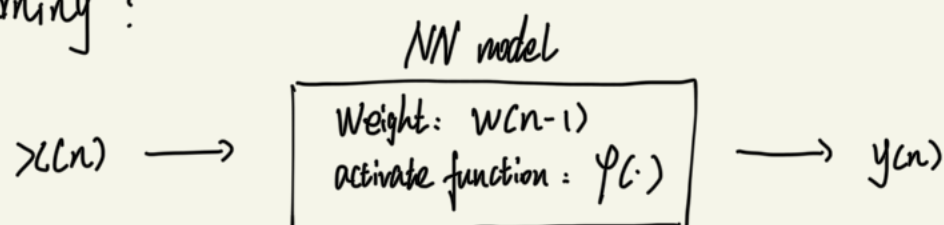


① What is Learning?



Neural Network accepts input and perform mapping through **Weights** to generate output.

When the model perform poorly . we adjust the weight

large $\mathcal{E}(n) = \sum_j \mathcal{E}_{y(n)}$

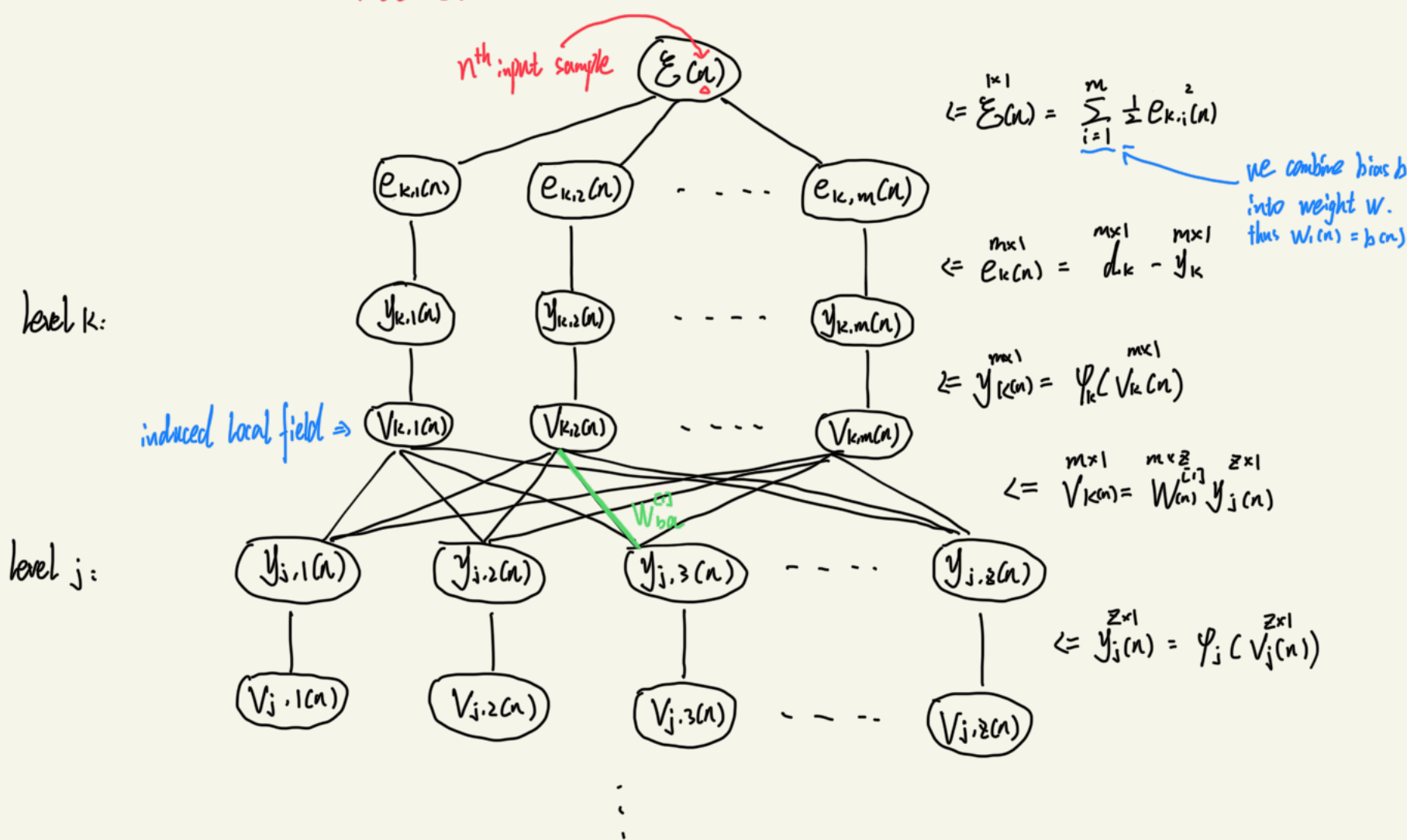
gradient descent

$$w(n+1) := w(n) + \Delta w$$

$$= w(n) - \eta \frac{\partial \mathcal{E}(n)}{\partial w}$$

② How we calculate $\frac{\partial \mathcal{E}(n)}{\partial w(n)}$?

measure how changes in weights $w(n)$ affect the training error $\mathcal{E}(n)$



I. for level "j-to-k": $\frac{\partial \mathcal{E}(n)}{\partial w_{ba}^{c,j}(n)}$, i.e. $y_a(n) \xrightarrow{w_{ba}^{c,j}(n)} y_b(n)$

Apply chain rule, we have:

$$\frac{\partial \mathcal{E}(n)}{\partial w_{ba}^{c,j}(n)} = \frac{\partial \mathcal{E}(n)}{\partial e_k(n)} \cdot \frac{\partial e_k(n)}{\partial y_k(n)} \cdot \frac{\partial y_k(n)}{\partial v_k(n)} \cdot \frac{\partial v_k(n)}{\partial w_{ba}^{c,j}(n)}$$

i. $\frac{\partial \mathcal{E}(n)}{\partial e_k(n)} = \frac{\partial \frac{1}{2} e_k(n)^T e_k(n)}{\partial e_k(n)} = e_k(n)$

ii. $\frac{\partial e_k(n)}{\partial y_k(n)} = \begin{bmatrix} \frac{\partial e_{k,1}(n)}{\partial y_{k,1}(n)} & \frac{\partial e_{k,1}(n)}{\partial y_{k,2}(n)} & \dots & \frac{\partial e_{k,1}(n)}{\partial y_{k,m}(n)} \end{bmatrix} = \begin{bmatrix} -1 & 0 & 0 & \dots & 0 \end{bmatrix} = -I_{m \times m}$

$$\begin{bmatrix} \frac{\partial e_{k,1}(n)}{\partial y_{k,1}(n)} & \frac{\partial e_{k,2}(n)}{\partial y_{k,2}(n)} & \dots & \frac{\partial e_{k,m}(n)}{\partial y_{k,m}(n)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial e_{k,m}(n)}{\partial y_{k,1}(n)} & \frac{\partial e_{k,m}(n)}{\partial y_{k,2}(n)} & \dots & \frac{\partial e_{k,m}(n)}{\partial y_{k,m}(n)} \end{bmatrix} \begin{bmatrix} 0 & -1 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & -1 & 0 \\ 0 & \dots & \dots & 0 & -1 \end{bmatrix}$$

Jacobian Matrix

$$\text{iii. } \frac{\partial y_k(n)}{\partial V_{k,c}(n)} = J = \begin{bmatrix} \phi'_k(V_{k,1}(n)) & 0 & \dots & 0 \\ 0 & \phi'_k(V_{k,2}(n)) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \phi'_k(V_{k,m}(n)) \end{bmatrix}_{m \times m}$$

$$\text{iv. } \text{know that } V_{k,c}(n) = W^{c,j}_k(n) y_{j,c}(n)$$

$$\text{i.e. } \begin{bmatrix} V_{k,1}(n) \\ V_{k,2}(n) \\ \vdots \\ V_{k,m}(n) \end{bmatrix}_{m \times 1} = \begin{bmatrix} W^{c,1}_k(n) \\ W^{c,2}_k(n) \\ \vdots \\ W^{c,m}_k(n) \end{bmatrix}_{m \times m} \begin{bmatrix} y_{j,1}(n) \\ y_{j,2}(n) \\ \vdots \\ y_{j,m}(n) \end{bmatrix}_{m \times 1} = \begin{bmatrix} [W^{c,1}_k(n)]^T y_{j,c}(n) \\ W^{c,1}_k(n) y_{j,1}(n) + \dots + W^{c,a}_k(n) y_{j,a}(n) + \dots + W^{c,m}_k(n) y_{j,m}(n) \\ \vdots \\ [W^{c,m}_k(n)]^T y_{j,c}(n) \end{bmatrix}$$

$$\text{we have } \frac{\partial V_{k,c}(n)}{\partial W^{c,j}_k(n)} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ y_{j,a}(n) \\ \vdots \\ 0 \end{bmatrix}$$

In summary.

$$\frac{\partial \mathcal{E}(n)}{\partial W^{c,j}_k(n)} = \begin{bmatrix} -e_{k,1}(n) \phi'_k(V_{k,1}(n)) & -e_{k,2}(n) \phi'_k(V_{k,2}(n)) & \dots & -e_{k,m}(n) \phi'_k(V_{k,m}(n)) \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ y_{j,a}(n) \\ \vdots \\ 0 \end{bmatrix} \quad (\text{index} = b)$$

define " $-\frac{\partial \mathcal{E}(n)}{\partial V_{k,b}(n)}$ " the local gradient for neuron b. Denote as $S_{k,b}(n)$

$$\text{for change on } W^{c,j}_k(n), \text{ we have: } \Delta W^{c,j}_k(n) = -\eta \frac{\partial \mathcal{E}(n)}{\partial W^{c,j}_k(n)}$$

$$= \eta S_{k,b}(n) y_{j,a}(n)$$

II. Local gradient for $y_a(n)$

$$S_a(n) = \frac{\partial \mathcal{E}(n)}{\partial y_{j,a}(n)} = \frac{\partial \mathcal{E}(n)}{\partial e_k(n)} \cdot \frac{\partial e_k(n)}{\partial y_k(n)} \cdot \frac{\partial y_k(n)}{\partial V_{k,c}(n)} \cdot \frac{\partial V_{k,c}(n)}{\partial y_{j,c}(n)} \cdot \frac{\partial y_{j,c}(n)}{\partial y_{j,a}(n)}$$

$$= [-e_{k,1}(n) \phi'_k(V_{k,1}(n)), +e_{k,2}(n) \phi'_k(V_{k,1}(n)) \dots, +e_{k,m}(n) \phi'_k(V_{k,m}(n))] \begin{bmatrix} W^{c,1}_k(n) & W^{c,2}_k(n) & \dots & W^{c,z}_k(n) \\ W^{c,1}_k(n) & W^{c,2}_k(n) & \dots & W^{c,z}_k(n) \\ \vdots & \vdots & \ddots & \vdots \\ W^{c,1}_k(n) & W^{c,2}_k(n) & \dots & W^{c,z}_k(n) \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \phi'_k(V_{j,a}(n)) \\ \vdots \\ 0 \end{bmatrix}$$

$$= \sum_{i=1}^m \underbrace{e_{k,i}(n) \psi'_k(V_{k,i}(n))}_{\text{local gradient}} W_{i,a}(n) \cdot \psi'_j(V_j, a(n))$$

$$= \left[\sum_{i=1}^m \delta_i(n) W_{i,a}(n) \right] \cdot \psi'_j(V_j, a(n))$$

③ Conclusion



$$\Delta W_{ji}(n) = \boxed{\eta} \cdot \boxed{\delta_j(n)} \cdot \boxed{y_i(n)}$$

2. for output neuron:

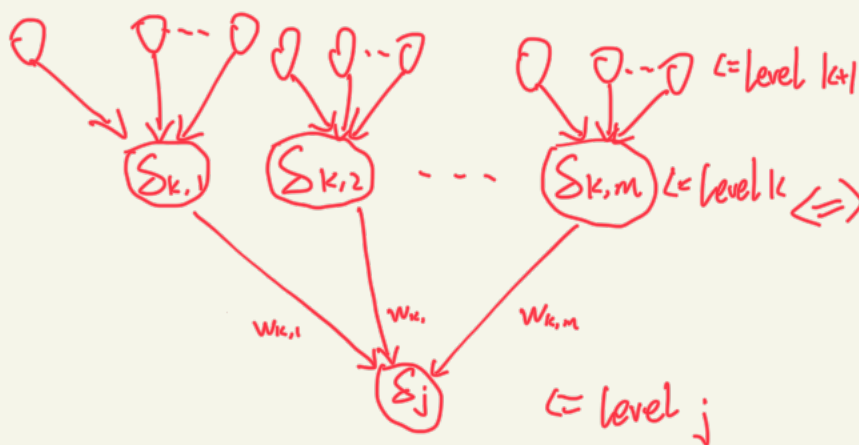
$$\delta_j(n) = e_j(n) \cdot \psi'_j(V_j(n))$$

1. for hidden neuron:

$$\delta_j(n) = \left[\sum_k \delta_k(n) W_{kj}(n) \right] \cdot \psi'_j(V_j(n))$$

all local gradient of the previous layer

if k is still not the output neuron



$$\delta_j(n) = \sum_k \left[\sum_{c(k,n)} \delta_{k,c}(n) W_{c(k,n),k}(n) \right] \cdot \psi'_j(V_j(n)) \Leftrightarrow$$

↑
high-level layer for k



Compute the local gradient of output neurons,
the the local gradients of hidden neurons layer-by-layer.