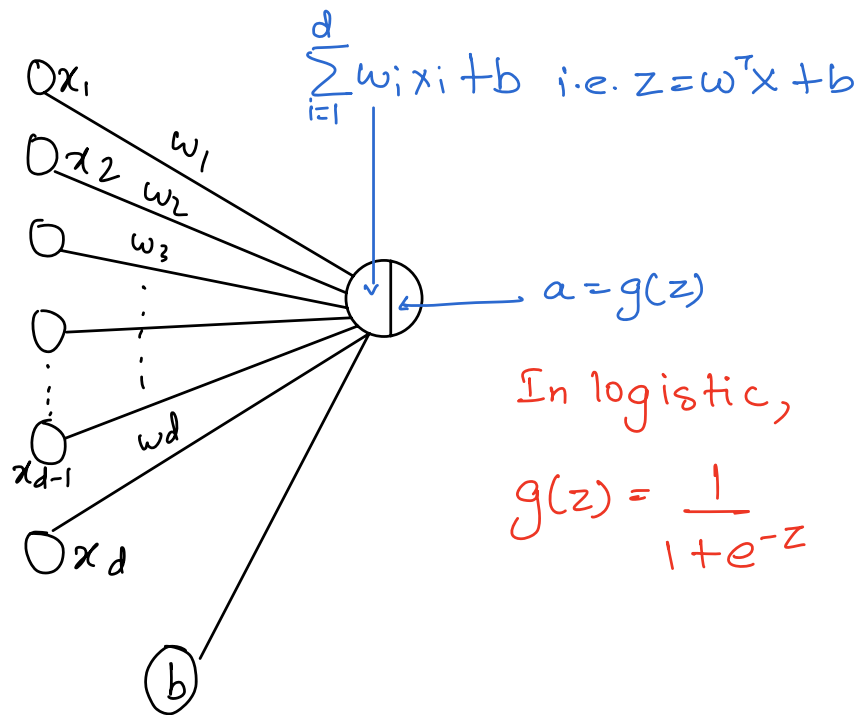


Deep Learning

Neural Networks learn features themselves.
Kernel methods require intuition on the part of the designer

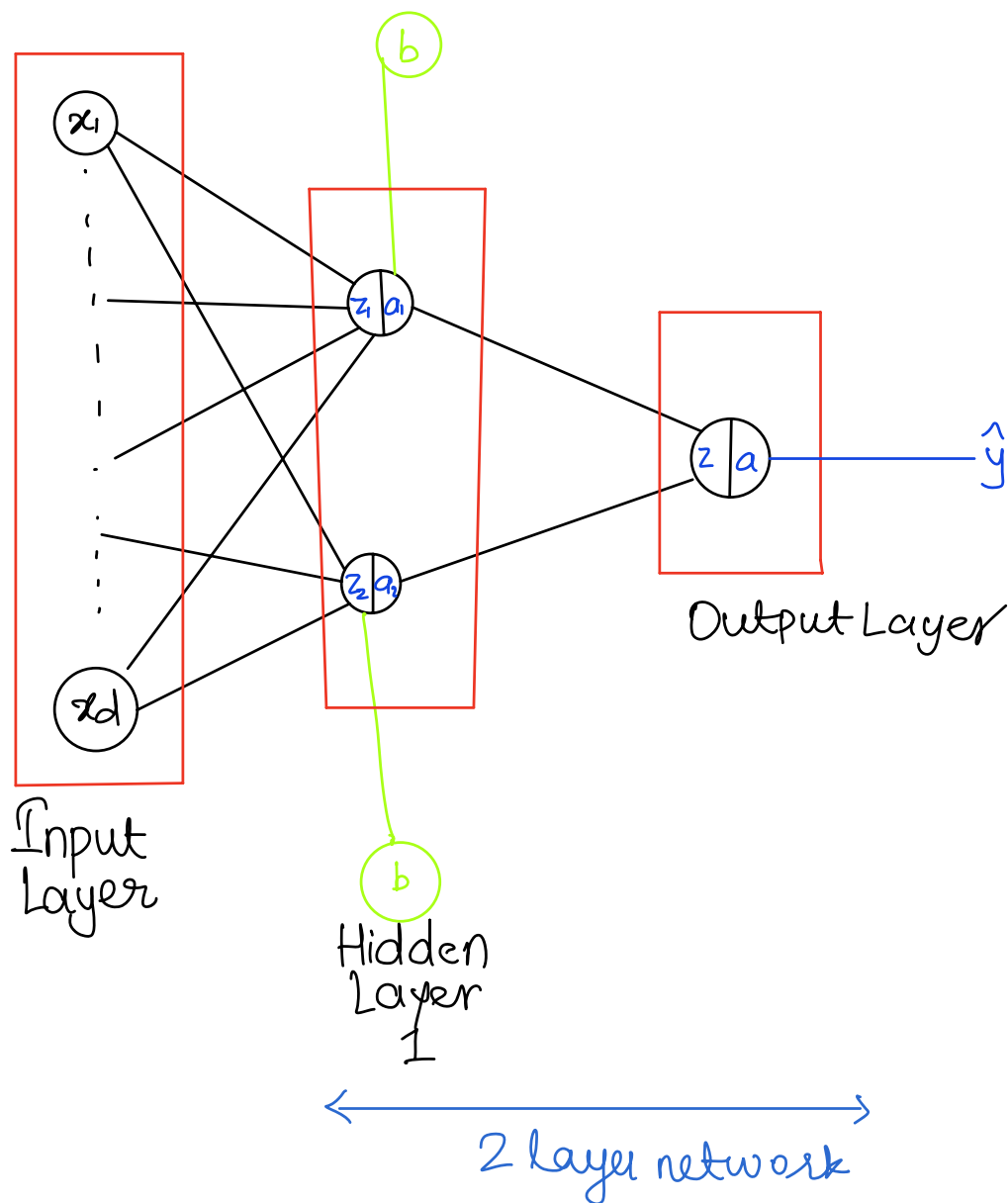
$$x \in \mathbb{R}^d$$



θ replaced by w

$$\hat{y} = a$$

$$l(y, \hat{y}) = y \log \hat{y} + (1 - y) \log (1 - \hat{y})$$



We don't count input layer

$w_{ij}^{[k]}$ connection
 $b_i^{[k]} \rightarrow$ bias
 $z_i^{[k]} \rightarrow w \cdot () + b$
 $a_i^{[k]} \rightarrow g(z)$

$w^{[k]}$ is the weight matrix

$\left[\begin{array}{c} \text{\#neurons} \\ \text{in } k^{\text{th}} \\ \text{layer} \end{array} \right] w^{[k]} \left[\begin{array}{c} \text{\#neurons} \\ \text{in } (k+1)^{\text{th}} \\ \text{layer} \end{array} \right]$

neurons in L -th layer.

$$x = a^{[0]}$$

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

no. of columns \rightarrow no. of neurons in input layer

$$z_2^{[1]} = \sum_j w_{2j}^{[1]} a_j^{[0]} + b_2 = w_2^{[1]T} a^{[0]} + b_2$$

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

$$z_1^{[2]} = \sum_j w_{1j}^{[2]} a_j^{[1]} + b_1 = w_1^{[2]T} a^{[1]} + b_1$$

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

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

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

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

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

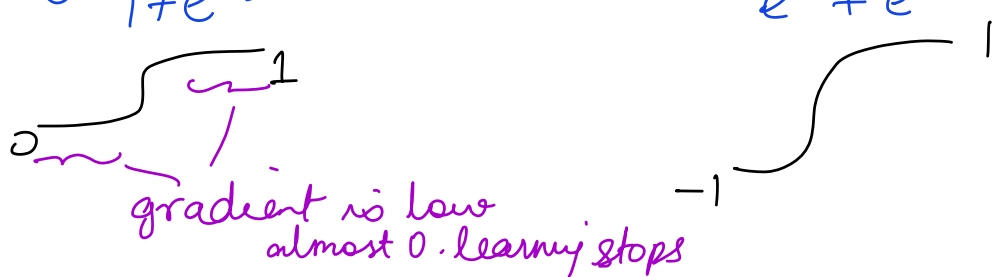
$$a^{[3]} = g(z^{[3]})$$

(g) What if $g(z) = z$?

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

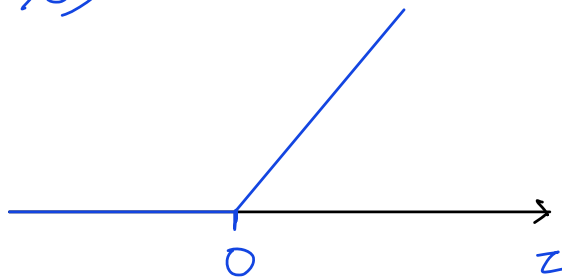
$g(z)$ can be non-linear function only

We need it to be monotonic.

$$g = \frac{1}{1+e^{-z}} \quad g = \tanh = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$


gradient is low
almost 0. learning stops

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



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

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

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

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

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

⋮

$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$$

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

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

For l in $1, 2 \dots L$

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

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

} Stochastic
Gradient

$\begin{matrix} \circ \\ \circ \\ \circ \\ \circ \\ \circ \\ \circ \end{matrix}$

Initialized randomly.

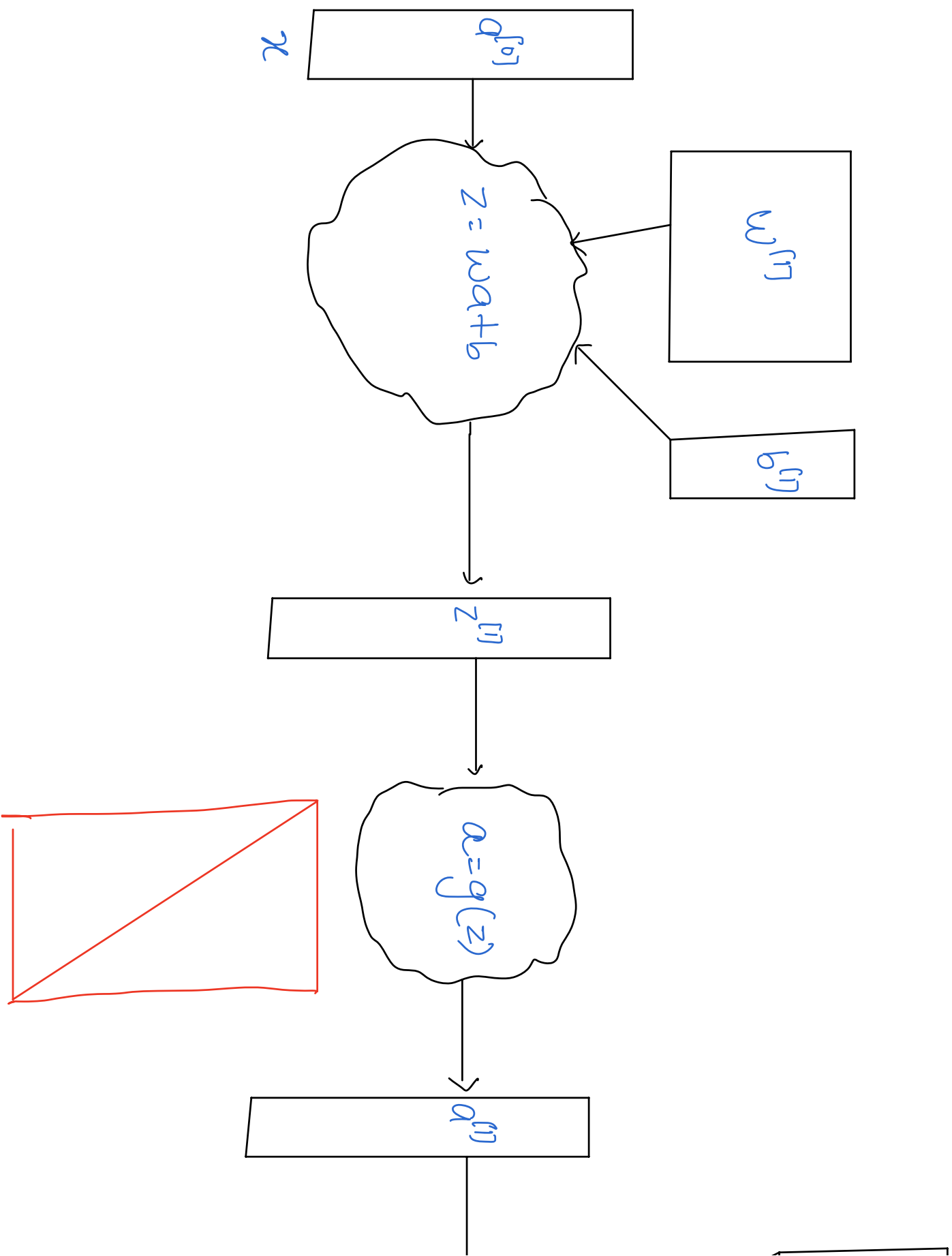
If initialized with all neurons will be same.

By random initialization, we ensure we reach diff.

local minima as loss is not convex

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

$$\sim \text{Unif}[-0.1, 0.1]$$



$$w^{[2]}$$

$$b^{[2]}$$

$$Z = W a + b$$

$$Z^{[2]}$$

$$a = g(Z)$$

$$a^{[2]}$$

$$w^{[3]}$$

$$Z =$$

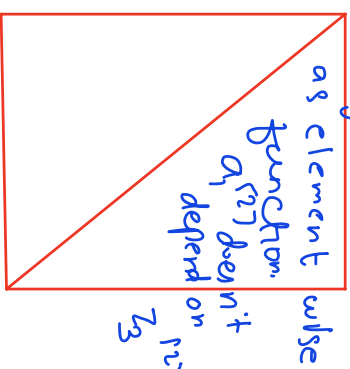
→ Jacobian

→ to make computation easy, we take derivative element wise, else it would be 3D, and very time consuming and it is very sparse

Diagonal matrix

as element wise function $a^{[2]}$ doesn't depend on z_3

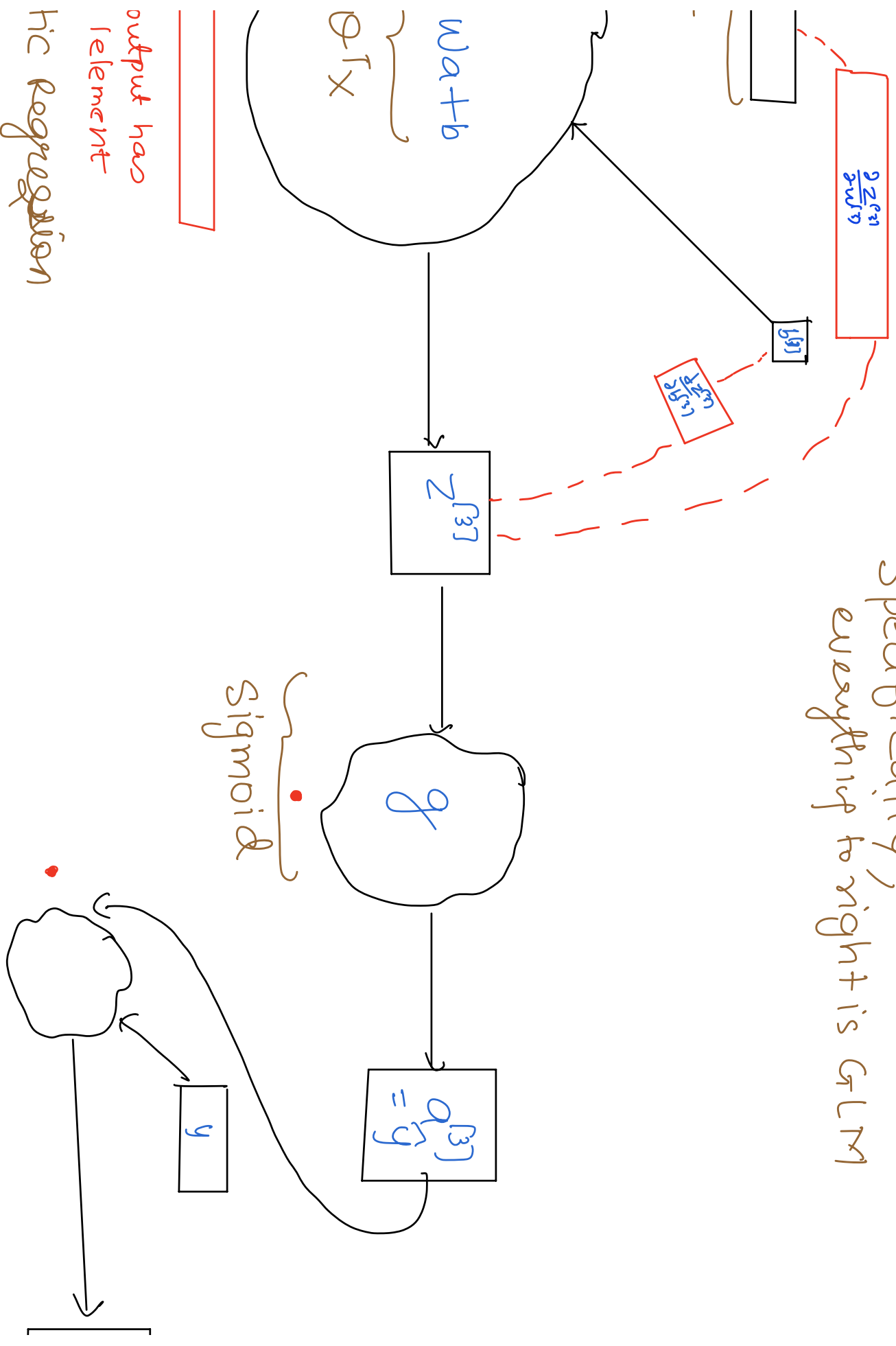
$$\phi_{w,b}(x)$$



→ logis.

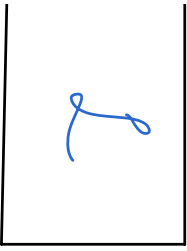
since only

Specifically, everything to right is GLM



output has
element

the regression



$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[2]}} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_{11}^{[2]}} & \dots & \frac{\partial \mathcal{L}}{\partial w_{13}^{[2]}} \\ \vdots & \ddots & \vdots \\ \vdots & \ddots & \vdots \\ \vdots & \ddots & \frac{\partial \mathcal{L}}{\partial w_{23}^{[2]}} \end{bmatrix}$$

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

$$= \frac{\partial}{\partial z^{[3]}} \left[-y \log \sigma(z^{[3]}) - (1-y) \log (1-\sigma(z^{[3]})) \right]$$

$$\vdots$$

$$= \sigma(z^{[3]}) - y$$