# Lecture 10. Backpropagation Algorithm. (BP)

1. BP is an approach to calculate the gradients for NNs.

BP is not an optimization algorithm.

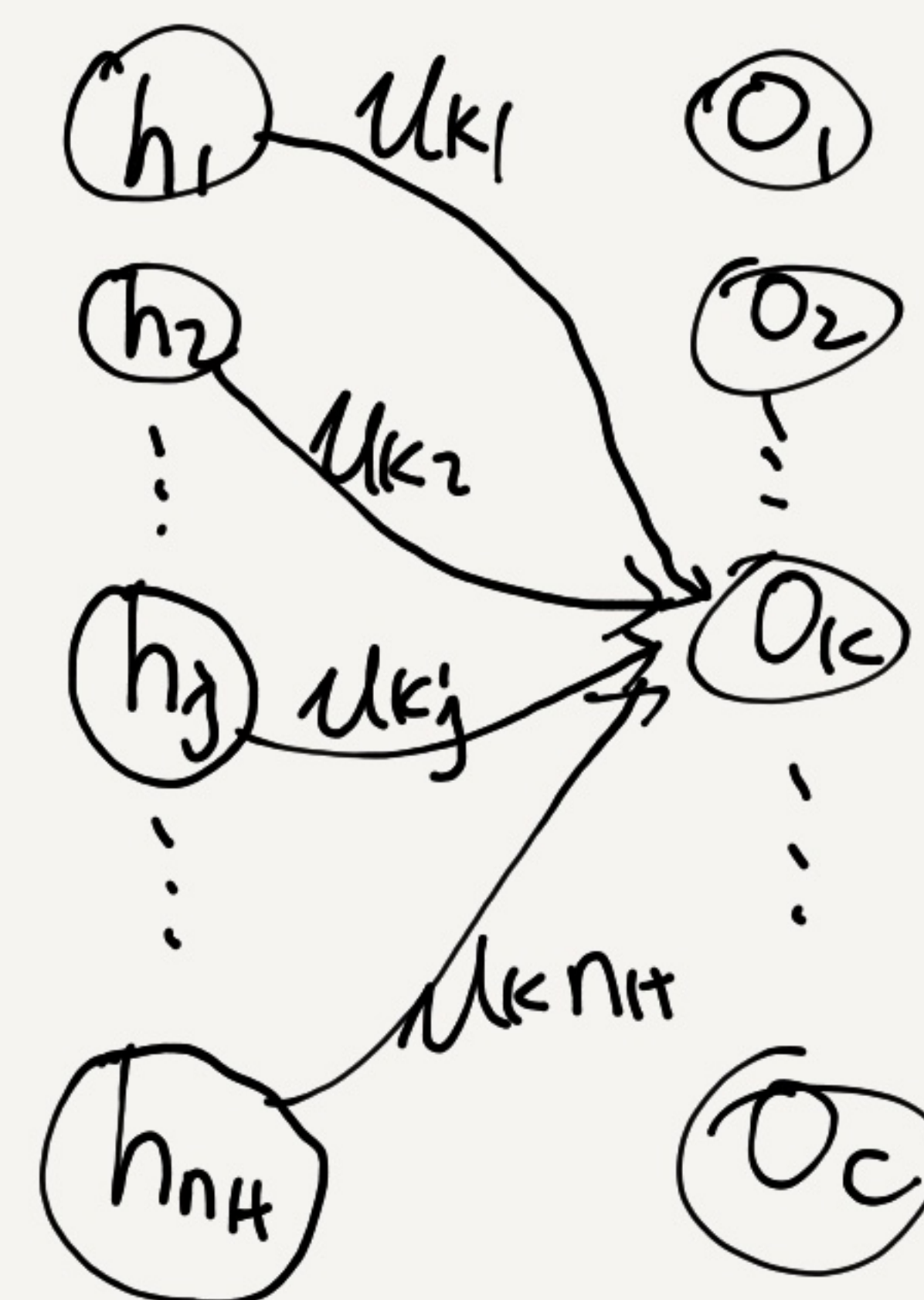In our 3-layer NN, BP will

(1) Calculate gradients for each weight/parameter between the output layer and the hidden layer.

$$\mathcal{U} = \begin{bmatrix} u_{11}, u_{12}, \cdots . u_{1n_H} \\ \vdots \\ u_{k1}, u_{k2}, \cdots u_{kn_H} \\ \vdots \\ u_{c1}, u_{c2}, \cdots , u_{cn_H} \end{bmatrix}_{c \times n_H} \begin{matrix} \longrightarrow O_1 \\ \\ \longrightarrow \text{all weights of } O_k \\ \\ \longrightarrow O_c \end{matrix}$$

For $k = 1, 2, \cdots, C$ (# of output node)

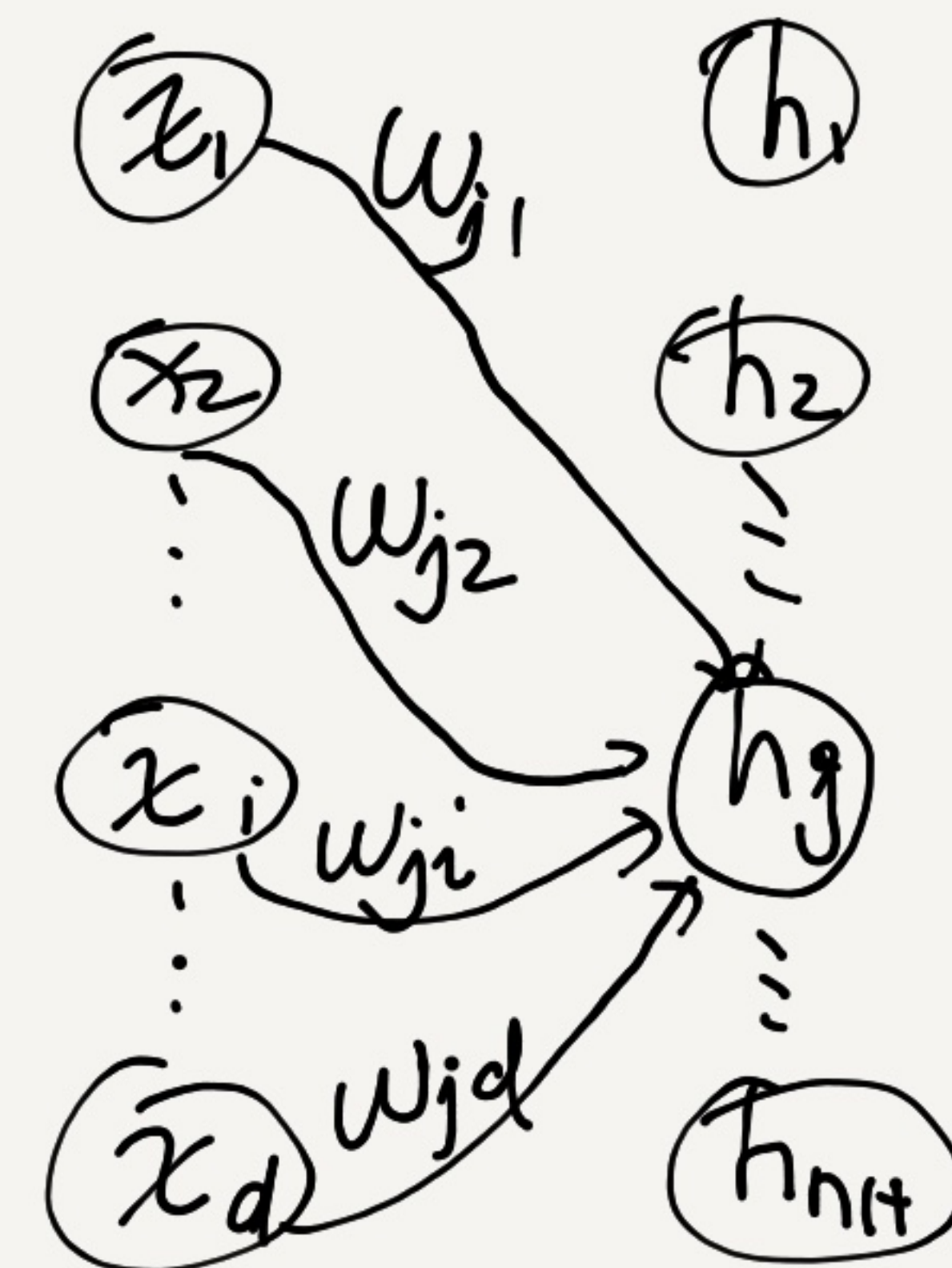$$\nabla_{u_{k1}} L, \nabla_{u_{k2}} L, \boxed{\nabla_{u_{k3}} L}, \cdots, \nabla_{u_{kn_H}} L$$

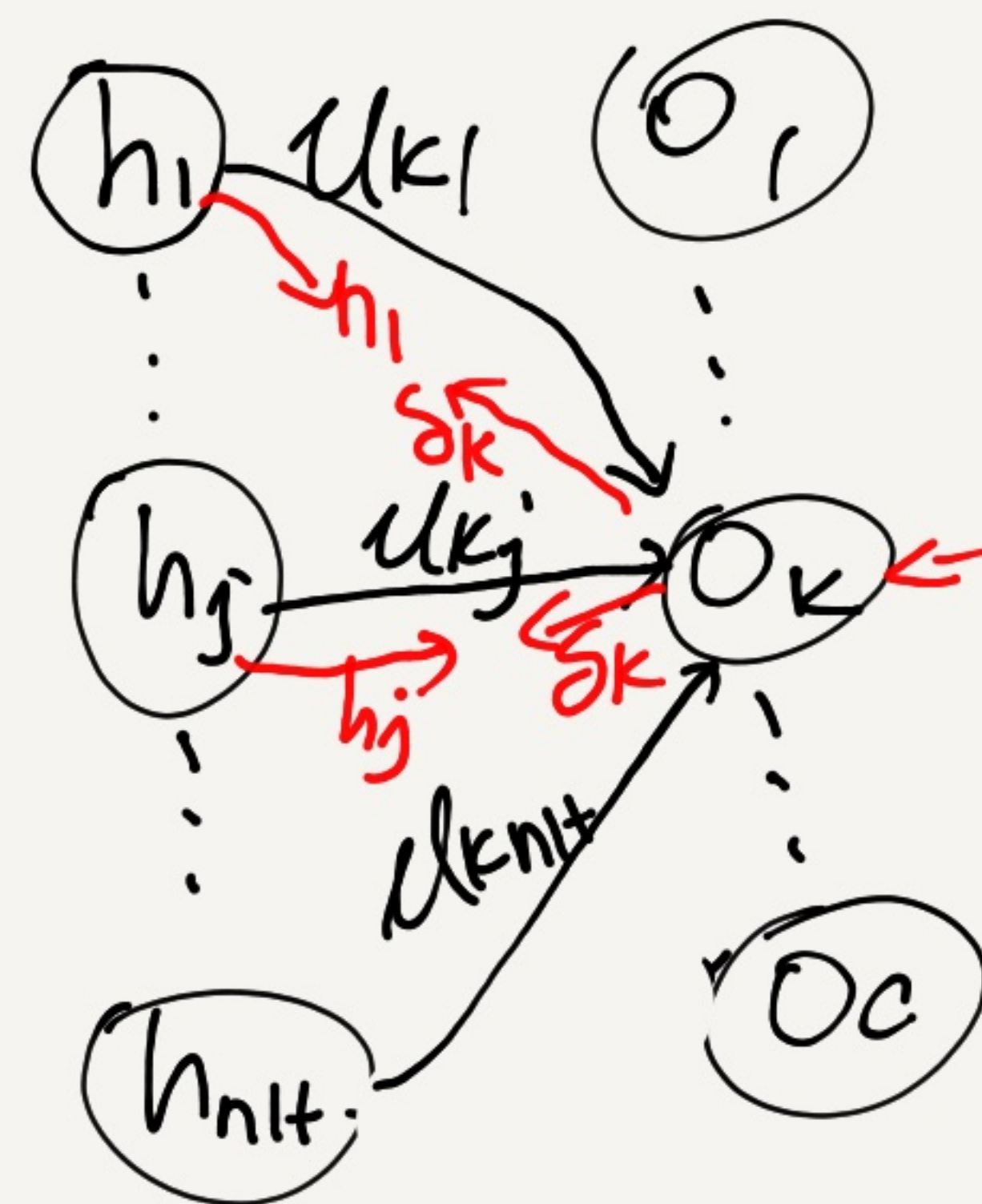(2) Calculate gradient for each weight between hidden and input layer.

$$W = \begin{bmatrix} w_{11}, w_{12}, \cdots, w_{1d} \\ \vdots \quad \vdots \quad\quad \vdots \\ w_{j1}, w_{j2}, \cdots, w_{jd} \\ \vdots \\ w_{n_{H}1}, w_{n_{H}2}, \cdots, w_{n_{H}d} \end{bmatrix} \begin{array}{l} \rightarrow h_1 \\ \\ \rightarrow h_j \\ \rightarrow h_{n_{H}} \end{array}$$



For $j = 1, 2, \cdots, n_{H}$

$$\nabla_{w_{j1}} L, \quad \boxed{\nabla_{w_{j2}} L}, \quad \cdots, \quad \nabla_{w_{jd}} L.$$

## 2. $\nabla_{u_{kj}} \mathcal{L}$ , $\quad K = 1, 2, \cdots, C.$
$\quad j = 1, 2, \cdots, n_H$



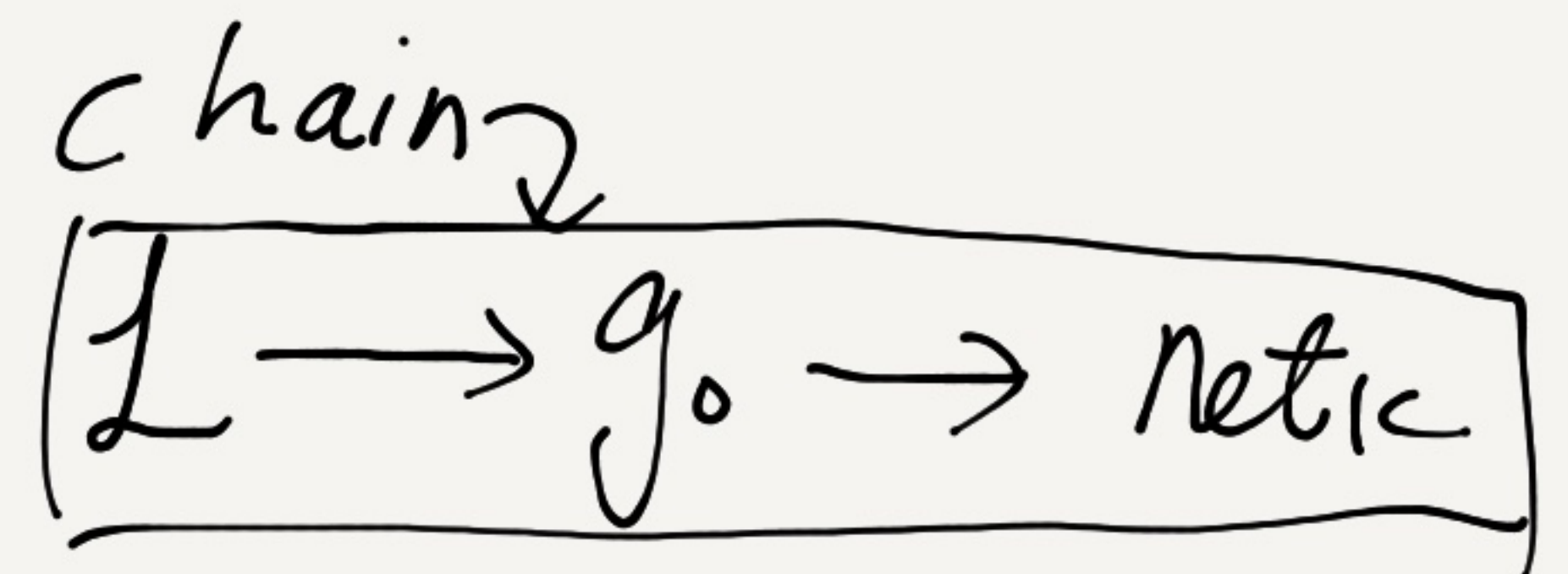Loss: $\mathcal{L}(w, u) = \left(\frac{1}{2}\right) \sum_{k=1}^{C} (y_k - o_k)^2$    prediction

model   $o_k = g_0(net_k)$    $net_k = \sum_{j=1}^{n_H} u_{kj} \cdot h_j$

$$\nabla_{u_{kj}} \mathcal{L} = \frac{\partial L}{\partial o_k} \cdot \frac{\partial o_k}{\partial net_k} \cdot \frac{\partial net_k}{\partial u_{kj}}$$

chain: $\boxed{L \longrightarrow g_0 \longrightarrow net_k}$

$$= \frac{1}{2} \cdot 2 \cdot (o_k - y_k) \cdot g_0' \cdot h_j$$

$$= \boxed{(o_k - y_k) \cdot g_0' \cdot} h_j = \delta_k \cdot h_j$$

$\delta_k$ (feed back from the $o_k$)

## 3. $\nabla_{w_{ji}} \mathcal{L}$



Loss: $\mathcal{L}(w, u) = \frac{1}{2} \sum_{k=1}^{C} (y_k - o_k)^2$

$\boxed{L \to g_0 \to net_k \to g_h \to net_j}$

model: $o_k = g_0(net_k)$ ,   $net_k = \sum_{j=1}^{n_H} u_{kj} \cdot h_j$

$h_j = g_h(net_j)$ ,   $net_j = \sum_{i=1}^{d} w_{ji} \cdot x_i$

$$\nabla_{w_{ji}} \mathcal{L} = \sum_{k=1}^{C} \left( \frac{\partial L}{\partial o_k} \cdot \frac{\partial o_k}{\partial net_k} \cdot \frac{\partial net_k}{\partial h_j} \right) \cdot \frac{\partial h_j}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ji}}$$

$\delta_j \cdot x_i$

$$= \left\{ \sum_{k=1}^{C} \boxed{(o_k - y_k) \cdot g_0' \cdot u_{kj}} \right\} \cdot g_h' \cdot x_i = \boxed{\sum_{k=1}^{C} (\delta_k \cdot u_{kj}) \cdot g_h'} \cdot x_i$$

$\delta_j \parallel$

## 4. summary of BP.

Output — hidden :　$\nabla_{u_{kj}} \mathcal{L} = \delta_k \cdot h_j$

hidden — input layer:　$\nabla_{w_{ji}} \mathcal{L} = \delta_j \cdot x_i$

$$\delta_k = (O_k - y_k) \cdot g_o'$$

$$\delta_j = \sum_{k=1}^{C} \delta_k \cdot u_{kj}$$

If we have multiple hidden layers ( more than 1 ), how can we extend the BP Algorithm?



4-layer NN.