

X-train shape (60000 x 28 x 28)

↓ Flatten

X-train shape (60000, 784)

X-test shape (10000, 784)

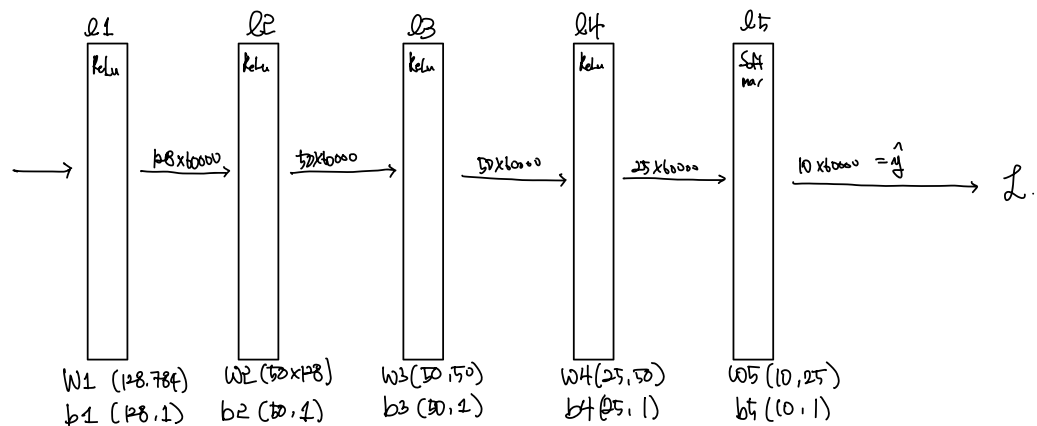
y-train shape (10 x 60000)

y-test shape (10 x 10000)

<input layer>

X-train.T,

(784, 60000)



$$W1 = X_{train.T} + b1$$

$$(128, 784) \quad (784, 60000) + (128, 1)$$

$$\Rightarrow (128, 60000) \Rightarrow Z1$$

$$ReLU(128, 60000) \Rightarrow A1$$

$$(10, 60000)$$

$$(10, 60000)$$

$$(10, 25) \quad (10, 1)$$

$$A5, Z5, L, W5, b5$$

$$\Rightarrow (10, 60000)$$

$$Z_5 = W_5 \cdot A_4 + b_5$$

$$W_5 = W_5 - \frac{dL}{dW_5} \cdot \alpha$$

$$\frac{dL}{dW_5} = \frac{dZ_5}{dW_5} \times \frac{dA_5}{dZ_5} \times \frac{dL}{dA_5}$$

$$L = -y \ln(A_5)$$

$$\frac{d}{dA} L = \frac{y}{A_5}$$

$$Z = \begin{pmatrix} z_0 \\ z_1 \\ z_2 \\ \vdots \\ z \end{pmatrix} \Rightarrow \sigma(Z) = \begin{pmatrix} \frac{e^{z_0}}{e^{z_0} + 1} \\ \vdots \\ \frac{e^{z_n}}{e^{z_n} + 1} \end{pmatrix}$$

$$\frac{dA_5}{dz_5} = \begin{pmatrix} a_0(1-a_0) & -a_0 a_1 & \dots & -a_0 a_9 \\ -a_0 a_1 & & & \vdots \\ \vdots & & & \vdots \\ a_9 a_0 & \dots & \dots & a_9(1-a_9) \end{pmatrix}$$

$$\frac{dL}{dA_5} = \begin{pmatrix} -\frac{y_0,1}{a_0,1} & & & -\frac{y_1,10000}{a_0,10000} \\ -\frac{y_1}{a_1} & 10000 & & \vdots \\ \vdots & & \ddots & \vdots \\ -\frac{y_9}{a_9} & & & -\frac{y_9,10000}{a_9,10000} \end{pmatrix}$$

$$\frac{dA_5}{dz_5} \times \frac{dL}{dA_5} = \begin{pmatrix} -y_0(1-a_0) + a_0 y_1 + a_0 y_2 + \dots + a_0 y_9 \\ -y_1(1-a_1) + a_1 \end{pmatrix}$$

$$\begin{aligned} & -y_0 + a_0 y_1 + a_0 (y_2 + \dots + y_9) \\ & -y_1 + a_1 \end{aligned}$$

$$\frac{dL}{dz_5} = A_5 - y \quad (10 \times 10000)$$

$$A_5 \times (A_5 - y) \quad (25 \times 10000)$$

here is forest.
think I should learn
matrix more.
still don't know

$$\frac{d}{dW_5} L = (A_5 - y) \times A_4^T \quad (10 \times 25) \times \frac{1}{m}$$

(m=60000)

$$\frac{d}{db_5} L = \frac{dZ_5}{db_5} \times \frac{dA_5}{dZ_5} \times \frac{d}{dA_5} L$$

$$= \frac{dZ_5}{db_5} \times (A_5 - y)$$

$$= (A_5 - y)$$

$$\frac{d}{dW_4} L = \frac{dZ_4}{dW_4} \times \frac{dA_4}{dZ_4} \times \frac{d}{dA_4} L$$

$$\frac{d}{dA_4} L = \frac{dZ_5}{dA_4} \times \frac{d}{dZ_5} L$$

$$= \frac{dZ_5}{dA_4} \times (A_5 - y)$$

$$= W_5^T \times (A_5 - y) \Rightarrow (25, 60000)$$

$$\frac{d}{dW_4} L = \frac{dZ_4}{dW_4} \times \frac{dA_4}{dZ_4} \times W_5^T \times (A_5 - y)$$

(25, 60000)

$$A_4 = \text{Relu}(Z_4)$$

$$Z_4 = W_4 \times A_3 + b_4$$

$A_3^T \Rightarrow (60000 \times 25)$ $A_4 = (25, 60000)$

$$\frac{d}{db_4} L = \frac{dZ_4}{db_4} \times \frac{dA_4}{dZ_4} \times \frac{dZ_5}{dA_4} \times dZ_5$$

why

$$A_4 \times (A_5 - y)$$

\Downarrow

$$(A_5 - y) \times A_4^T$$

(I saw a blog about

why $A_4 \rightarrow A_4^T$, but

I don't know if they

can apply commutative property

in this matrix chain rule.

+) why do I have to divide in 60000

$$\frac{d}{dw_3} L = \frac{dz_3}{dw_3} \times \frac{dA_3}{dz_3} \times \frac{d}{dA_3} L$$

$$\frac{d}{dA_3} L = \frac{dz_4}{dA_3} \times \frac{d}{dz_4} L$$

$$z_4 = w_4 \cdot A_3 + b_4$$

$$\begin{array}{c} \downarrow \\ w_4^T \\ (10 \times 1) \end{array} \quad \begin{array}{c} \downarrow \\ (25, 60000) \end{array}$$

$$\frac{d}{dw_3} L = \frac{dz_3}{dw_3} \times \frac{dA_3}{dz_3} \times \frac{d}{dA_3} L$$

$$\begin{array}{c} \downarrow \\ A_2^T \\ (60000 \times 50) \end{array}$$

$$(50 \times 60000)$$

$$(50 \times 60000)$$

$$z_3 = w_3 \times A_2 + b_3$$

Back propagation.

$$① \frac{\partial \mathcal{L}}{\partial W_5}$$

$$\underbrace{\frac{\partial Z_5}{\partial W_5}}_{\downarrow A_4^T (60000 \times 25)} \times \underbrace{\frac{\partial \mathcal{L}}{\partial Z_5}}_{\hookrightarrow A_5 - Y, (10 \times 60000)} = \text{np.dot}(A_5 - Y, A_4^T) = dW_5, (10 \times 25).$$

cache dZ_5 (10×60000)

$$② \frac{\partial \mathcal{L}}{\partial W_4}$$

$$\frac{\partial Z_4}{\partial W_4} \times \frac{\partial \mathcal{L}}{\partial Z_4}$$

$$\frac{\partial \mathcal{L}}{\partial Z_4} = \underbrace{\frac{\partial A_4}{\partial Z_4}}_{\text{ReLU back.} (25, 60000)} \times \underbrace{\frac{\partial Z_5}{\partial A_4}}_{W_5^T (25 \times 10)} \times \boxed{\frac{\partial \mathcal{L}}{\partial Z_5}}_{(10 \times 60000)} = \text{ReLU back.}(\text{np.dot}(W_5^T, dZ_5))$$

$\hookrightarrow dZ_4, (25, 60000)$.
cache.

$$\frac{\partial Z_4}{\partial W_4} = A_3^T, (60000, 50)$$

$$dW_4 = \text{np.dot}(dZ_4, A_3^T)$$

$$③ \frac{\partial \mathcal{L}}{\partial W_3}$$

$$\frac{\partial Z_3}{\partial W_3} \times \frac{\partial \mathcal{L}}{\partial Z_3}$$

$$\frac{\partial \mathcal{L}}{\partial Z_3} = \frac{\partial A_3}{\partial Z_3} \times \frac{\partial Z_4}{\partial A_3} \times dZ_4 = \text{ReLU back.}(\text{np.dot}(W_4^T, dZ_4))$$

$\hookrightarrow dZ_3 (50 \times 60000)$, cache.

$\downarrow W_4^T$

$$dW_3 = \text{np.dot}(dZ_3, A_2^T)$$

④ $\frac{\partial}{\partial W_2} L$

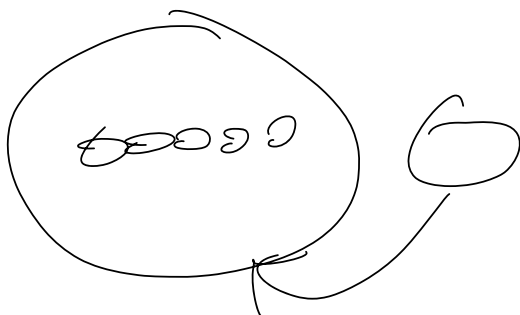
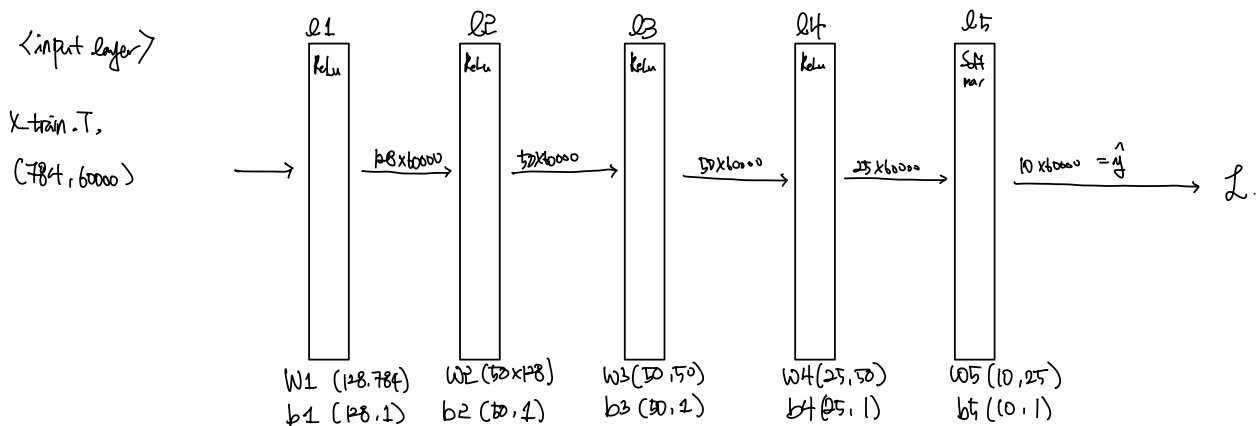
$\frac{\partial Z_2}{\partial W_2} \times \frac{\partial}{\partial Z_2} L$

$\frac{\partial}{\partial Z_2} L = \frac{\partial A_2}{\partial Z_2} \times \frac{\partial Z_1}{\partial A_2} \times dZ_1 = \text{relu_back}(np.dot(W_3^T, dZ_1))$
 \downarrow
 W_3^T
 $\hookrightarrow dZ_2, (10 \times 60000), \text{ cache.}$
 $dW_2 = np.dot(dZ_2, A_1^T)$

⑤ $\frac{\partial}{\partial W_1} L$

$\frac{\partial Z_1}{\partial W_1} \times \frac{\partial}{\partial Z_1} L$

$\frac{\partial}{\partial Z_1} L = \frac{\partial A_1}{\partial Z_1} \times \frac{\partial Z_2}{\partial A_1} \times dZ_2 = \text{relu_back}(np.dot(dZ_2, W_2^T))$
 $\hookrightarrow dZ_1, (128 \times 60000) \text{ cache.}$
 $dW_1 = np.dot(dZ_1, X_{\text{train}})$



200

$$(28 \times 184) \times (184 \times 1) = 28 \times 1$$

$$(50 \times 128) \times (128 \times 1) = 50 \times 1$$

$$50 \times 1 \times 50 \times 1 = 50 \times 1$$

$$25 \times 50 \times 50 \times 1 = 25 \times 1$$

$$10 \times 25 \times 25 \times 1 = 10 \times 1$$

$$(10 \times 1) - (10 \times 1) = (10 \times 1)$$

$$(10 \times 1) \times (1 \times 25) = 10 \times 25$$