

Layer

$$\begin{aligned} \text{input : } x \in \mathbb{R}^m &\rightarrow Wx + b \in \mathbb{R}^n \\ &\rightarrow f(Wx + b) \in \mathbb{R}^n : \text{output} \end{aligned} \quad (1)$$

Network

$$\text{input : } x_1 \rightarrow y_1 = f_1(W_1x_1 + b_1) \quad (2)$$

$$x_2 = y_1 \rightarrow y_2 = f_2(W_2x_2 + b_2) \quad (3)$$

\vdots

$$x_N = y_{N-1} \rightarrow y_N = f_N(W_Nx_N + b_N) : \text{output} \quad (4)$$

Training

Training data: a collection of pairs (x, y_{true})

$$\text{Loss function: } L = \frac{1}{N_{\text{samples}}} \sum_{x \in \text{samples}} \|y_N(x) - y_{\text{true}}\|^2 \quad (\text{for example}) \quad (5)$$

Parameters repeatedly updated according to

$$W \rightarrow W - r \left. \frac{\partial L}{\partial W} \right|_{\text{batch}} \quad (6)$$

$$b \rightarrow b - r \left. \frac{\partial L}{\partial b} \right|_{\text{batch}} \quad (7)$$

where "batch" is a randomly selected subset of the training data.

Gradient

Let λ_n denote a parameter (W or b) in the n -th layer. Then

$$\frac{\partial L}{\partial \lambda_n} = \frac{\partial L}{\partial y_N^i} \frac{\partial y_N^i}{\partial x_N^j} \frac{\partial y_{N-1}^j}{\partial x_{N-1}^k} \dots \frac{\partial y_n^l}{\partial \lambda_n} \quad (8)$$

and introducing the notation

$$(J_{f_m})^i_j := \left. \frac{\partial f_m^i(x)}{\partial x^j} \right|_{x=W_mx_m+b_m} \quad (9)$$

the derivatives in Eq. (8) are given by

$$\frac{\partial y_m^i}{\partial x_m^j} = (J_{f_m})^i_k (W_m)^k_j \quad (10)$$

$$\frac{\partial y_n^i}{\partial (W_n)^j_k} = (J_{f_n})^i_j x_n^k \quad (11)$$

$$\frac{\partial y_n^i}{\partial b_n^j} = (J_{f_n})^i_j \quad (12)$$