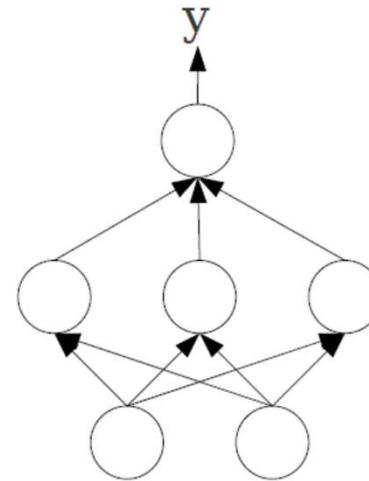


Training of Multi-Layer NN: Back-Propagation

How to Train Multi-Layer NN

Define sum-squared error:

$$E = \frac{1}{2} \sum_p (d^p - y^p)^2$$

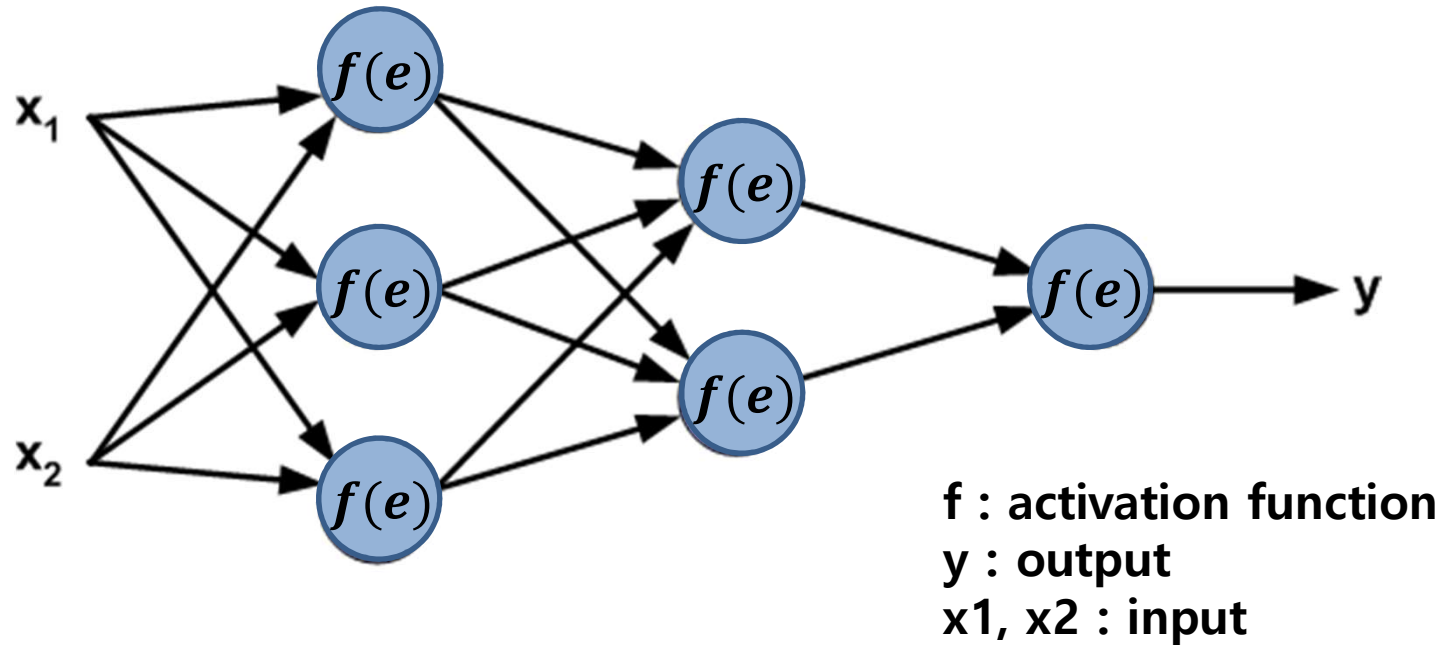


Use gradient descent error minimization:

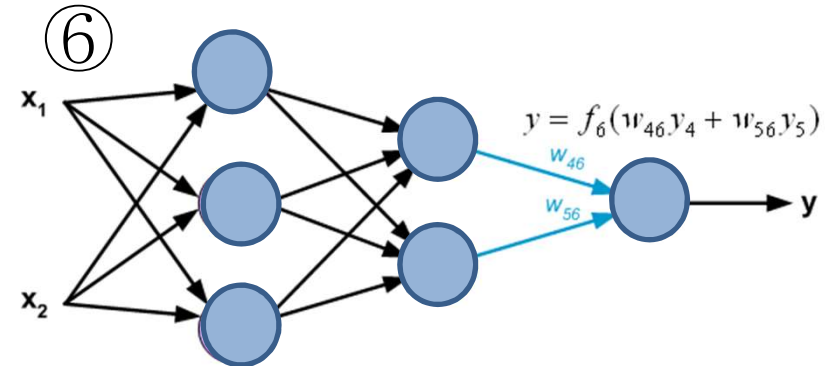
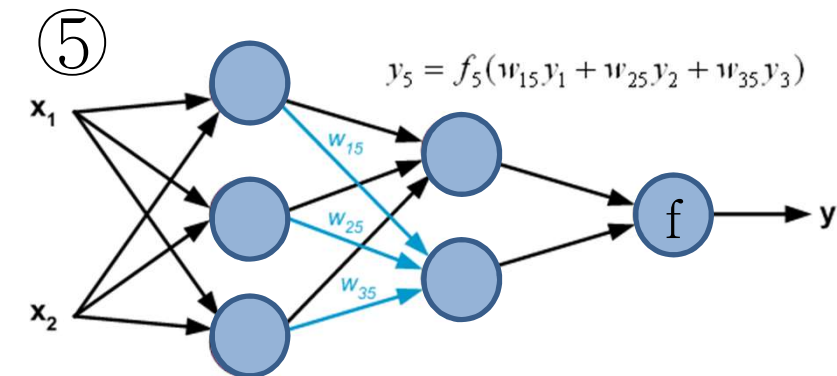
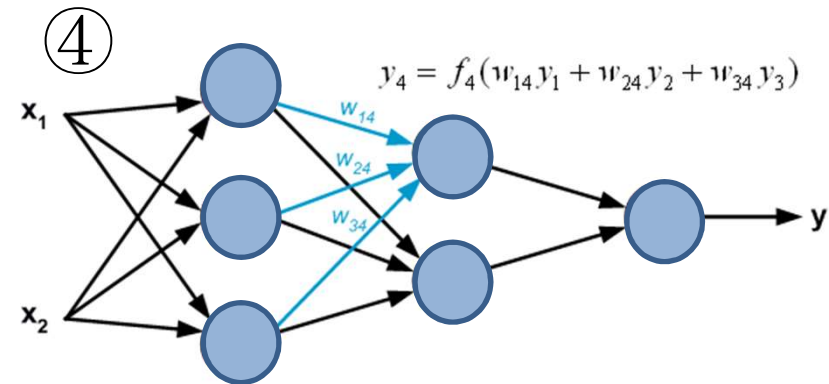
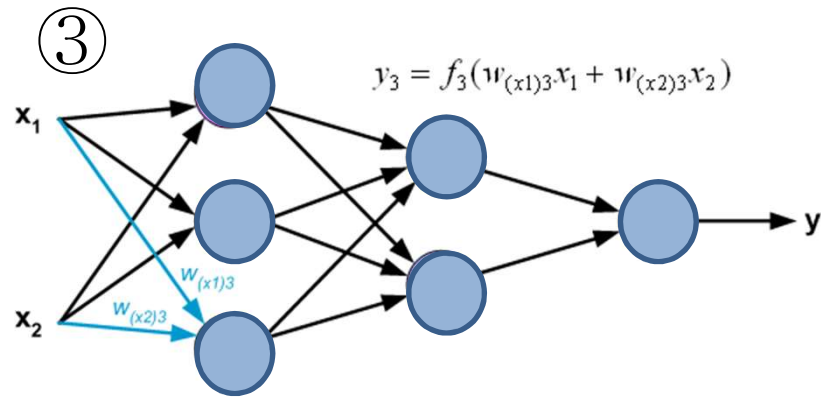
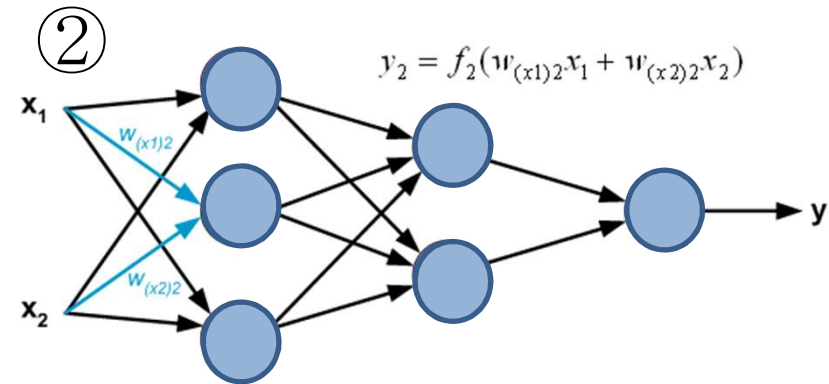
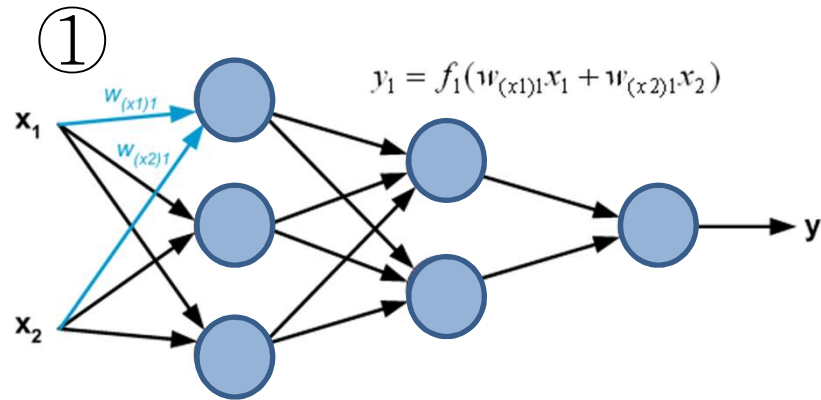
$$\Delta w_{ij} = -\eta \frac{\partial E}{\partial w_{ij}}$$

Works if the nonlinear transfer function is differentiable.

Multilayer Neural Network

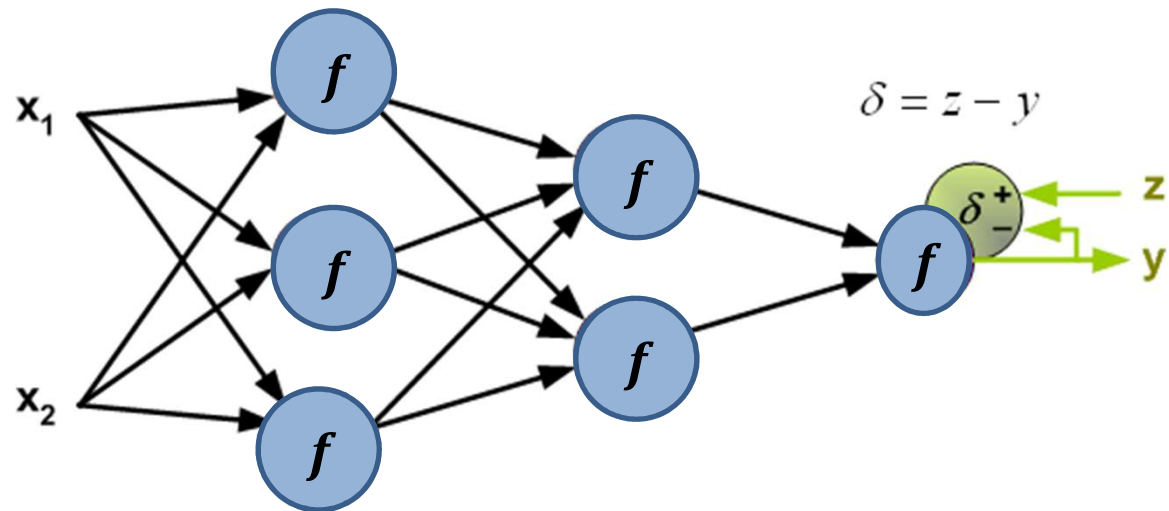


Forward-Propagation



Back-Propagation

Minimize error(δ) by finding the weights (w)

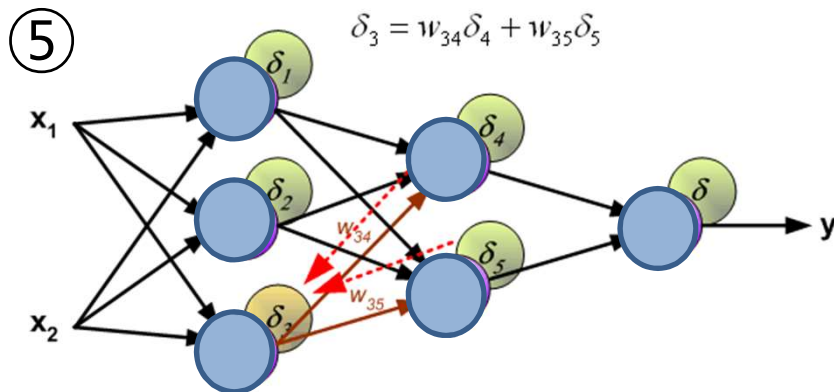
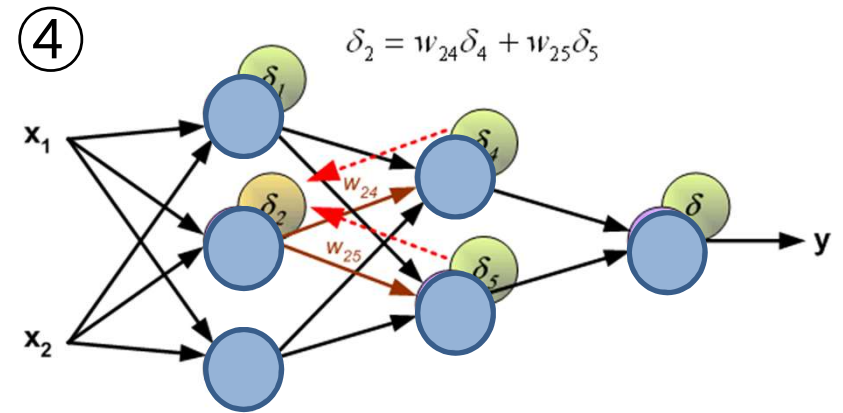
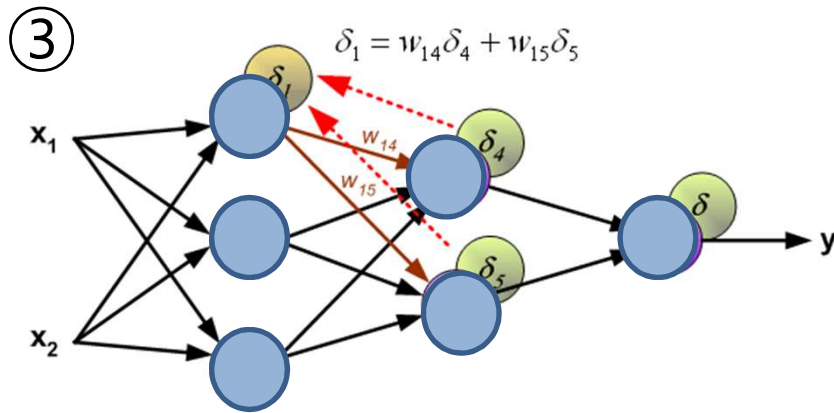
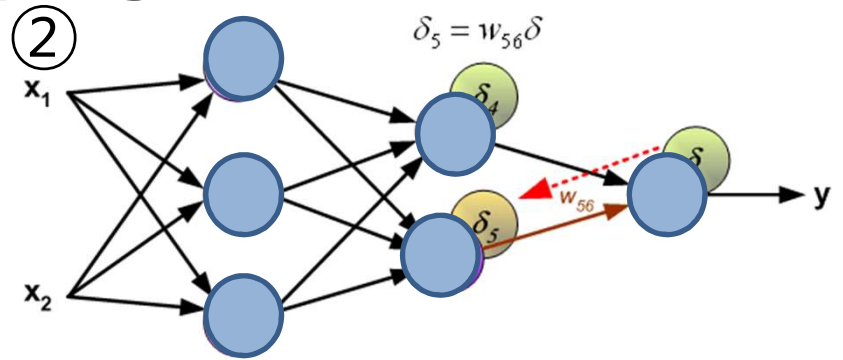
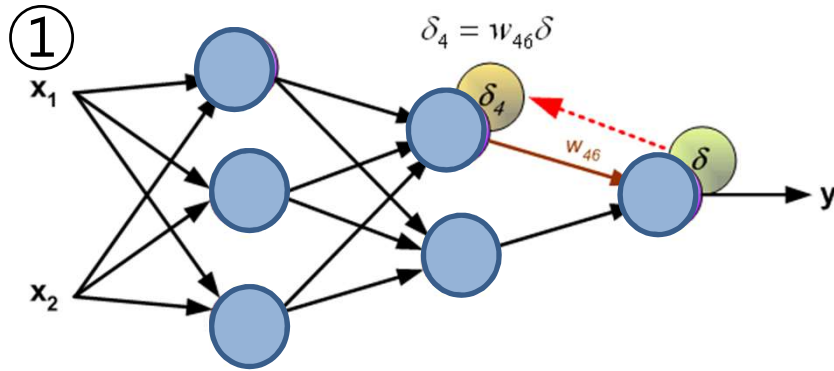


δ = error

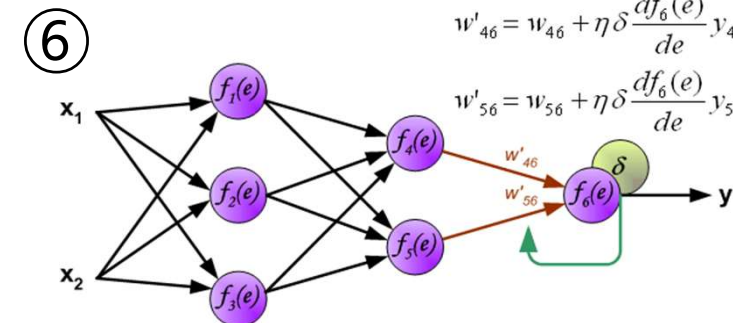
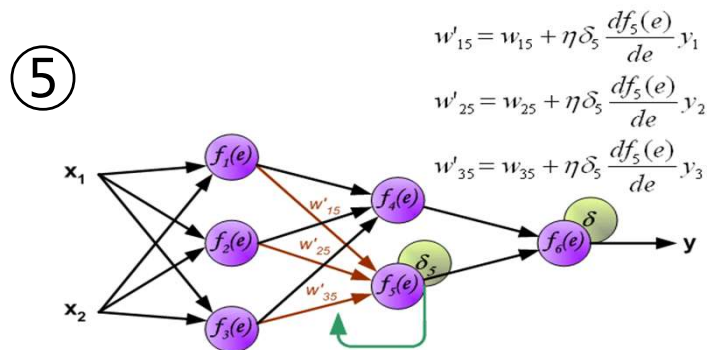
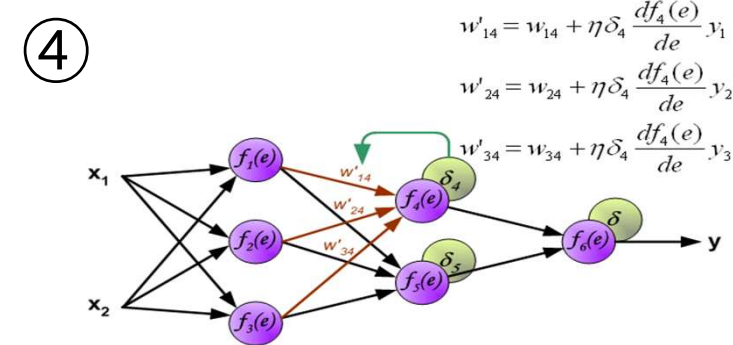
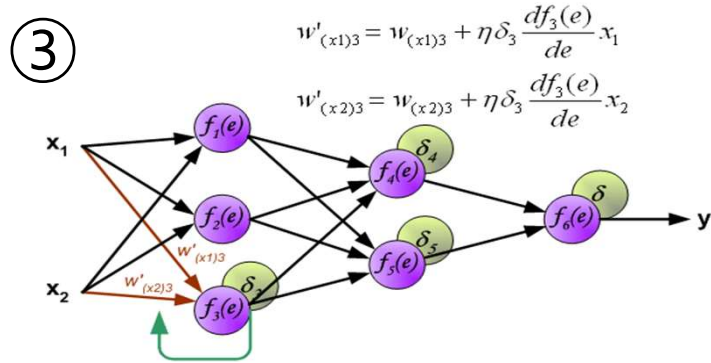
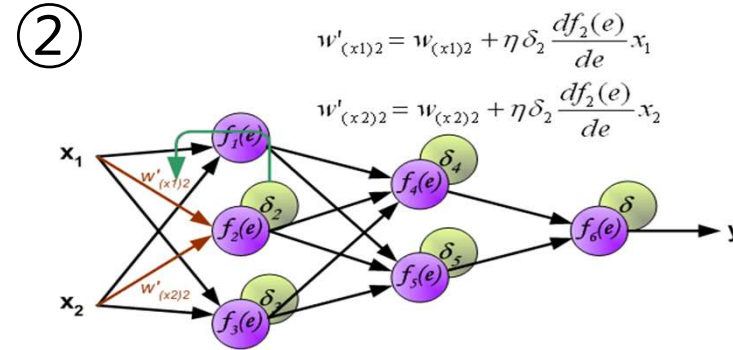
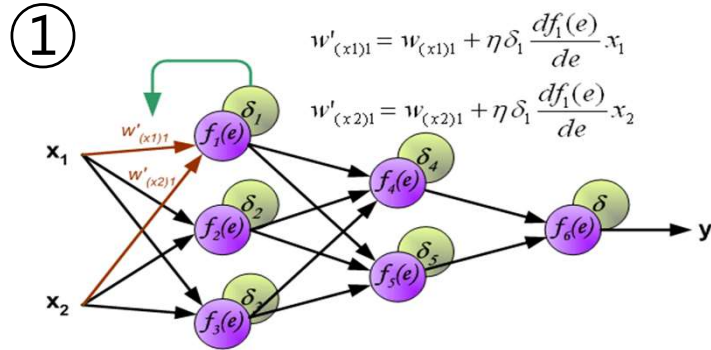
z = desired output

y = actual output

Back-Propagation



Update the weights, w



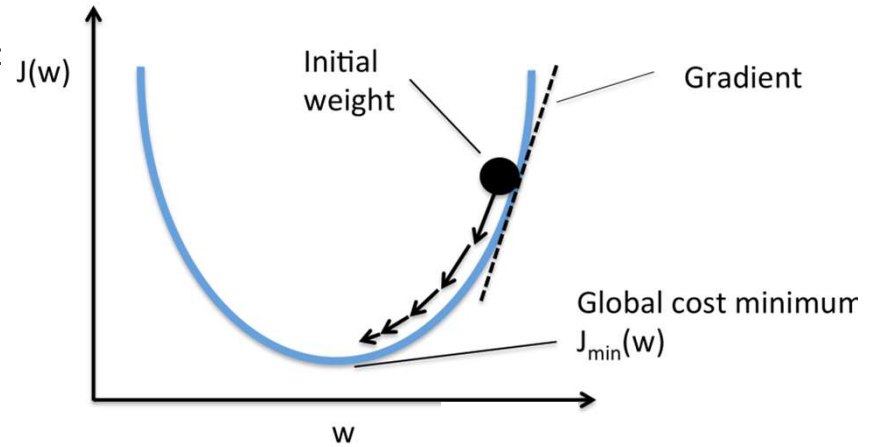
Gradient Descent

The cost function $J(\cdot)$, the sum of squared errors (SSE), can be written as:

$$J(\mathbf{w}) = \frac{1}{2} \sum_i (\text{target}^{(i)} - \text{output}^{(i)})^2$$

$$SSE = \sum_i (\text{target}^{(i)} - \text{output}^{(i)})^2$$

$$MSE = \frac{1}{n} \times SSE$$



The magnitude and direction of the weight update is computed by taking a step in the opposite direction of the cost gradient

$$\Delta w_j = -\eta \frac{\partial J}{\partial w_j},$$

where η is the learning rate. The weights are then updated after each epoch via the following update rule:

$$\mathbf{w} := \mathbf{w} + \Delta \mathbf{w},$$

where $\Delta \mathbf{w}$ is a vector that contains the weight updates of each weight coefficient w , which are computed as follows:

$$\begin{aligned} \Delta w_j &= -\eta \frac{\partial J}{\partial w_j} \\ &= -\eta \sum_i (\text{target}^{(i)} - \text{output}^{(i)}) (-x_j^{(i)}) \\ &= \eta \sum_i (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}. \end{aligned}$$

<http://sebastianraschka.com/faq/docs/closed-form-vs-gd.html>

Example

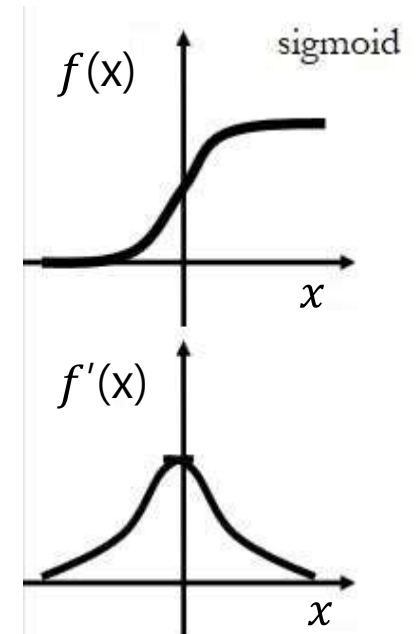
If activation function is a sigmoid function

$$E \text{ or } J = \frac{1}{2} \sum_{i=1}^N (z_i - f(x_i w_i))^2$$

$$\frac{\partial E}{\partial w} = - \sum_{i=1}^N (z_i - f(x_i w_i)) f'(x_i w_i) x_i$$

$$= - \sum_{i=1}^N (z_i - f(x_i w_i)) f(x_i w_i) (1 - f(x_i w_i)) x_i$$

$$f(x) = \frac{1}{1 + e^{-x}}$$
$$f'(x) = \frac{-e^{-x}}{(1 + e^{-x})^2}$$



Training rule :

$$w_i \leftarrow w_i + \Delta w_i$$
$$\Delta w_i = \eta (z_i - y_i) f(x_i w_i) (1 - f(x_i w_i)) x_i$$

η = learning rate

as in the textbook.