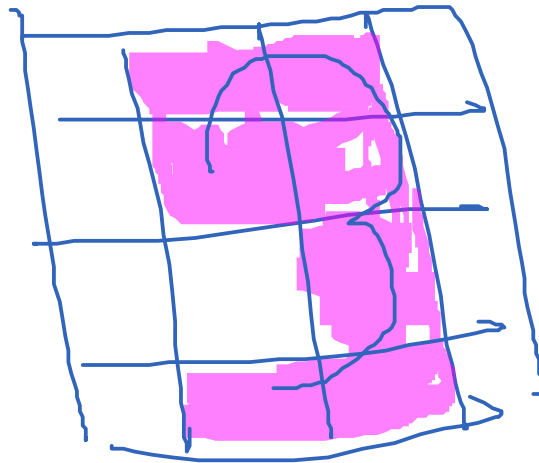# Neural Network Foundation

Tianchu.Zhao@uts.edu.au

# MNIST

- MNIST (Modified National Institute of Standards and Technology database)
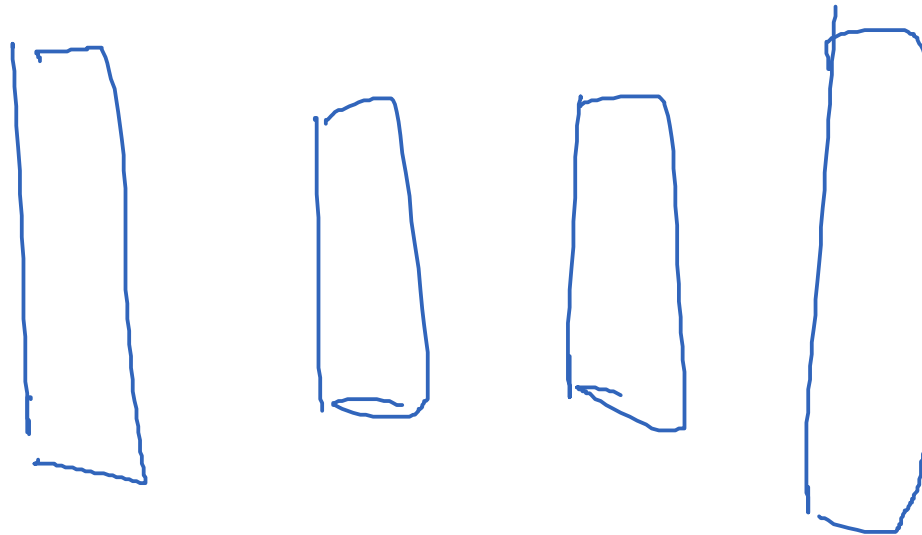- is a large database that contains handwritten digits images

- we can see our eye's retina as a matrix and when we see 3
- the corresponding locations lights up

- Now we want to build a computer program that is able to do the same
- Neural Network

28 {
| 0.1 | 0.1 | 0.1 | 0 |
| 0.1 | 0.1 | | |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 |

28

flatten

→ 784 {
| 0.1 |
| 0.1 |
| 0.1 |
| 0 |
| 0.1 |

(input)

10 {
9
1
2
⋮
8
9

⟶

10 {
| 0.1 |
| 0.2 |
| 0.1 |
| 0.5 |

(output)

- the neurons in hidden layer and the number of hidden layer can be any arbitrary number

- you need experimentation to determine what's the best for the task

# An explanation of how we recognise 3

- we hope that the hidden layer will do the same

$a_1$

$W_1$

strength of
___
in image

- we can compute weighted sum

$$w_1 a_1 + w_2 a_2 + \cdots$$

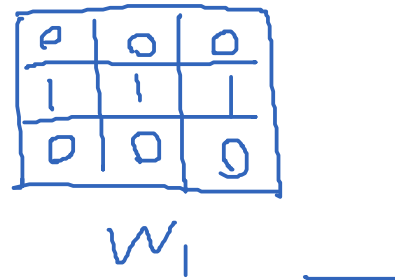- we can also compare weighted sum, so that we know which is more related to the desire number

- to do this we use a function called activation function
  - eg sidmoid

$$\sigma : \frac{1}{1+e^{-x}}$$

$$x \in R \quad \sigma \in (0, 1)$$

- then we have

$$\sigma(w_1 a_1 + w_2 a_2 + w_3 a_3 \ldots)$$

- we can also say that we only want this to be activated when weighted sum is > 5

$$\sigma(w_1 a_1 + w_2 a_2 + w_3 a_3 \ldots - 5)$$

- This represents the compute flow from 1 layer to anther layer

# Parameters(weights)

- In total we have the following parameters



weights: 784×4    4×4    4×10        = 3192
bias:       4        4      10        = 18
                                      = 3210

- Manually configure all 3210 parameters is are not feasible
- Thus we have machine learning
- The learning part is to let the computer find the appropriate weight and bias

- We first need to compute the output
- We first initialise all weights randomly and feed the data layer by layer to obtain output, this is called forward pass

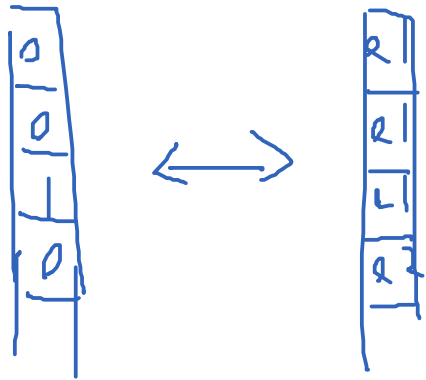$$\boxed{3} \rightarrow \quad a^1 \quad \rightarrow \quad a^2 \quad \longrightarrow \quad \text{output}$$

$$\text{output}: \quad \begin{array}{|c|} \hline 2.1 \\ \hline 0.1 \\ \hline 0.1 \\ \hline 0.2 \\ \hline \end{array} \quad \rightarrow \quad 4$$

- the output is random, because random weights
- to let the model learn the parameters, we need to know how much the current network is not doing right

$$\Sigma \ (\hat{y} - y)^2 \ \rightarrow \ \text{cost of 1 training example}$$

# Simple relation between weight and cost



$C(w)$

$w$

weight that optimise $C(w)$

find

- we can find out the derivative of C with respect to w

- derivative which measures the how much given parameters affect steepness of the function

- a partial derivative, which measures the how much one parameter affect steepness of the function holding other parameters constant

- the minimal Cost has the following properties

$$\frac{dC}{dw} C(w) = 0$$

- In reality the Cost vs weights graph is a lot more complex



- We can't calculate the minimum value of C straight away
- We use a strategy called gradient descent

- The idea is
- shift left if the slope is positive
- shift right if the slope is negative
- These will lower the cost, and once the Cost gets smaller, this means we are closer to the real prediction
- repeat this step will reach the local minima

$$\vec{W} \begin{bmatrix} \\ \\ \\ \end{bmatrix} \implies$$

negative
gradient

$$-\nabla C(\vec{w}) = \begin{bmatrix} \ell \\ -0.5 \end{bmatrix}$$

each weight has different
gradient

- The algorithm for computing this gradient efficiently is called back propagation

$$\text{forward pass} \quad a^0 \rightarrow a^1 \rightarrow a^2 \rightarrow a^3$$

$$\text{backprop} \quad a^{L-3} \underset{w^{L-2} b^{L-2}}{\longleftarrow} a^{L-2} \underset{w^{L-1} b^{L-1}}{\longleftarrow} a^{L-1} \underset{w^L b^L}{\longleftarrow} a^L$$

$$\text{Cost} \quad C = (\underline{a^L} - y)^2$$

$$a^L = \sigma(z^L)$$

$$z^L = w^L \underline{a^{L-1}} + b^L$$

$$w^L b^L \cdot a^{L-1} \rightarrow z^L \rightarrow a^L \rightarrow C$$

$$\frac{dC}{dw^L} = \frac{dz^L}{dw^L} \frac{da^L}{dz^L} \frac{dC}{da^L}$$

$$C = (a^L - y)^2$$

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

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

$$\frac{dC}{da^L} = 2(a^L - y) \qquad \left[\frac{1}{u(x)}\right]' = \frac{u'(x)}{u(x)^2}$$

$$\frac{da^L}{dz^L} = \sigma'(z^L) \longrightarrow \frac{d}{dz}\left[\frac{1}{1+e^{-z}}\right] = \frac{-\frac{d}{dz}[1+e^{-z}]}{(1+e^{-z})^2}$$

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

$$\frac{dz^L}{dw^L} = a^{L-1} \qquad\qquad = \frac{e^z}{(e^z+1)^2}$$

we have

$$\frac{dC}{dw^{(L)}} = a^{L-1} \cdot \sigma'(z^L) \cdot 2(a^L - y)$$

- this is just for 1 training example,
- such that this image want to change this parameter this way, another image want to change this parameter another way
- to change weight so that it is suitable for all training examples we average them

$$\frac{\partial C}{\partial w^L} = \frac{1}{n} \sum_{k=1}^{n} \frac{\partial C_k}{\partial w^L}$$

- Put everything together, we have

$$\text{forward pass} \quad a^0 \rightarrow a^1 \rightarrow a^2 \rightarrow a^3$$

$$\text{backprop} \quad a^{L-3} \leftarrow a^{L-2} \leftarrow a^{L-1} \leftarrow \dot{a}^L$$
$$\qquad\qquad w^{L-2} b^{L-2} \qquad w^{L-1} b^{L-1} \qquad w^L b^L$$

$$X \nearrow z^{L-2} = a^{L-2} \nearrow z^{L-1} = a^{L-1} \nearrow z^L = a^L \nearrow C$$
$$\quad w^{L-2} \qquad\qquad w^{L-1} \qquad\qquad w^L \qquad\qquad Y \nearrow$$
$$\quad b^{L-2} \qquad\qquad b^{L-1} \qquad\qquad b^L$$