



w (weight) is for different edges (each pair)

b (bias)

$$a_j^L = \sigma \left(\sum_k w_{jk}^L a_k^{(L-1)} + b_j^L \right)$$

σ = activation function

$$C_1 = \frac{1}{2} \sum_j (y_j - a_j^L)^2$$

$$C = \frac{1}{n} \sum_n \frac{1}{2} \sum_j (y_j - a_j^L)^2 = \frac{1}{n} \sum_n C_1 = \frac{1}{2n} \sum_j (y_j - a_j^L)^2$$

n = # of sample for training

Gradient Descent

$$w = w - m \frac{\partial C}{\partial w}$$

m = learning rate

too high = fail to converge

too low = inefficient, slow to converge