# Multi-layer perceptron



$h_\ell = \gamma_\ell(s_\ell)$
element wise
(coordinatewise)

① $s_\ell = W_\ell h_{\ell-1} + b_\ell$

$\ell = 1, \dots, L$

① $s_{\ell,k} = \langle w_{\ell,k}, h_{\ell-1} \rangle + b_{\ell,k}$

② $h_{\ell,k} = \text{sigmoid}(s_{\ell,k}) = \dfrac{e^{s_{\ell,k}}}{1 + e^{s_{\ell,k}}}$

or $\text{ReLU}(s_{\ell,k}) = \max(0, s_{\ell,k})$

leaky ReLU

$\theta = (w_\ell, b_\ell, \ell = 1, \dots, L)$

Forward process · how to predict y based on x

$$X \longrightarrow h_1 \longrightarrow h_2 \longrightarrow \dots \longrightarrow h_{\ell-1} \longrightarrow h_\ell \longrightarrow \dots \longrightarrow h_L \longrightarrow \hat{y}$$

$$w_1, b_1 \qquad w_2, b_2 \qquad w_{\ell-1}, b_{\ell-1} \quad w_\ell, b_\ell \qquad w_L, b_L$$

Backward process · error back propagation via chain rule
enables training of NN by gradient descent, calc derivative

$$\leftarrow h_1 \leftarrow \dots \leftarrow h_{\ell-1} \overset{①}{\leftarrow} h_\ell \overset{①}{\leftarrow} \dots \overset{①}{\leftarrow} h_L \overset{①}{\leftarrow} e$$

$$w_1, b_1 \qquad\qquad w_{\ell-1}, b_{\ell-1} \quad w_\ell, b_\ell \qquad w_L, b_L$$

$\dfrac{\partial \text{Loss}}{\partial h_{\ell-1}} \overset{①}{\longleftarrow} \dfrac{\partial \text{Loss}}{\partial h_\ell}$

② $\dfrac{\partial \text{Loss}}{\partial w_\ell}$

$$\frac{\partial L}{\partial h_{\ell-1}^T} = \sum_{k=1}^{d} \frac{\partial L}{\partial h_{\ell,k}} \frac{\partial h_{\ell,k}}{\partial s_{\ell,k}} \frac{\partial s_{\ell,k}}{\partial h_{\ell-1}^T}$$

$h_{\ell-1}$   col vector

$h_{\ell-1}^T$   row vector

$$= \sum_{k=1}^{d} \frac{\partial L}{\partial h_{\ell,k}} r_\ell'(s_{\ell,k}) \, w_{\ell,k}$$

$$= \begin{pmatrix} \frac{\partial L}{\partial h_{\ell,k}} r_\ell'(s_{\ell,k}) \\ 1, 2, \dots \quad k, \quad \dots \quad d \end{pmatrix} \begin{pmatrix} w_{\ell,k} \end{pmatrix} \begin{matrix} 1 \\ 2 \\ \vdots \\ k \\ \vdots \\ d \end{matrix}$$

$$= \begin{pmatrix} \frac{\partial L}{\partial h_{\ell,k}} \\ 1, 2, \dots \quad k \quad \dots \quad d \end{pmatrix} \begin{pmatrix} r_\ell'(s_{\ell,k}) w_{\ell,k} \end{pmatrix} \begin{matrix} 1 \\ 2 \\ \vdots \\ k \\ \vdots \\ d \end{matrix}$$

row wise multiplication

$$\frac{\partial L}{\partial h_{\ell-1}^T} = \frac{\partial L}{\partial h_\ell^T} \, r_\ell' \odot w_\ell \qquad ①$$

$$r_\ell' = \begin{pmatrix} r_\ell'(s_{\ell,k}) \end{pmatrix} \begin{matrix} 1 \\ 2 \\ \vdots \\ k \\ \vdots \\ d \end{matrix}$$

$$\frac{\partial L}{\partial h_{\ell-1}^T} = \frac{\partial L}{\partial h_\ell^T} \, \text{diag}\left(r_\ell'\right) w_\ell$$
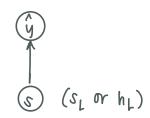
$$② \quad \frac{\partial L}{\partial w_{\ell,k}} = \frac{\partial L}{\partial h_{\ell,k}} \frac{\partial h_{\ell,k}}{\partial s_{\ell,k}} \frac{\partial s_{\ell,k}}{\partial w_{\ell,k}}$$

$$= \frac{\partial L}{\partial h_{\ell,k}} \, r_\ell'(s_{\ell,k}) \, h_{\ell-1}^T$$

$$\frac{\partial L}{\partial w_\ell} = \begin{pmatrix} \frac{\partial L}{\partial w_{\ell,k}} \end{pmatrix} \begin{matrix} 1 \\ 2 \\ \vdots \\ k \\ \vdots \\ d \end{matrix} = \begin{pmatrix} \frac{\partial L}{\partial w_{\ell,k}} r_\ell'(s_{\ell,k}) \end{pmatrix} \begin{matrix} 1 \\ 2 \\ \vdots \\ k \\ \vdots \\ d \end{matrix} h_{\ell-1}^T = r_\ell' \odot \frac{\partial L}{\partial h_\ell} h_{\ell-1} = \text{diag}\left(r_\ell'\right) \frac{\partial L}{\partial h_\ell} h_{\ell-1}$$

Top layer

linear model
$$s = x^T \beta$$

neural network

$$x \rightarrow h_1 \rightarrow \dots \rightarrow h_{\ell-1} \rightarrow s \rightarrow \hat{y}$$
$$\qquad\qquad\qquad\qquad (h_L)$$

$\hat{y}$

$s$ $(s_L$ or $h_L)$

linear regression

$$y \sim N(s, \sigma^2) \quad Loss = \tfrac{1}{2}(y-s)^2 \quad \frac{dL}{ds} = -(y-s) = -e$$

logistic regression

$$y \sim Bernoulli(p = sigmoid(s)) \quad Loss = -(y_i s_i - \log(1 + e^s)) \qquad \frac{dL}{ds} = -(y-p) = -e$$

$\hat{y}$

$s$

regression:
$$Loss = \tfrac{1}{2}|y-s|^2$$
$$\frac{dL}{ds} = -(y-s) = -e$$

logistic: multiclass classification

C classes $\qquad (1, 2, \dots, c, \dots, C)$

$$y = \text{one-hat} =
\begin{array}{c|c}
0 & 1 \\
0 & 2 \\
\vdots & \vdots \\
1 & c \\
\vdots & \vdots \\
0 & C
\end{array}$$

soft max
$$P(y = c \mid s) = \frac{e^{s_c}}{\sum_{c'=1}^{C} e^{s_{c'}}}$$

$$s = \text{logit } s =
\begin{array}{c|c}
s_1 & 1 \\
s_2 & 2 \\
\vdots & \vdots \\
s_c & c \\
\vdots & \vdots \\
s_C & C
\end{array}$$