# 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 ..? )
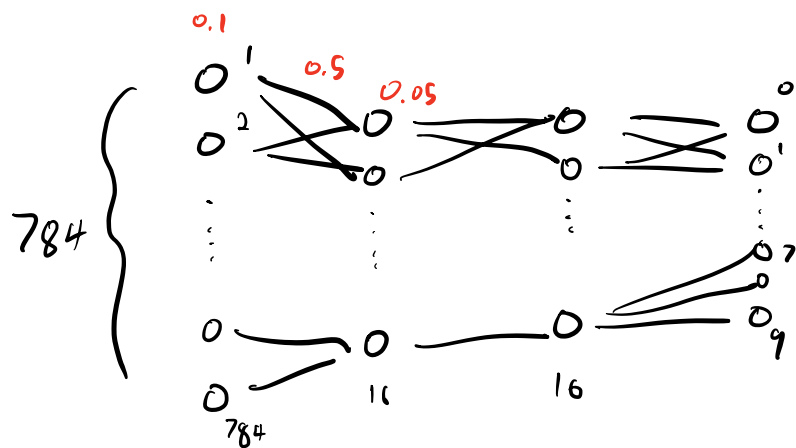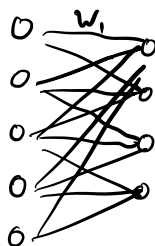
▷ Data processing :
    • Use Mnist or USPS ?

▷ Implementation : no hidden layers first ?

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

$28 \{$ $\boxed{7}$ $\longleftrightarrow$ 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)

$a^{(0)}$ → $\mathbb{R}^{784}$

0.1    784
1 O   $\underline{W_{1,1}}$   O 1    $0.1 \cdot W_{1,1} +$
$\overline{W_{1,2}}$     $0.2 \cdot W_{1,2} + \cdots$
0.2
2 O      O 2       $+ 0.8 \cdot W_{1,784}$
⋮
4
0.8   $W_{1,784}$   ⋮
784 O
      O 16
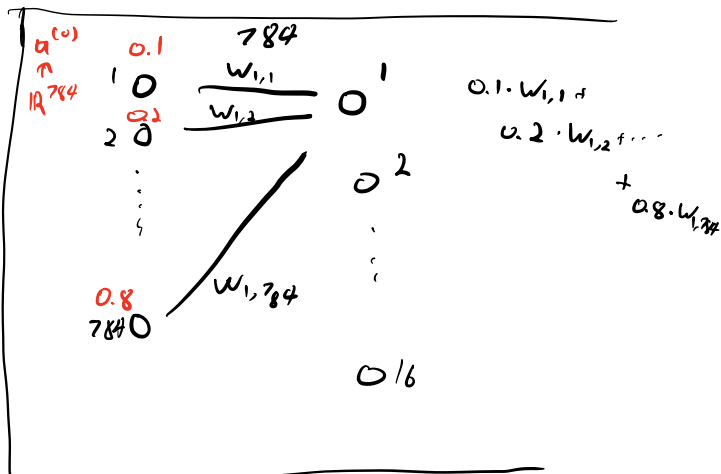
First layer : 784 neurons.

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

Final layer : 10 neurons

2 inner layers, 16 neurons each.

Layer 0        Layer 1

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

$a^{(0)} \longmapsto \sigma\left( W^{(1)} a^{(0)} + b^{(1)} \right)$    $\underline{\sigma}(x) = \dfrac{1}{1+e^{-x}}$

                      Sigmoid.

$$\begin{pmatrix} W_{1,1} & \cdots & W_{1,784} \\ W_{2,1} & \cdots & W_{2,784} \\ & \vdots & \\ W_{16,1} & \cdots & W_{16,784} \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ b_{16} \end{pmatrix} \qquad a^{(0)} \begin{pmatrix} a^{(0)}_1 \\ \vdots \\ a^{(0)}_{784} \end{pmatrix}$$
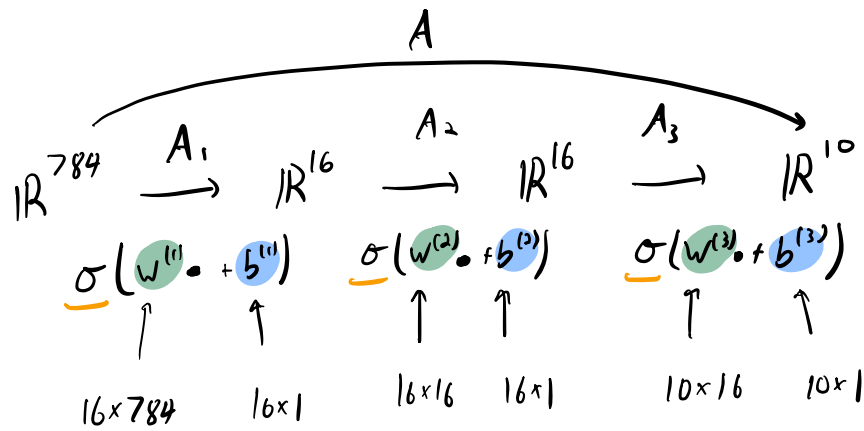
$W_1 \in M_{16 \times 784}(\mathbb{R})$    weight matrix.

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

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

$a^{(1)}_i = \sigma\left( \left( \displaystyle\sum_{j=0}^{784} W_{i,j} a^{(0)}_j \right) + b^{(1)}_i \right)$

Layer 0         Layer 1         Layer 2         Layer 3

$$A$$

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

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

$$16\times784 \qquad 16\times1 \qquad\quad 16\times16 \quad 16\times1 \qquad\quad 10\times16 \quad 10\times1$$

Overall   activation

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

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

activation   function :

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

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

| Choices | | # parameters : | Choices | | # parameters |
|---------|--|---------------|---------|--|-------------|
| $W^{(1)} \in M_{784\times16}(\mathbb{R})$ | | $784\times16 = 12{,}544$ | $b^{(1)} \in \mathbb{R}^{16}$ | | 16 |
| $W^{(2)} \in M_{16\times16}(\mathbb{R})$ | | $16\times16 = 256$ | $b^{(2)} \in \mathbb{R}^{16}$ | | 16 |
| $W^{(3)} \in M_{16\times10}(\mathbb{R})$ | | $16\times10 = 160$ | $b^{(3)} \in \mathbb{R}^{10}$ | | 10 |

total # parameters :   13,002 .

## Training:

Cost function built from a training 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

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

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

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 $\iff$ the weights $\underline{w}, \underline{b}$ do a good job guessing the correct answer.

$\underset{\underline{w}^{(1)}, \underline{w}^{(2)}, \underline{w}^{(3)}}{\overset{\|}{}} \quad b^{(1)}, b^{(2)}, b^{(3)}$

$\implies$ $\boxed{\text{Goal}}$ 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:**

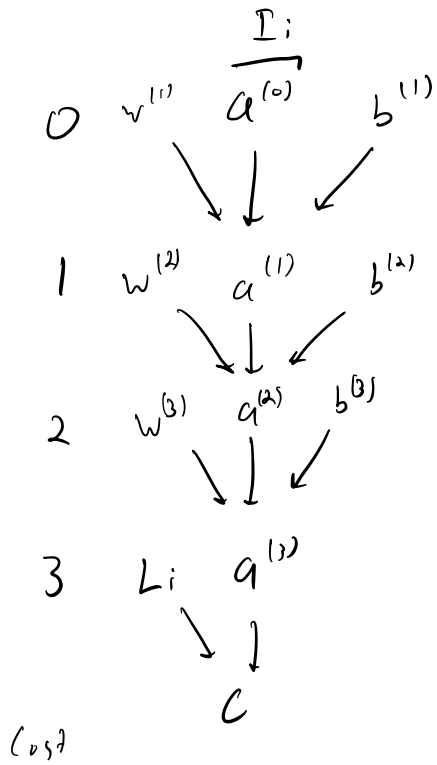For ②, use $-\nabla C(\underline{w}, \underline{b})$.

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

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

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

- - - - - - - -
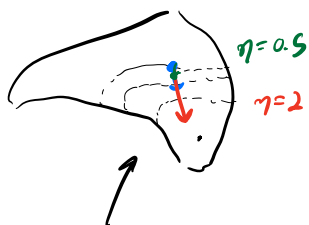
☆ To do: details for Gradient descent ☆

$$= \frac{1}{n} \sum_{r=1}^{n} \| \underbrace{A_3 A_2 A_1 (a^{(0)})}_{a^L} - L$$

$$\underline{I_i}$$

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

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

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

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

$C$

Cost

July 3$^{rd}$

$\eta$

$$\underline{w_0} - \eta \nabla \cdot \mathcal{L}(\underline{w_0})$$

$\eta = 0.5$
$\eta = 2$

$y = \sigma(z)$

$\sigma(z - 2)$

$1$

$\frac{1}{2}$

$2$

$z$

$\boxed{\text{Recalling} \quad \text{activation}}$

Layer 0      Layer 1      Layer 2      Layer 3
(input layer)                                    (output layer)

$j-1 \xrightarrow{\;w^{(i)}\;} i$



\# neurons     784       16       16       10

$\boxed{\text{activations:}}$

Notation: $\quad w_{jk}^{(i)}$

i) **weight** by which neuron $k$ in layer $i-1$ influences neuron $j$ in layer $i$
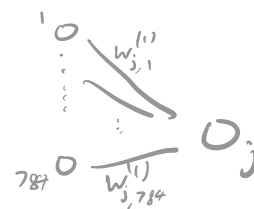
## From Layer 0 to layer 1:

**initial activation:**
(input to layer zero)

$$a^{(0)} = \begin{pmatrix} a_1^{(0)} \\ \vdots \\ a_{784}^{(0)} \end{pmatrix} \in \mathbb{R}^{784}.$$

activation of $j^{th}$ neuron in $1^{st}$ layer:     $(j = 1, \cdots, 16)$

$$a_j^{(1)} = \sigma\left( \left[ \sum_{k=1}^{784} w_{jk}^{(1)} a_k^{(0)} \right] + b_j^{(1)} \right)$$

define: $\quad z_j^{(1)} = \quad$

Keep track of all neurons at once using vectors:

$$a^{(1)} := \begin{pmatrix} a_1^{(1)} \\ \vdots \\ a_{16}^{(1)} \end{pmatrix} \in \mathbb{R}^{16} \quad \& \quad z^{(1)} := \begin{pmatrix} z_1^{(1)} \\ \vdots \\ z_{16}^{(1)} \end{pmatrix} \in \mathbb{R}^{16}$$

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

$j,k^{th}$ entry $\quad w_{jk}^{(1)}$

**All together:** $a^{(1)} = \sigma(z^{(1)})$,      $z^{(1)} = \underbrace{W^{(1)}}_{\substack{16 \times 784 \\ \text{matrix}}} \cdot \underbrace{a^{(0)}}_{\substack{784 \times 1 \\ \text{vector}}} + \underbrace{b^{(1)}}_{\substack{16 \times 1 \\ \text{vector}}}$

$\underbrace{\phantom{a^{(1)}}}_{\substack{16 \times 1 \\ \text{vector}}}$ $\underbrace{\phantom{z^{(1)}}}_{\substack{16 \times 1 \\ \text{vector}}}$

## Layer 1 to layer 2

$a^{(2)} = \sigma(z^{(2)})$,      $z^{(2)} = \underbrace{W^{(2)}}_{16 \times 16} \cdot \underbrace{a^{(1)}}_{16 \times 1} + \underbrace{b^{(2)}}_{16 \times 1}$

$\underbrace{\phantom{a}}_{16 \times 1}$ $\underbrace{\phantom{z}}_{16 \times 1}$

## From layer 2 to layer 3

$a^{(3)} = \sigma(z^{(3)})$,      $z^{(3)} = \underbrace{W^{(3)}}_{10 \times 16} \cdot \underbrace{a^{(2)}}_{16 \times 1} + \underbrace{b^{(3)}}_{10 \times 1}$

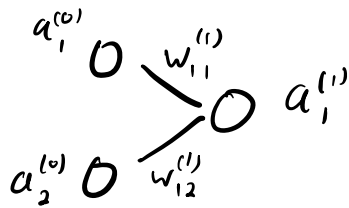$\underbrace{\phantom{a}}_{16 \times 1}$ $\underbrace{\phantom{z}}_{10 \times 1}$

---

## Final cost:

$$C = \underbrace{\frac{1}{\#\ x\ in}}_{\substack{\text{training} \\ \text{data}}} \sum_{\substack{x\ in\ training \\ data}} C_x ,$$

$x \in \mathbb{R}^{784}$

$$C_x = \| \underbrace{\overbrace{A_{\underline{w},\underline{b}}(x)}^{10 \times 1}}_{\substack{\text{guess based on} \\ \text{given weights \&} \\ \text{biases.}}} - \overbrace{\text{actual value}(x)}^{10 \times 1} \|^2$$

need to find      $\nabla C = \frac{1}{\#\dots} \nabla \cdot C_x$

## Simple example:



$$a^{(0)} = \begin{pmatrix} a_1^{(0)} \\ a_2^{(0)} \end{pmatrix}, \qquad a^{(1)} = \sigma(z^{(1)}),$$

$$z^{(1)} = \underset{1 \times 2}{\left( w_{11}^{(1)} \quad w_{12}^{(1)} \right)} \underset{2 \times 1}{\begin{pmatrix} a_1^{(0)} \\ a_2^{(0)} \end{pmatrix}} + \underset{1 \times 1}{b_1^{(1)}}$$

$$= w_{11}^{(1)} a_1^{(0)} + w_{12}^{(1)} a_2^{(0)} + b_1^{(1)}$$

variables in cost function

---

**Exercise:** If $w^{(1)} = (2 \quad \ln(2))$ and $a^{(0)} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$,

& $b^{(1)} = 0$,

find the activation of the neuron in the output layer.

---

$\forall : \frac{1 - e^{-z}}{1} = \frac{1+1}{1} = \frac{1}{3}$

Suppose our training data is one single

data point:

$$\left( \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \quad 1 \right)$$

(This means the input of $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$ to the first layer

should output the "vector" 1 (rember there is only one nerron).

## Cost function:

$$\| \sigma( v_{11}^{(1)} \cdot 1 + w_{12}^{(1)} \cdot 2 + b_1^{(1)} ) - 1 \|^2$$

$w_1 = w_{11}^{(1)}$

$w_2 = v_{12}^{(1)}$

$b = b_1^{(1)}$

Cost function

$$C(w_1, w_2, b) := \left( \sigma(w_1 + 2w_2 - b) - 1 \right)^2$$

## Exercise: Find the Gradient of the cost function.

$$\frac{\partial C}{\partial w_1} = 2 \left( \sigma(w_1 + 2w_2 - b) - 1 \right) \cdot \sigma'(w_1 + 2w_2 - b) \cdot \underline{1}$$

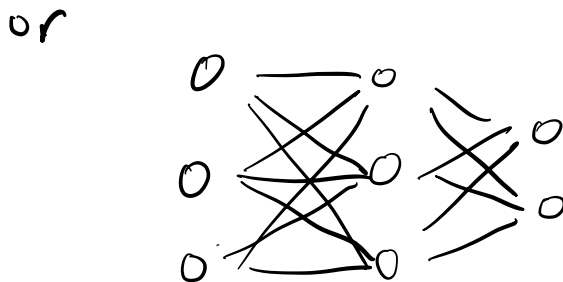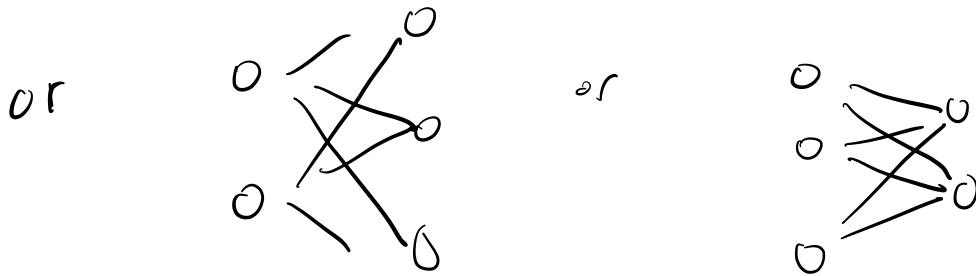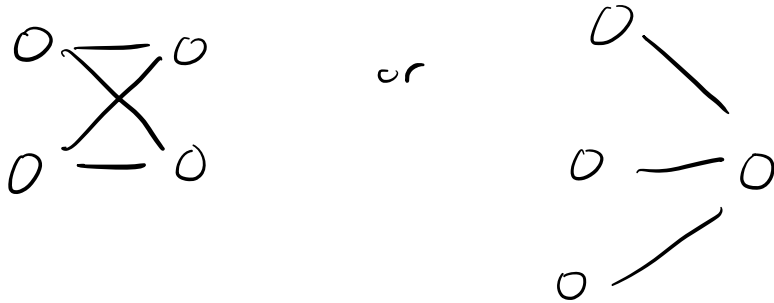$$\frac{\partial w_1}{\partial w_1} = 1$$

$$\frac{\partial C}{\partial w_2} = 2 \left( \sigma(w_1 + 2w_2 - b) - 1 \right) \cdot \sigma'(v_1 + 2w_2 - b) \cdot 2$$

$$\frac{\partial(w_1 + 2w_2 - b)}{\partial w_2} = 2.$$

$$\frac{\partial C}{\partial b} = 2 \left( \sigma(w_1 + 2w_2 - b - 1 \right) \cdot \sigma'(w_1 + 2w_2 - b) \cdot (-1)$$

$$\nabla C = \begin{pmatrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial w_2} \\ \frac{\partial C}{\partial b} \end{pmatrix}$$

## Exercise: Go through chapter 2 in the book, and for all the formulas displayed there, write what they are in this simple case.

Exercise: Do the same for other simple neural networks,

e,g.



or



or



or



or

"Assignment":

▷ Reading: Ch1 & Ch2 of the book.
(More theoretical)
math!

▷ Programming: ( Recommended as group work! )
- Run the code in the book for MNIST digits
( experiment with using different structures of neural
networks, e.g. of shape [784, 16, 16, 10] )

- Run the code in the book for USPS digits.
experiment to find good parameters.

- If there is some other data set you want
to analyze, go for it!

Next meeting w/ me: July 17 - 21
( in person or zoom... )

Students only meeting next week encouraged!