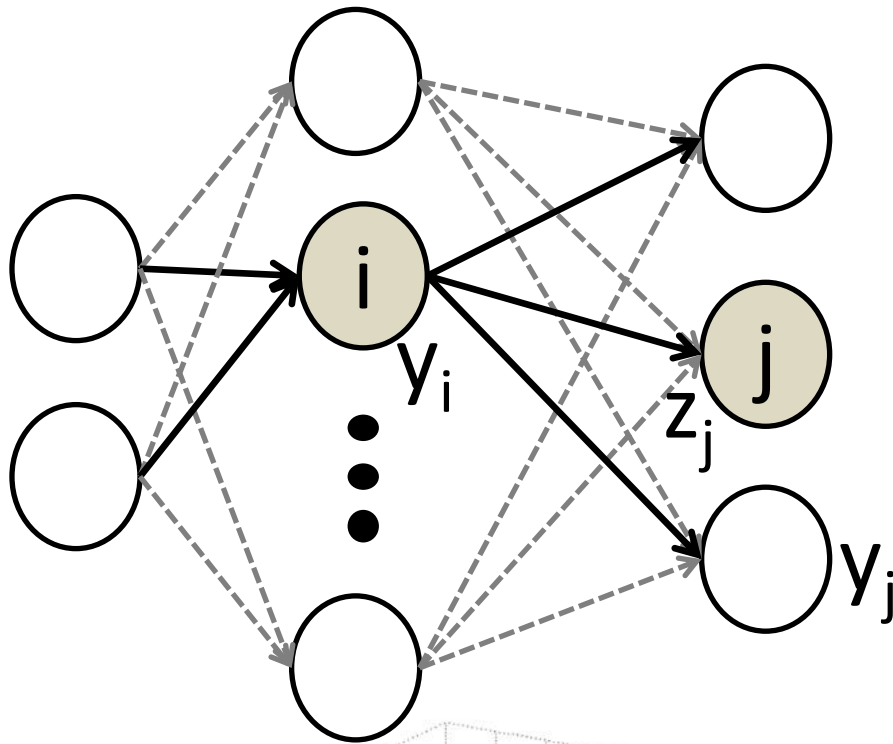
# Forward Propagation



$$z_i = \theta_{0,i}^{(1)} + \sum_{j=1}^m x_j \theta_{j,i}^{(1)} \quad \hat{y}_i = \theta_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j \theta_{j,i}^{(2)}$$

<http://introtodeeplearning.com/index.html>

# Back Propagation

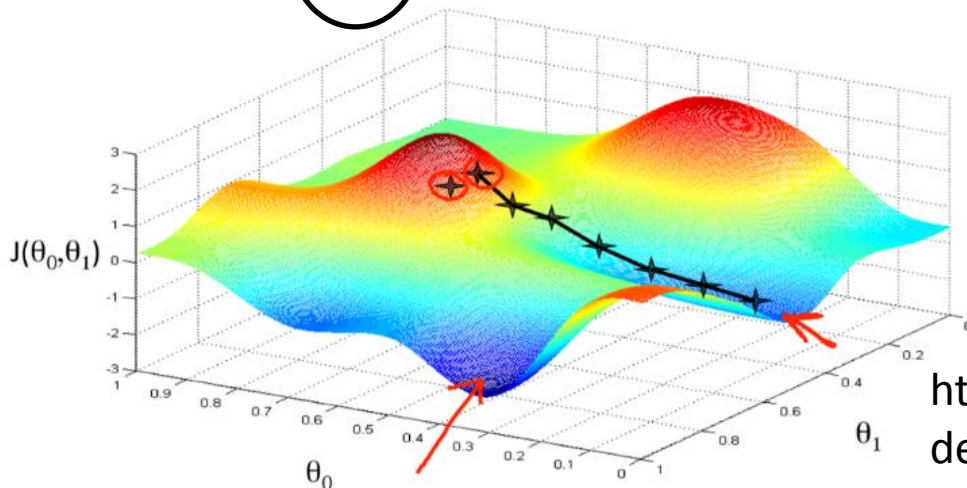


$$\frac{\partial E}{\partial z_j} = \frac{dy_j}{dz_j} \frac{\partial E}{\partial y_j} = y_j (1 - y_j) \frac{\partial E}{\partial y_j}$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial z_j}{\partial w_{ij}} \frac{\partial E}{\partial z_j} = y_i \frac{\partial E}{\partial z_j}$$

$$w_{ij}^l(t+1) = w_{ij}^l(t) - \lambda \frac{\partial E(W(t))}{\partial w_{ij}^l}$$

$$\frac{\partial E}{\partial y_i} = \sum_j \frac{dz_j}{dy_i} \frac{\partial E}{\partial z_j} = \sum_j w_{ij} \frac{\partial E}{\partial z_j}$$

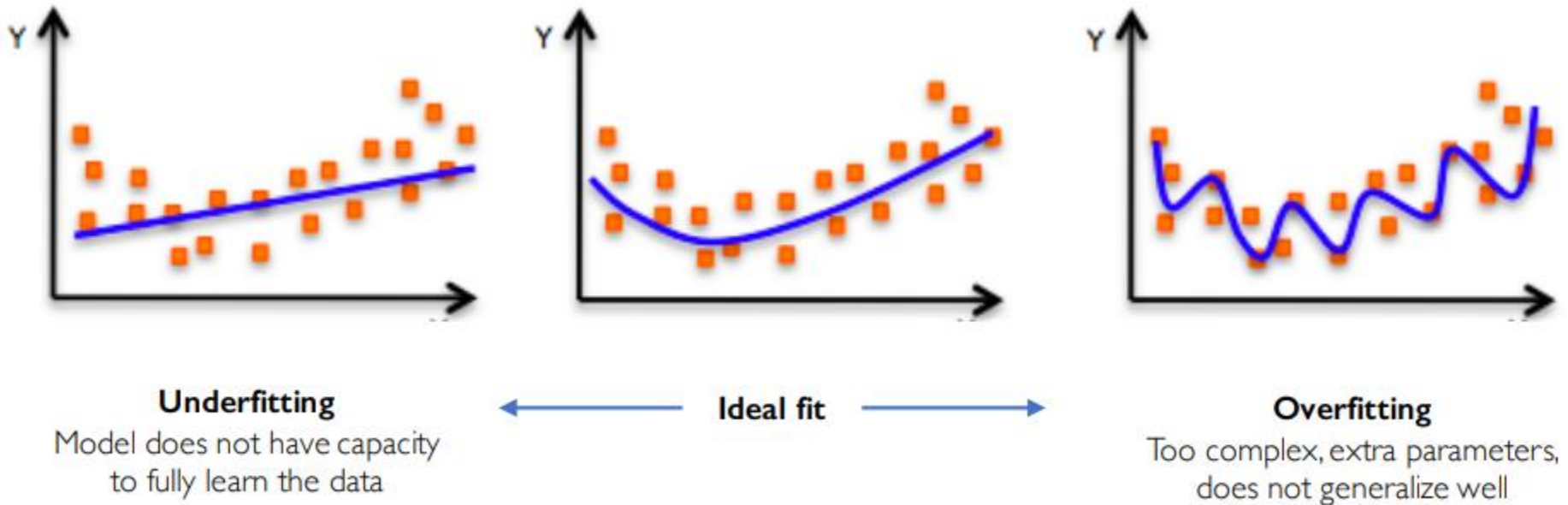


<https://mc.ai/a-dummies-guide-to-gradient-descent-and-backpropagation/>

Single Layer Neural Networks (Back Propagation)

# Regularization

## Overfitting Problem



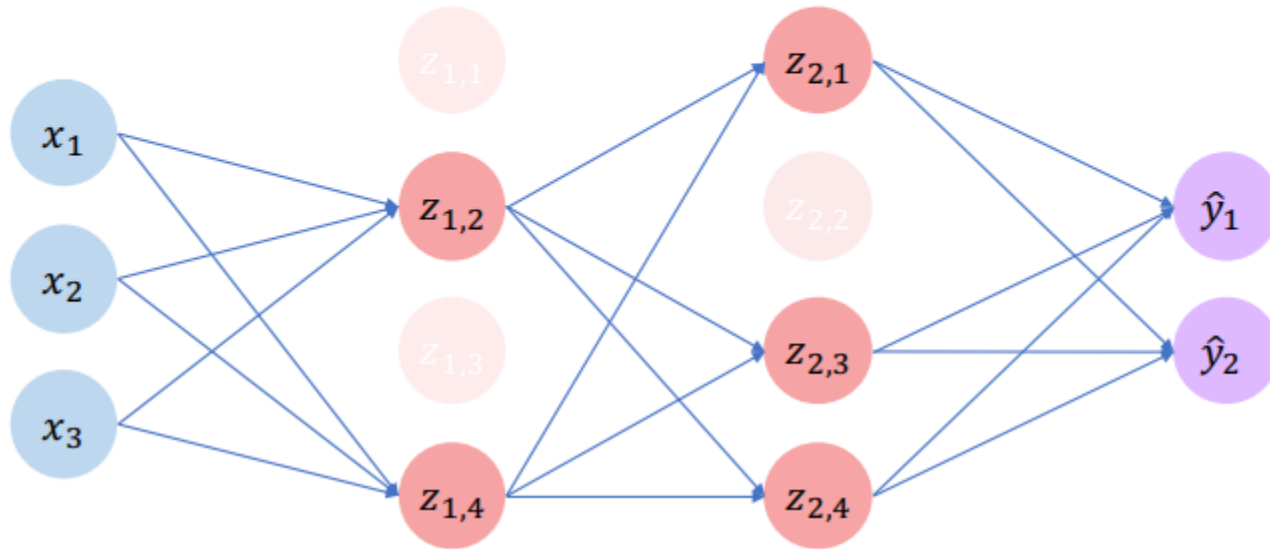
<http://introtodeeplearning.com/index.html>

**Solution** —————→ **Regularization**

Helps to generalize  
our network

# Regularization: Dropout

- During training, randomly set some activations to 0
  - Typically 'drop' 50% of activations in layer
  - Forces network to not rely on any 1 node



<http://introtodeeplearning.com/index.html>



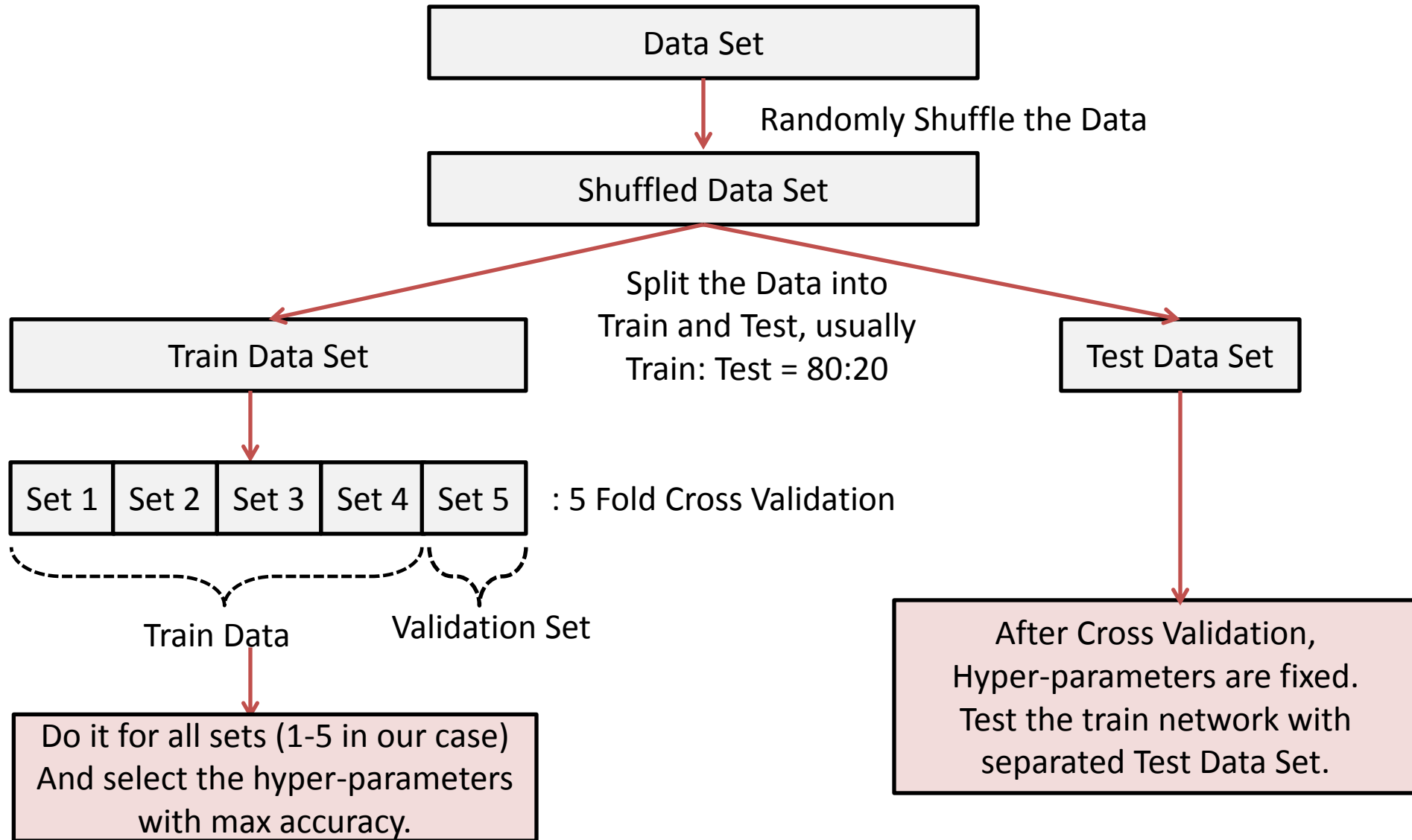
# Regularization: Early Stopping

- Stop training before we have a chance to overfit



<http://introtodeeplearning.com/index.html>

# Training and Validation



Cross Validation