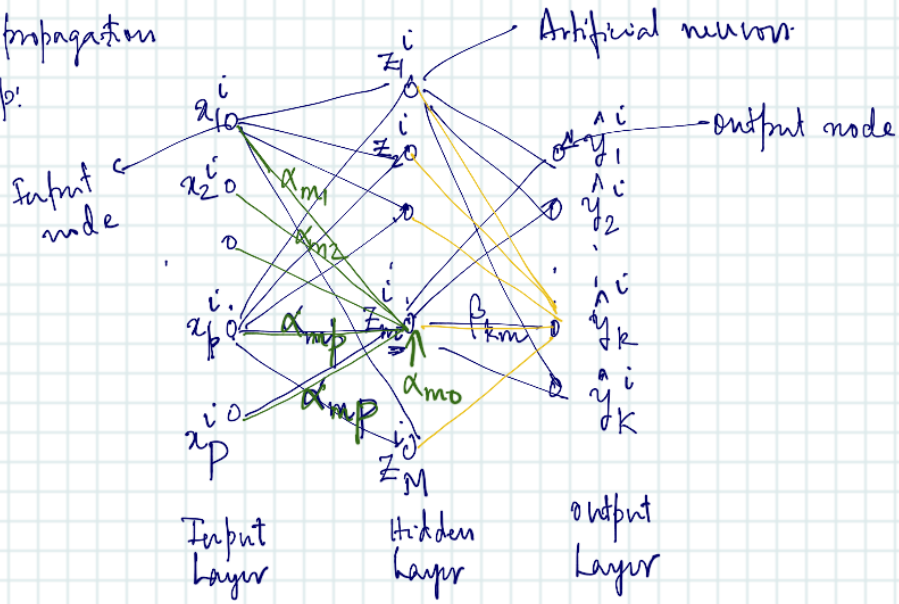


- Recap
- ANN
- Gradient descent
- Backpropagation
- Recap:



Multilayer Perceptron

- Goals:
- Find the input output relation ✓
 - In the supervised learning setting, describe & define loss $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
 - Find a method to minimize loss

Let $\underline{x}^i = [1, x_1^i, \dots, x_p^i]^T$,
 $\underline{z}^i = [1, z_1^i, \dots, z_M^i]^T$
 $\underline{y}^i = [\hat{y}_1^i, \dots, \hat{y}_K^i]^T$

• Input output relation:

Input-hidden layer relation: \rightarrow non-linearity

$$z_m^i = \sigma \left(\langle \underline{x}^i, \underline{\alpha}_m \rangle \right)$$

Annotations for the equation above:
- σ : node index
- m : data point
- i : data point
- $\underline{\alpha}_m$: weight vector

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

Annotations for the equation above:
- α_{m0} : bias

$$z_m^i = \sigma \left(\sum_{j=0}^P x_j^i \cdot \alpha_{mj} \right)$$

$$z_m^i = \sigma \left(\underbrace{1 \cdot \alpha_{m0}}_{\text{bias term}} + x_1^i \alpha_{m1} + \dots + x_{\cancel{p}}^i \alpha_{m\cancel{p}} \right) \quad - (1)$$

$$\underline{z}^i = [1, z_1^i, z_2^i, \dots, \underbrace{z_m^i}, \dots, z_M^i]^T$$

$$\underline{z}^i = [1, \sigma(\langle \underline{x}_1^i, \underline{\alpha}_1 \rangle), \sigma(\langle \underline{x}_2^i, \underline{\alpha}_2 \rangle), \dots, \sigma(\langle \underline{x}_m^i, \underline{\alpha}_m \rangle)]^T \quad - (2)$$

parameters in the input-hidden layer node connections = $M \cdot (P+1)$

• hidden layer - output layer relation.

$$\hat{y}_k^i = \text{softmax}(\langle \underline{z}^i, \underline{\beta}_k \rangle) \quad \underline{\beta}_k = [1, \beta_{k1}, \dots, \beta_{kM}]^T$$

$$\hat{y}_k^i = \frac{e^{\langle \underline{z}^i, \underline{\beta}_k \rangle}}{\sum_{j=1}^K e^{\langle \underline{z}^i, \underline{\beta}_j \rangle}} \quad - (3)$$

$$\hat{y}_k^i = \text{softmax}(\langle [1, \sigma(\langle \underline{x}_1^i, \underline{\alpha}_1 \rangle), \dots, \sigma(\langle \underline{x}_m^i, \underline{\alpha}_m \rangle)]^T, \underline{\beta}_k \rangle) \quad - (4)$$

$\Rightarrow \hat{y}_k^i = f(\underline{x}^i; \underline{\theta})$ where $f(\cdot)$ is defined in (4),

$$\underline{\theta} = \{ \underline{\alpha}_1, \underline{\alpha}_2, \dots, \underline{\alpha}_M, \underline{\beta}_1, \underline{\beta}_2, \dots, \underline{\beta}_K \}$$

• Risk or cost : $R(\underline{\theta}) = \sum_{i=1}^n d(\underline{y}^i, \hat{\underline{y}}^i)$