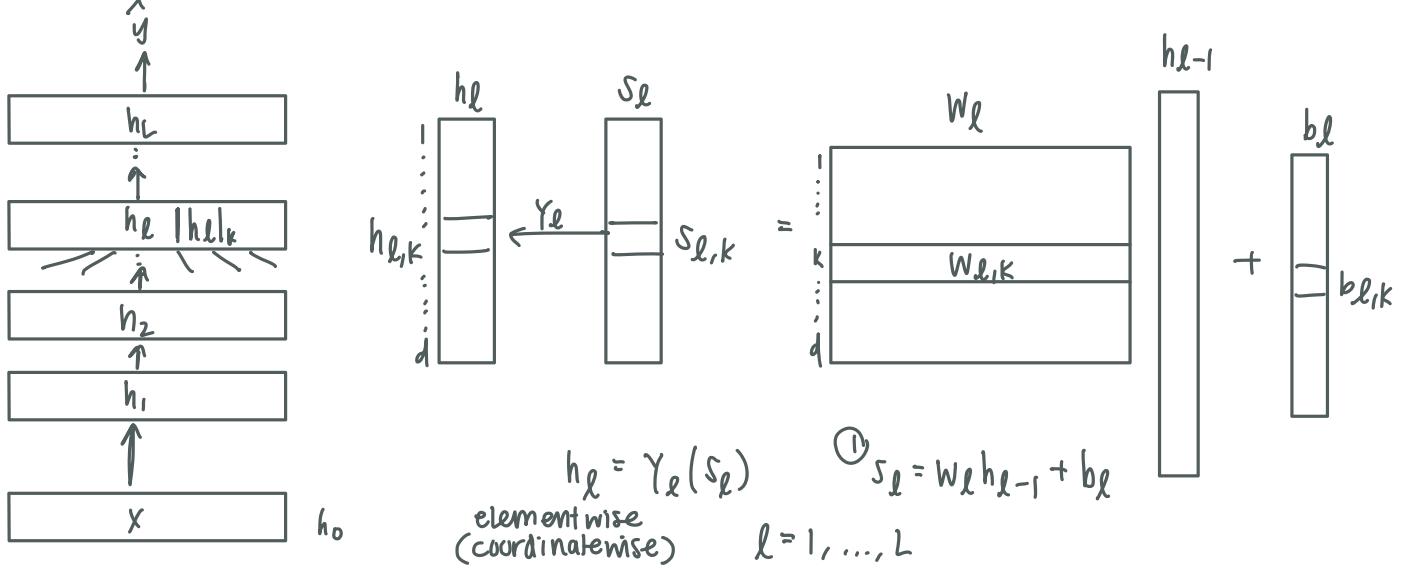


Multilayer perceptron



$$\textcircled{1} \quad s_l = w_l h_{l-1} + b_l$$

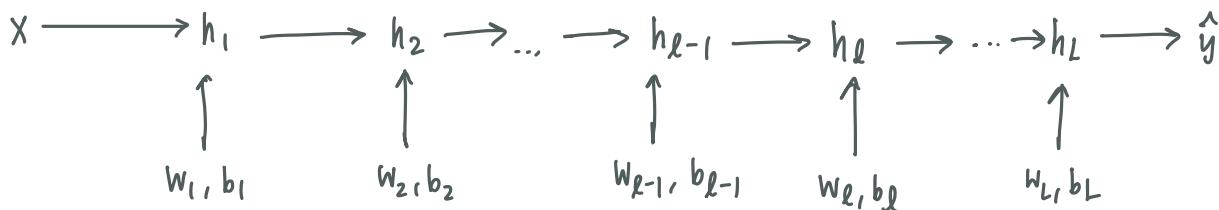
$$\textcircled{2} \quad s_{l,k} = \langle w_{l,k}, h_{l-1} \rangle + b_{l,k}$$

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

$$\theta = (w_l, b_l, l=1, \dots, L)$$



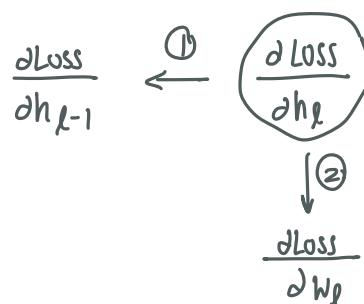
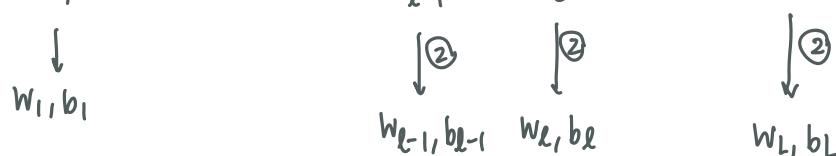
Forward process · how to predict y based on x



Backward process · error back propagation via chain rule

enables training of NN by gradient descent, calc derivative

$$e \leftarrow h_L \leftarrow \dots \leftarrow h_l \leftarrow \dots \leftarrow h_1 \leftarrow e$$



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

h_{l-1} col vector
 h_{l-1}^T row vector

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

$$= \left(\begin{array}{c} \frac{\partial L}{\partial h_{l,k}} r'_l(s_{l,k}) \\ \vdots \\ 1, 2, \dots, k, \dots, d \end{array} \right) \left(\begin{array}{c} 1 \\ \vdots \\ w_{l,k} \\ \vdots \\ d \end{array} \right)$$

$$= \left(\begin{array}{c} \frac{\partial L}{\partial h_{l,k}} \\ \vdots \\ 1, 2, \dots, k, \dots, d \end{array} \right) \left(\begin{array}{c} r'_l(s_{l,k}) w_{l,k} \\ \vdots \\ k \\ \vdots \\ d \end{array} \right)$$

row wise multiplication

$$\frac{\partial L}{\partial h_{l-1}^T} = \frac{\partial L}{\partial h_l^T} r'_l \odot w_l \quad \textcircled{1}$$

$$r'_l = \left(\begin{array}{c} r'_l(s_{l,k}) \\ \vdots \\ k \\ \vdots \\ d \end{array} \right)$$

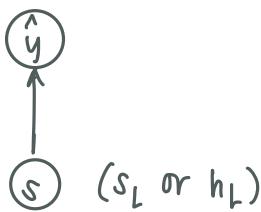
$$\frac{\partial L}{\partial h_{l-1}^T} = \frac{\partial L}{\partial h_l^T} \text{diag}(r'_l) w_l$$

$$\textcircled{2} \quad \frac{\partial L}{\partial w_{l,k}} = \frac{\partial L}{\partial h_{l,k}} \frac{\partial h_{l,k}}{\partial s_{l,k}} \frac{\partial s_{l,k}}{\partial w_{l,k}}$$

$$= \frac{\partial L}{\partial h_{l,k}} r'_l(s_{l,k}) h_{l-1}^T$$

$$\frac{\partial L}{\partial w_l} = \left(\frac{\partial L}{\partial w_{l,k}} \right)_{k=1}^d = \left(\frac{\partial L}{\partial w_{l,k}} r'_l(s_{l,k}) \right)_{k=1}^d h_{l-1}^T = r'_l \odot \frac{\partial L}{\partial h_l} h_{l-1} = \text{diag}(r'_l) \frac{\partial L}{\partial h_l} h_{l-1}$$

Top layer



linear model

$$s = x^T \beta$$

neural network

$$x \rightarrow h_1 \rightarrow \dots \rightarrow h_{L-1} \rightarrow s \rightarrow \hat{y}$$

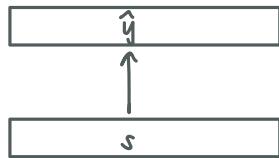
(h_L)

linear regression

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

logistic regression

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



regression:

$$\text{Loss} = \frac{1}{2} |y-s|^2$$

$$\frac{\partial \text{Loss}}{\partial s} = -(y-s) = -e$$

logistic: multiclass classification

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

$$y = \text{one-hot} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ c \\ 0 \\ \vdots \\ 0 \end{bmatrix}^T$$

soft max

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

$$s = \text{logit } s = \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_c \\ \vdots \\ s_C \end{bmatrix}^T$$