# Intros

Structure: Weekly meetings, more like office hours.
  ↳ also, show me what has been done in the code.

Go over rough idea:

## To do:

▷ Watch 3b1b videos

▷ Find meeting time for next week
( July 3ʳᵈ only day I am available )

▷ Code implementation
( probably best to use python. Google collab .. ? )
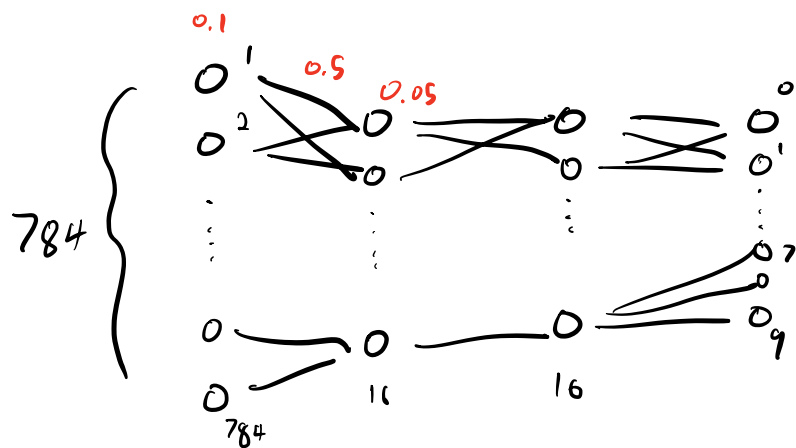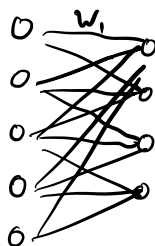
▷ Data processing:
• Use Mnist or USPS ?

▷ Implementation: no hidden layers first ?

▷ Recommended: Students meet once per week for
a couple hours to try to code the neural network.

$28\{$ $\boxed{7}$ $\longleftarrow\longrightarrow$ vector w/ $28 \cdot 28 = \underline{784}$ components.

$\underbrace{\phantom{xxxx}}_{28}$

$\uparrow$
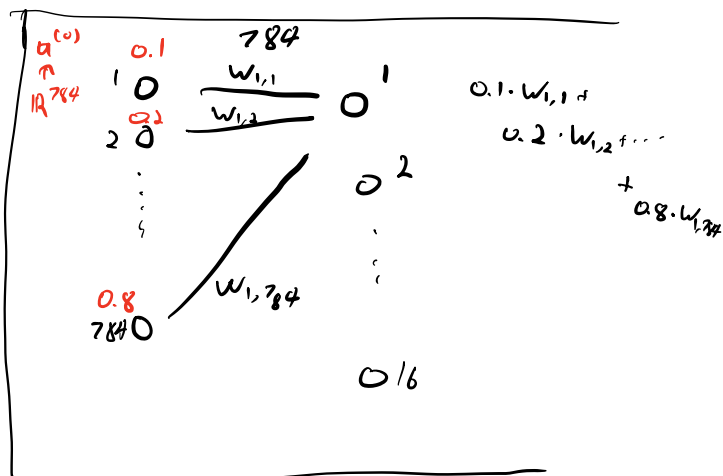
$\mathbb{R}^{784}$



784 {

0.1 0.5 0.05

(multilayer perceptron)

First layer: 784 neurons.

input space: $[0,1]^{784}$

Final layer: 10 neurons

2 inner layers, 16 neurons each.

Layer 0        Layer 1

$$\mathbb{R}^{784} \longrightarrow \mathbb{R}^{16}$$

$$a^{(0)} \longmapsto \sigma\left(W^{(1)}a^{(0)} + b^{(1)}\right)$$

or ReLu

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

Sigmoid.

$$\begin{pmatrix} W_{1,1} & \cdots & W_{1,784} \\ W_{2,1} & \cdots & W_{2,784} \\ & \vdots & \\ W_{16,1} & \cdots & W_{16,784} \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ b_{16} \end{pmatrix} \qquad a^{(0)} \begin{pmatrix} a_1^{(0)} \\ \vdots \\ a_{784}^{(0)} \end{pmatrix}$$
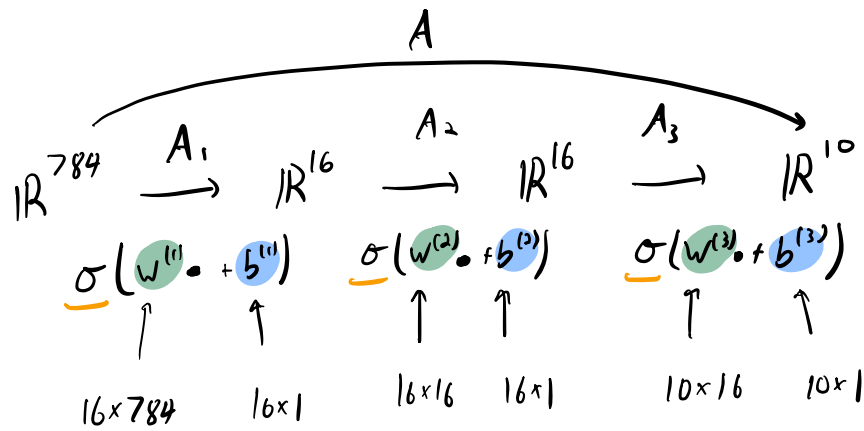
$W_1 \in M_{16 \times 784}(\mathbb{R})$    weight matrix.

$b_1 \in M_{16 \times 16}(\mathbb{R})$    bias matrix.

So $a^{(1)} = \sigma\left(W^{(1)}a^{(0)} + b^{(1)}\right)$

$$a_i^{(1)} = \sigma\left(\left(\sum_{j=0}^{784} W_{i,j}a_j^{(0)}\right) + b_i^{(1)}\right)$$

Layer 0　　　Layer 1　　　Layer 2　　　Layer 3

$$\overset{\displaystyle A}{\overbrace{\qquad\qquad\qquad\qquad\qquad\qquad}}$$

$$\mathbb{R}^{784} \xrightarrow{A_1} \mathbb{R}^{16} \xrightarrow{A_2} \mathbb{R}^{16} \xrightarrow{A_3} \mathbb{R}^{10}$$

$$\underline{\sigma}\left( W^{(1)} \bullet + b^{(1)} \right) \quad \sigma\left( W^{(2)} \bullet + b^{(1)} \right) \quad \sigma\left( W^{(3)} \bullet + b^{(3)} \right)$$

$$16 \times 784 \qquad 16 \times 1 \qquad 16 \times 16 \quad 16 \times 1 \qquad 10 \times 16 \quad 10 \times 1$$

Overall    activation

$$A = \sigma\left( W^{(3)} \cdot \sigma\left( W^{(2)} \cdot \sigma\left( W^{(1)}(\bullet) + b^{(1)} \right) + b^{(2)} \right) + b^{(3)} \right)$$

$$A_3 \circ A_2 \circ A_1 (\bullet)$$

Tactivation    function:

$$A : \mathbb{R}^{784} \longrightarrow \mathbb{R}^{10} \qquad A = A_{W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, W^{(3)}, b^{(3)}}$$

$$( \text{notation:} \quad A_{\underline{W}, \underline{b}} )$$

| Choices | | # parameters : | Choices | | # parameters |
|---|---|---|---|---|---|
| $W^{(1)} \in M_{784 \times 16}(\mathbb{R})$ | | $784 \times 16 = 12,544$ | $b^{(1)} \in \mathbb{R}^{16}$ | | 16 |
| $W^{(2)} \in M_{16 \times 16}(\mathbb{R})$ | | $16 \times 16 = 256$ | $b^{(2)} \in \mathbb{R}^{16}$ | | 16 |
| $W^{(3)} \in M_{16 \times 10}(\mathbb{R})$ | | $16 \times 10 = 160$ | $b^{(3)} \in \mathbb{R}^{10}$ | | 10 |

$\underline{\text{total}}$ # $\underline{\text{parameters}}$ :     13,002 .

## Training:

Cost   function   built   from   a   training   $\underline{\text{Set}}$ .

Collection   of   pairs

$$TS = \{ (I_1, L_1), (I_2, L_2), \dots, (I_n, L_n) \}$$

where   $I_k \in \mathbb{R}^{784}$   is an   input,   $L_1 \in \mathbb{R}^{10}$   is   the   correct output.

$$C = C_{TS} : \mathbb{R}^{13,002} \longrightarrow \mathbb{R},$$

really

$$C: M_{784 \times 16} \times \mathbb{R}^{16} \times M_{16 \times 16} \times \mathbb{R}^{16} \times M_{16 \times 10} \times \mathbb{R}^{10} \longrightarrow \mathbb{R}$$

$$\left( W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, W^{(3)}, b^{(3)} \right) \longmapsto \frac{1}{n} \sum_{i=1}^{n} \| A_{\underline{w},\underline{b}} (I_i) - L_i \|^2$$

Average over all images in the training data.

distance   of   the   guess   $\left( A_{\underline{w},\underline{b}} (I_i) \right)$   from   the   answer   $L_i$

# Intuition of Cost function:

$C(\underline{W}, \underline{b})$ is small $\Longleftrightarrow$ the weights $\underline{W}, \underline{b}$ do a good job guessing the correct answer.

$\underset{W^{(1)}, W^{(2)}, W^{(3)}}{\|} \quad b^{(1)}, b^{(2)}, b^{(3)}$

$\Longrightarrow$ $\boxed{\text{Goal}}$ Should be to finds weights

$\underline{W}, \underline{b}$ minimizing $C$.

# Gradient Descent

Idea: ① Start with Random inputs $(\underline{w}_0, \underline{b}_0)$ to $C$.

② Find the direction of $(\underline{w}_0, \underline{b}_0)$ that leads the the larges decrease in $C$.

③ Change $(\underline{w}_0, \underline{b}_0)$ to $(\underline{w}_1, \underline{b}_1)$ according to the direction from ②.

More precisely:

For ②, use $-\nabla C(\underline{w}, \underline{b})$.

$$C: \mathbb{R}^n \longrightarrow \mathbb{R}, \qquad \nabla C: \mathbb{R}^n \longrightarrow \mathbb{R}$$

$$\nabla C = \begin{pmatrix} \frac{\partial C}{\partial x_1} \\ \vdots \\ \frac{\partial C}{\partial x_n} \end{pmatrix}$$

$$\left( \text{e.g.} \quad \begin{matrix} C: \mathbb{R}^2 \longrightarrow \mathbb{R} \\ (x,y) \longmapsto x^2 y + y \end{matrix} \quad , \quad \nabla C = \begin{pmatrix} \frac{\partial C}{\partial x} \\ \frac{\partial C}{\partial y} \end{pmatrix} = \begin{pmatrix} 2xy \\ x^2 + 1 \end{pmatrix} \right)$$

⁑ To do: details for Gradient descent ⁑

$$= \frac{1}{n} \sum_{r=1}^{2} \left\| \underbrace{A_3 A_2 A_1 (a^{(i)})}_{a^L} - L \right.$$

$$\overline{\underset{a^{(0)}}{\Gamma_i}}$$

$0 \quad v^{(1)} \quad a^{(0)} \quad b^{(1)}$

$1 \quad w^{(2)} \quad a^{(1)} \quad b^{(2)}$

$2 \quad w^{(3)} \quad a^{(2)} \quad b^{(3)}$

$3 \quad Li \quad a^{(3)}$

$C$

Cost