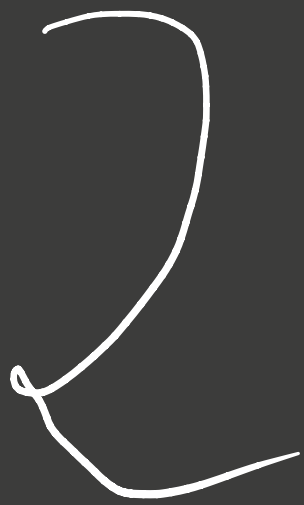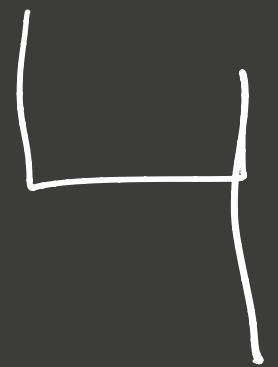# NEURAL NETWORKS

* ORIGINALLY BIOLOGICALLY INSPIRED

* STATE-OF-THE ART IN MOST FIELDS

* TURING AWARD WINNERS — BENGIO, LECONN, HINTON

Q). what is the following?

2 3 44 4

* EASY FOR US TO RECOGNIZE

* WHAT ABOUT COMPUTERS.?

\* COMPUTER SEES 16 pixels

\* NON-TRIVIAL TO WRITE PROGRAM!

$\Rightarrow$ LEARNING!

# Perceptron

* Artificial neuron developed in 1950s/60s by Rosenblatt inspired by McCulIoh & Pitts

(Binary I/P)
$x_1$

$x_2$
:
$x_3$

(Binary)
Output

Perceptron

# Perceptron



$x_1 \xrightarrow{w_1 \ (weights)}$

$x_2 \xrightarrow{w_2}$

$x_3 \xrightarrow{w_3}$

$\rightarrow O/p$

$$O/p = \begin{cases} 0 \; ; \; \text{if} \; \sum w_i x_i \leq \text{threshold} \\ 1 \; ; \quad \sum w_i x_i > \text{threshold} \end{cases}$$

# Perceptron

$x_1$ —— $w_1$ (weights)

$x_2$ —— $w_2$ ——→ O/P

$x_3$ —— $w_3$

1 —— bias (b)

$$O/p = \begin{cases} 0 & ; \text{ if } \sum w_i x_i + b \leq 0 \\ 1 & ; \quad \sum w_i x_i + b > 0 \end{cases}$$

# Perceptron

Q). For 2 i/p learn perceptron for binary AND.

| $x_1$ | $x_2$ | Out |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

# Perceptron

(Q). For 2 i/p learn perceptron for binary AND.

$$w_1 = 1; w_2 = 1; b = -1.5$$

| $x_1$ | $x_2$ | Out |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

# Perceptron

Q). For 2 i/p learn perceptron for binary OR

| $x_1$ | $x_2$ | Out |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

$x_1$ $w_1$

$x_2$ $w_2$

$b$

$1$

$\rightarrow O/p$

# Perceptron

Q). For 2 i/p learn perceptron for

binary OR

$$w_1 = 1 \quad w_2 = 1 \quad ; b = -0.5$$

| $x_1$ | $x_2$ | Out |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

# Perceptron

Q). For 1 i/p learn perceptron for binary NOT

| $x_1$ | Out |
|-------|-----|
| 0     | 1   |
| 1     | 0   |

# Perceptron

Q). For 1 i/p learn perceptron for
binary NOT

$$w_1 = -1; \quad b = 0.5$$

| $x_1$ | Out |
|-------|-----|
| 0     | 1   |
| 1     | 0   |

# Perceptron

Q). For 2 i/p learn perceptron for binary XOR

| $x_1$ | $x_2$ | Out |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$x_1 \xrightarrow{w_1}$

$x_2 \xrightarrow{w_2}$

$1 \xrightarrow{b}$

$\rightarrow 0/p$

AND

OR

$x_2$

$x_1$

XOR

NON-LINEARLY SEPARABLE

# Perceptron

Cost funcⁿ

$$J(w_1, w_2, b) = \frac{1}{4} \sum_{i=1}^{4} (y_i - \hat{y_i})^2$$

Optimum

$$w_1 = w_2 = 0$$

$$b = \frac{1}{2}$$



← Prediction

# Perceptron

Let's add more neurons to learn XOR

(Bias implicit)



Can this network of perceptrons learn XOR?

$$\hat{y} = w_1 X_1 + w_2 X_2 + b$$

$$X_1 = \mu_1 x_1 + \mu_2 x_2 + \mu_0$$

$$X_2 = \lambda_1 x_1 + \lambda_2 x_2 + \lambda_0$$

$$\therefore \boxed{\hat{y} = \gamma_1 x_1 + \gamma_2 x_2 + \gamma_0}$$ ← Still linear !!

Need some non-linearity.

How?

Activation functions!

# Activation Functions

1) Add non-linearity

2) Ensures small change in weights/bias
$$\Rightarrow \text{Small change in o/p}$$

Desirable for learning

3) In some cases,
maps $[-\infty, \infty] \rightarrow [a, b]$

$$x_1$$

$$x_2$$

$$\vdots$$

$$x_d$$

$$\sum w_i^0 x_i^0 + b$$

$$z$$

$$\sigma(z)$$

$$a$$

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

Activation

# Perceptron Algebraic Form

Step function

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

0

Small change in w or b can lead to large change in o/p.

# Sigmoid Neuron

* "Smoothed" out step function
* $\sigma(z) = 1/{1+e^{-z}}$



Small change in $w$ $\Rightarrow$ Small change in Output

# Activation Function

$x_1$ $w_1$

$x_2$ $w_2$

$+$

$w_1 x_1 + w_2 x_2 + b$

$\boxed{\text{Act}}$

$f(w_1 x_1 + w_2 x_2 + b)$

o/p is $f(wx+b)$

$f$ is non-linear

# SIGMOID UNIT

$x_1$   $w_1$

$x_2$   $w_2$

$\vdots$

$wd$

$x_d$   $b$

$1$

$O/p = \text{Activation} \left( \Sigma w_p x_i + b \right)$

$\quad = \sigma \left( \Sigma w_i x_i + b \right)$

$\sigma(z) = \dfrac{1}{1+e^{-z}}$

# Other activation functions

## Rectified Linear Unit (ReLU)

$f(z) = \max\{0, z\}$

z

## Tanh

$f(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$

1

−1

Simplest Neural Networks
MULTI LAYER PERCEPTRON (MLP)
( Layers of Sigmoid units)

I/P

Hidden Layers
( Neither o/p or I/p)

O/P

$x_1$

$x_2$

$\vdots$

$x_d$

Simplest Neural Networks

MULTI LAYER PERCEPTRON (MLP)

I/P

Hidden Layers
( Neither o/p or I/p )

O/p

$x_1$

$x_2$

$\vdots$

$x_d$

I/p also
a layer

# XOR USING MLP



$x_1$

$x_2$

$\hat{y}$

Hidden unit 1, Layer 1

Hidden unit 2, Layer 1

# XOR USING MLP



$$\hat{y} = \text{Activation of Layer 2 o/p}$$

$$= a^{[2]}$$

$$\text{I/p} = \text{Activation of layer 0} = a^{[0]} = [x_1 \; x_2]$$

1x2

# XOR USING MLP

Layer

Node

$x_1$

$x_2$

$z_1^{[1]}$

$a^{[2]}$

$\hat{y}$

$= a^{[2]}$

$$z_1^{[1]} = [x_1 \; x_2] \begin{bmatrix} \bullet \\ \bullet \end{bmatrix} + b_1^{[1]}$$

$a^{[0]}$

$w_1^{[1]}$

# XOR USING MLP

Layer

Node

$x_1$

$z_1^{[1]}$

$z_2^{[1]}$

$a^{[2]}$

$\hat{y}$

$= a^{[2]}$

$x_2$

$$z_1^{[1]} = [x_1 \; x_2] \begin{bmatrix} \bullet \\ \bullet \end{bmatrix} + b_1^{[1]}$$

$a^{[0]}$

$w_1^{[1]}$

$$z_2^{[1]} = [x_1 \; x_2] \begin{bmatrix} \bullet \\ \bullet \end{bmatrix} + b_2^{[1]}$$

$a^{[0]}$

$w_2^{[1]}$

# XOR USING MLP

$$x = [x_1 \ x_2] = a^{[0]}$$

$$z^{[1]}_1 \ _{1\times1} = a^{[0]}_{1\times2} \ w^{[1]}_1 \ _{2\times1} + b^{[1]}_1 \ _{1\times1}$$

$$z^{[1]}_2 \ _{1\times1} = a^{[0]}_{1\times2} \ w^{[1]}_2 \ _{2\times1} + b^{[1]}_2 \ _{1\times1}$$

$$\left[ z^{[1]}_1 \ z^{[1]}_2 \right]_{1\times2} = a^{[0]}_{1\times2} \left[ w^{[1]}_1 \ \ w^{[1]}_2 \right]_{2\times2} + \left[ b^{[1]}_1 \ \ b^{[1]}_2 \right]_{1\times2}$$

# XOR USING MLP

$$Z^{[1]} = \underset{1\times2}{a^{[0]}} \; \underset{2\times2}{w^{[1]}} + \underset{1\times2}{b^{[1]}}$$

$1\times2$

# XOR USING MLP



$$a_1^{[1]} = \sigma\left(z_1^{[1]}\right)$$

$$a_2^{[1]} = \sigma\left(z_2^{[1]}\right)$$

## XOR USING MLP

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} & a_2^{[1]} \end{bmatrix} = \sigma\left(z^{[1]}\right)$$

$$a_1^{[1]} = \sigma\left(z_1^{[1]}\right)$$

$$a_2^{[1]} = \sigma\left(z_2^{[1]}\right)$$

# XOR USING MLP



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

$$a^{[2]} = 6\left(z^{[2]}\right)$$

# XOR USING MLP

ONLY 1 POINT

$$x_1 = 1; \; x_2 = 1; \; y = 0$$

$$6 = RELU$$

$$w^{[1]} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}; \; b^{[1]} = [0 \; -1]; \; w^{[2]} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

$$b^{[2]} = 0$$

## XOR USING MLP

$$x_1 = 1; \quad x_2 = 1; \quad y = 0 \implies a^{[0]} = [1 \quad 1]$$

$$W^{[1]} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}; \quad b^{[1]} = [0 \quad -1]$$

$$z^{[1]} = a^{[0]} W^{[1]} + b^{[1]} = [1 \quad 1] \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + [0 \quad -1]$$

$$= [2 \quad 2] + [0 \quad -1]$$

$$= [2 \quad 1]$$

$$a^{[1]} = \sigma(z^{[1]}) = \max\{[0,0], [2,1]\} = [2,1]$$

# XOR USING MLP

$x_1 = 1; x_2 = 1; y = 0 \Rightarrow a^{[0]} = [1 \quad 1]$

$a^{[1]} = [2 \quad 1] \qquad\qquad w^{[2]} = [1 \quad -2] \; ; \; b^{[2]} = 0$

$$z^{[2]} = [2 \quad 1] \begin{bmatrix} 1 \\ -2 \end{bmatrix} + 0 = 0$$

$a^{[2]} = \sigma(z^{[2]}) = 0$

$\therefore \hat{y}(1,1) = a^{[2]} = 0$

# XOR USING MLP

let's redo for $x_1 = 0$; $x_2 = 1$; $y_{TRUE} = 1$

$$W^{[1]} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} ; \quad b^{[1]} = \begin{bmatrix} 0 & -1 \end{bmatrix} ; \quad W^{[2]} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} ; \quad b^{[2]} = 0$$

# XOR USING MLP

let's redo for $x_1 = 0$ ; $x_2 = 1$ ; $y_{TRUE} = 1$

$$W^{[1]} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \; ; \; b^{[1]} = \begin{bmatrix} 0 & -1 \end{bmatrix} \; ; \; w^{(2)} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \; ; \; b^{[2]} = 0$$

$$Z^{[1]} = a^{[0]} w^{[1]} + b^{[1]} = \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & -1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & -1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$a^{[1]} = \sigma(Z^{[1]}) = MAX\left\{ \begin{bmatrix} 0,0 \end{bmatrix}, \begin{bmatrix} 1,0 \end{bmatrix} \right\} = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

# XOR USING MLP

let's redo for $x_1 = 0; x_2 = 1 ; y_{TRUE} = 1$

$$w^{[1]} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} ; \quad b^{[1]} = \begin{bmatrix} 0 & -1 \end{bmatrix} ; \quad w^{[2]} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} ; \quad b^{[2]} = 0$$

$$a^{[1]} = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$z^{[2]} = a^{[1]} w^{[2]} + b^{[2]} = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix} + 0$$

$$= 1$$

$$a^{[2]} = 6 \left[ z^{[2]} \right] = 1$$

$$\therefore \quad \hat{y} (0,1) = a^{[2]} = 1 = y_{TRUE}$$

## COMPUTATION FOR 'M' INSTANCES

$$X = \begin{bmatrix} - & x_1 & - \\ - & x_2 & - \\ & \cdots & \\ - & x_M & - \end{bmatrix} \quad \text{where } x_i \in R^d$$

$M \times d$

$$X = \begin{bmatrix} a^{[0](1)} \\ a^{[0](2)} \\ a^{[0](M)} \end{bmatrix}$$

$(i)$ denotes instance #

# COMPUTATION FOR 'M' instances

$$X = \begin{bmatrix} - & x_1 & - \\ - & x_2 & - \\ & \cdots & \\ - & x_M & - \end{bmatrix} \quad \text{where } x_i \in R^d$$

$M \times d$

$$X = \begin{bmatrix} a^{[0](1)} \\ a^{[0](2)} \\ a^{[0](M)} \end{bmatrix}$$

$(i)$ denotes instance #

$$z^{[1](1)} = a^{[0](1)} w^{[1]} + b^{[1]}$$
$$z^{[1](2)} = a^{[0](2)} w^{[1]} + b^{[1]} \implies z^{[1]} = A^{[0]} w^{[1]} + b^{[1]}$$

$\cdots \cdots$

$$A^{[0]} = \begin{bmatrix} a^{[0](1)} \\ \vdots \\ a^{[0](m)} \end{bmatrix}$$

# MLP FOR XOR (OVER 'm' SAMPLES)

$$X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \quad ; \quad y_{TRUE} = \begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix}^T$$

$$W^{[1]} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad ; \quad b^{[1]} = \begin{bmatrix} 0 & -1 \end{bmatrix}$$

$$\therefore Z^{[1]} = \underset{4\times2}{A^{[0]}} \; \underset{2\times2}{W^{[1]}} + \underset{1\times2}{b^{[1]}} \quad \leftarrow \text{Broadcasted}$$

$$Z^{[1]} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + \begin{bmatrix} 0 & -1 \end{bmatrix}$$

# MLP FOR XOR (OVER 'm' SAMPLES)

$$X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \quad ; \quad y_{TRUE} = \begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix}^T$$

$$Z^{[1]} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + \begin{bmatrix} 0 & -1 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

$$A^{[1]} = \sigma(Z^{[1]}) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

# MLP FOR XOR (OVER 'm' SAMPLES)

$$X = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \quad ; \quad y_{TRUE} = \begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix}^T$$

$$A^{[1]} = \sigma\left(Z^{[1]}\right) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

$$w^{[2]} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \quad ; \quad b^{[2]} = 0$$

$$\therefore Z^{[2]} = A^{[1]} w^{[2]} + b^{[2]} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix} + 0 = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$= y_{TRUE}$$

# FORWARD & BACKWARD PROPAGATION



Backward propagation
Compute derivatives
&
update
weights

Forward propagation
Compute $\hat{y}$

# COMPUTATION GRAPH

$$\hat{y} = 6(2x_1 + 3x_2) \; ; \; LOSS = L(\hat{y}, y)$$



$$u = 2 * x_1 \; ; \; v = 3 * x_2$$
$$w = u + v \; ; \; a = 6(w)$$

# COMPUTATION GRAPH

$$\hat{y} = \sigma(w_1 x_1 + w_2 x_2) \; ; \; LOSS = L(\hat{y}, y)$$



$$\frac{\partial L}{\partial w_1} = ?$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial w} \frac{\partial w}{\partial v} \frac{\partial v}{\partial w_1}$$

# Derivative of Activation Functions

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

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z))$$

$$\sigma(z) = \text{MAX}\{0, z\}$$

$$\frac{\partial \sigma(z)}{\partial z} = \begin{cases} 0 & ; \ z < 0 \\ 1 & ; \ z > 0 \\ \text{undefined} & ; \ z = 0 \end{cases}$$

# Backpropagation for XOR Network

$x$

$W^{[1]}$

$b^{[1]}$

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

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

$W^{[2]}$
$b^{[2]}$

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

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

$L(a^{[2]}, y)$

# Backprogation for XOR Network

$x$

$w^{[1]}$

$b^{[1]}$

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

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

$w^{[2]}$
$b^{[2]}$

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

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

$$L(a^{[2]}, y)$$

$$L(a^{[2]}, y) = -y \log a^{[2]} - (1-y) \log(1-a^{[2]})$$

$$Q: \frac{\partial L(a^{[2]}, y)}{\partial w^{[2]}} = ?$$

# Backprogation for XOR Network

$$x$$

$$w^{[1]}$$

$$b^{[1]}$$

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

$$a^{[1]} = 6(z^{[1]})$$

$$w^{[2]}$$
$$b^{[2]}$$

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

$$a^{[2]} = 6(z^{[2]})$$

$$L(a^{[2]}, y)$$

$$L(a^{[2]}, y) = -y \log a^{[2]} - (1-y) \log (1 - a^{[2]})$$

$$\frac{\partial L(a^{[2]}, y)}{\partial z^{[2]}} = \frac{\partial L(a^{[2]}, y)}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}}$$

$$\frac{\partial L(a^{[2]}, y)}{\partial a^{[2]}} = -\frac{y}{a^{[2]}} + \frac{(1-y)}{1 - a^{[2]}} \quad \cdots \textcircled{1}$$

# Backprogation for XOR Network

$$x$$

$$w^{[1]}$$

$$b^{[1]}$$

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

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

$$w^{[2]}$$
$$b^{[2]}$$

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

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

$$L(a^{[2]}, y)$$

$$\frac{\partial a^{[2]}}{\partial z^{[2]}} = \sigma(z^{[2]})(1 - \sigma(z^{[2]})) \quad \cdots \quad \text{②}$$

$$\text{①} \& \text{②}$$

$$\frac{\partial L(a^{[2]}, y)}{\partial z^{[2]}} = \left\{ -\frac{y}{a^{[2]}} + \frac{(1-y)}{1 - a^{[2]}} \right\} \left\{ \sigma(z^{[2]})(1 - \sigma(z^{[2]})) \right\} \quad \cdots \quad \text{③}$$

# Backprogation for XOR Network

$x$

$w^{[1]}$

$b^{[1]}$

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

$$a^{[1]} = \sigma(z^{[1]})$$
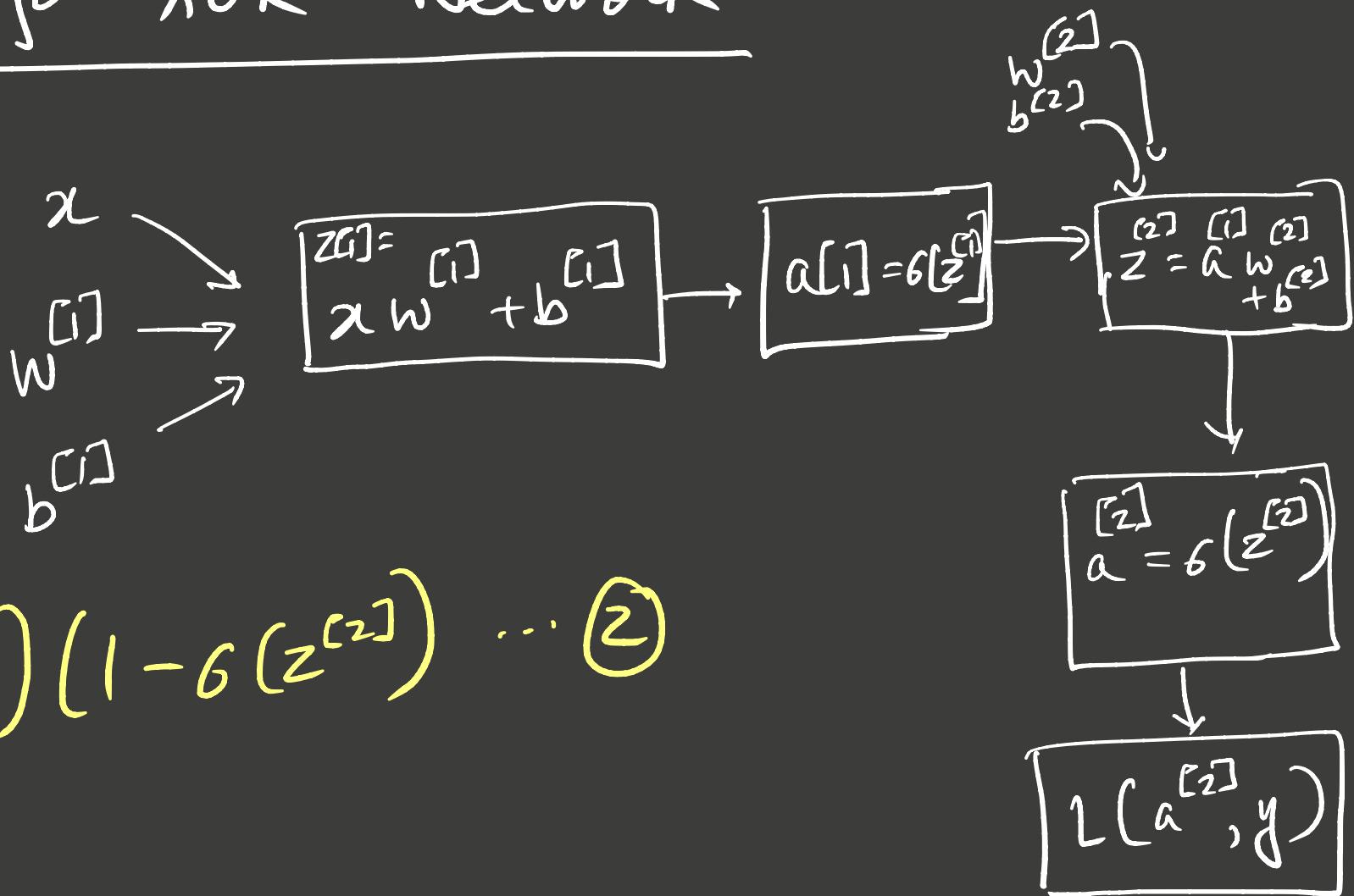
$w^{[2]}$, $b^{[2]}$

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

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

$$L(a^{[2]}, y)$$

$$\frac{\partial L(a^{[2]}, y)}{\partial w^{[2]}} = \frac{\partial L(a^{[2]}, y)}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial w^{[2]}}$$

$$= a^{[1]} \left\{ -\frac{y}{a^{[2]}} + \frac{(1-y)}{1-a^{[2]}} \right\} \left\{ \sigma(z^{[2]})(1-\sigma(z^{[2]})) \right\} \quad \ldots ④$$

Similarly

$$\frac{\partial L(a^{[2]}, y)}{\partial b^{[2]}} = \left\{ -\frac{y}{a^{[2]}} + \frac{(1-y)}{1-a^{[2]}} \right\} \left\{ \sigma(z^{[2]})(1-\sigma(z^{[2]})) \right\} \quad \ldots ⑤$$

# Backprogation for XOR Network

$x$

$w^{[1]}$

$b^{[1]}$

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

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

$w^{[2]}$
$b^{[2]}$

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

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

$$L(a^{[2]}, y)$$

$$\frac{\partial L(a^{[2]}, y)}{\partial z^{[1]}} = \frac{\partial L(a^{[2]}, y)}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}}$$

$w^{[2]}$

$$\sigma(z^{[1]})(1 - \sigma(z^{[1]}))$$

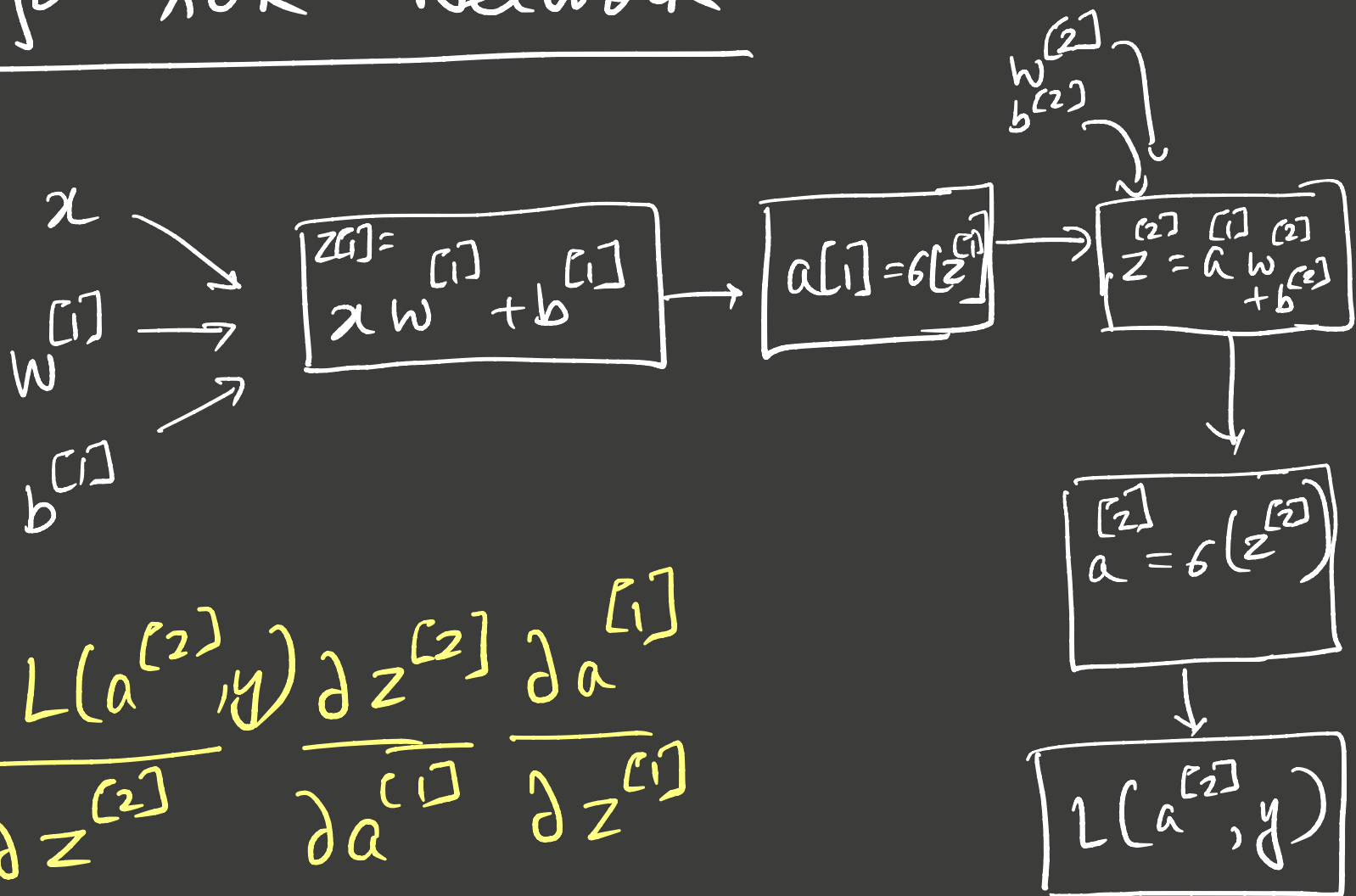$$\frac{\partial L(a^{[2]}, y)}{\partial w^{[1]}} = \frac{\partial L(a^{[2]}, y)}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial w^{[1]}} = \frac{\partial L(a^{[2]}, y)}{\partial z^{[1]}} * x$$

# WORKED OUT EXAMPLE

$x_1 = 1; x_2 = 1; y_{TRUE} = 0$ ; $\sigma^{(1)} = RELU$

$\sigma^{(2)} = RELU$

$w^{[1]}, b^{[1]}$

$w[2], b[2]$



$$W^{[1]} = \begin{bmatrix} w_1^{[1]} & w_2^{[1]} & \cdots \\ & & \end{bmatrix}$$

No. of I/ps = Dimensionality of $a^{[0]}$

$\underbrace{\qquad\qquad}$ # Hidden units in layer 1

$b^{[1]} = [\quad]$

$\longleftarrow$ # hidden units in Layer 1

# WORKED OUT EXAMPLE

$x_1 = 1; x_2 = 1; y_{TRUE} = 0$ ; $\sigma^{(1)} = RELU$
$\sigma^{(2)} = RELU$

$$w^{[1]} \, b^{[1]}$$

$$W[2], b[2]$$



$x_1$

$x_2$

$$W^{[1]} = \begin{bmatrix} .1 & .2 \\ .0 & .1 \end{bmatrix}$$

$$b^{[1]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

RANDOM INIT

$$W^{[2]} = \begin{bmatrix} .1 \\ .2 \end{bmatrix}$$

$$b^{[2]} = \begin{bmatrix} 0 \end{bmatrix}$$

# WORKED OUT EXAMPLE

$$x_1 = 1; x_2 = 1 ; y_{TRUE} = 0 \quad ; \quad \sigma^{(1)} = RELU$$
$$\sigma^{(2)} = RELU$$

$$a^{[1]} = \sigma\left(a^{[0]} w^{[1]} + b^{[1]}\right) = \sigma\left\{\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} .1 & .2 \\ .0 & .1 \end{bmatrix} + \begin{bmatrix} 0 & 0 \end{bmatrix}\right\} = \begin{bmatrix} .1 & .2 \end{bmatrix}$$

$$a^{[2]} = \sigma\left[a^{[1]} w^{[2]} + b^{[2]}\right] = \sigma\left[\begin{bmatrix} .1 & .2 \end{bmatrix} \begin{bmatrix} .1 \\ .2 \end{bmatrix} + \begin{bmatrix} 0 \end{bmatrix}\right] = .05$$

$$w^{[1]} = \begin{bmatrix} .1 & .2 \\ .0 & .1 \end{bmatrix} \qquad b^{[1]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$w^{[2]} = \begin{bmatrix} .1 \\ .2 \end{bmatrix} \qquad b^{[2]} = \begin{bmatrix} 0 \end{bmatrix}$$

# WORKED OUT EXAMPLE

$$x_1 = 1; x_2 = 1; y_{TRUE} = 0 \qquad ; \qquad \sigma^{(1)} = RELU$$
$$\sigma^{(2)} = RELU$$

$$a^{[1]} = \sigma\left(a^{[0]} w^{[1]} + b^{[1]}\right) = \sigma\left\{\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} .1 & .2 \\ .0 & .1 \end{bmatrix} + \begin{bmatrix} 0 & 0 \end{bmatrix}\right\} = \begin{bmatrix} .1 & .2 \end{bmatrix}$$

$$a^{[2]} = \sigma\left[a^{[1]} w^{(2)} + b^{(2)}\right] = \sigma\left[\begin{bmatrix} .1 & .2 \end{bmatrix} \begin{bmatrix} .1 \\ .2 \end{bmatrix} + \begin{bmatrix} 0 \end{bmatrix}\right] = .05$$

$$\text{Let } L(a^{[2]}, y) = \frac{1}{2}\left\{a^{[2]} - y\right\}^2$$

$$\therefore \frac{\partial L(a^{[2]}, y)}{\partial a^{[2]}} = a^{[2]} - y = .05 - 0 = .05$$

# WORKED OUT EXAMPLE

$$x_1 = 1; \; x_2 = 1; \; y_{TRUE} = 0 \qquad ; \qquad \sigma^{(1)} = RELU$$
$$\sigma^{(2)} = RELU$$

$$a^{[1]} = \sigma\left(a^{[0]} w^{[1]} + b^{[1]}\right) = \sigma\left\{\begin{bmatrix} 1 & 1 \end{bmatrix}\begin{bmatrix} .1 & .2 \\ .0 & .1 \end{bmatrix} + \begin{bmatrix} 0 & 0 \end{bmatrix}\right\} = \begin{bmatrix} .1 & .2 \end{bmatrix}$$

$$a^{[2]} = \sigma\left[a^{[1]} w^{[2]} + b^{[2]}\right] = \sigma\left[\begin{bmatrix} .1 & .2 \end{bmatrix}\begin{bmatrix} .1 \\ .2 \end{bmatrix} + \begin{bmatrix} 0 \end{bmatrix}\right] = .05$$

$$\therefore \boxed{\frac{\partial L(a^{[2]}, y)}{\partial a^{[2]}} = a^{[2]} - y = .05 - 0 = .05}$$

$$\frac{\partial L(a^{[2]}, y)}{\partial w^{[2]}} = \frac{\partial L(a^{[2]}, y)}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial w^{[2]}} = .05 * 1 * a^{[1]}$$

$$= \begin{bmatrix} .005 & .01 \end{bmatrix}$$

# WORKED OUT EXAMPLE

$$x_1 = 1; x_2 = 1; y_{TRUE} = 0 \qquad ; \qquad \sigma^{(1)} = RELU$$
$$\sigma^{(2)} = RELU$$

$$a^{[1]} = \sigma\left(a^{[0]} w^{[1]} + b^{[1]}\right) = \sigma\left\{ \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} .1 & .2 \\ .0 & .1 \end{bmatrix} + \begin{bmatrix} 0 & 0 \end{bmatrix} \right\} = \begin{bmatrix} .1 & .2 \end{bmatrix}$$

$$a^{[2]} = \sigma\left[ a^{[1]} w^{(2)} + b^{(2)} \right] = \sigma\left[ \begin{bmatrix} .1 & .2 \end{bmatrix} \begin{bmatrix} .1 \\ .2 \end{bmatrix} + \begin{bmatrix} 0 \end{bmatrix} \right] = .05$$

$$\frac{\partial L(a^{[2]}, y)}{\partial w^{[2]}} = \begin{bmatrix} .005 & .01 \end{bmatrix}$$

$$\frac{\partial L(a^{[2]}, y)}{\partial b^{(2)}} = \frac{\partial L(a^{[2]}, y)}{\partial a^{(2)}} * \frac{\partial a^{[2]}}{\partial z^{(2)}} * \frac{\partial z^{(2)}}{\partial b^{(2)}} = .05 * 1 = .05$$

$$\therefore \text{ update rule: } w^{(2)} = w^{(2)} - \text{learning} * \begin{bmatrix} .005 & .01 \end{bmatrix}$$
$$\text{Rate} \qquad --$$

# Digit Classifier using MLP

* 64 x 64 grayscale image

  * 0 $\longrightarrow$ Black

    1 $\longrightarrow$ white

    [0-1] $\longrightarrow$ B/w Black & white

* I/p layer : 64 x 64 = 4096

Question1 : Is new digit 9 or not?

$\Rightarrow$ O/p size = __1__ neurons

If Hidden layer sizes are

(100, 20, 1 (o|p layer)

what is # params?

If Hidden layer sizes are

$$(100, 20, 1 (o/p \, layer))$$

what is # params?

$$a^{[0]} = [ \, \text{------} \, ]_{1 \times 4096}$$

$$w^{[1]} = \begin{bmatrix} \uparrow \\ 4096 \\ \downarrow \\ \text{-----} \end{bmatrix}_{4096 \times 100}$$

$$\leftarrow 100 \longrightarrow$$

$$b^{[1]} = [ \quad ]_{1 \times 100}$$

If Hidden layer sizes are

$$(100, 20, 1 (\text{o/p layer}))$$

what is # params?

$$a^{[0]} = 1 \times 7096 \quad ; \quad w^{[1]} = 4096 \times 100 \quad ; \quad b^{[1]} = 1 \times 100$$

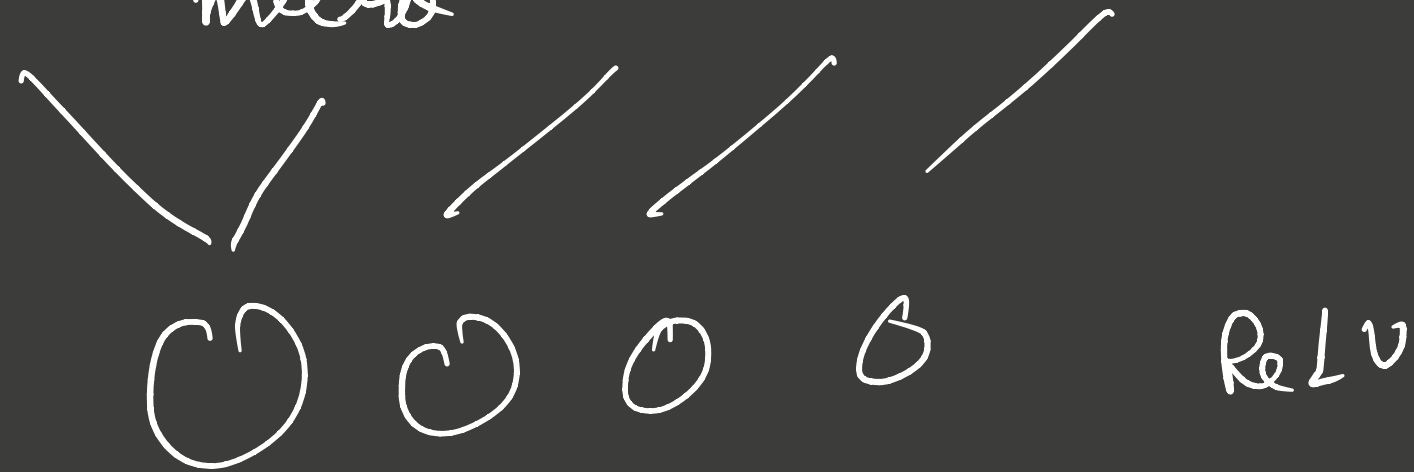$$z^{[1]} = a^{[0]} w^{[1]} + b^{[1]} = 1 \# 100$$

$$a^{[1]} = 1 \times 100$$

$$w^{[2]} = 100 \times 20 \quad ; \quad b^{(2)} = 1 \times 20$$
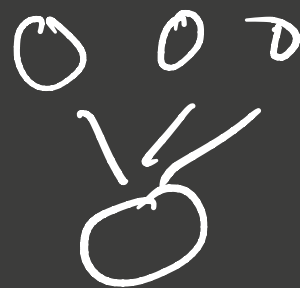
$$w^{[3]} = 1 \times 20 \qquad b^{[3]} = 1 \times 1$$

$$\text{Total params} = \sum_{i=1}^{3} \text{Size}(w^{[i]}) + \text{Size}(b^{[i]})$$

# Case Study I : Housing price prediction

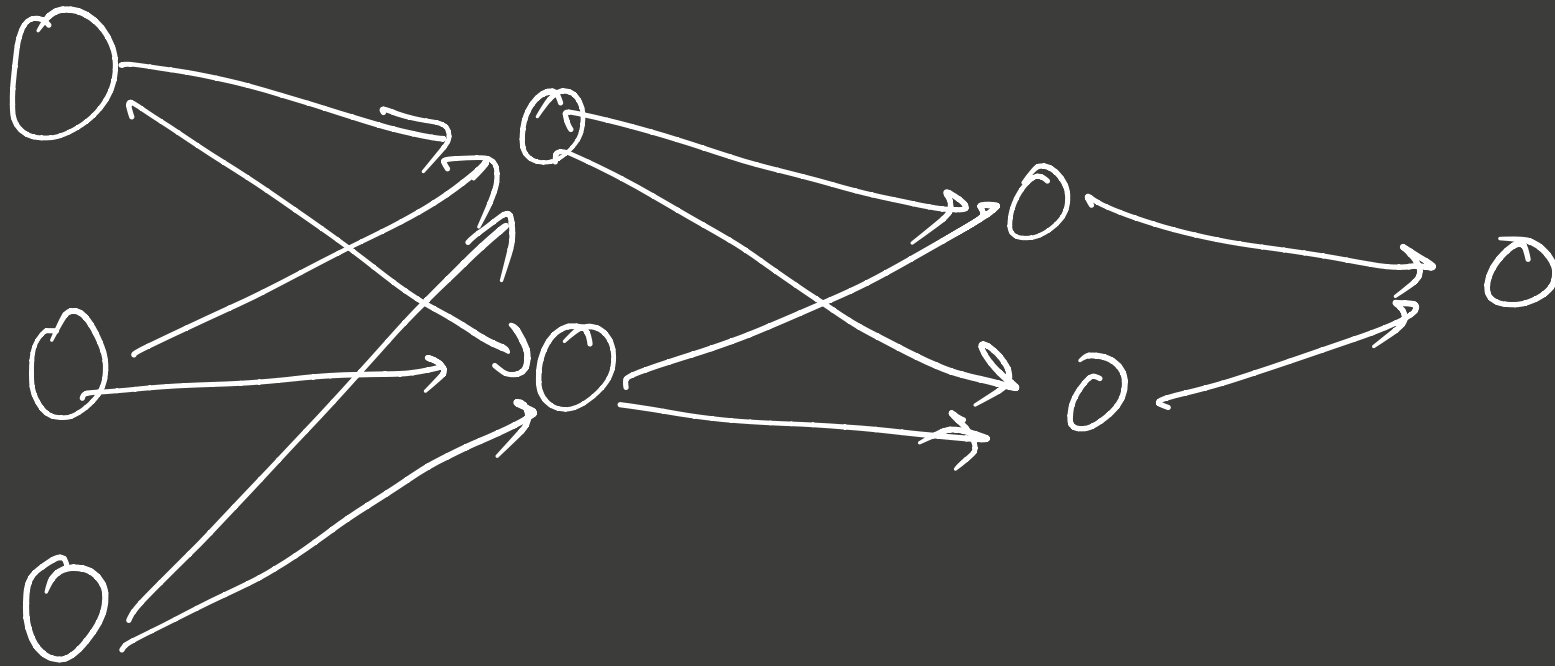$$x_i = \{ \text{Area}, \text{Distance to}, \text{\# schools}, \ldots \}$$

metro



ReLU

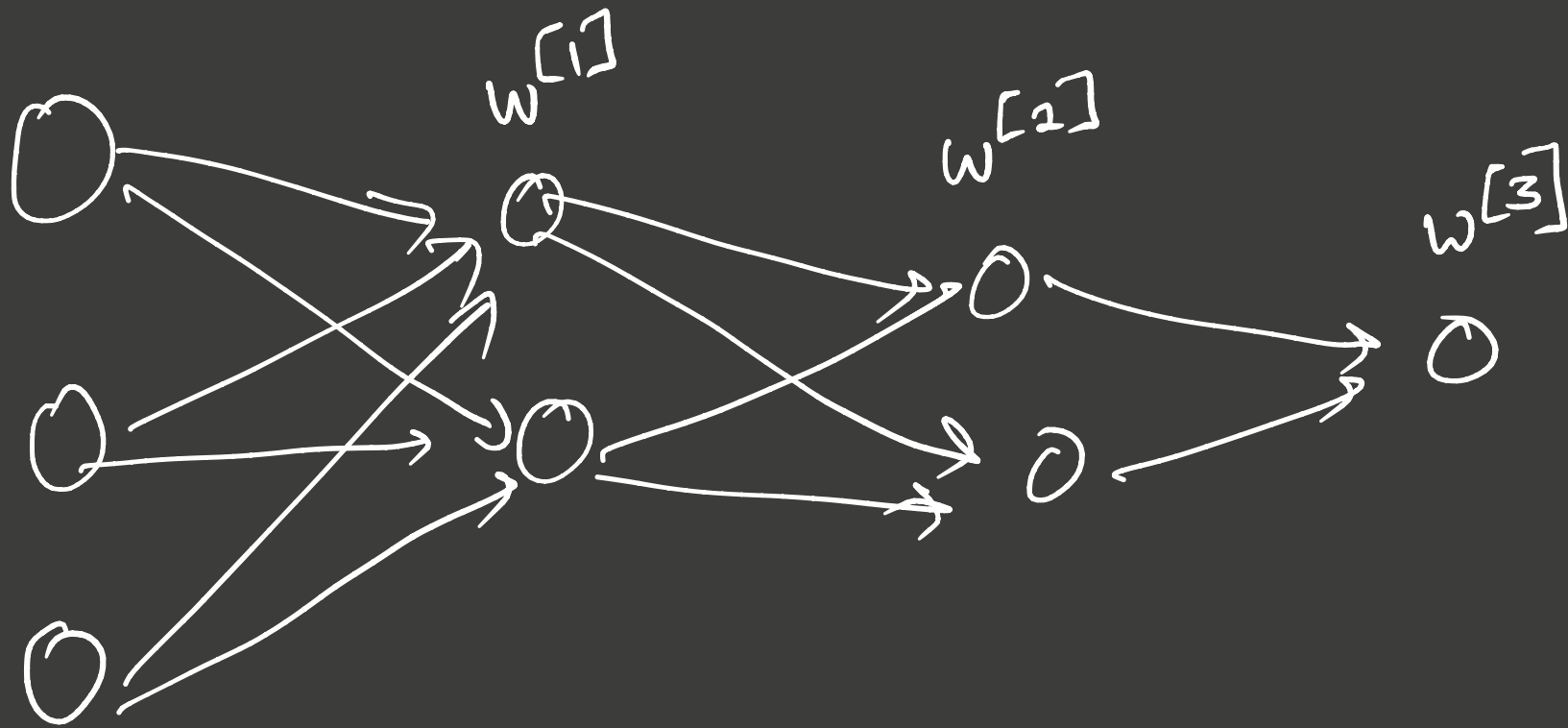Linear or ReLU activation

# REGULARIZATION / PREVENT OVERFITTING

① Dropout

② Penalty Terms.

③ Data augmentation

④ Early stopping
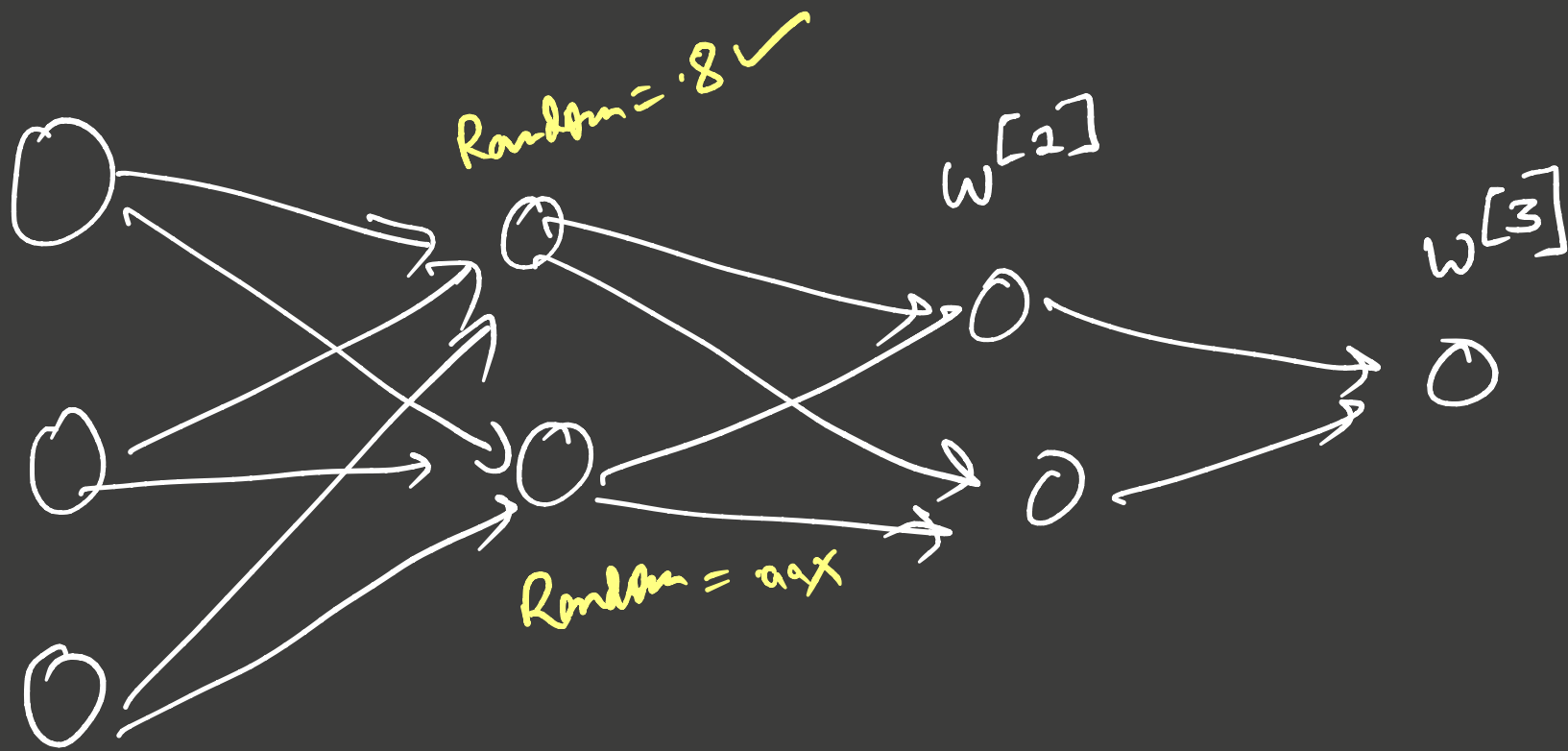
# Dropout



\* with a probability 'p' keep a node...

# Dropout



* with a probability $p$ (eg. 0.9) keep = a node
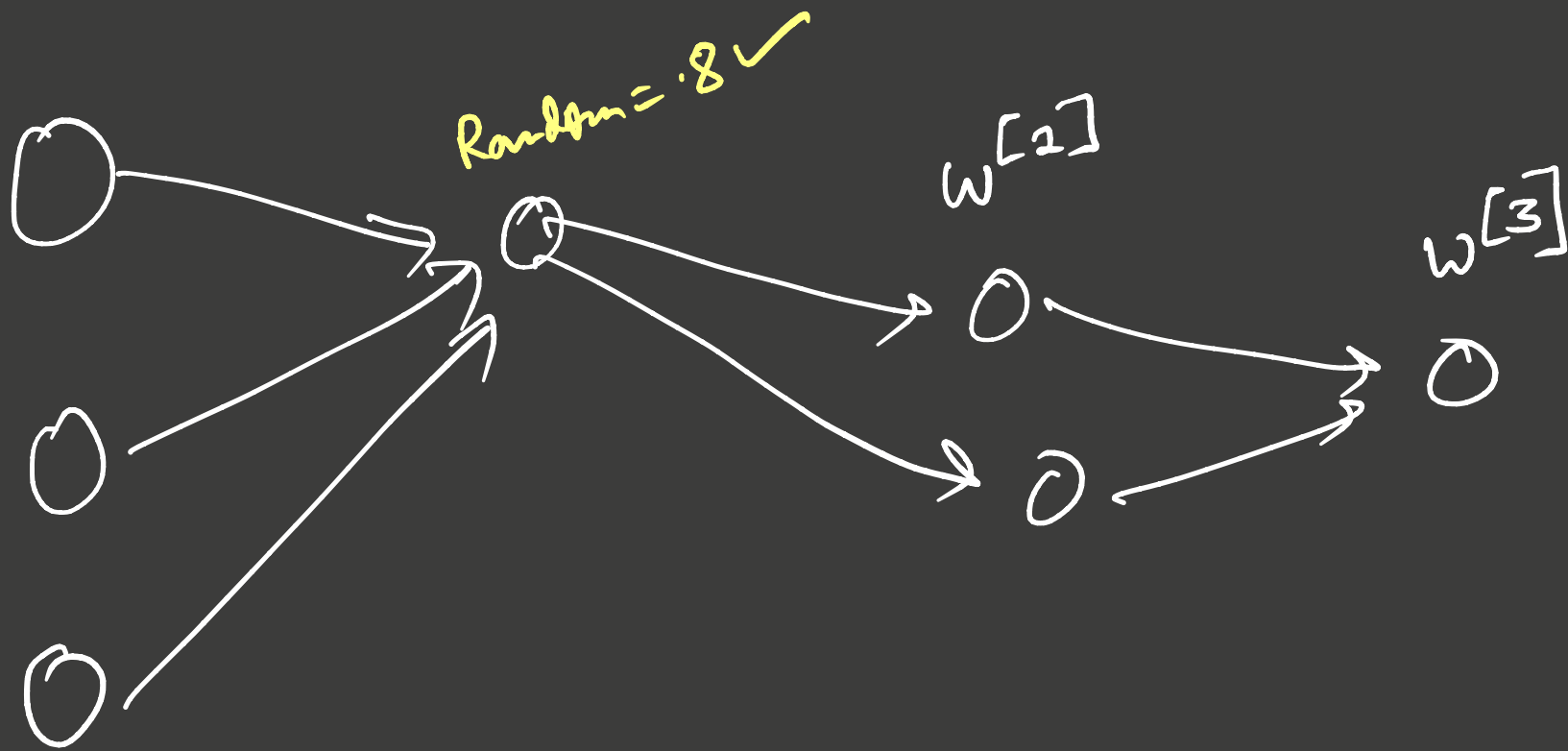
* Random() $< p$ : Keep ; Else : Drop

# Dropout



Random = .8 ✓

$w^{[2]}$

$w^{[3]}$

Random = aax

* with a probability $\underline{P}$ (eg. 0.9) keep a node

* Random() < $p$ : Keep; Else: Drop

# Dropout



Random $= .8$ ✓

$w^{[2]}$

$w^{[3]}$
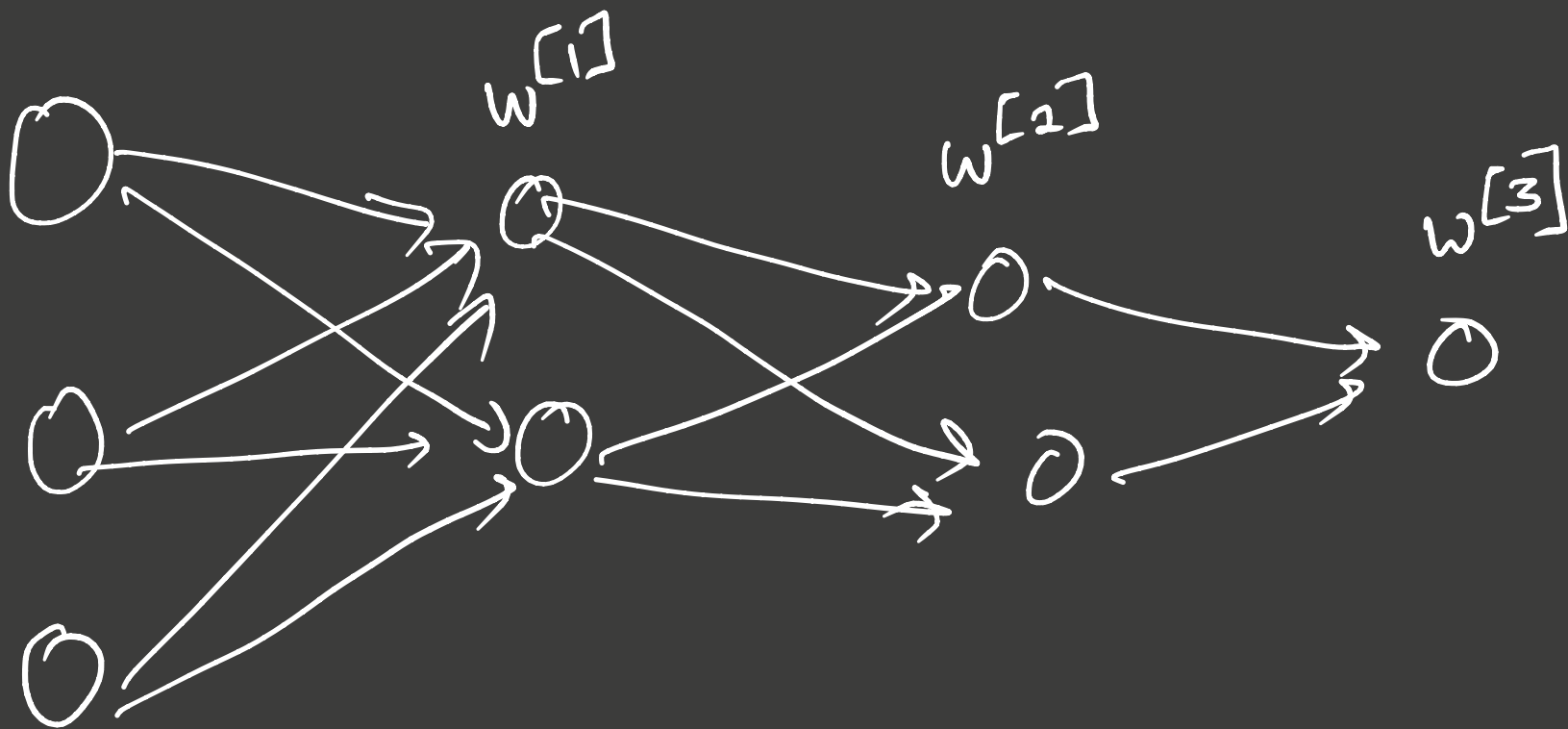
* with a probability $\underline{P}$ (eg. 0.9) keep a node

* Random ( ) $< P$ : Keep ; Else : Drop

# Dropout



$$A^{[1]} = A^{[1]} * MASK$$

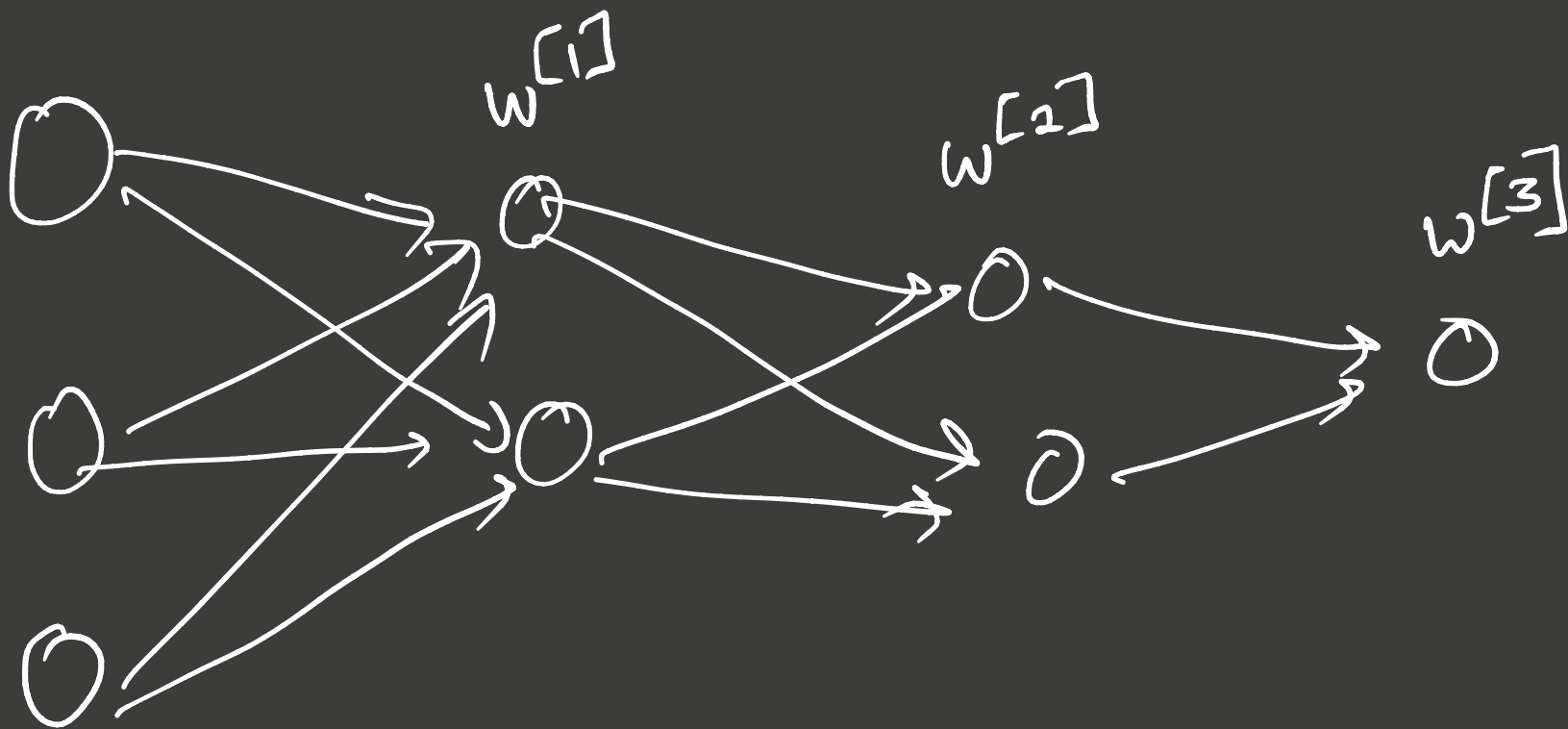$\left\{ \begin{array}{l} \text{Element-} \\ \text{wise} \\ \text{multiplication} \end{array} \right\}$

$$Mask = RANDOM(A^{[i]}.shape[0], A^{[i]}.shape[1])$$

$$< p$$

$$= \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

# Dropout



$W^{[1]}$  $W^{[2]}$  $W^{[3]}$

$$A^{[1]} = A^{[1]} * MASK$$

$$A^{[1]} = A^{[1]} / P$$

( why?

$\therefore$ we removed some nodes,

$E(A^{[1]})$ would reduce ...

$$Mask = RANDOM(A^{[1]}.shape(0), A^{[1]}.shape(1))$$

$$< P$$

$$= \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

# why Dropout works

① Smaller nets $\Rightarrow$ less overfitting

② Since nodes can be "Shut" at
   random,
   weight spread across nodes
   $$\Downarrow$$
   Shrinkage (akin $L_2$)

# REGULARISATION USING L1/L2 PENALTY

$$J\left(w^{[1]}, b^{[1]}, \ldots \, \omega^{[e]}, b^{[e]} \ldots \right)$$

$$= \sum_i L\left(\hat{y}_i, y_i\right) + \lambda \sum_{\ell=1}^{L} \|w^{[\ell]}\|^2$$

LOSS

L2
REGULARISATION

# Data Augmentation

Example: Add transformations of images
to make train set "bigger"