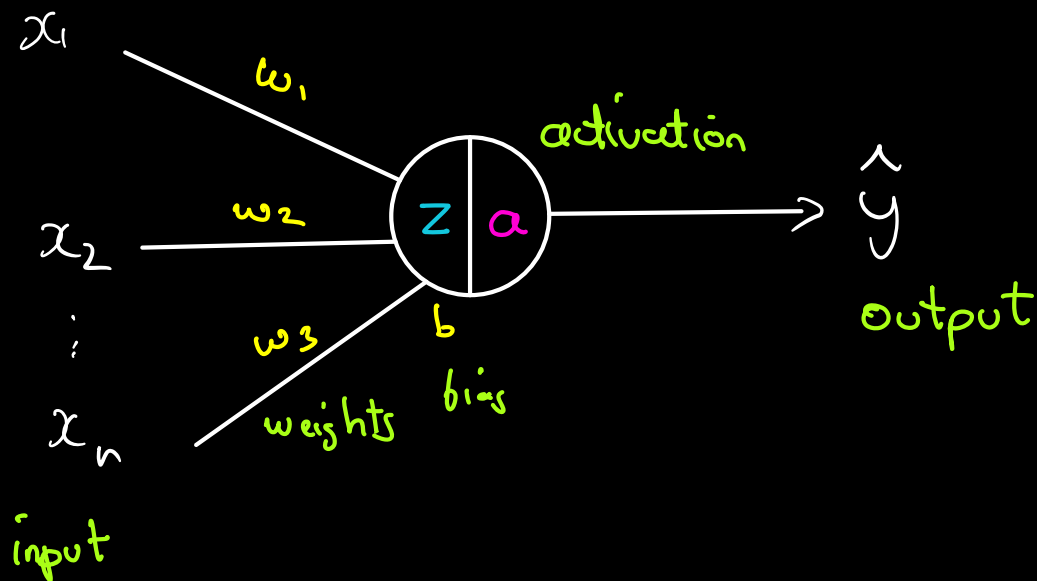


Forward Propagation & Loss

[Neural Networks]

Recap

Neuron



$$\hat{y} = a(w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b)$$

↑

Generalized: if we change this activation function it behaves like different models

For example:

$a = \text{step} \rightarrow \text{perceptron}$

$a = \text{linear} \rightarrow \text{linear regression}$

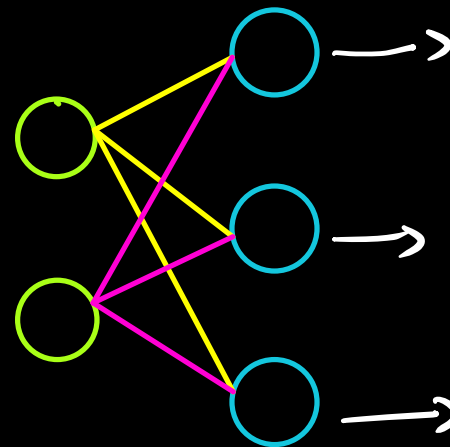
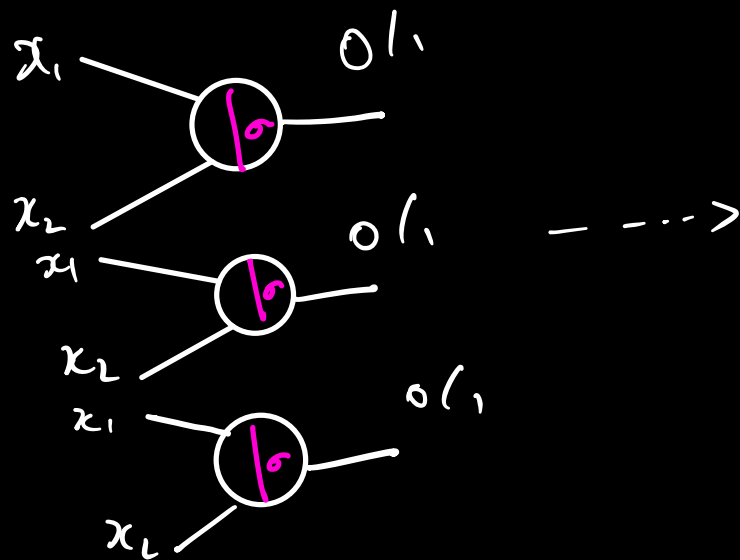
$a = \text{sigmoid} \rightarrow \text{logistic regression}$

... and so on.

Let's build a small Network

Using one vs rest approach with binary

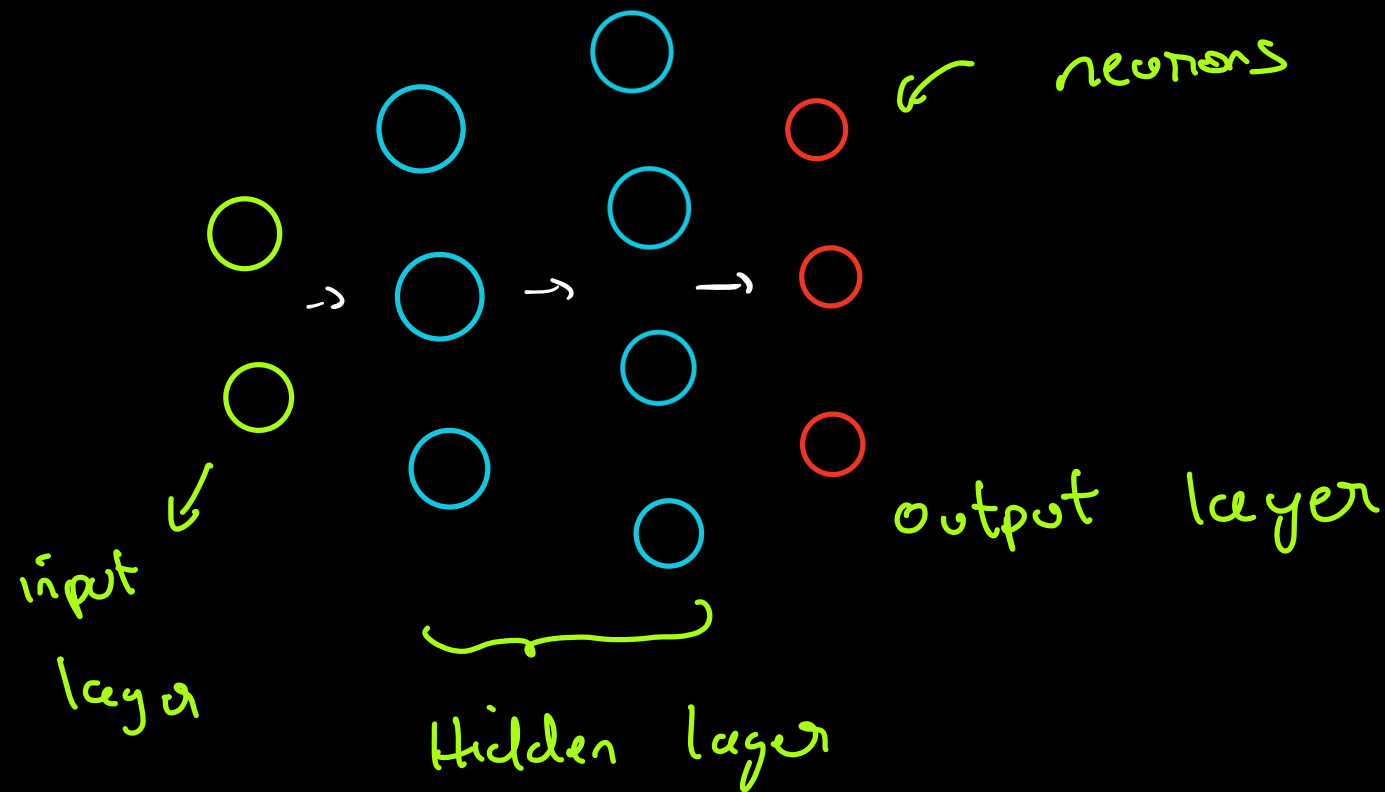
LRU, (log Reg Units) ; i.e. $\rightarrow a = \leftarrow \text{sigmoid}$



Now this is a single model which will be trained (All LRUs simultaneously)

Obviously we need to fix the probabilities which we will do later...

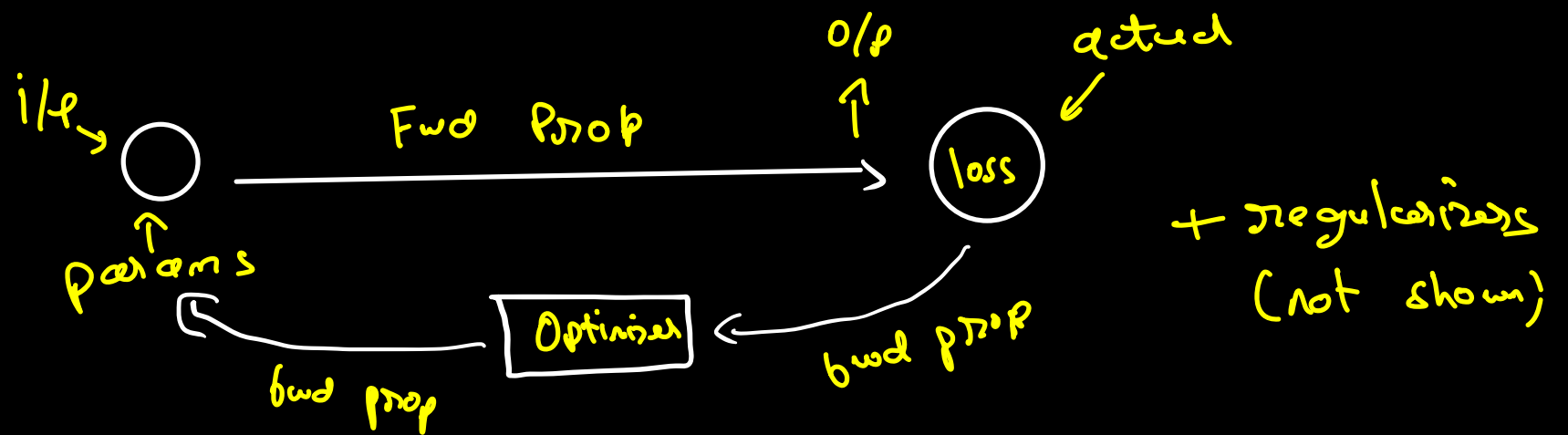
Components of a neural network



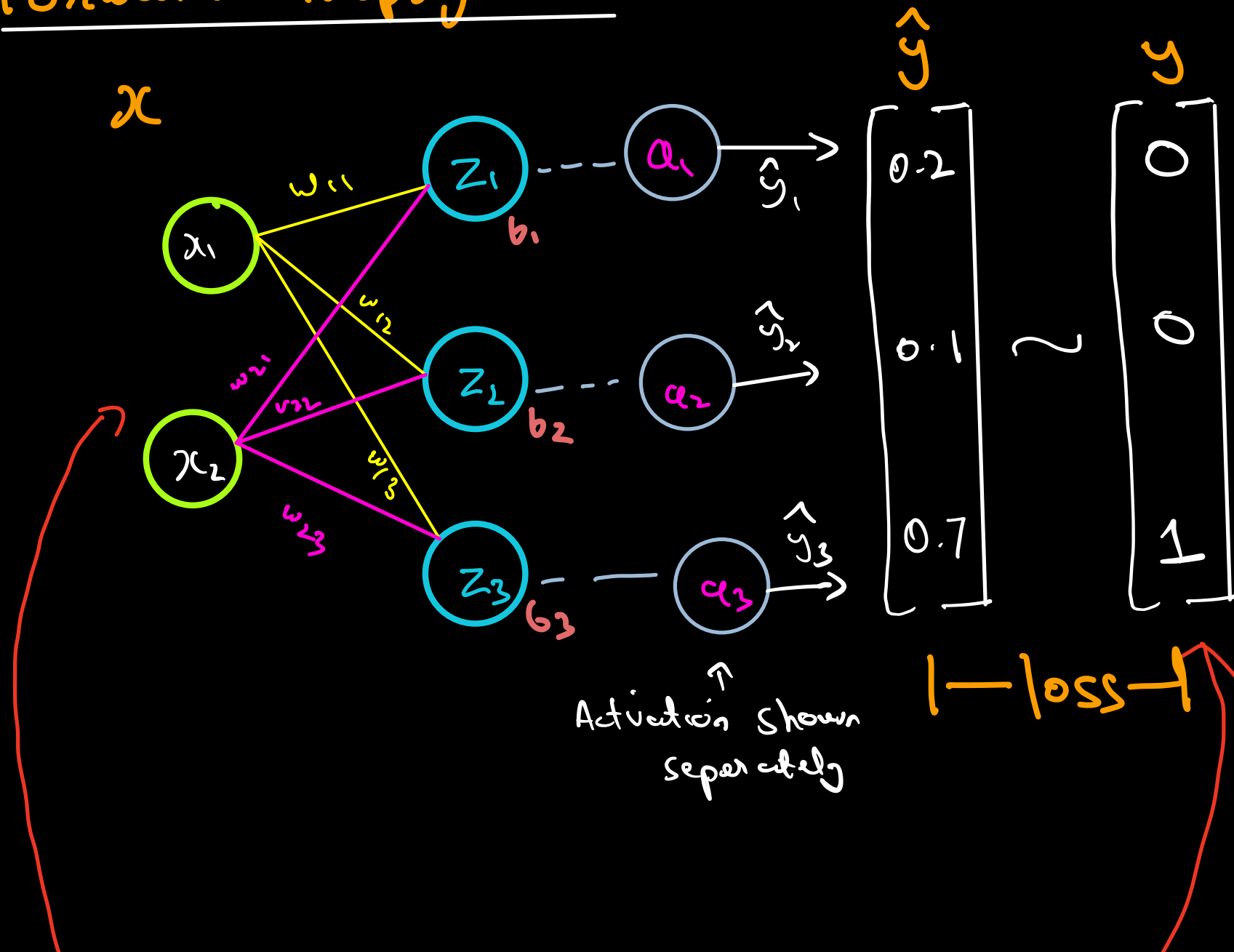
Components

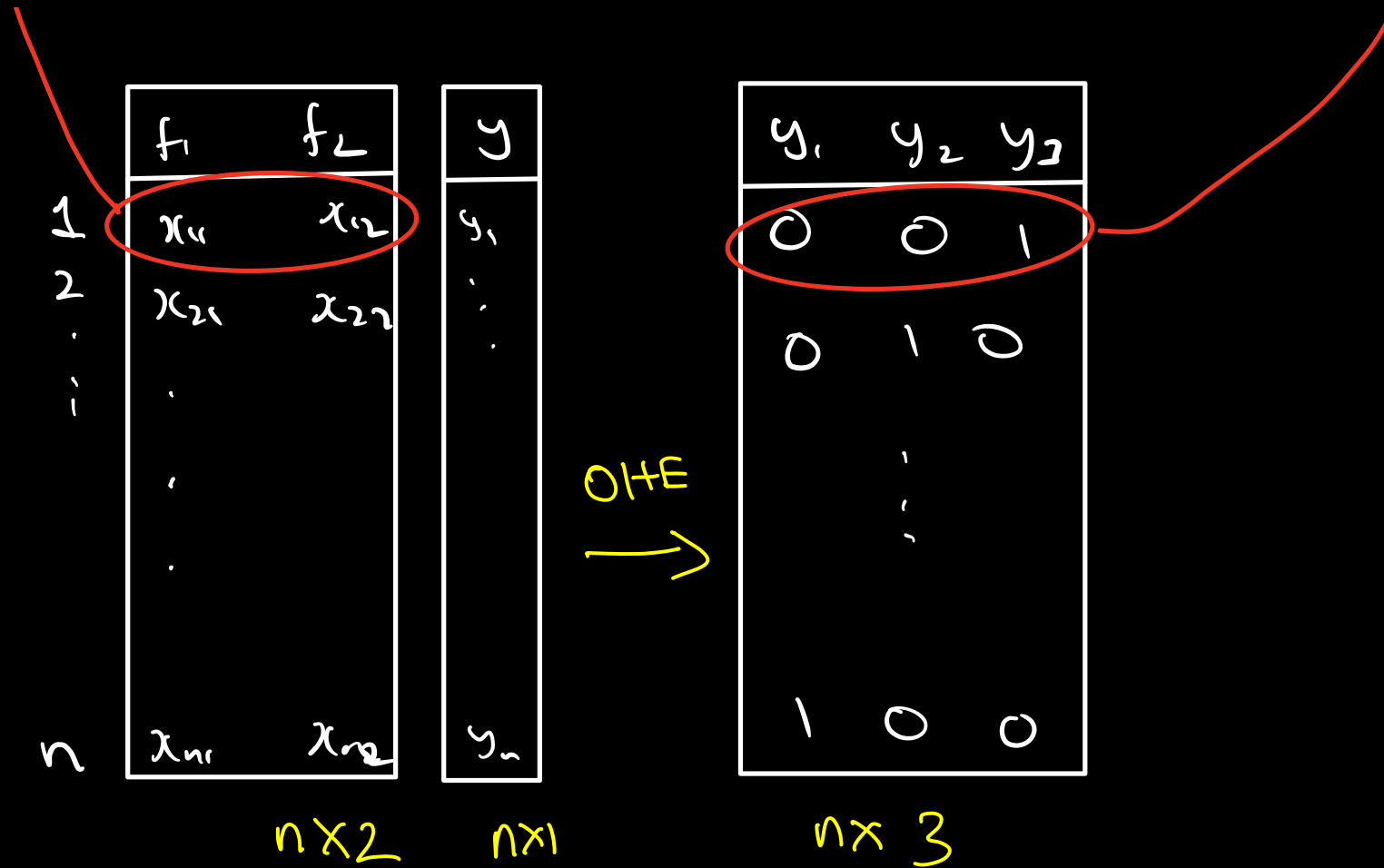
- Neurons + layers

- input / output
- Fwd Prop
- loss
- bwd Prop
- optimizers
- regularisers



Forward Propagation





Data Matrix
 $n \times d$
 # points \leftarrow # features

Output Matrix
 $n \times K$
 \hookrightarrow # classes

Weights
 $d \times K$

$$a(w_{11}x_1 + w_{21}x_2 + b_1) = \hat{y}_1$$

$$a(w_{12}x_1 + w_{22}x_2 + b_2) = \hat{y}_2$$

$$a(w_{13}x_1 + w_{23}x_2 + b_3) = \hat{y}_3$$

$$\begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nd} \end{bmatrix} \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} = \begin{bmatrix} \phantom{x_{11}} \\ \phantom{x_{11}} \\ \phantom{x_{11}} \end{bmatrix}$$

$d \times K$
 eg: 2×3
 $= 6 \checkmark$

$n \times d$ $n \times K$

But we also need to add biases

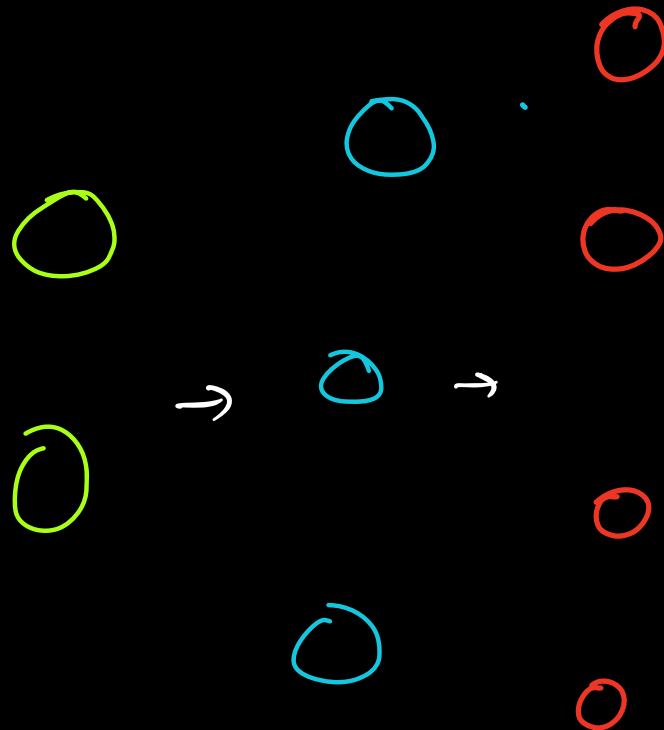
$$\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ \vdots & \vdots \end{bmatrix}_{n \times d} \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ \vdots & \vdots & \vdots \end{bmatrix}_{d \times k} + \begin{bmatrix} b_1 & b_2 & b_3 \\ b_1 & b_2 & b_3 \\ \vdots & \vdots & \vdots \\ b_1 & b_2 & b_3 \end{bmatrix}_{n \times k} = \begin{bmatrix} \underline{a_{11}} & a_{12} & a_{13} \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \underline{a_{n3}} \end{bmatrix}_{n \times k}$$

$$a(X \cdot W + b) = \hat{y} \quad n \times k$$

Q: What is the total number of learnable params in this model?

→ 6 weights + 3 biases = 9 params

Q: Num of params \rightarrow ?



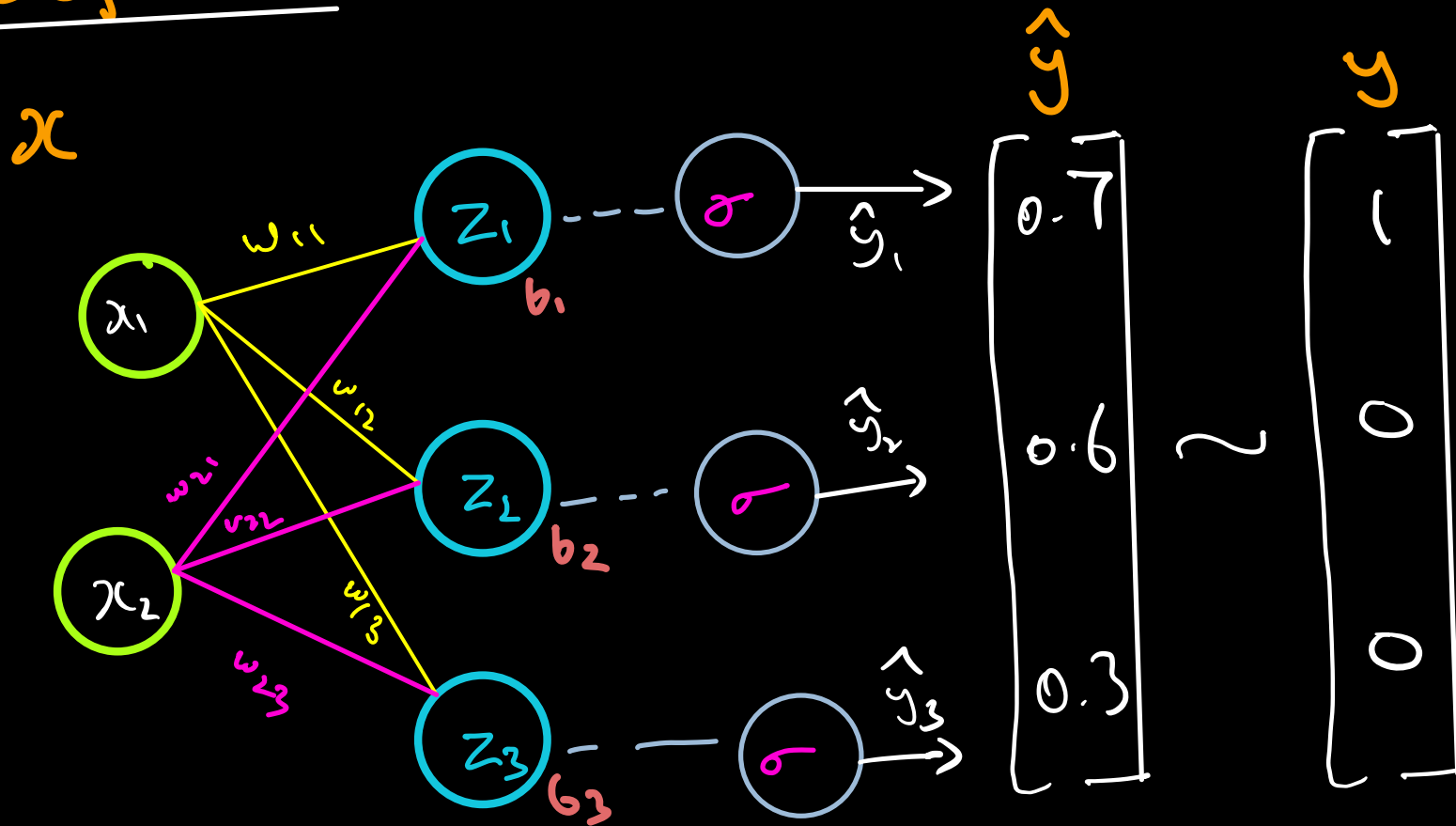
$$\text{Ans} \rightarrow (2 \times 3) + 3 + (3 \times 4) + 4$$

$$= 6 + 3 + 12 + 4$$

$$= 9 + 16 =$$

$$\boxed{25}$$

Softmax



In our model the activation is sigmoid

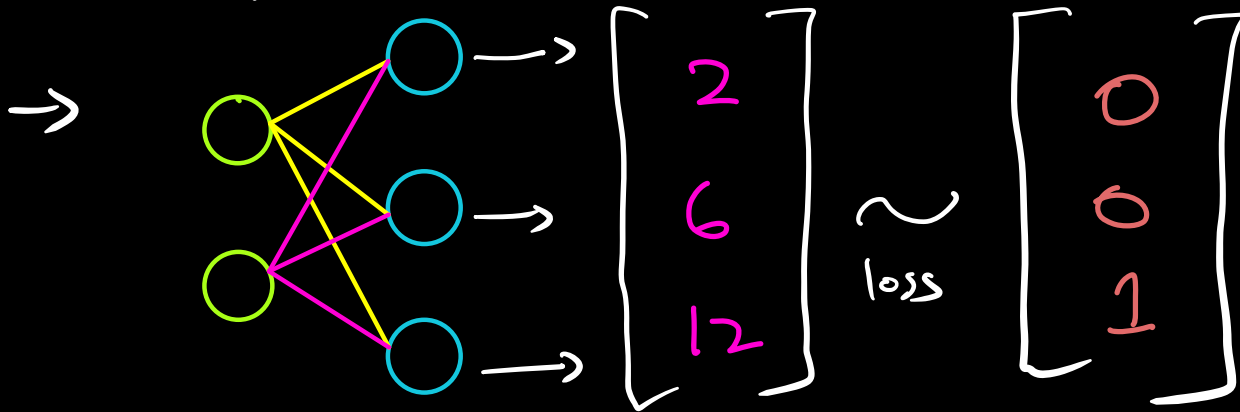
So there is no guarantee that $\sum_{\text{probability}} \hat{y}_i = \underline{\underline{1}}$

Hence we need to normalise, we have already seen that

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} ; K = \# \text{ classes}$$

why e^x ?

→ easy to differentiate



direct normalize:

$$\frac{2}{2+6+12}, \frac{6}{2+6+12}, \frac{12}{2+6+12}$$

z 0.1 0.3 0.6

$$= \begin{bmatrix} 0.1 \\ 0.3 \\ 0.6 \end{bmatrix} \sim \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \rightarrow \text{High difference}$$

Softmax

$$\frac{e^2}{e^2 + e^6 + e^{12}}$$

$$\dots = \begin{bmatrix} 0.0004 \\ 0.002 \\ 0.997 \end{bmatrix} \sim \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Low difference

Loss

Q: Can we use rmse to compare \hat{y}, y ?

→ No, better choice:

For 2 classes

$$\min_{w, b} J(w, b) : -\frac{1}{N} \sum_{i=1}^N y_i \log(\hat{p}_i) + (1 - y_i) (\log(1 - \hat{p}_i))$$

A.K.A → logloss → binary cross entropy loss

Similarly for K classes

$$\min_{w, b} J(w, b) = \sum_{i=1}^N \sum_{j=1}^K y_{ij} \log(p_{ij})$$

aka → categorical cross entropy loss

Back Propagation

To update the param w, b I need,

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

But J is not a direct function,
hence

