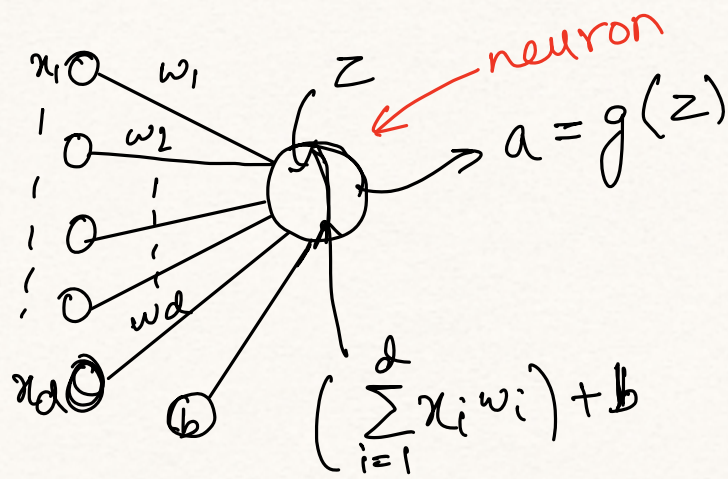


$$x \in \mathbb{R}^d$$



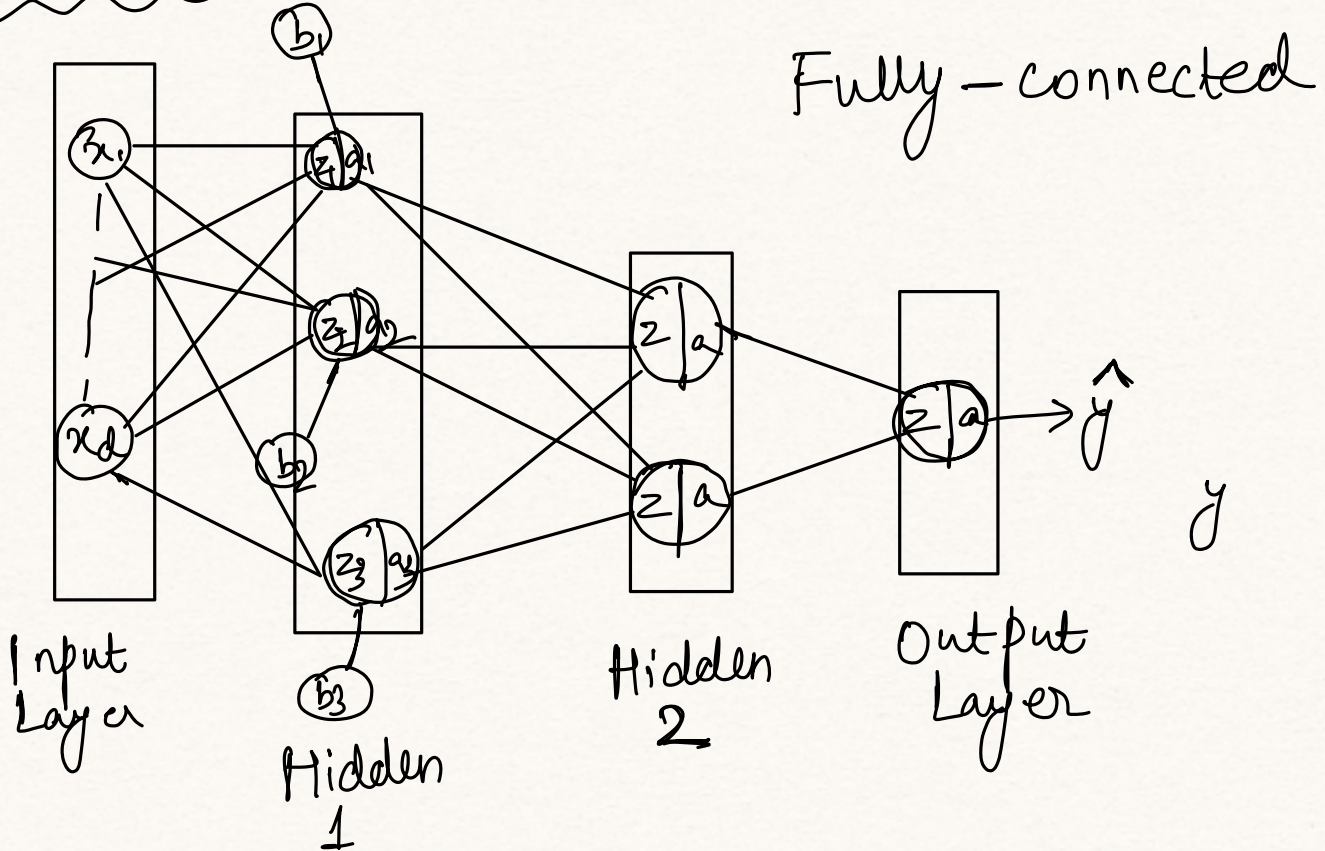
$$z = x^T w + b$$

For logistic reg.

$$g(z) = \frac{1}{1 + e^{-z}}$$

activation function

Neural Network



Then calculate loss for \hat{y} .

Terminologies → layer

$w_{ij}^{[l]}$ — connection

$b_i^{[l]}$ — bias

$z_i^{[l]}$ — $w \cdot () + b$

$a_i^{[l]}$ — $g(z)$

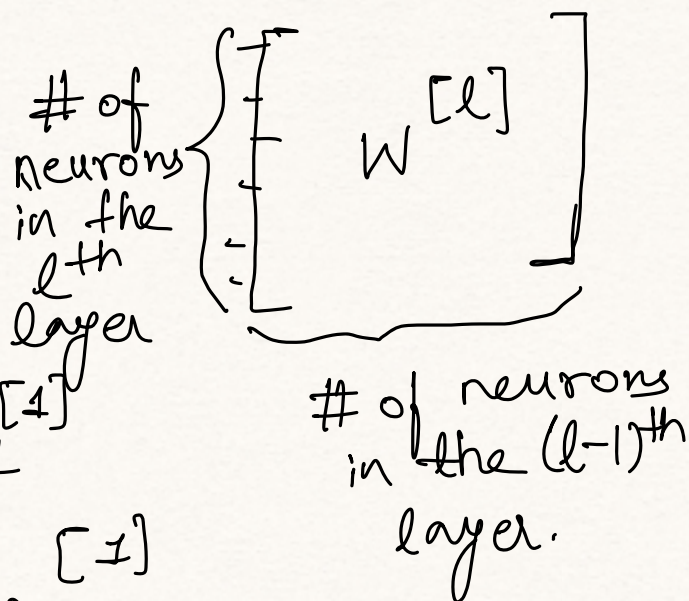
$x = a^{[0]}$

$$z_1^{[1]} = \sum_j w_{1j}^{[1]} a_j^{[0]} + b_1^{[1]}$$

$$z_2^{[1]} = w_2^{[1]T} a^{[0]} + b_2^{[1]}$$

$$a_1^{[1]} = g(z_1^{[1]})$$

$w^{[l]}$ is weight matrix



$z^{[1]}$	$=$	$W^{[1]}$	$a^{[0]}$	$+ b^{[1]}$
3		3x2	2	3

$$a^{[1]} = g(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = g(z^{[2]})$$

$$z^{[3]} = W^{[3]} a^{[2]} + b^{[3]}$$

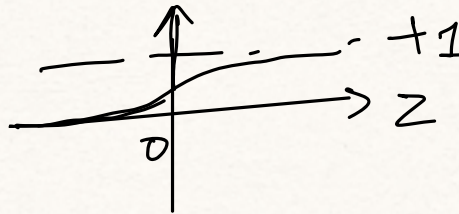
$$a^{[3]} = g(z^{[3]})$$

⋮

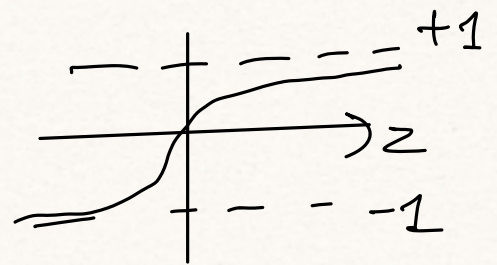
What if $g(z) = z$, $a = z$

$$\begin{aligned}
 a^{[3]} &= W^{[3]} a^{[2]} \\
 &= W^{[3]} z^{[2]} \\
 &= W^{[3]} W^{[2]} a^{[1]} \\
 &= W^{[3]} W^{[2]} z^{[1]} = W^{[3]} W^{[2]} W^{[1]} x \\
 &= \tilde{W} x
 \end{aligned}$$

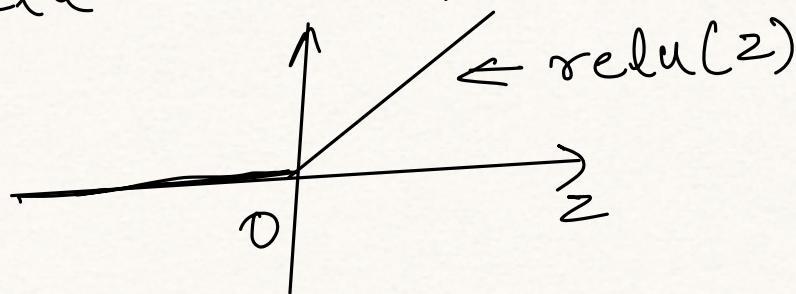
$$g(z) = \frac{1}{1+e^{-z}}$$



$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



$$g(z) = \text{relu} = \max(z, 0)$$



$$a^{[0]} = x^{(i)}$$

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

$$a^{[1]} = g(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = g(z^{[2]})$$

⋮

$$a^{[L]} = g(z^{[L]})$$

$$\hat{y}^{(i)} = a^{[L]}$$

$$L(y^{(i)}, \hat{y}^{(i)}) = -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})]$$

For l in $1, 2, \dots, L$

$$W^{[l]} := W^{[l]} - \alpha \cdot \frac{\partial L}{\partial W^{[l]}}$$

$$b^{[l]} := b^{[l]} - \alpha \cdot \frac{\partial L}{\partial b^{[l]}}$$

⌈ not to be initialized 0.
Initialize w, b

Xavier/Kle
init.

$$w_{ij}^{[l]} \sim N\left(0, \sqrt{\frac{2}{n^{[l]} + n^{[l-1]}}}\right)$$

For t in $1, \dots$

$$\theta := \theta - \alpha ()$$

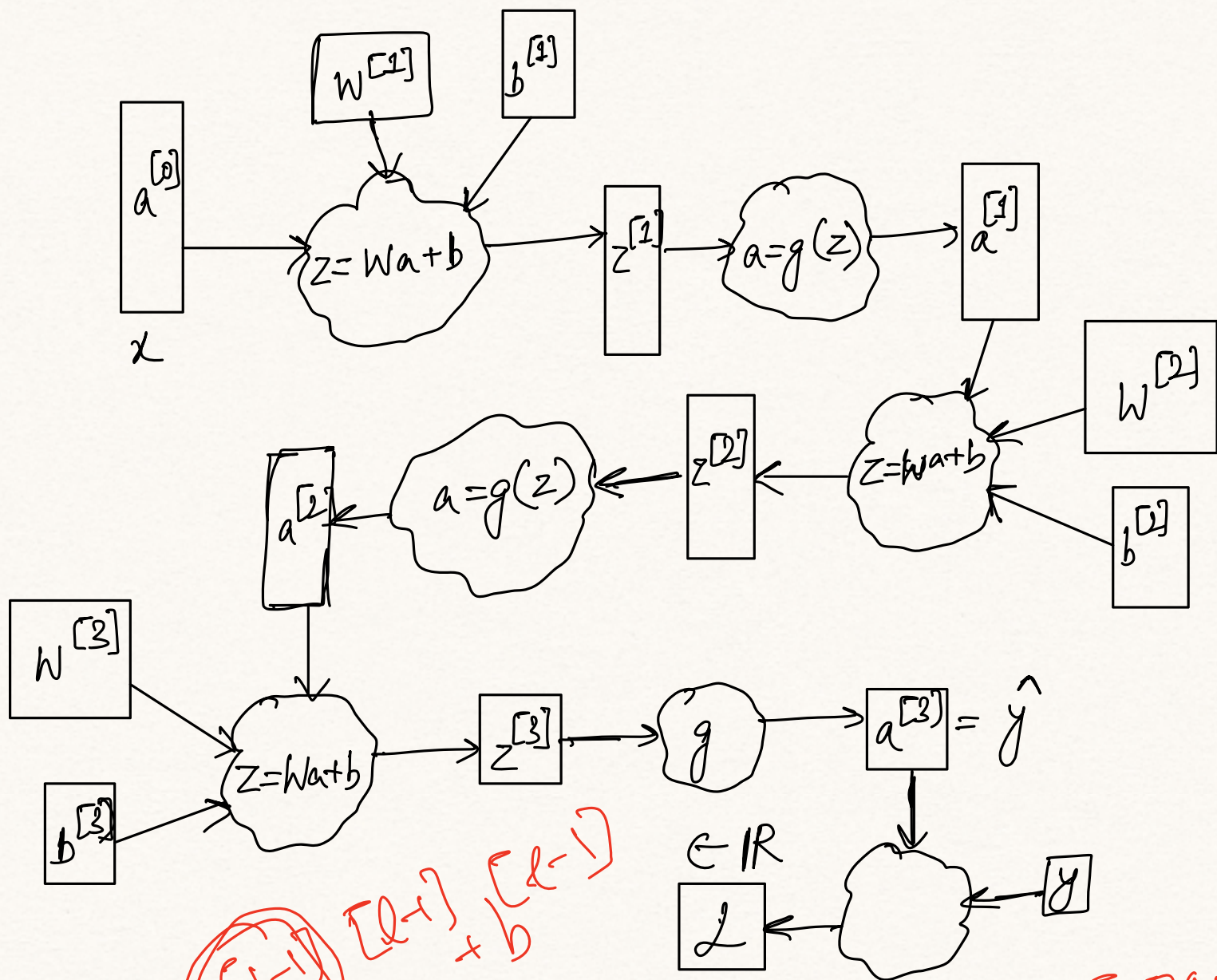
$$w_{ij}^{[l]} \sim U([-0.1, 0.1])$$

For t in $1, 2, \dots$

For l in $1, 2, \dots, L$

$$W^{[l]} := W^{[l]} - \alpha \frac{\partial L}{\partial W^{[l]}}$$

$$b^{[l]} := \dots$$



$z^{[l]} = W^{[l-1]} \cdot a^{[l-1]} + b^{[l]}$

$L \in \mathbb{R}$

loss funcⁿ.

We are considering binary classification.

$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial W_{11}^{[2]}} & \dots & \frac{\partial \mathcal{L}}{\partial W_{13}^{[2]}} \\ \vdots & & \vdots \\ \frac{\partial \mathcal{L}}{\partial W_{33}^{[2]}} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial z^{[3]}} &= \frac{\partial}{\partial z^{[3]}} [-y \log \hat{y} - (1-y) \log(1-\hat{y})] \\ &= \frac{\partial}{\partial z^{[3]}} \left[-y \log(\sigma(z^{[3]})) - (1-y) \log(1-\sigma(z^{[3]})) \right] \end{aligned}$$

$$= a^{[3]} - y$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_{ij}^{[2]}} &= \frac{\partial \mathcal{L}}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial W_{ij}^{[2]}} \\ &= \frac{\partial \mathcal{L}}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial z^{[3]}} \cdot \frac{\partial z^{[3]}}{\partial W_{ij}^{[2]}} \end{aligned}$$

diag[1]

$$= \underbrace{\frac{\partial \mathcal{L}}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial z^{[2]}}}_{(a^{[3]} - y)} \cdot \frac{\partial z^{[3]}}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}}$$

$$= (a^{[3]} - y) W^{[3]} \cdot \frac{\partial z^{[2]}}{\partial W_{ij}^{[2]}}$$

BACK PROPAGATION

i^{th}
 $\text{pos.} \rightarrow \begin{bmatrix} 0 \\ \vdots \\ a_j^{[1]} \\ \vdots \end{bmatrix}$

$$= \underbrace{(a^{[3]} - y) W^{[3]} \circ g'(z^{[2]})}_{1 \times 2} \begin{bmatrix} 0 \\ \vdots \\ a_j^{[1]} \\ \vdots \end{bmatrix}_{2 \times 1}$$

$$= \left[(a^{[3]} - y) W^{[3]} \circ g'(z^{[2]}) \right]_i a_j^{[1]}$$

$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = (a^{[3]} - y) W^{[3]} \circ g'(z^{[2]}) a_j^{[1]T}$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\frac{\partial f}{\partial x} \in \mathbb{R}^{n \times m}$$

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{[2]}} = \underbrace{\frac{\partial \mathcal{L}}{\partial a^{[L]}} \cdot \frac{\partial a^{[L]}}{\partial z^{[L]}} \cdot \frac{\partial z^{[L]}}{\partial a^{[L-1]}}}_{(y - a^{[L]}) \cdot W^{[L]}} \cdot \underbrace{\frac{\partial a^{[L-1]}}{\partial z^{[L-1]}} \cdots \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial W_{ij}^{[2]}}}_{\text{diag. matrix} \uparrow W^{[L-1]} \uparrow \text{diag}(g')}$$

$$\frac{\partial a^{[L-1]}}{\partial z^{[L-1]}} =$$

$$\begin{matrix} & \leftarrow z \rightarrow \\ \begin{matrix} \uparrow \\ a \\ \downarrow \end{matrix} \left[\begin{array}{cccc} g'(z) & 0 & 0 & 0 \\ 0 & g'(z) & & 0 \\ & & \ddots & \\ 0 & & & g'(z) \end{array} \right] \end{matrix}$$

↑
diagonal matrix

$$a_i = g(z_i)$$

$$\frac{\partial a_j}{\partial z_i} = 0 ; i \neq j$$

$$\boxed{\frac{\partial z^{[L-1]}}{\partial a^{[L-2]}} = W^{[L-1]}}$$

$z^{[L-1]} = W^{[L-1]} \cdot a^{[L-2]} + b^{[L-1]}$

$$\frac{\partial z^{[2]}}{\partial W_{ij}^{[2]}} =$$

$$\begin{bmatrix} 0 \\ \vdots \\ a_j \\ \vdots \\ 0 \end{bmatrix} \leftarrow i\text{th}$$

$$J(w, b) = \sum_{i=1}^B L(y^{(i)}, \hat{y}^{(i)})$$

B : size of minibatch

$$\underbrace{z}_{m \times 1}^{[1]} = \underbrace{W}_{m \times d}^{[1]} \underbrace{x^{(i)}}_{d \times 1} + \underbrace{b}_{m \times 1}^{[1]}$$

$m \times 1$

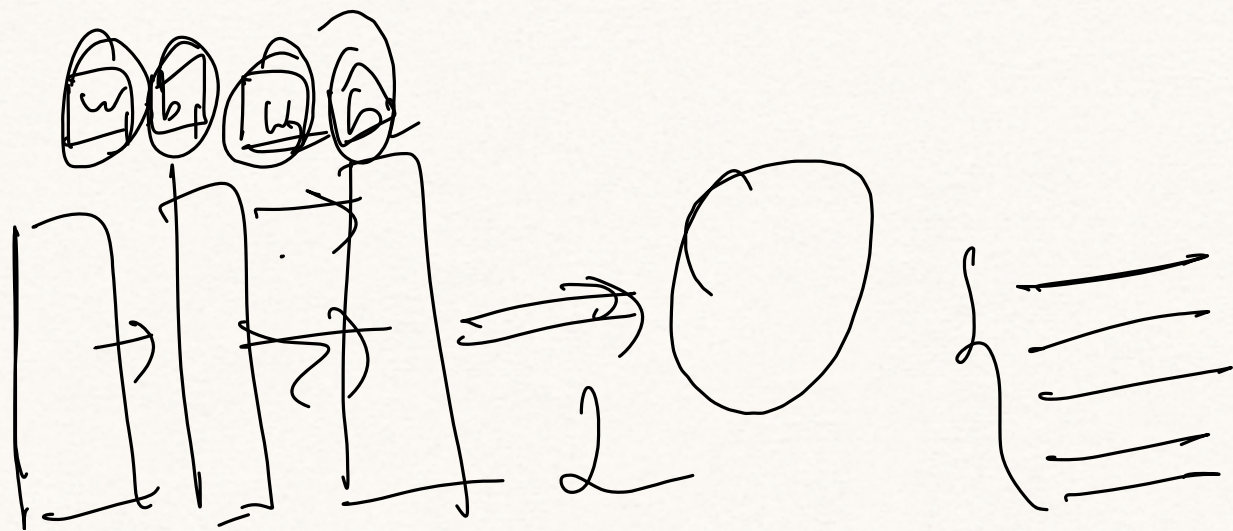
$$= \underbrace{W}_{m \times d}^{[1]} \cdot \underbrace{[x^{(1)} x^{(2)} \dots x^{(B)}]}_{d \times B} + \underbrace{b}_{m \times 1}^{[1]}$$

$$= m \times B + m \times 1 \quad (1 \times B)$$

("Broadcasting") $\Rightarrow \left\{ \begin{bmatrix} 1 \end{bmatrix} \rightarrow \underbrace{\begin{bmatrix} 1 & 1 & 1 \end{bmatrix}}_B \right\}$

B
copies
of col.

$N.N =$ Learnable feature map
+
Linear model



$$\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial b_1} \text{ --- }$$

R^{10}

ϕ

$$y = \begin{bmatrix} 1 & 0 & \dots & 0 \end{bmatrix}^T$$

$$J = K \left\{ \begin{array}{l} \text{---} \\ \text{---} \\ \text{---} \end{array} \right\} \rightarrow \textcircled{i}$$

$$= \sum_{i=1}^n y_i \log y_i$$

$\textcircled{0} 1$

\hookrightarrow