

Neural Networks

Sumohana S. Channappayya

Review

- Machine Learning
 - **Supervised**
 - Unsupervised
 - Reinforcement

Review: Supervised Learning

| Input (\mathbf{x}) | Label (y) |
|--|---------------|
|  | Cat |
|  | Bird |
|  | Dog |
|  | Fish |
|  | ? |

Review: Supervised Learning

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

Review: Supervised Learning

- How to learn parameters θ ?
 - Quantify error: $d(y, \hat{y})$
 - Find **error** between estimate \hat{y} 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)

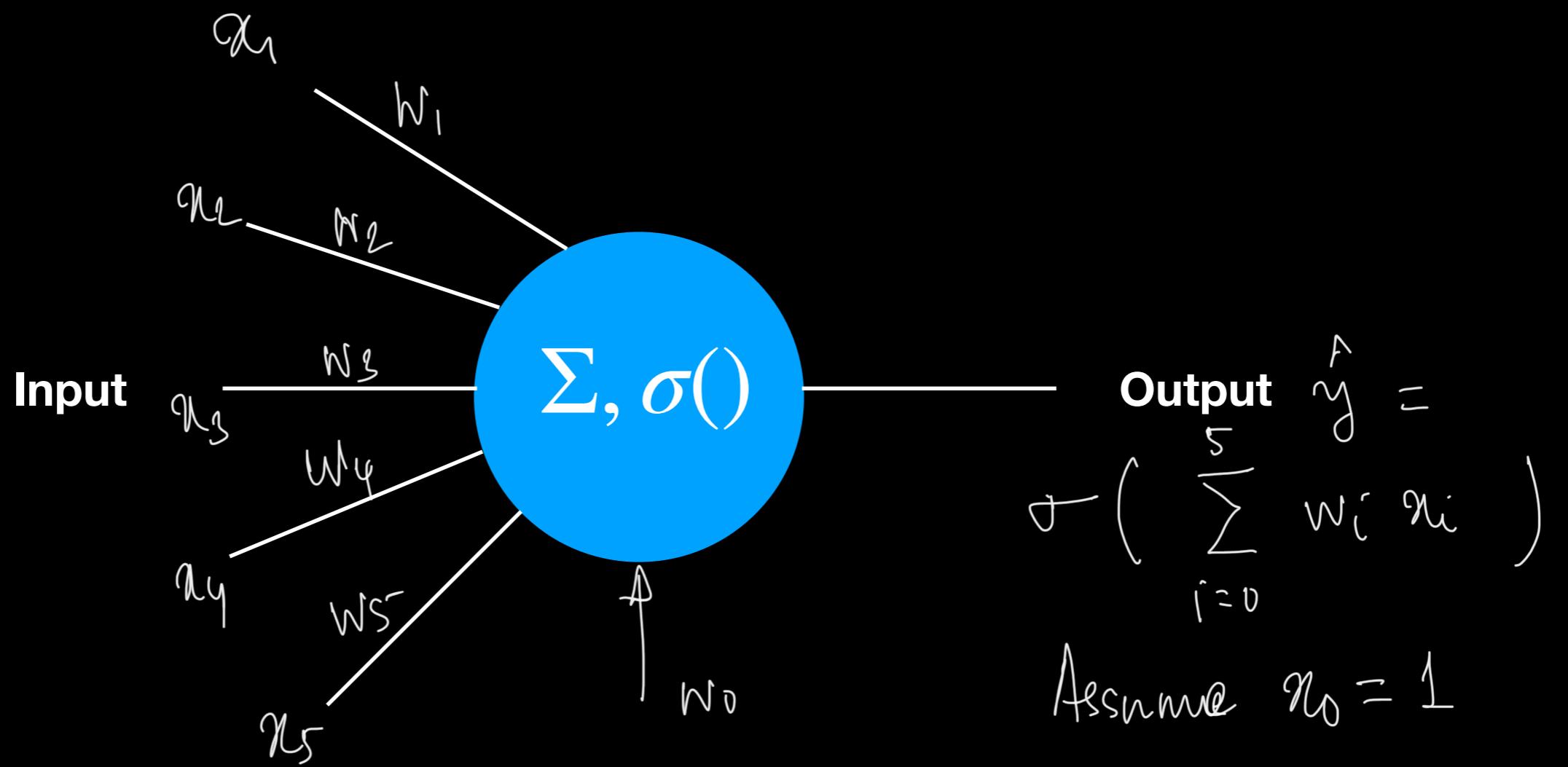
Motivation for Neural Networks

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



| Input (\mathbf{x}) | | Label (y) |
|------------------------|---|---------------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

A Neuron



Key function: $\sigma()$

- Recall motivation: **nonlinear** input output relationship

- Properties of $\sigma()$

- Nonlinear

- Continuous

- Differentiable?

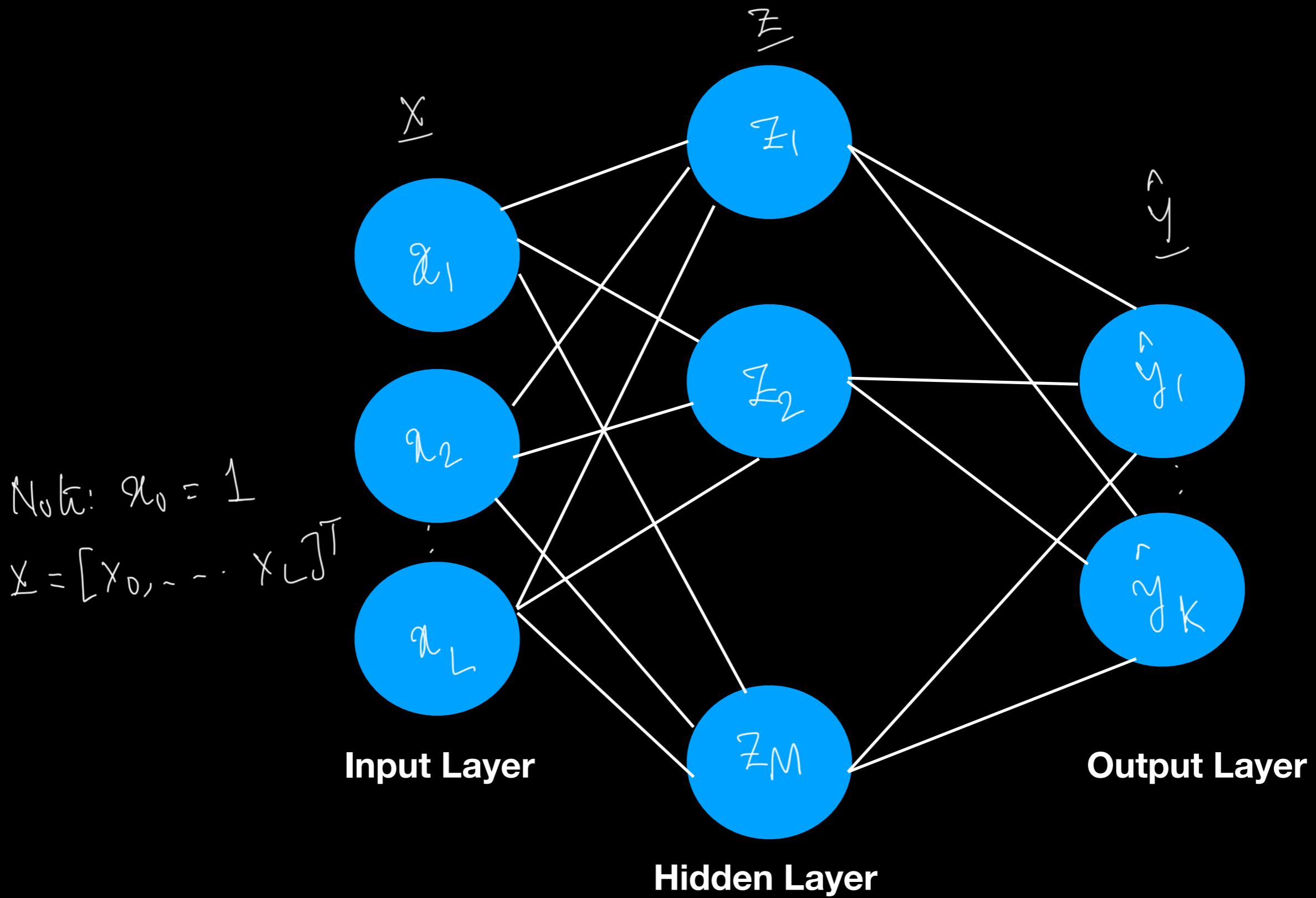
- Examples: tanh, sigmoid, ReLU etc.

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma(x) = \max\{0, x\}$$

A Neural Network



Feedforward Neural Network

$$\underline{\alpha}_m = [\alpha_{m0}, \alpha_{m1}, \dots, \alpha_{mL}]^T$$

$$\underline{\beta}_k = [\beta_{k0}, \dots, \beta_{kM}]^T$$

$$\underline{\theta} = [\underline{\alpha}_1^T, \dots, \underline{\alpha}_M^T, \underline{\beta}_1^T, \dots, \underline{\beta}_K^T]^T$$

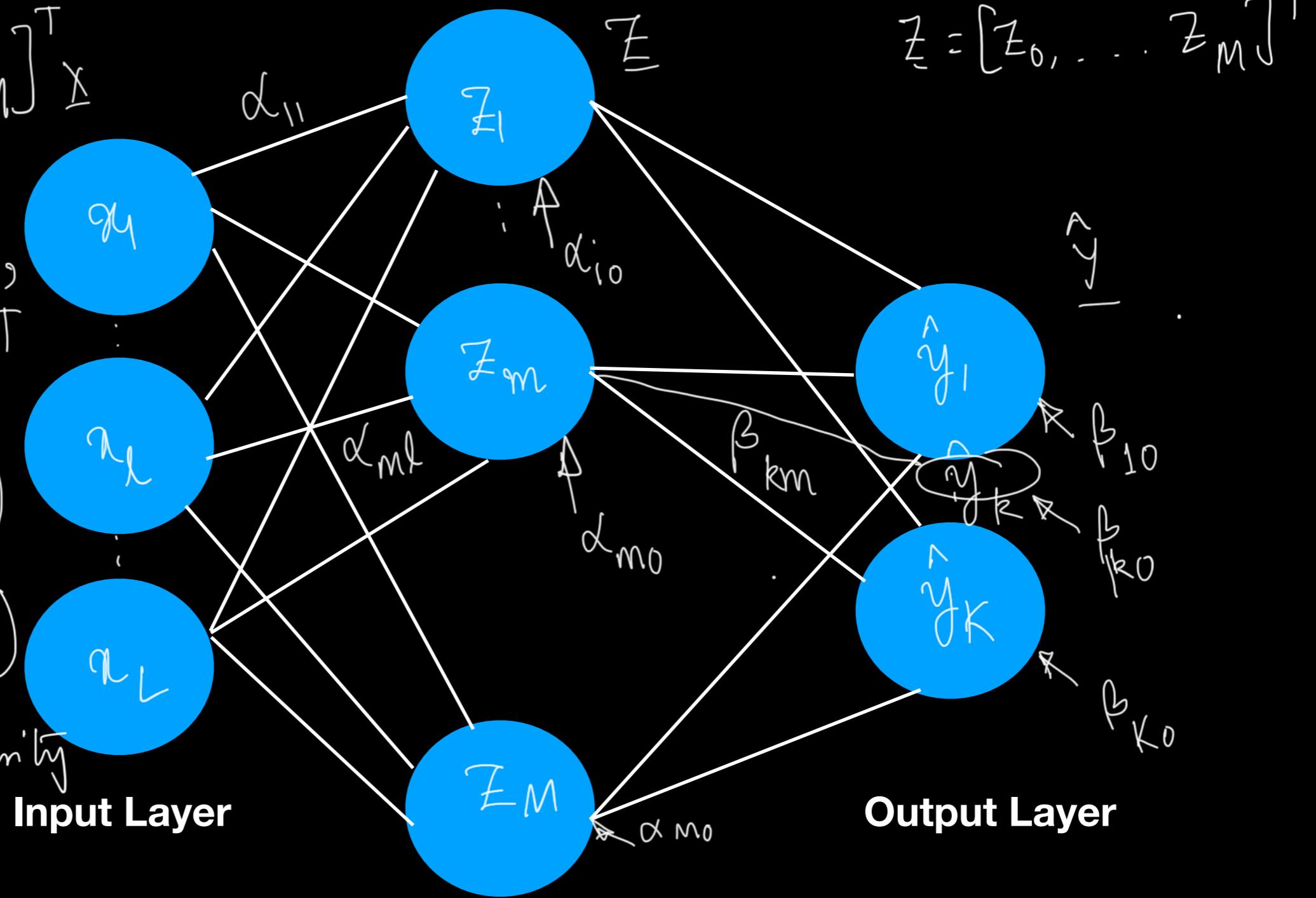
$$z_m = f(\underline{\alpha}_m^T \cdot \underline{x})$$

$$\hat{y}_k = g(\underline{\beta}_k^T \cdot \underline{z})$$

g : output nonlinearity

Input Layer

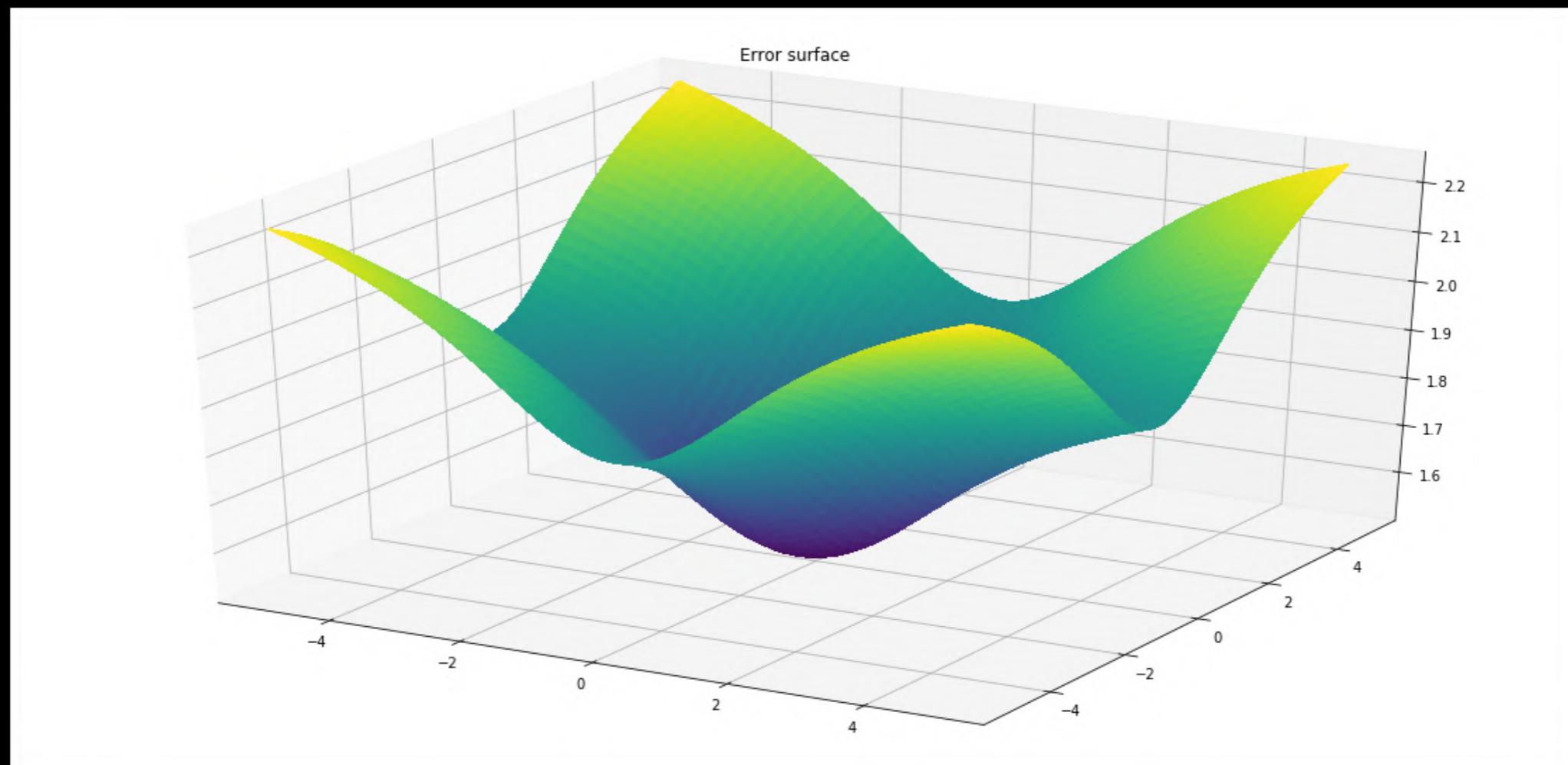
Hidden Layer



$$\text{Note: } z_0 = 1$$

$$\underline{z} = [z_0, \dots, z_m]^T$$

Error Surface



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

S.S.E: $d(y, \hat{y}) = \sum_{i=1}^K (y_i - \hat{y}_i)^2$; $d_2(y, \hat{y}) = -\sum_{i=1}^K y_i \log \hat{y}_i$

- Derivation:

$$R(\underline{\theta}) = \sum_{i=1}^n \sum_{k=1}^K (y_{ki} - \hat{y}_{ki})^2$$

Let's find $\frac{\partial R(\underline{\theta})}{\partial \beta_{km}}$. Observation: β_{km} affects only \hat{y}_{ki} .
 \Rightarrow we can ignore the contribution of other outputs to $R(\underline{\theta})$.

$$\begin{aligned} \frac{\partial}{\partial \beta_{km}} \sum_{i=1}^n \sum_{j=1}^K (y_{ji} - \hat{y}_{ji})^2 &= \frac{\partial}{\partial \beta_{km}} \sum_{i=1}^n (y_{ki} - \hat{y}_{ki})^2 \\ &= \frac{\partial}{\partial \beta_{km}} \sum_{i=1}^n (y_{ki} - g(\beta_{k \cdot}^T \cdot \underline{\xi}_i))^2 \end{aligned}$$

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

- Derivation: $\Rightarrow \frac{\partial R(\theta)}{\partial \beta_{km}} = -\sum_{i=1}^n 2(y_{ki} - \hat{y}_{ki}) \cdot g'(\beta_k^T \cdot z_i) \circ z_{mi} - \textcircled{1}$

$$\frac{\partial}{\partial \beta_{km}} [\beta_k^T \cdot z_i] = \frac{\partial}{\partial \beta_{km}} [\beta_{k0} \cdot z_{m0} + \beta_{k1} z_{m1} + \dots + \beta_{km} z_{mi} + \dots + \beta_{RM} z_{Mi}] = z_{mi}$$

$\therefore \boxed{\frac{\partial R(\theta)}{\partial \beta_{km}} = \sum_{i=1}^n \delta_{ki} z_{mi}}$; where $\delta_{ki} = -2(y_{ki} - \hat{y}_{ki}) \cdot g'(\beta_k^T \cdot z_i)$

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

$$\begin{aligned}
 & \underline{\alpha}_m^T \cdot \underline{x}_{li} \\
 &= \alpha_{m0} \cdot x_{0i} + \alpha_{m1} x_{1i} + \dots \\
 &+ \alpha_{ml} \circled{a_{li}} \\
 &+ \dots + \alpha_{mL} x_{Li} \\
 &= \sum_{i=1}^n \sum_{k=1}^K \frac{\partial}{\partial \alpha_{ml}} (y_{ki} - g(\beta_k^T \underline{z}_{li}))^2 \\
 &= \sum_{i=1}^n \sum_{k=1}^K \frac{\partial}{\partial \alpha_{ml}} (y_{ki} - g(\beta_k^T \underline{z}_{li}))^2 = \sum_{i=1}^n \sum_{k=1}^K \frac{\partial}{\partial \alpha_{ml}} (y_{ki} - g(\beta_k^T \cdot \underline{z}_{li}))^2 \\
 &= \sum_{i=1}^n \sum_{k=1}^K \frac{\partial}{\partial \alpha_{ml}} (y_{ki} - g(\dots + \beta_{km} \cdot \underbrace{\tau(\underline{\alpha}_m^T \cdot \underline{x}_{li})})) \\
 &\frac{\partial R(\theta)}{\partial \alpha_{ml}} = \sum_{k=1}^K \sum_{i=1}^n -2 \cdot (y_{ki} - g(\dots)) \cdot g'(\dots) \cdot \beta_{km} \\
 &+ \tau(\underline{\alpha}_m^T \cdot \underline{x}_{li}) \cdot \underline{x}_{li} \quad \text{--- (2)}
 \end{aligned}$$

Back propagation

$$\frac{\partial R(\underline{\theta})}{\partial \alpha_{ml}} = \sum_{i=1}^n s_{mi} \alpha_{li}$$

where $s_{mi} = \sigma^1(\underline{\alpha}_m^T \cdot \underline{x}_i) \cdot \sum_{k=1}^K \beta_{km} \cdot \delta_{ki}$.

back propagation of error.

Back propagation

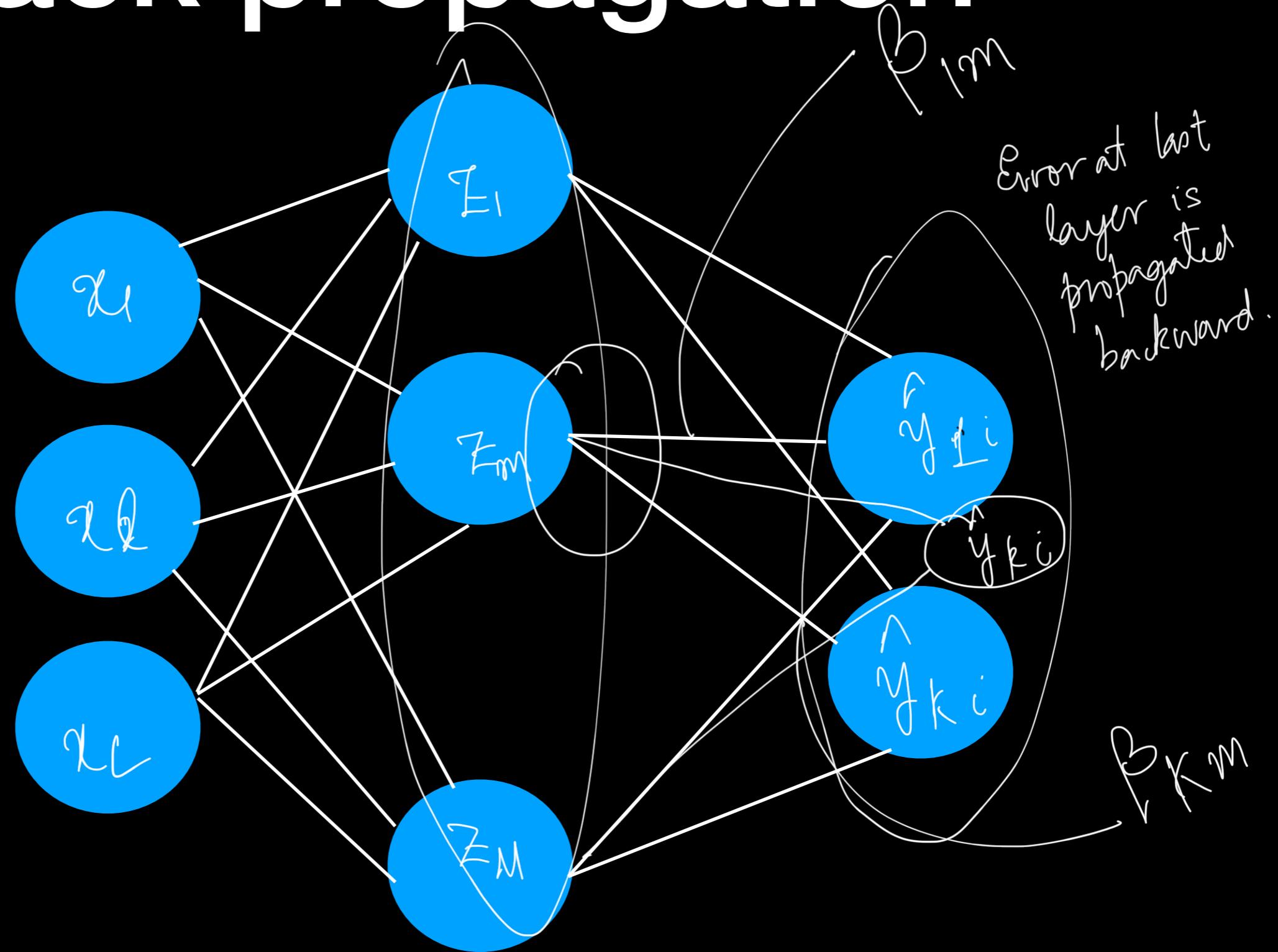
$$\beta_{km}^{(r+1)} = \beta_{km}^{(r)} - \gamma_r \cdot \frac{\partial R(\theta)}{\partial \beta_{km}}$$

$$\alpha_{ml}^{(r+1)} = \alpha_{ml}^{(r)} - \gamma_r \cdot \frac{\partial R(\theta)}{\partial \alpha_{ml}}$$

r: iteration count

γ_r : learning rate

Back propagation



Challenges

- Initialization
- Overfitting
- Input range
- Network architecture
- Multiple minima