

# Neural Networks

Sumohana S. Channappayya

# Review

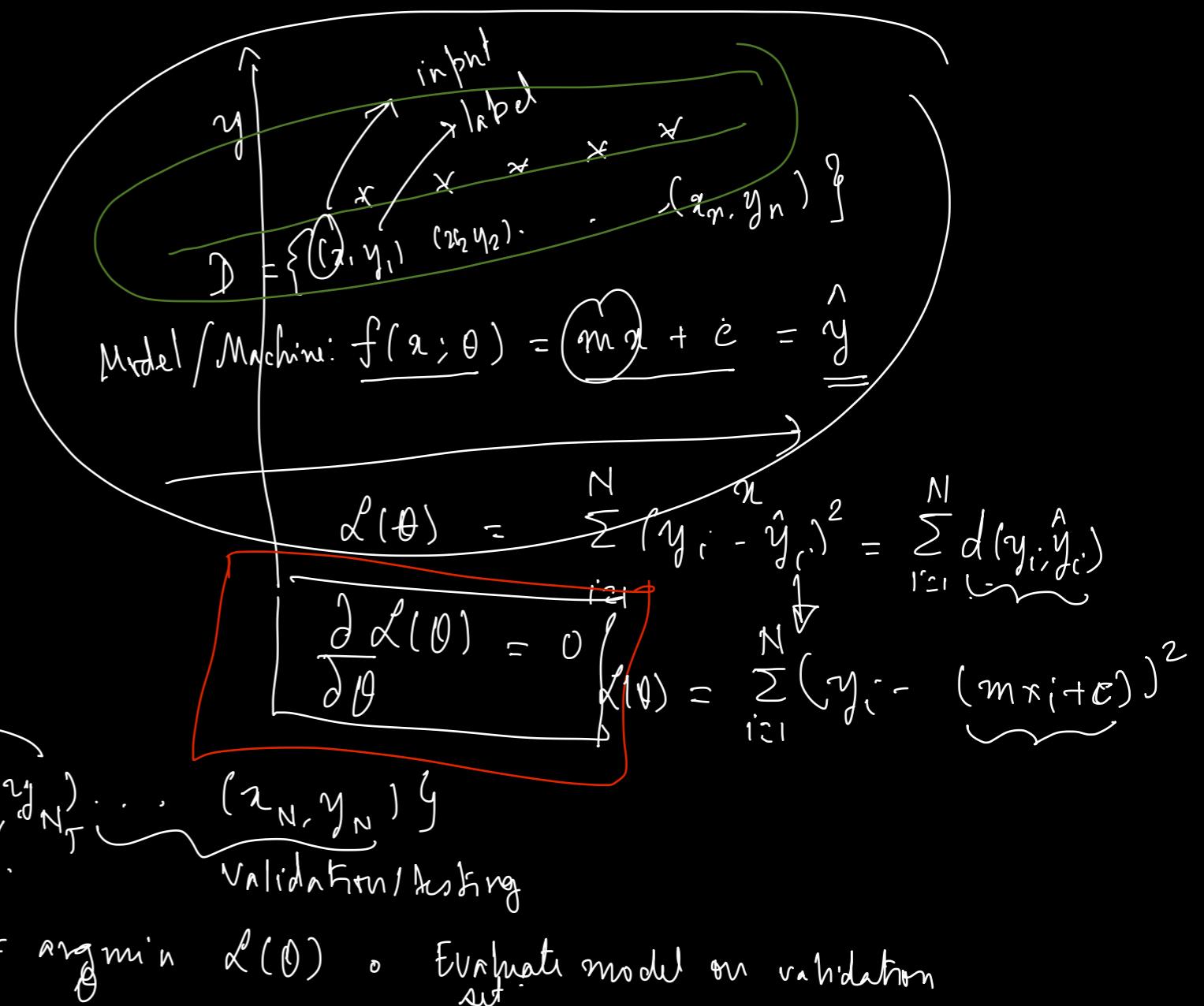
- Machine Learning

- Supervised**

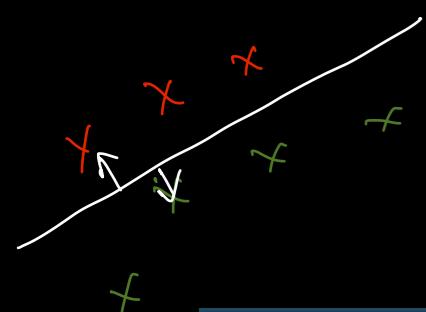
- Unsupervised

- Reinforcement

- SL:  $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$  training
- $\hat{y} = f(x; \theta)$  ← machine model.
- $\ell(\theta) = \sum_{i=1}^N d(y_i, \hat{y}_i)$ ;  $\theta^* = \arg \min_{\theta} \ell(\theta)$  Evaluate model on validation set



# Review: Supervised Learning



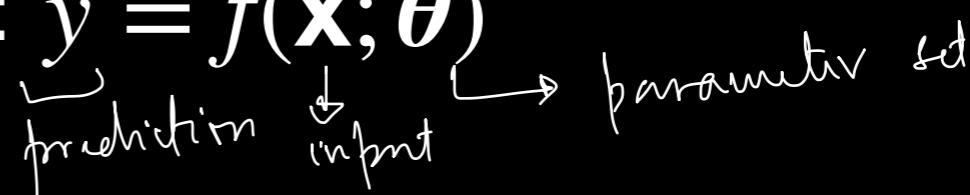
Input ( $\mathbf{x}$ )	Label ( $y$ )
	$\hat{y} = f(\mathbf{x}; \theta)$ ? Cat 0
	$\hat{y} = \sum_{i=1}^N \sum_{j=1}^M w_{i,j} x_{i,j}$ parameters Bird 1
	$\hat{y} = f(\mathbf{x}; \theta)$ Dog 2
	$= \text{NLO Linear}$ $= \sigma\left(\sum_{i=1}^N w_i x_i + w_0\right)$ Fish 3
	?

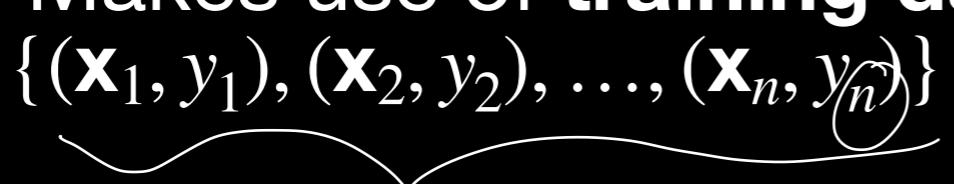
# Review: Supervised Learning

- Given a set of **input points** and corresponding **ground truth labels**:  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$

 labels

- Discrete labels: **classification**
- Continuous labels: **regression**
- Learn a **parameterised ( $\theta$ ) functional relationship** to estimate label from input:  $\hat{y} = f(\mathbf{x}; \theta)$
- Makes use of **training data** to learn parameters  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$

 parametric set



# Review: Supervised Learning

- How to learn parameters  $\theta$ ?

- Quantify error:  $d(y, \hat{y})$

- Find **error** between estimate and ground truth over **training data points**:  $R(\theta) = \sum_{i=1}^n d(y_i, \hat{y}_i)$

- Find parameters that **minimise error**

- How **good** is our **model**?

- Evaluate model's performance on **test data**  
 $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$

# Review: Supervised Learning

- Models covered so far:
  - Linear regression
  - k-NN
  - Naive Bayes
  - Logistic regression
  - Support Vector Machines (SVM) ✓
  - Neural Network

# Motivation for Neural Networks

- Modelling the neuron



- Modelling **nonlinear** relation between  $\{\mathbf{x}_i\}$  and  $\{y_i\}$

XOR

$$\begin{aligned}f(x_1, x_2) \\= \boxed{\text{XOR}(x_1, x_2)}\end{aligned}$$

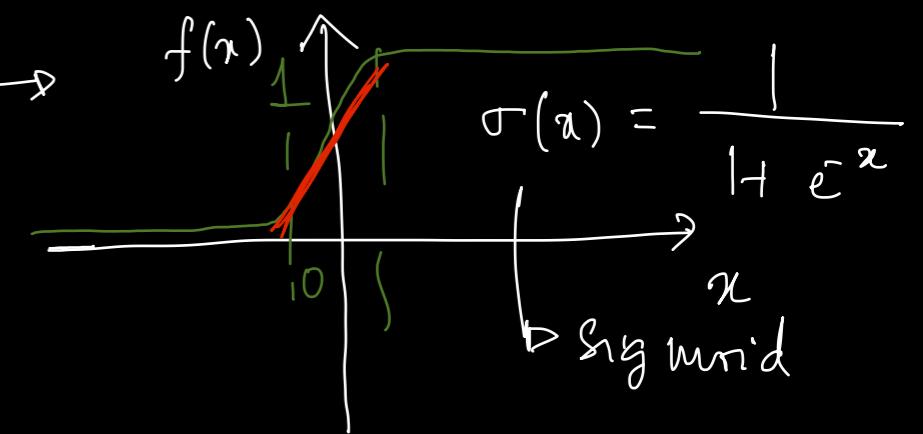
$$\begin{aligned}\hat{y} &= w_1 x_1 + \\&= \underline{w_2 x_2} \\d(y, \hat{y}) &\uparrow\end{aligned}$$

Input ( $\mathbf{x}$ )	Label ( $y$ )
✓ 0	✓ 0
0 ✓	1
1 ✓	0
1 ✓	1

Can we  
find a  
mod. for  
the XOR  
gate?

Yes. Neural Network

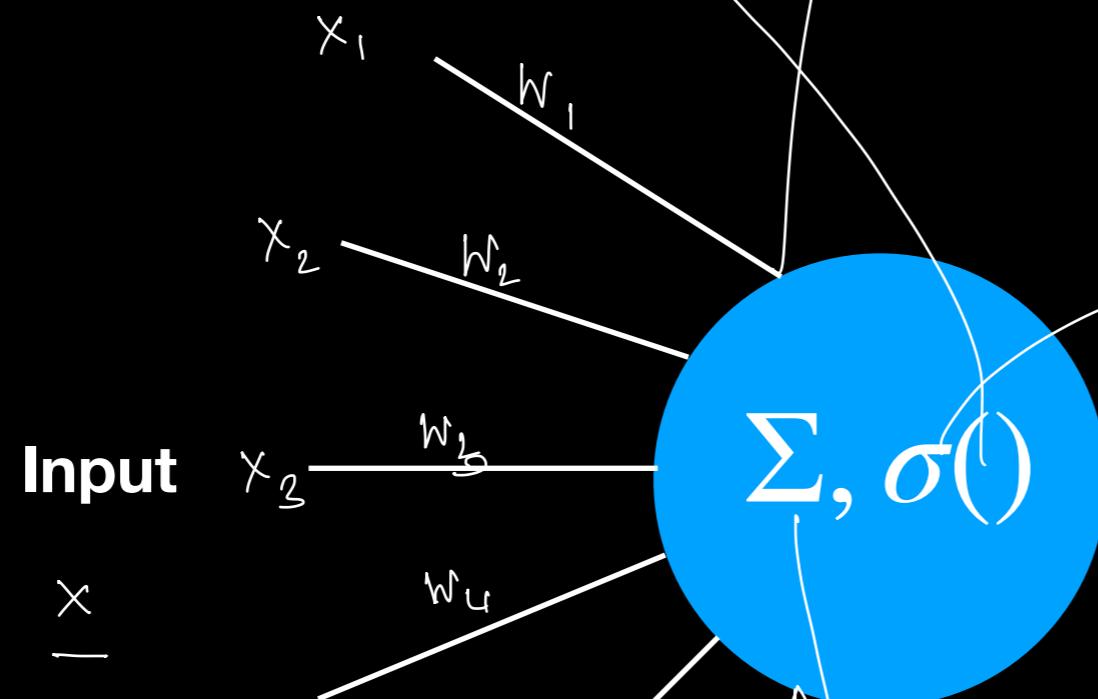
# A Neuron



$$y \in \{0, 1\}$$

$$f = g \circ h$$

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_5 \end{bmatrix}$$



$$\underline{x}' = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_s \end{bmatrix}$$

$$\hat{y} = \sigma(\underline{w}', \underline{x}')$$

$$\underline{w}' = \begin{bmatrix} w_0 \\ \vdots \\ w_s \end{bmatrix}$$

non-linearity

$$d(\theta) = \sum_{i=1}^N d(y_i, \hat{y}_i)$$

Output

$$\hat{y} = \sum_{i=1}^L w_i x_i + w_0$$

$$\hat{y} = \sigma(\hat{y}_L)$$

$$\hat{y} = \sigma\left(\sum_{i=1}^5 w_i x_i + w_0\right)$$

bias

$$\hat{y} = \sigma\left(\sum_{i=0}^s w_i x'_i\right)$$

# Key function: $\sigma()$

- Recall motivation: nonlinear input output relationship

- Properties of  $\sigma()$  nonlinearity

- Nonlinear ✓

- Continuous ✓

- Differentiable? ✓

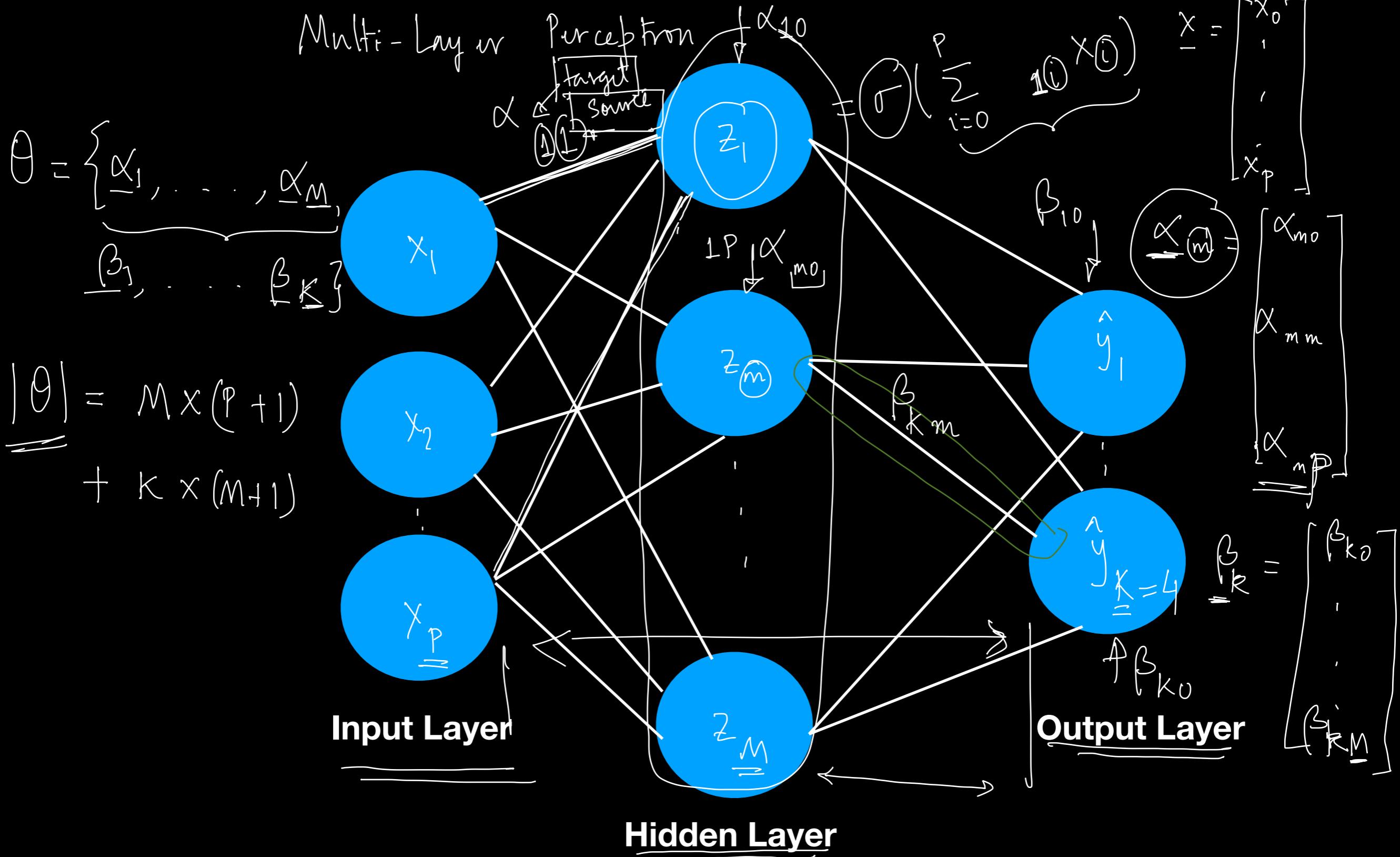
- Examples: tanh, sigmoid, ReLU etc.

$$\hat{y} = \sigma\left(\sum_{i=0}^P w_i x_i\right)$$

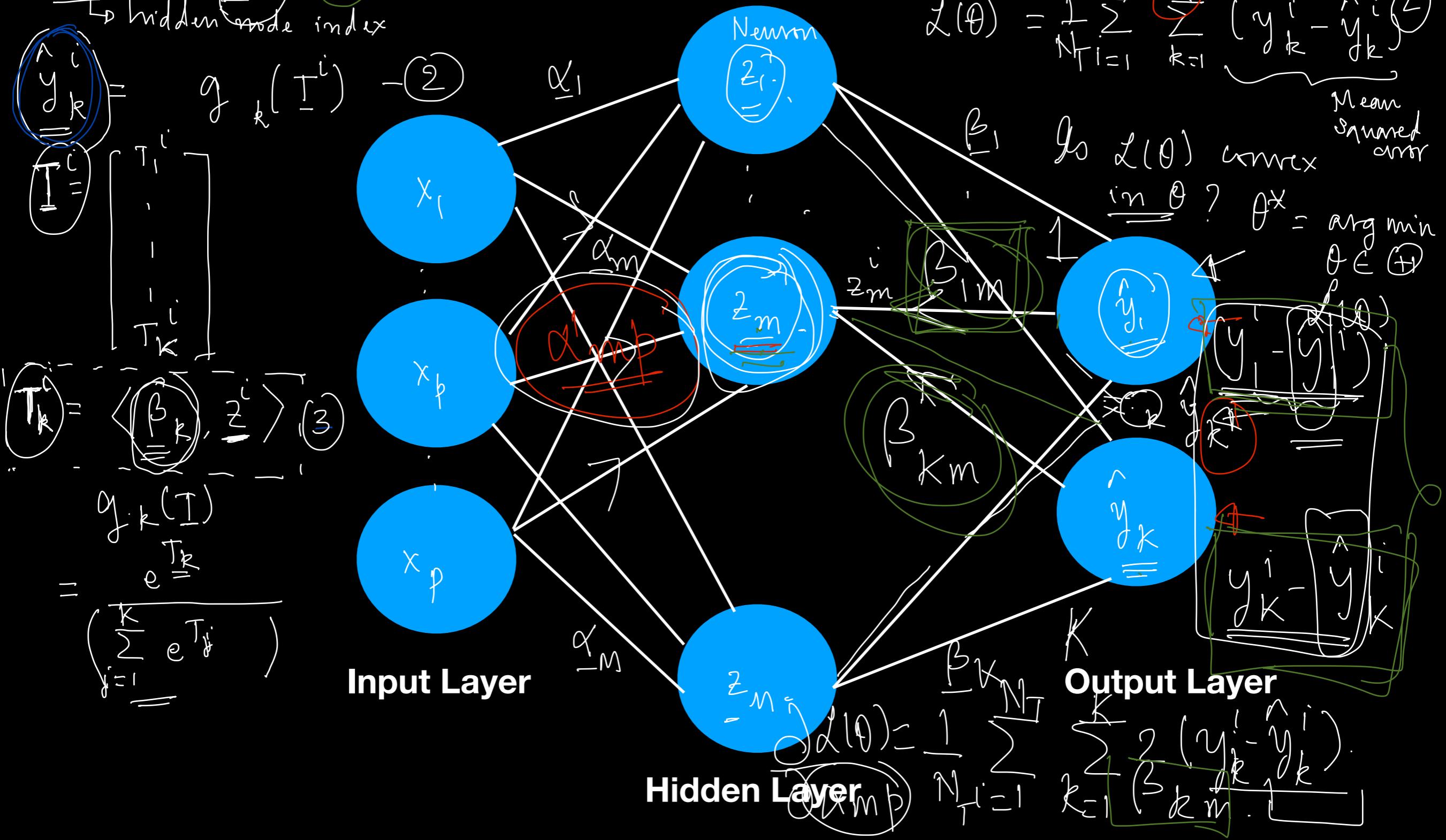
$\underline{x} = [1, x_1, \dots, x_p]^T$

$\underline{w} = [w_0, w_1, \dots, w_p]^T$

# A Neural Network



# Feedforward Neural Network



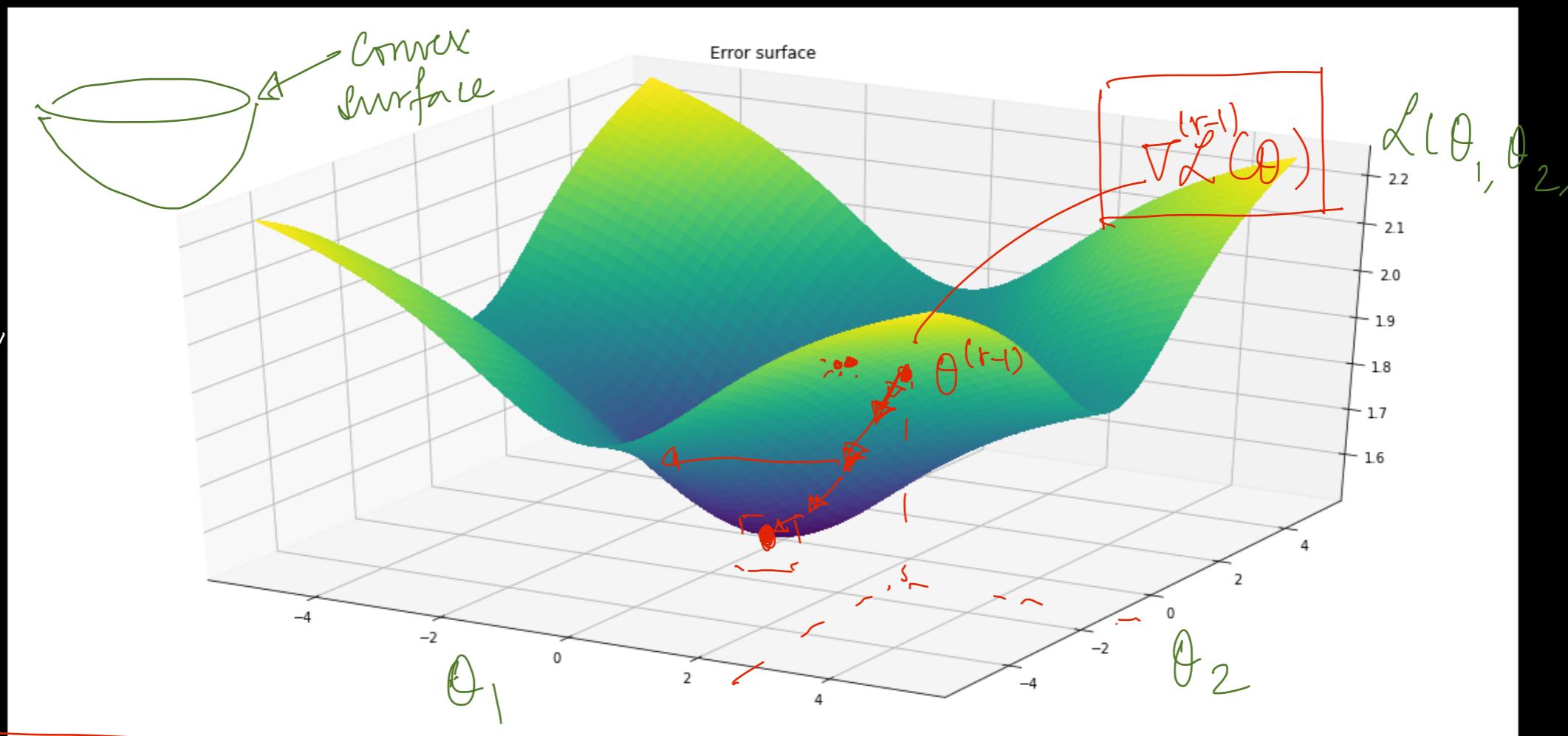
# Error Surface

- $\ell(\theta)$  is not convex in  $\theta \Rightarrow$  no closed form solution ✓  
This is due to the non-linearity present in our model.

- $\theta^*$   
is not  
going  
to global

- Iterative,  
gradient  
based  
solution

$$\begin{aligned} \theta^{(r)} &= \theta^{(r-1)} \\ &\quad - \eta \nabla \ell(\theta) \end{aligned}$$



$$\boxed{\theta^* = \hat{A}^{-1} \hat{X}}$$

$\nabla \ell(\theta)$  ?

# Minimising Error: $\frac{\partial R(\theta)}{\partial \theta_i}$

- Find local minimum ✓
- Iterative procedure ✓
- Gradient descent ✓
- Let's derive the gradient update expression

# Minimising Error: $\frac{\partial R(\theta)}{\partial \theta_i}$

$$\boxed{\begin{array}{l} \nabla \mathcal{L}(\theta) \\ = \left[ \begin{array}{l} \frac{\partial \mathcal{L}(\theta)}{\partial \theta_0} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial \theta_{10}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial \theta_{M,P}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial \theta_{km}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial \theta_{KM}} \end{array} \right] \end{array}}$$

• Derivation:

$$\frac{\partial \mathcal{L}(\theta)}{\partial \beta_{km}}$$

$$\mathcal{L}(\theta) = \frac{1}{N_T} \sum_{i=1}^{N_T} \sum_{k=1}^K (y_k^{(i)} - \hat{y}_k^{(i)})^2 = \frac{1}{N_T} \sum_{i=1}^{N_T} \left( \sum_{k=1}^K (y_k^{(i)} - g_k(\underline{T}^{(i)}))^2 \right)$$

$$= \frac{1}{N_T} \sum_{i=1}^{N_T} \left( \sum_{k=1}^K (y_k^{(i)} - g_k(\langle \beta_k, z^{(i)} \rangle))^2 \right)$$

$(\because \beta_{km} \text{ affects only } \hat{y}_k^{(i)})$

$$= \frac{1}{N_T} \sum_{i=1}^{N_T} -2(y_k^{(i)} - \hat{y}_k^{(i)})$$

$$\text{Now, let's find } \frac{\partial \hat{y}_k^{(i)}}{\partial \beta_{km}} = \frac{\partial}{\partial \beta_{km}} g_k(\underline{T}^{(i)})$$

$$= g_k'(\underline{T}^{(i)}) \cdot \frac{\partial \langle \beta_k, z^{(i)} \rangle}{\partial \beta_{km}} = g_k'(\underline{T}^{(i)}) \cdot \frac{\partial (\beta_{km} z_m^{(i)})}{\partial \beta_{km}}$$

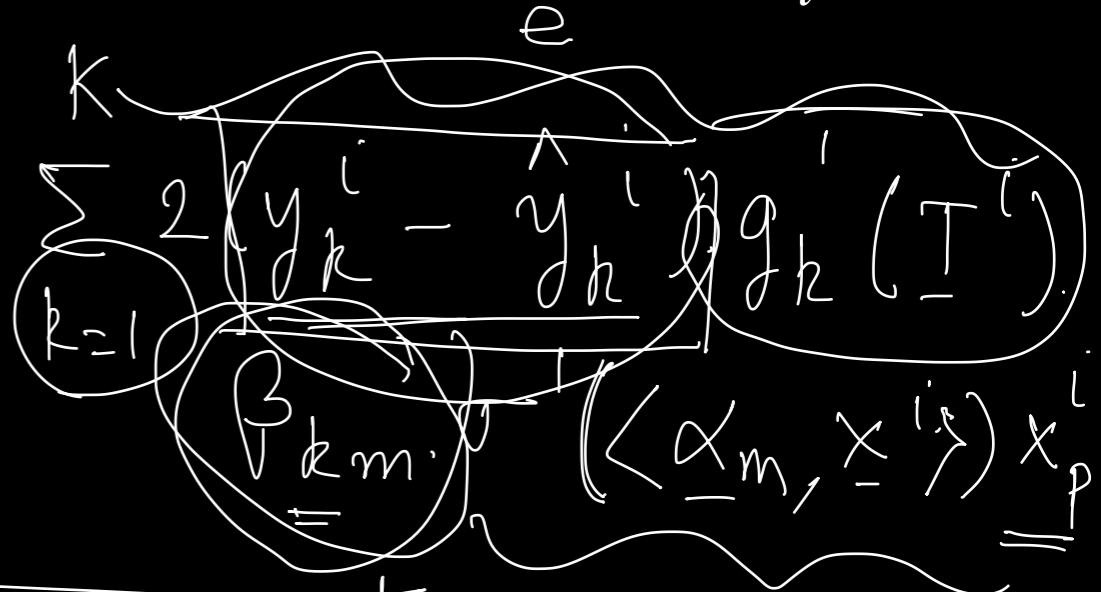
Plug (2) into (1) gives us

$$\frac{\partial \mathcal{L}(\theta)}{\partial \beta_{km}} = g_k'(\underline{T}^{(i)}) \cdot z_m^{(i)}$$

# Minimising Error: $\frac{\partial R(\theta)}{\partial \theta_i}$

- Derivation:

$$\frac{\partial \mathcal{L}(\theta)}{\partial \alpha_m} = -\frac{1}{N_T} \sum_{i=1}^{N_T} \sum_{k=1}^K \frac{1}{2} (y_k^i - \hat{y}_k^i)^2 g_k(T^i)$$



$$\underline{\theta}^{(r)} = \underline{\theta}^{(r-1)} - \eta \nabla \mathcal{L}(\theta)$$

found using  
back propagation.

# Back propagation

- o  $\boxed{\text{Neuron}}$   $\hat{y} = \sigma(\langle \underline{w}, \underline{x} \rangle)$

- o A multi-layer perceptron

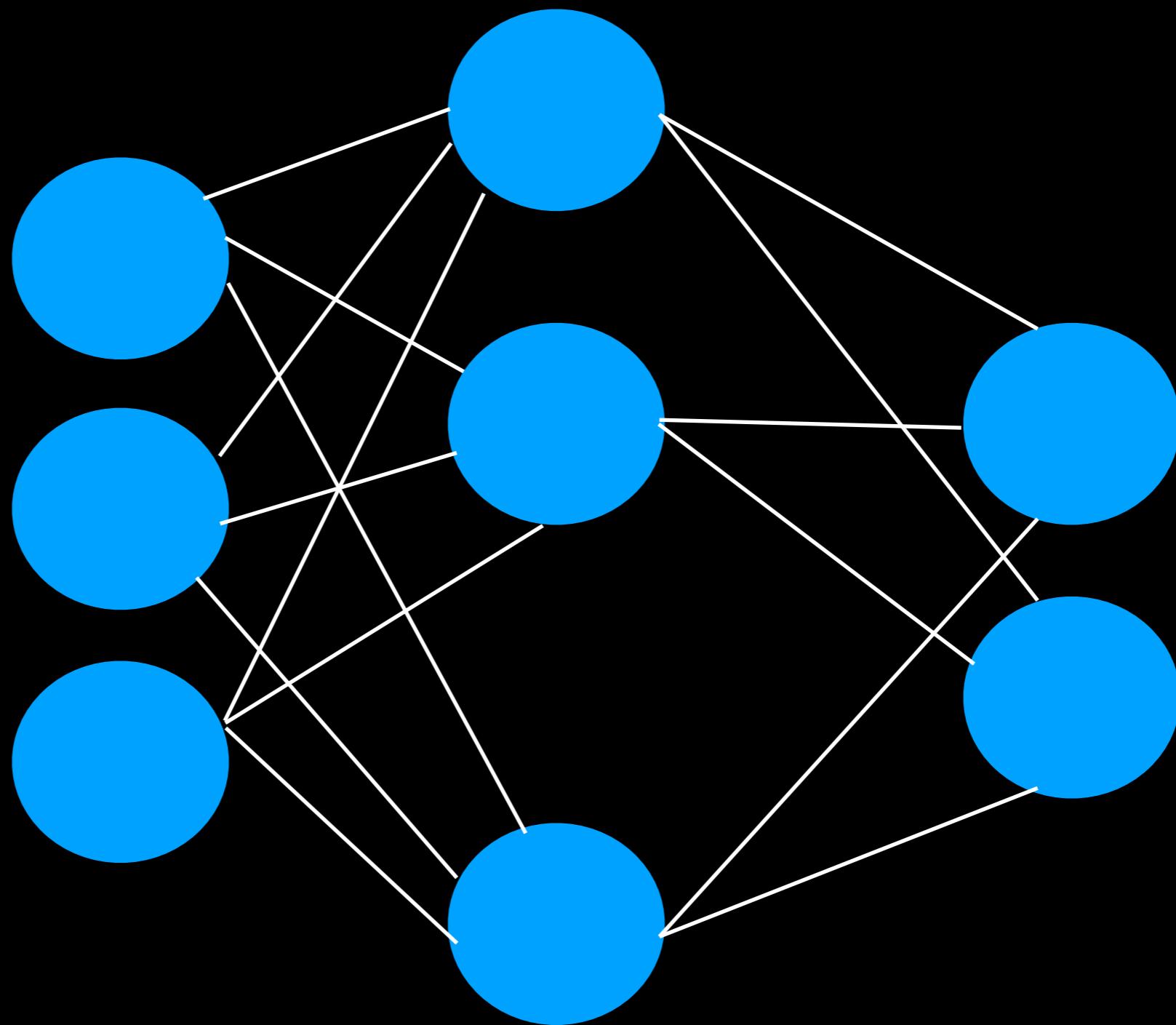
- o  $\hat{y} = f(\underline{x}; \theta)$

- o  $\mathcal{L}(\theta)$

- o  $\nabla \mathcal{L}(\theta)$ : back propagation

- o  $\underline{\theta}^{(r)} = \underline{\theta}^{(r-1)} - \eta \nabla^{(r)} \mathcal{L}(\theta)$

# Back propagation



# Challenges

- Initialization : ensure weights are not 0. Random init.
- Overfitting :  $\|\theta\|$  grows with # neurons
- Input range : Large dynamic range
- Network architecture : choice of neurons not straightforward  
# data points,  $\exists$  relationship
- Multiple minima