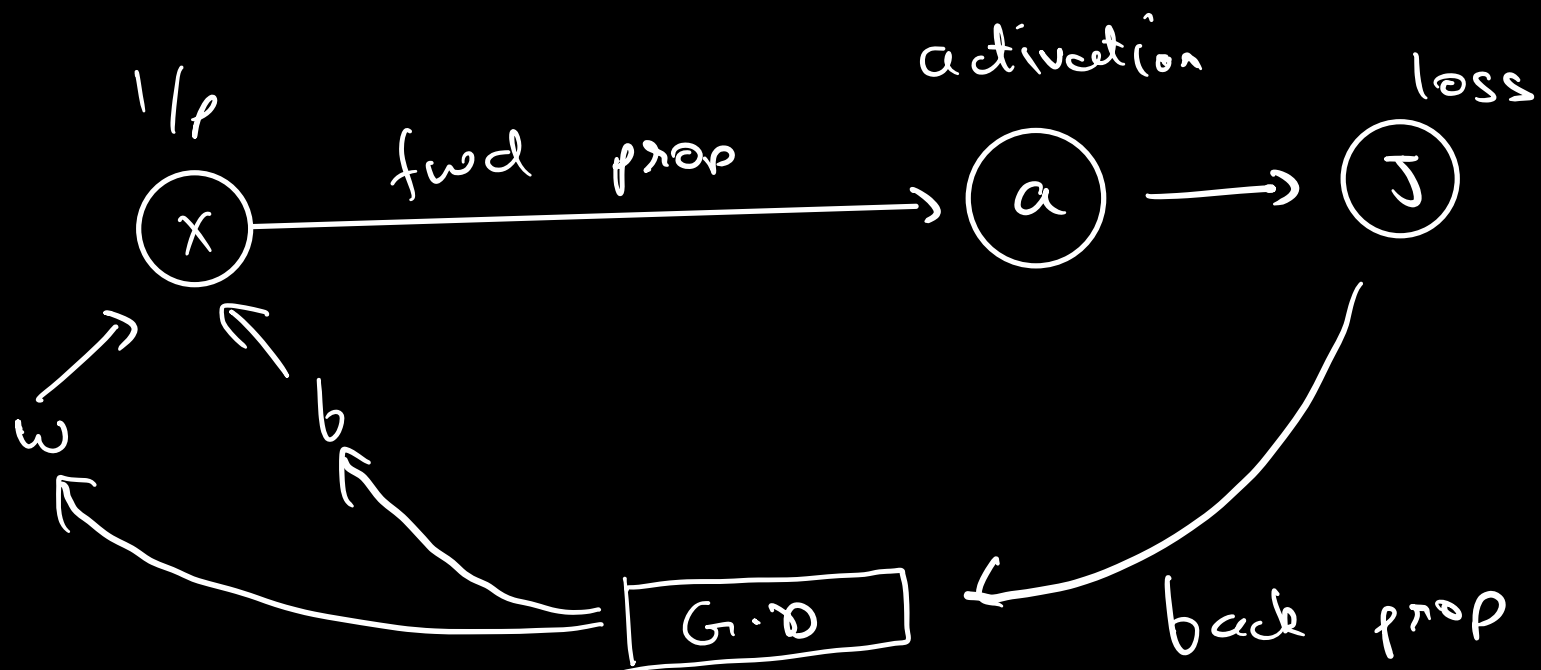


Backpropagation

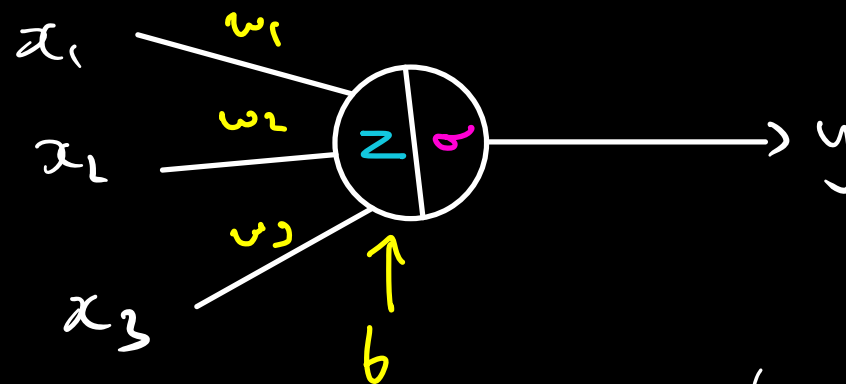
&

Multi Layer Neural Networks

Training a Neuron



Let's consider a single Neuron



L.R.U

Shapes

$$w = (1, d) \quad b = (1, 1)$$

$$x = (d, 1) \quad a = (1, 1)$$

$$L = \text{scalar}$$

Fwd

$$z = \vec{w} \cdot \vec{x} + b$$

$$\hat{y} = a = \sigma(z)$$

$$J = -(y \log(a) + (1-y) \log(1-a))$$

Gr.D

$$w' = w - \alpha \frac{\partial J}{\partial w}, \quad b' = b - \alpha \frac{\partial J}{\partial b}$$

Bwd

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial b}$$

$$1) \frac{\partial J}{\partial a} = - \left(y \frac{\partial \log(a)}{\partial a} + (1-y) \frac{\partial \log(1-a)}{\partial a} \right)$$

$$= - \left(y \cdot \frac{1}{a} + (1-y) \frac{1}{a-1} \right)$$

$$2) \frac{\partial a}{\partial z} = \frac{\partial}{\partial z} \left(\frac{1}{1+e^{-z}} \right)$$

$$\frac{\partial}{\partial z} [(1 + e^{-z})^{-1}]$$

$$= -\frac{1}{(1 + e^{-z})^2} \cdot (-e^{-z})$$

$$= \frac{e^{-z}}{(1 + e^{-z})^2}$$

$$= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}}$$

$$= \sigma(z) \cdot \left[\frac{e^{-z} + 1 - 1}{1 + e^{-z}} \right]$$

$$= \sigma(z) [1 - \sigma(z)]$$

$$= a(1-a)$$

$$3) \frac{\partial Z}{\partial w} = x$$

$$4) \frac{\partial Z}{\partial b} = 1$$

$$\Rightarrow \frac{\partial J}{\partial w} = - \left(y \cdot \frac{1}{a} + (1-y) \frac{1}{a-1} \right) (a)(1-a) \cdot x$$

$$= - \left(y \cdot \frac{1}{q} \cdot q(1-a) + (1-y) \frac{1}{q} a \cdot (\cancel{1-a}) \right) x$$

$$= - \left[y(1-a) - (1-y)a \right]$$

$$= - \left[y - \cancel{ay} - a + \cancel{ay} \right]$$

$$= \boxed{(a - y) \cdot \vec{x}}$$

$$\text{OR } (a - y) \cdot x^T$$

[because w shape is
transpose of x shape]

$$\Rightarrow \frac{\partial J}{\partial y} = (a - y)$$

\therefore G.V

$$w' = w - \alpha (a - y) x^T$$

$$b = b - \alpha (a - y)$$

$$\text{eg: } a = \hat{y} = 0.7$$

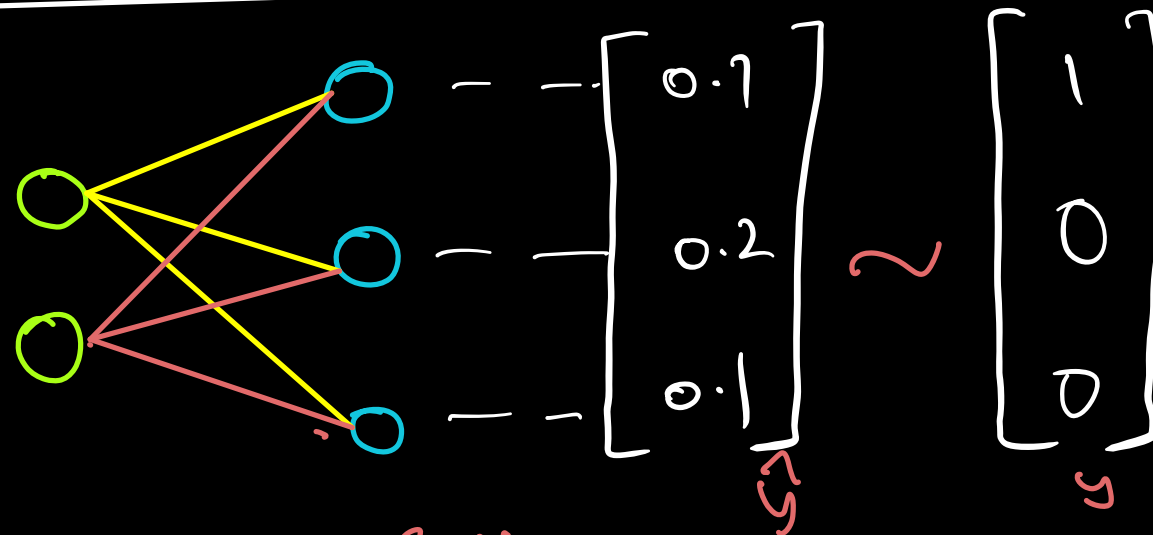
$$y = 1$$

$$\therefore \frac{\Delta w}{\Delta x} = 0.7 - 1 = -0.3$$

This is for a single point \vec{x}_i . When we have entire data X , then we will

use \rightarrow TODO

For a multiclass Network



Backprop through Softmax

Long mathematical proof!! → Post Read

Hack →

Sigmoid is a special case of
softmax

$$S_M(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

[softmax]

for 2 classes

$$S_M(z_A) = \frac{e^{z_A}}{e^{z_A} + e^{z_B}}$$

$$= \frac{1}{1 + e^{(z_B - z_A)}}$$

$$= \frac{1}{1 + e^{-(z_A - z_B)}}$$

$$\approx \frac{1}{1 + e^{-z'}} \quad \dots \text{Sigmoid}$$

This is loosely like a sigmoid, hence the derivative also comes out to be similar even for multiple classes:

$$\therefore \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial z} = (a - y)$$

$$\therefore \frac{\partial J}{\partial w} = (a - y) \cdot x^T, \quad \frac{\partial J}{\partial b} = (a - y)$$

Note that here

$$a = \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad b = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\Delta w = (a - y) \cdot x^T$$

$$= k \times 1, \quad 1 \times d$$

$$= k \times d$$

$$\text{ex: } 3 \times 2$$

Please note that shapes depend on initial settings and conventions.

When you have many points,

$$w' = w - \alpha \frac{1}{n} \sum_{i=1}^n (a_i - y_i) x_i^T$$

$$b = b - \alpha \frac{1}{n} \sum_{i=1}^n (a_i - y_i)$$

These calc can also be done in matrix form

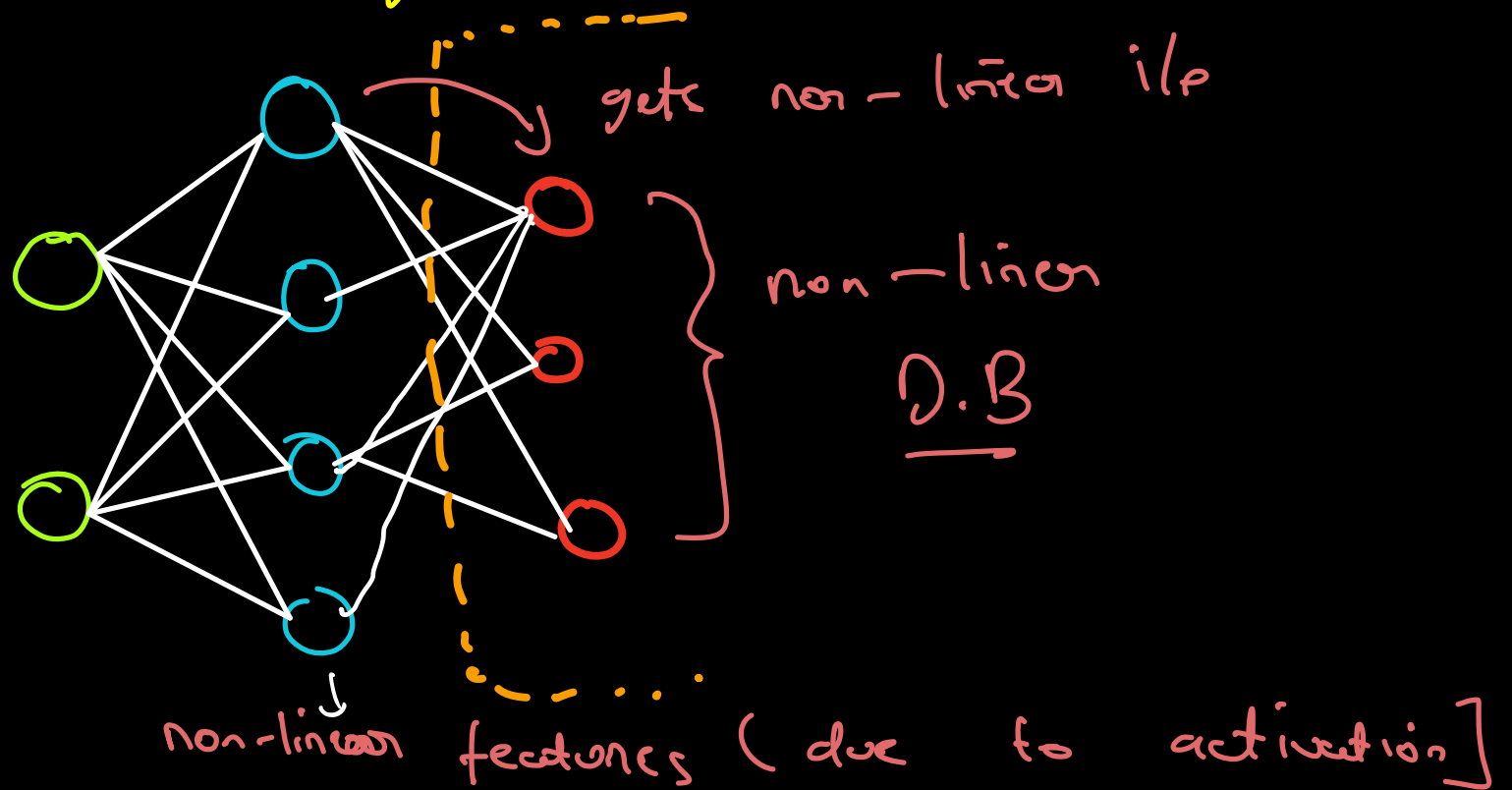
$$\text{eg: } \rightarrow w - \alpha \underbrace{(A - Y) \cdot X^T}_{\rightarrow \text{take mean in code}}$$

→ code

Multi-layer networks

We got a linear decision boundary. How to get non-linear?

→ Add another layer



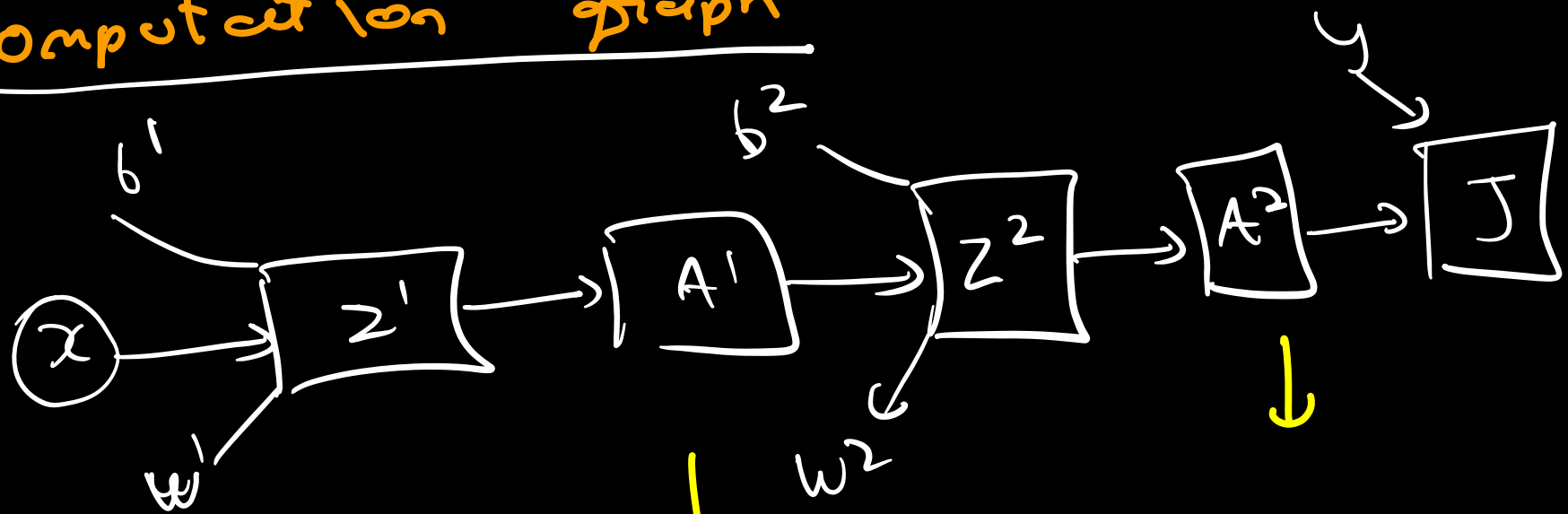
Do not confuse

$$\begin{aligned} z &= wx + b \\ &\quad \downarrow \\ &\quad (k \times d)(d \times n) + (k \times n) \\ &= \underline{\underline{k \times n}} \end{aligned}$$

$$\begin{aligned} z &= \cancel{x}w + b \\ &\quad \downarrow \\ &\quad (n \times d)(d \times k) + (n \times k) \\ &\quad \underline{\underline{n \times k}} \end{aligned}$$

Both are correct, depends on initial shape.

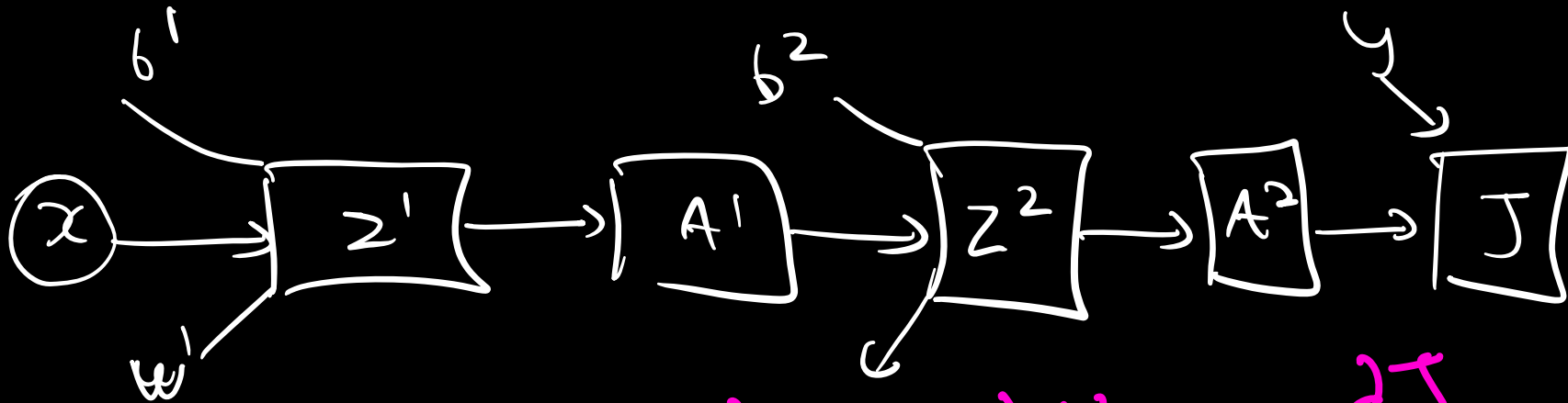
Computation graph



Fwd Prop $\rightarrow A^1 = a(xw^1 + b^1)$

$$A^2 = a(\underbrace{A^1 w^2 + b^2}_{\text{non linear}})$$

Back prop



$$\begin{aligned} \frac{\partial A^1}{\partial z^1} &\leftarrow \frac{\partial z^2}{\partial A^1} & \frac{\partial A^2}{\partial z^2} &\leftarrow \frac{\partial J}{\partial A^2} \\ &\swarrow \quad \searrow & \swarrow \quad \searrow \\ \frac{\partial z^1}{\partial w^1} & \quad \frac{\partial z^1}{\partial b} & \frac{\partial z^2}{\partial w^2} & \quad \frac{\partial z^2}{\partial b^2} \end{aligned}$$

$$\therefore \frac{\partial \mathcal{J}}{\partial \omega'} = \frac{\partial \mathcal{J}}{\partial A^2} \cdot \frac{\partial A^2}{\partial Z^2} \cdot \frac{\partial Z^2}{\partial A'} \cdot \frac{\partial A'}{\partial Z'} \cdot \frac{\partial Z'}{\partial \omega'}$$