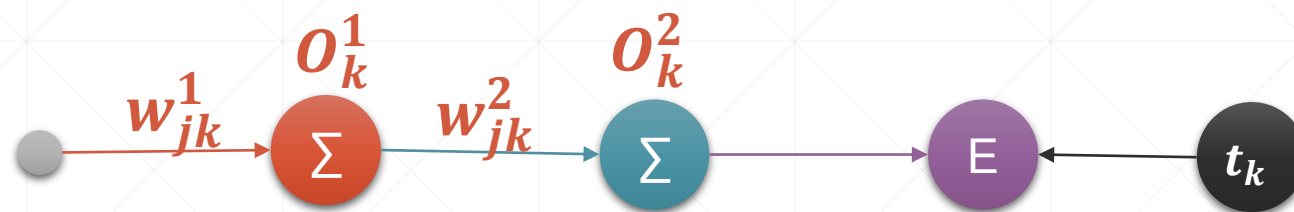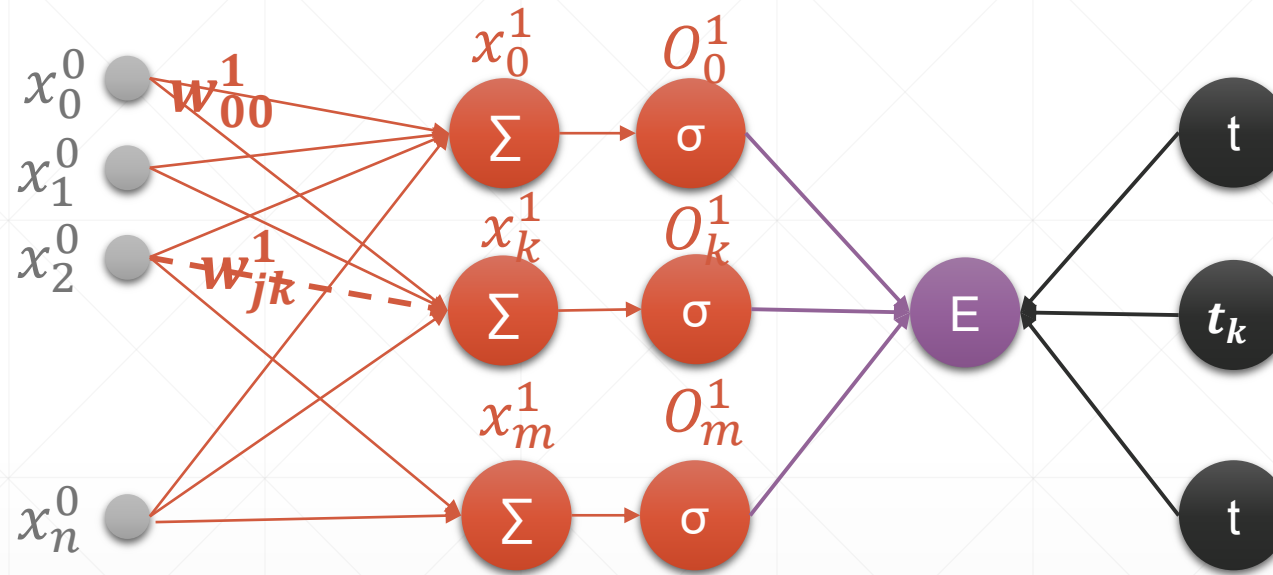# 多层感知机梯度

主讲：龙良曲

# Chain rule



$$\frac{\partial E}{\partial w_{jk}^1} = \frac{\partial E}{\partial O_k^1} \frac{\partial O_k^1}{\partial x} = \frac{\partial E}{\partial O_k^2} \frac{\partial O_k^2}{\partial O_k^1} \frac{\partial O_k^1}{\partial x}$$
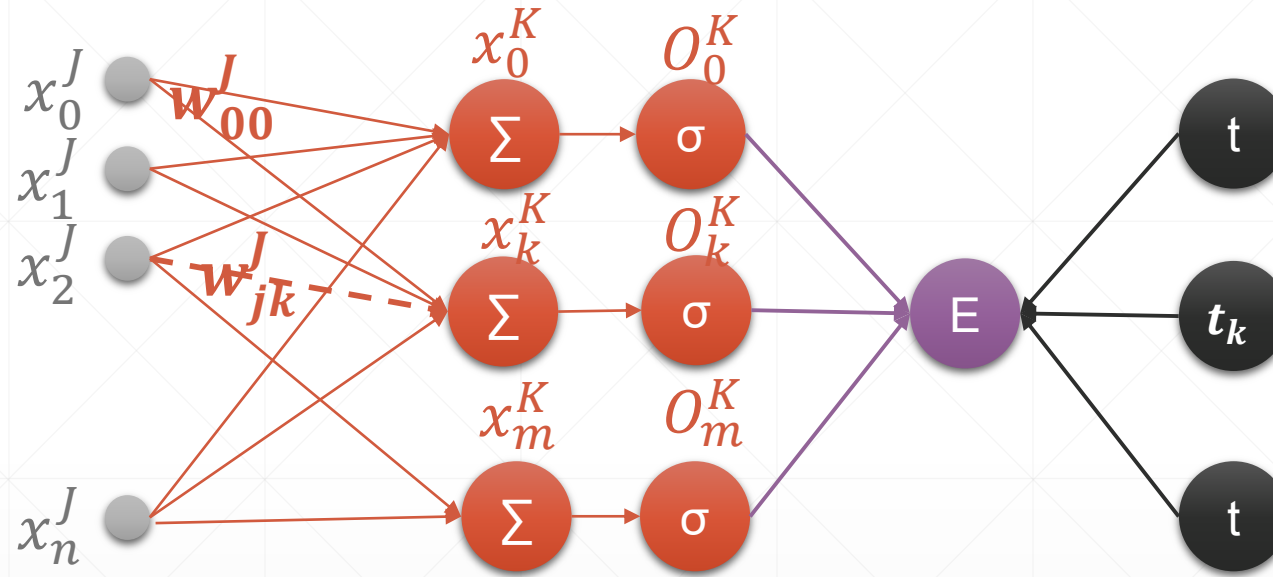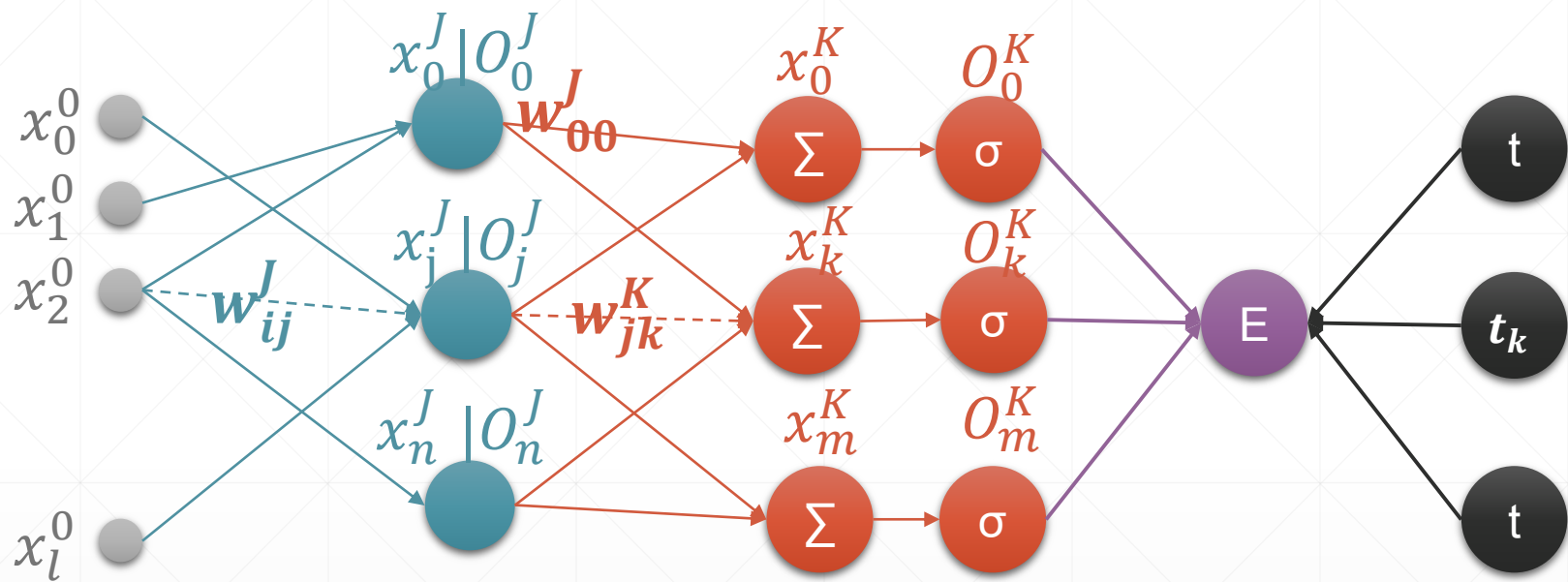
# Multi-output Perceptron



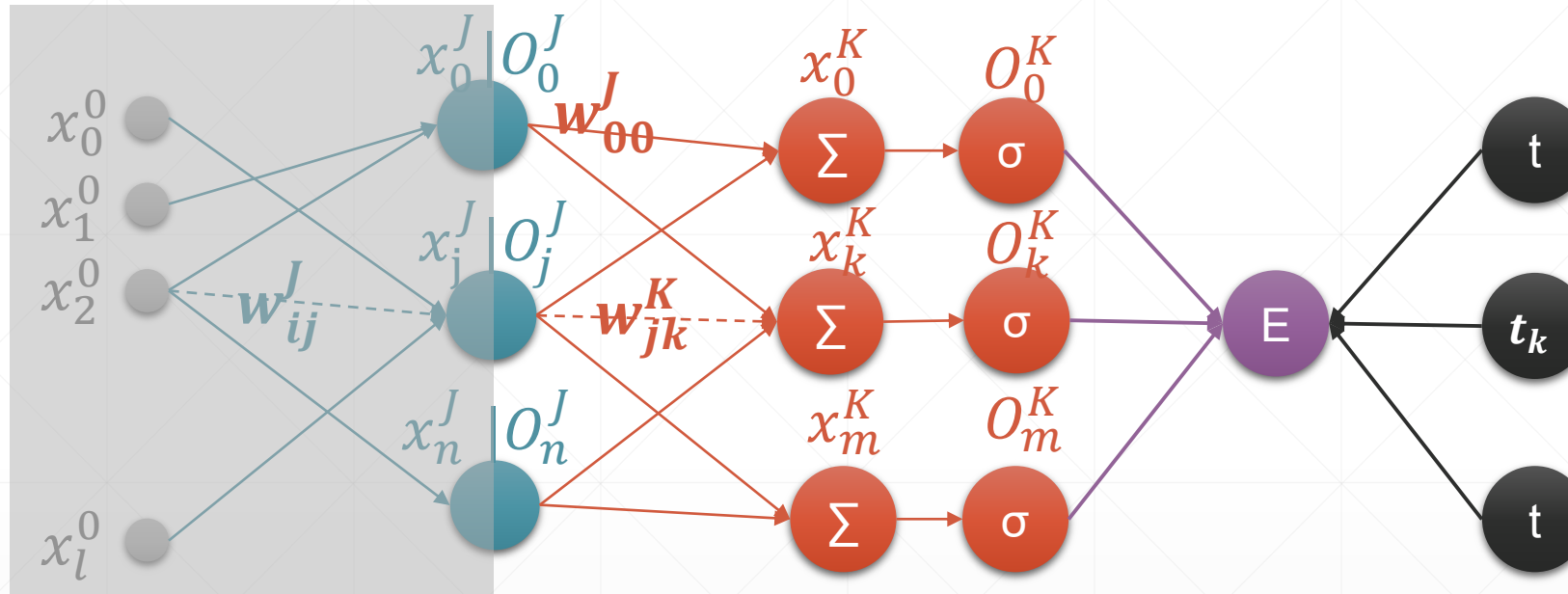$$\frac{\partial E}{\partial w_{jk}} = \left(O_k - t_k\right) O_k \left(1 - O_k\right) x_j^0$$

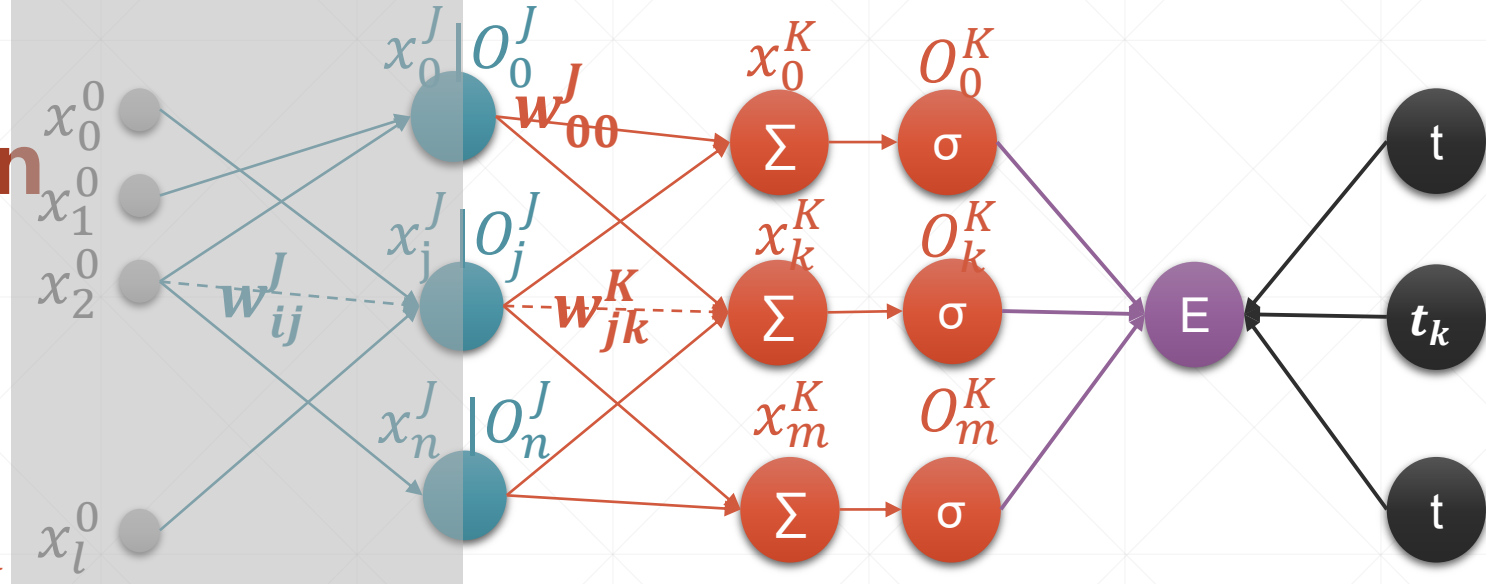# Multi-Layer Perceptron

# Multi-Layer Perceptron

# Multi-Layer Perceptron



$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k)O_k(1 - O_k)\, x_j^0$$

$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k)O_k(1 - O_k)\, O_j^J$$

# Multi-Layer Perceptron



$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k)O_k(1 - O_k)\, O_j^J$$

$$\frac{\partial E}{\partial w_{jk}} = \delta_k^K \qquad\qquad O_j^J$$

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial}{\partial W_{ij}} \frac{1}{2} \sum_{k \in K} (\mathcal{O}_k - t_k)^2$$

$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (\mathcal{O}_k - t_k) \frac{\partial}{\partial W_{ij}} \mathcal{O}_k$$

$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (\mathcal{O}_k - t_k) \frac{\partial}{\partial W_{ij}} \sigma(x_k)$$
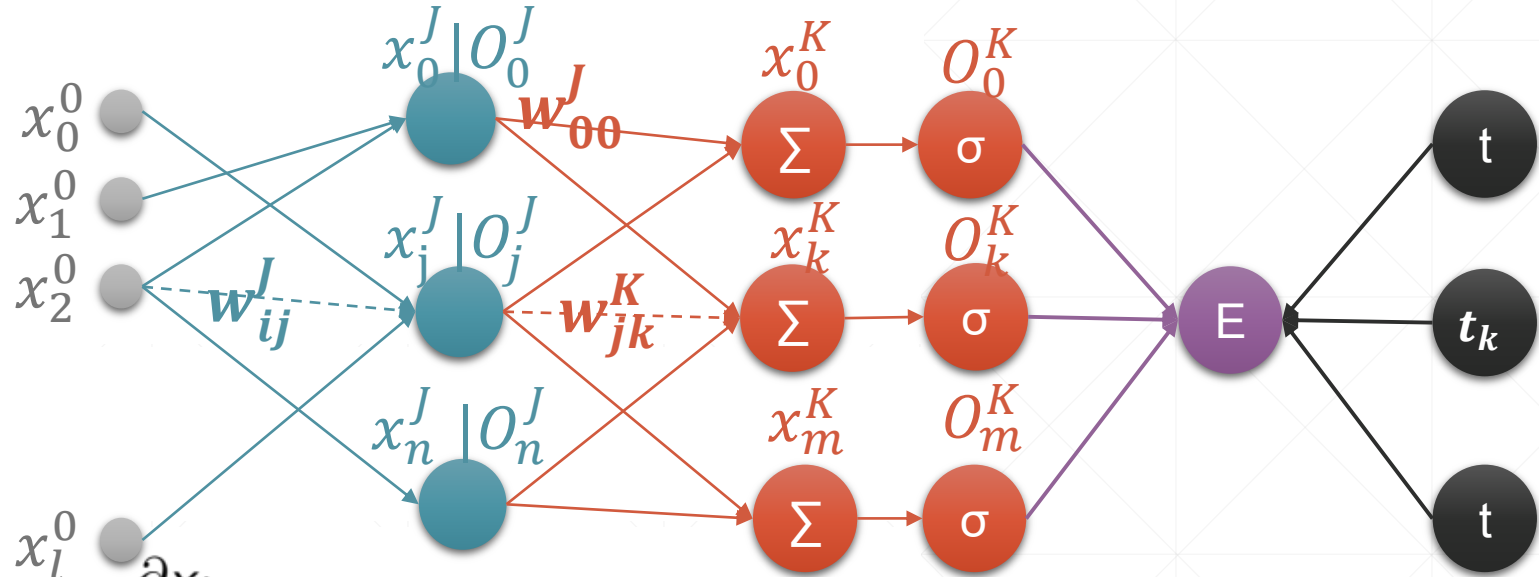
$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (\mathcal{O}_k - t_k) \sigma(x_k)(1 - \sigma(x_k)) \frac{\partial x_k}{\partial W_{ij}}$$

$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (\mathcal{O}_k - t_k) \mathcal{O}_k (1 - \mathcal{O}_k) \frac{\partial x_k}{\partial \mathcal{O}_j} \cdot \frac{\partial \mathcal{O}_j}{\partial W_{ij}}$$

$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (\mathcal{O}_k - t_k) \mathcal{O}_k (1 - \mathcal{O}_k) W_{jk} \frac{\partial \mathcal{O}_j}{\partial W_{ij}}$$
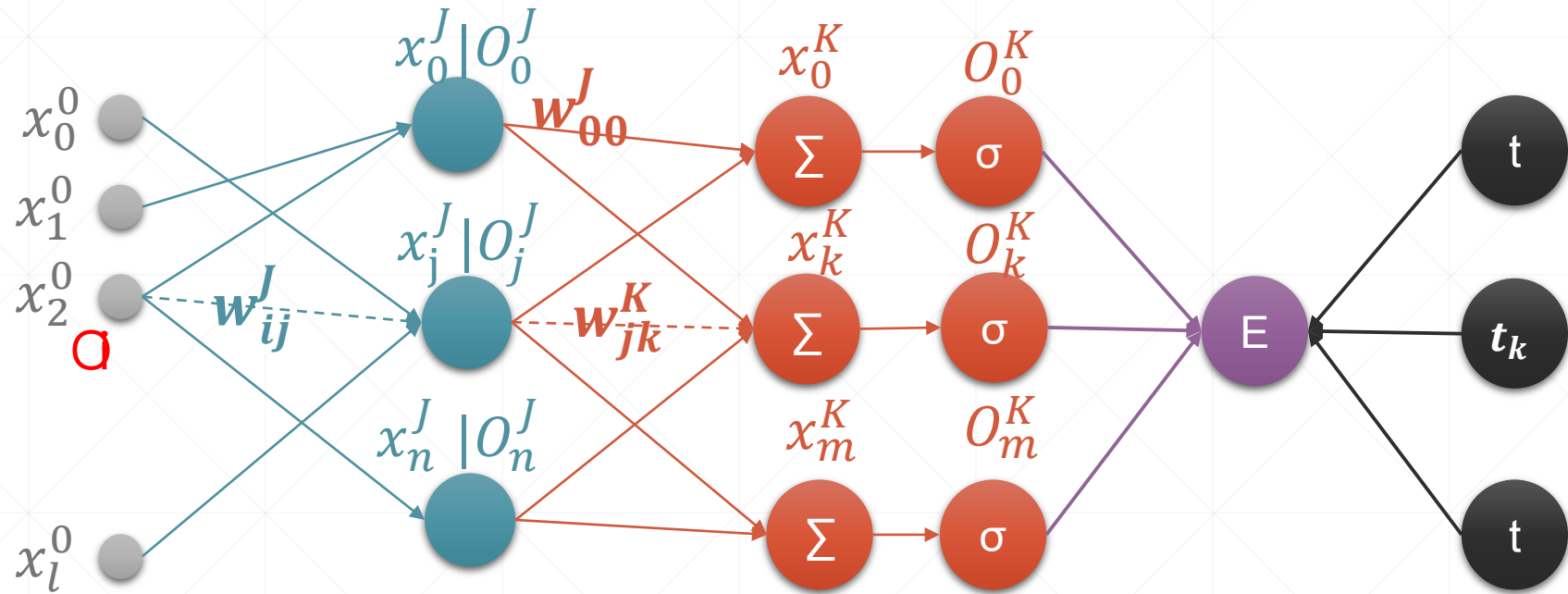
$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial \mathcal{O}_j}{\partial W_{ij}} \sum_{k \in K} (\mathcal{O}_k - t_k) \mathcal{O}_k (1 - \mathcal{O}_k) W_{jk}$$

$$\frac{\partial E}{\partial W_{ij}} = \mathcal{O}_j (1 - \mathcal{O}_j) \frac{\partial x_j}{\partial W_{ij}} \sum_{k \in K} (\mathcal{O}_k - t_k) \mathcal{O}_k (1 - \mathcal{O}_k) W_{jk}$$

$$\boxed{\frac{\partial E}{\partial W_{ij}} = \mathcal{O}_j (1 - \mathcal{O}_j) \mathcal{O}_i \sum_{k \in K} (\mathcal{O}_k - t_k) \mathcal{O}_k (1 - \mathcal{O}_k) W_{jk}}$$

$$\frac{\partial E}{\partial W_{ij}} = \mathcal{O}_j(1 - \mathcal{O}_j)\mathcal{O}_i \sum_{k \in K} (\mathcal{O}_k - t_k)\mathcal{O}_k(1 - \mathcal{O}_k)W_{jk}$$
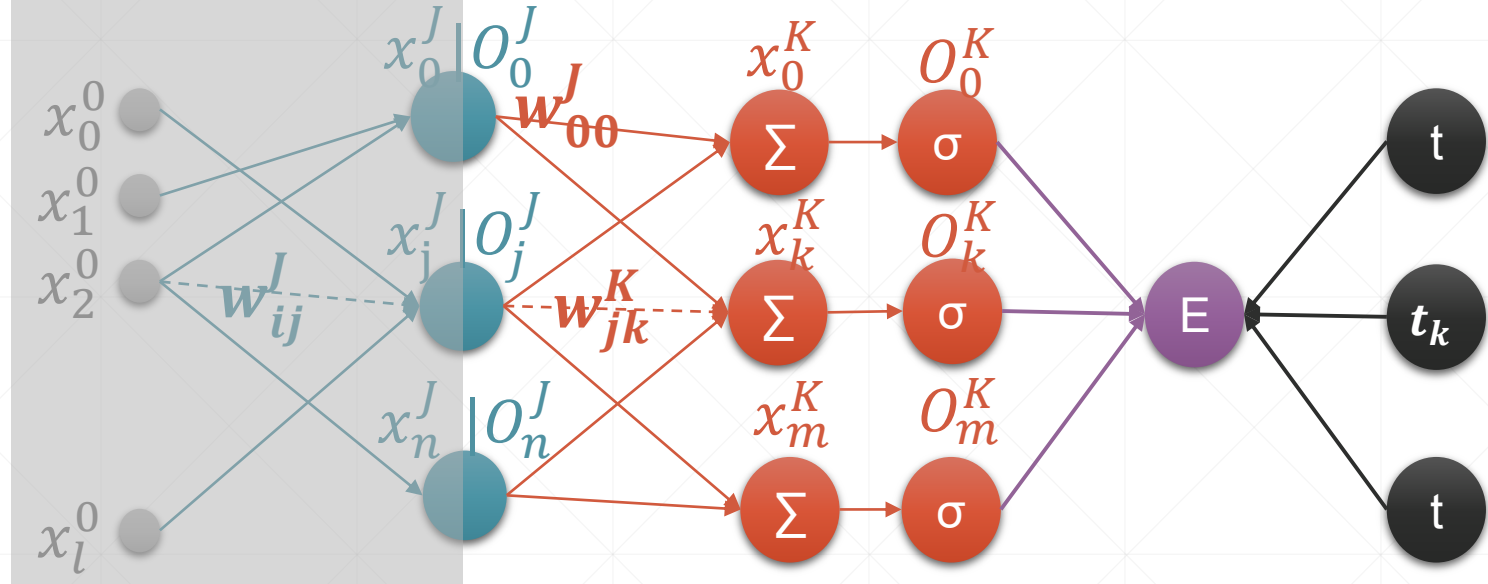
$$\frac{\partial E}{\partial W_{ij}} = \mathcal{O}_i\mathcal{O}_j(1 - \mathcal{O}_j) \sum_{k \in K} \delta_k W_{jk}$$

$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k)O_k(1 - O_k)\,O_j^J$$

$$\frac{\partial E}{\partial w_{jk}} = \delta_k^K \qquad O_j^J$$



$$\frac{\partial E}{\partial W_{ij}} = \mathcal{O}_j(1 - \mathcal{O}_j)\mathcal{O}_i \sum_{k \in K}(\mathcal{O}_k - t_k)\mathcal{O}_k(1 - \mathcal{O}_k)W_{jk}$$

$$\frac{\partial E}{\partial W_{ij}} = \mathcal{O}_i\mathcal{O}_j(1 - \mathcal{O}_j)\sum_{k \in K}\delta_k W_{jk}$$

$$\frac{\partial E}{\partial W_{ij}^{(l)}} = O_i^{(l-1)} \delta_j^{(l)}$$

For an output layer node $k \in K$

输出层: $\delta^{(l)} = O^{(l)} .* (1 - O^{(l)}) .* (O^{(l)} - t)$

隐藏层: $\delta^{(l)} = O^{(l)} .* (1 - O^{(l)}) .* \delta^{(l+1)} (W^{(l+1)})^T$

$$\frac{\partial E}{\partial W_{jk}} = O_j \delta_k$$

where

$$\delta_k = \underline{O_k(1 - O_k)}(O_k - t_k)$$
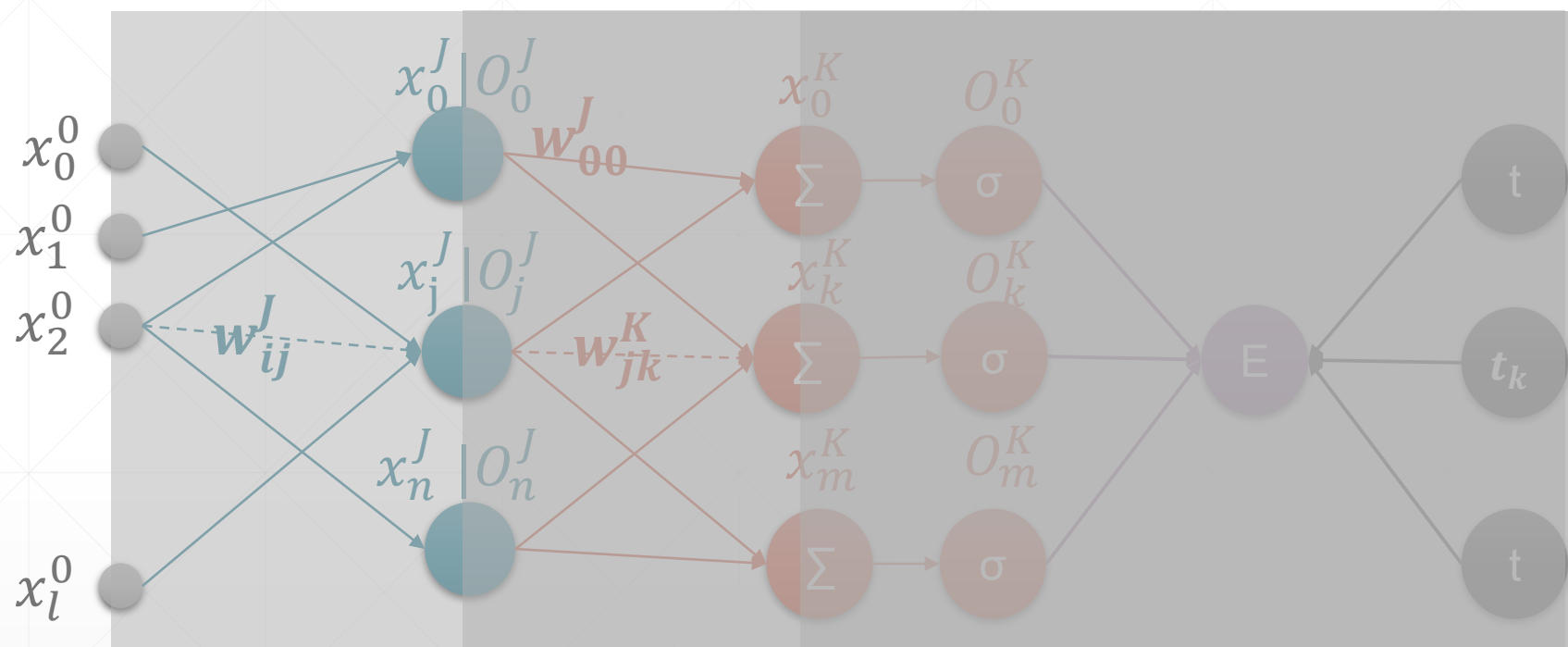
若最后一层无激活函数，则没有这一部分

For a hidden layer node $j \in J$

$$\frac{\partial E}{\partial W_{ij}} = O_i \delta_j$$
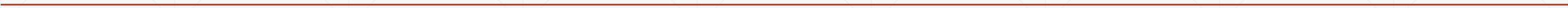
此处的参数矩阵与Ng课程中的互为转置，
因为Ng中每层为列向量，而此处为行向量

where

$$\delta_j = O_j(1 - O_j) \sum_{k \in K} \delta_k W_{jk}$$

# Congratulations!

下一课时

优化与训练

# Thank You.