

## Cross entropy loss.

—————, n examples

$$y : \downarrow \begin{bmatrix} y_{11} & y_{12} & \dots & y_{1n} \\ y_{21} \\ \vdots \\ y_{10,1} & & y_{10,n} \end{bmatrix}_{10 \times n}$$

Labels.

$$\hat{y} = \begin{bmatrix} \hat{y}_{11} & \dots & \hat{y}_{1n} \\ \vdots & \ddots & \vdots \\ \hat{y}_{10,1} & \dots & \hat{y}_{10,n} \end{bmatrix}_{10 \times n}$$

$$\text{Loss} = \left\{ -\sum y_{r1} \log(\lambda_1) \quad \dots \quad -\sum y_{rn} \log(\hat{\lambda}_{rn}) \right\}_{1 \times n}$$

$y_{r1} \rightarrow$  ~~the~~  $r^{\text{th}}$  row, 1<sup>st</sup> col.

$$-\sum y_{ri} \log(\hat{y}_{ri}) \rightarrow -\underbrace{(\sum y_{11} \log(\hat{y}_{11}) + y_{21} \log(\hat{y}_{21}) + \dots)}_{\text{multi cross entropy loss.}}$$

if we sum up the "loss" matrix's all columns  
then we get the cost of model.

def. forward Pass (input\_list)

\* inputs-list shape : (42,000, 784)

targets-list. shape : (42,000, 10)

px1 px2 px3 --  $\xrightarrow{\text{Pixel values}}$  px784

inputs-list:

Ex 1  
Ex 2

Example Ex. 42,000

]

target-list:

Example 2  
Example 3  
42,000

0 1 2  
0 0 1  
0 0 0  
0 0 0

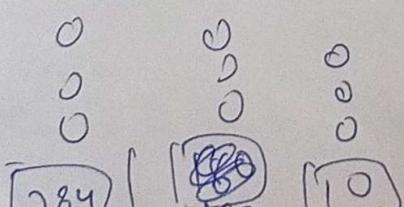
$\xrightarrow{\text{label}}$

this means  
first image has  
label = 2

lets take hidden neur = 160

input neur = 784

output " = 10



weight  
matrix

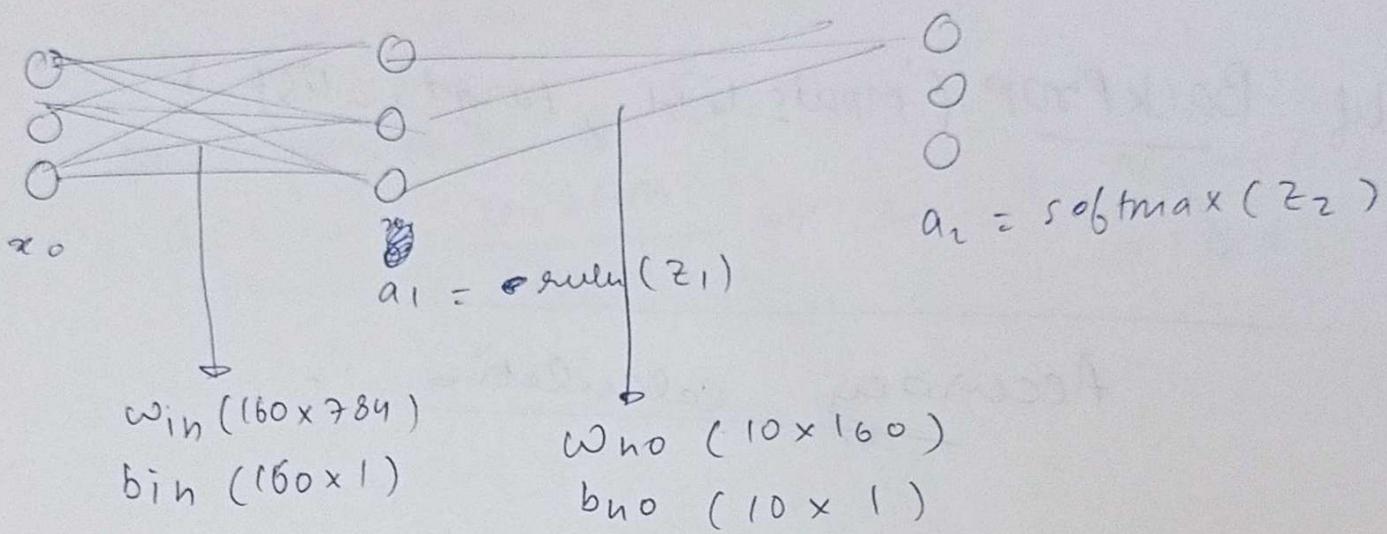
$$w_{in} = (160, 784)$$

$\xrightarrow{784}$   
 $160 \left[ w_{in} \right]$

$b_{in} = (160, 1) \quad [160] \quad [ ]$

inputs = (Input.list)<sup>T</sup>

$$\text{inputs} = (784, 42000) = x_0$$



$$z_1 = w_{ho} \cdot x_0 + b_{ho}$$

$$= \underbrace{\begin{bmatrix} 160 \times 784 \end{bmatrix}}_{w_{in}} \underbrace{\begin{bmatrix} 784 \times 42000 \\ x_0 \end{bmatrix}} + \begin{bmatrix} \text{bin} \\ \vdots \\ \vdots \end{bmatrix}_{160 \times 1}$$

$$= \begin{bmatrix} 160 \times 42000 \end{bmatrix} + \begin{bmatrix} 160 \times 1 \end{bmatrix}$$

$\downarrow$   
broadcast to  
 $(160, 42000)$

$$z \in \{160, 42000\}$$

$$q_1 = \text{relu } z_1 = 160, 42000$$

$$\begin{aligned} z_2 &= w_{ho} \cdot q_1 + b_{ho} \\ &= (10 \times 160) - (160 - 42000) + (10 \times 1) \\ &\quad a_2 = (10, 42000) \end{aligned}$$

# Backprop.

- ① input-list
- ② target list

- ① return loss ~~top~~
- ② and update params.

$$L = \text{self.CEL}(y, \hat{y})$$

$\hat{y}$  comes from forward pass.

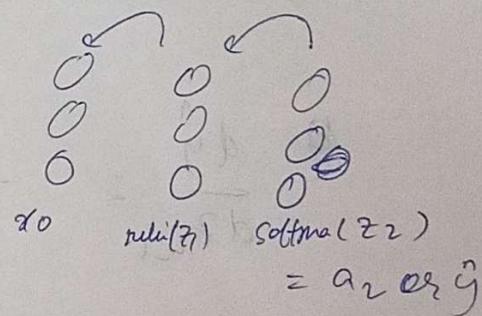
$$\hat{y} = \text{forward(inputs\_list)}$$

$$L = \checkmark$$

$$q_2 = \hat{y}$$

now for param update,

$$\frac{dL}{dw_{ho}} \quad \cancel{\frac{dL}{dy}} \quad \cancel{\frac{dL}{dz_1}} \times \cancel{\frac{dy}{dz_1}} \times \cancel{\frac{dL}{d\hat{y}}}$$



$$L = -\varepsilon(y * \log(\hat{y})) + \text{const}$$

$$\hat{y} = \text{softmax}(z_2)$$

$$z_1 = w_{ho} \cdot a_1 + b_{ho}$$

$$a_1 = \text{relu}(z_1)$$

$$z_1 = w_{in} \cdot x_0 + b_{in}$$

Terminologies

~~a~~  $x_0$  = input

~~a~~  $z_1 = z_h$   $a_2 = \hat{y}$

~~a~~  $z_2 = z_o$

$z_2 = z_o$

$$\frac{dL}{dw_{ho}} = \frac{dL}{d\hat{y}} \times \frac{d\hat{y}}{dz_2} \times \frac{dz_2}{d w_{ho}}$$

↳ same for  $b_{ho}$

$$\frac{dL}{dw_{in}} = \frac{dL}{d\hat{y}} \times \frac{d\hat{y}}{dz_2} \times \frac{dz_2}{da_1} \times \frac{da_1}{dz_1} \times \frac{dz_1}{d w_{in}}$$

↳ same for  $b_{in}$

$\frac{dL}{d\hat{y}} = \text{\$ cross-entropy-loss-derivative}$

$$= -\sum_i \frac{d(\hat{y}_i \log \hat{y}_i)}{d\hat{y}_i} - \sum_i -\hat{y}_i \frac{1}{\hat{y}}$$

$\hat{y}_i$  is  
the  $i$ -th  
label

$$\frac{d\hat{y}}{dz_2} = \frac{d\text{softmax}(z_2)}{dz_2} = \frac{d}{dz_2} \left( \frac{e^{z_i}}{\sum e^{z_i}} \right)$$

on solving,

$$\frac{dL}{dz_2} = \frac{dL}{d\hat{y}} \times \frac{d\hat{y}}{dz_2} = \hat{y} - y$$

$$\begin{aligned} \frac{dL}{dw_{h0}} &= (\hat{y} - y) \frac{d z_2}{d w_{h0}} = (\hat{y} - y) \cdot a_1^T \\ &\quad \downarrow \\ &\quad (10, 42000) \rightarrow (160, 42000)^T \\ &\quad = (42000, 160) \\ &\quad = (10, 160) \end{aligned}$$

$$\frac{dL}{dw_{ih}} = \frac{dL}{dz_2} \times \frac{d z_2}{da_1} \times \frac{da_1}{dz_i} \times \frac{dz_i}{dw_{ih}} = (\hat{y} - y) \times \frac{d \text{softmax}(z_2)}{d z_i}$$

$$\begin{aligned}
 \frac{dL}{dw_{in}} &= (\hat{y} - y) \times \omega_{ho} \times \frac{d(\text{relu } z_1)}{dz_1} \times x_0 \\
 &= (\hat{y} - y) \times \omega_{ho} \times \cancel{\frac{d(\text{relu } z_1)}{dz_1}} \times x_0 \\
 &\quad \xrightarrow{\text{relu-derivative}(z_1)} \\
 &= (\omega_{ho}^T \cdot \hat{y} - y) \cdot x_0^T \times \text{rd}(z) \\
 &\quad \downarrow \qquad \downarrow \\
 &\quad (10 \times 160)^T \quad (10,42,000) \xrightarrow{T} (784, 42,000)^T \\
 &= (160 \times 784)
 \end{aligned}$$

same for

$$\frac{dL}{db_{in}} \quad \text{and} \quad \frac{dL}{db_{ho}}$$

Fit

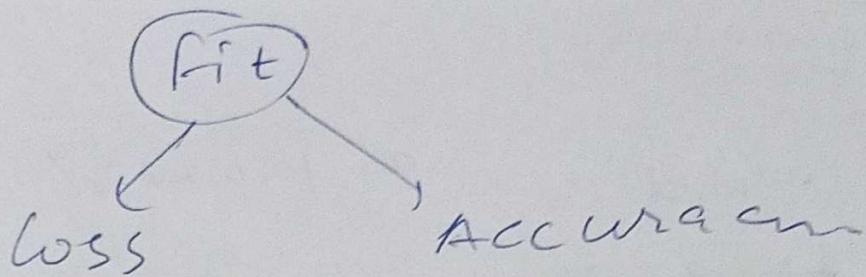
Accuracy (train & val)

$$(X_{\text{train}}) \text{ inputs\_list} = (42,000, 784) \quad } \rightarrow \text{train}$$

$$(y_{\text{train}}) \text{ target\_list} = (42,000, 10) \quad }$$

$$X_{\text{val}} \text{ validation\_labels} = (12,000, 784) \quad } \rightarrow \text{val}$$

$$(y_{\text{val}}) \text{ valid\_data} = (12,000, 10) \quad }$$



① loss.

$t\_loss = \text{self.backprop(inputs\_list, targets\_list)}$

Or ~~REPLACED~~

$tloss = \text{np.mean}(\text{CEL}(tr-a, tr-p))$

$tr-a = \text{targets\_list.T}$  ~~(\*)~~

$tr-p = \text{forward}(\text{inputs\_list})$

$vloss = \text{np.mean}(\text{CEL}(v-a, v-p))$

$v-a = \text{val\_label.T}$

$v-p = \text{forward}(\text{val-} \cancel{\text{label}} \text{data})$

## ② Accuracy

tr-pred = ~~for~~ forward( $X_{\text{train}}$ ).T

~~tr-a~~ =

---

tr-pred  $\rightarrow$  np.argmax(tr-pred, axis=1)

tr-a = np.argmax(tr-a, axis=1)

some for val.

## Momentum

$$\beta = 0.9.$$

$$\omega = \omega - \frac{dL}{dw} \times l \cdot r \rightarrow \text{Before momentum.}$$

$$w = \omega - \frac{v_w}{\beta} \times l \cdot r \rightarrow \text{after momentum.}$$

$$v_w = \beta \times w + (1 - \beta) \times \frac{dL}{dw}$$

$$\therefore v_b = \beta \times b + (1 - \beta) \times \frac{dL}{db}$$

$$i6 \quad \beta = 0$$

$$v_w = \frac{dL}{dw}, \quad v_b = \frac{dL}{db}.$$