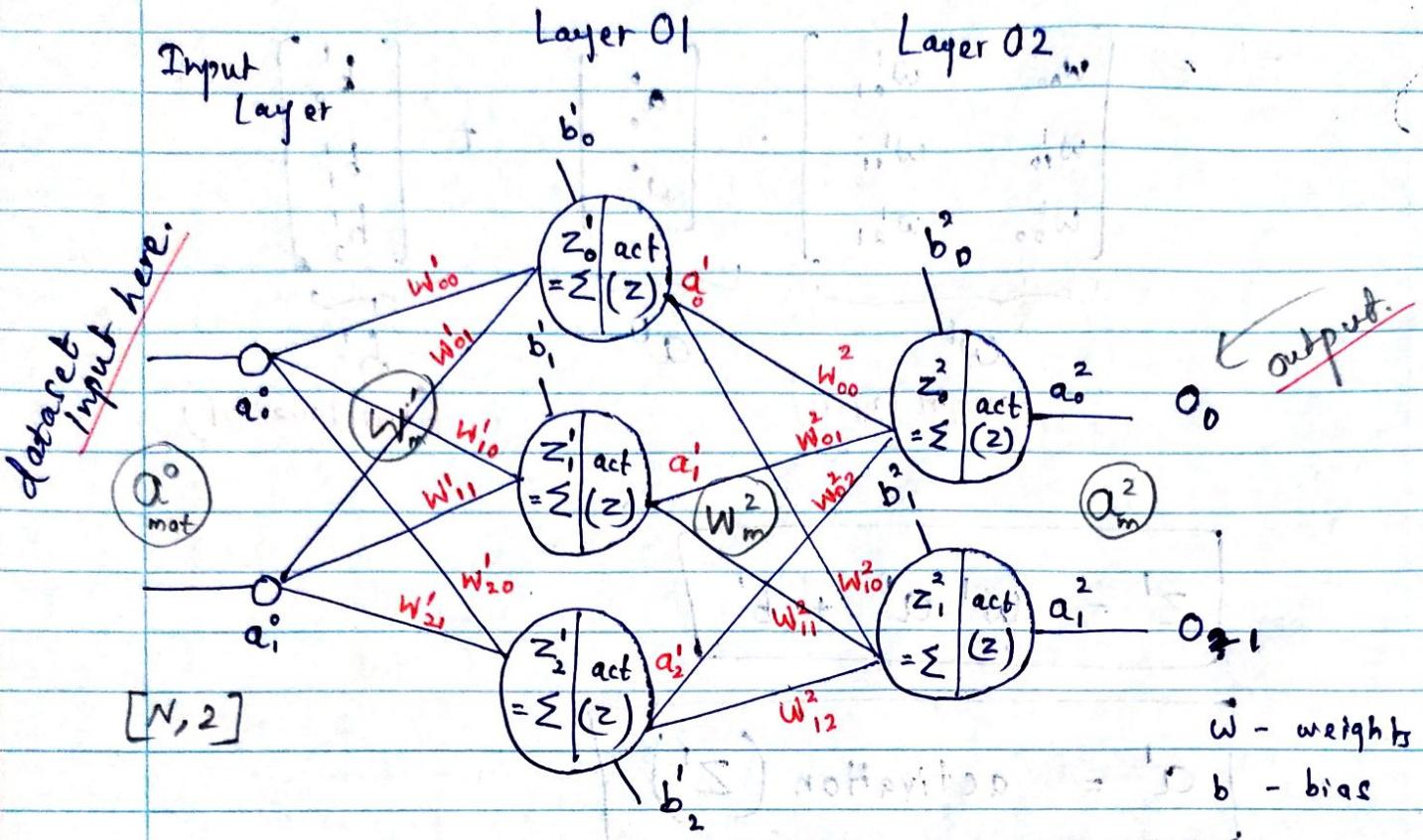


Mathematic of neural network.

No: _____

Date: _____

①



* $act()$ - activation function (ReLU, softmax...) the purpose of the activation function is to introduce the non-linearity to the network. Otherwise, the node will be just a linear representation of previous layer.

Layer 01 - (forward prop.)

$$z'_0 = \left| w'_{00} a^0_0 + w'_{01} a^0_1 \right| + b'_0 \quad a'_0 = act(z'_0)$$

$$z'_1 = \left| w'_{10} a^0_0 + w'_{11} a^0_1 \right| + b'_1 \quad a'_1 = act(z'_1)$$

$$z'_2 = \left| w'_{20} a^0_0 + w'_{21} a^0_1 \right| + b'_2 \quad a'_2 = act(z'_2)$$

$$z' = \begin{bmatrix} w'_{00} & w'_{01} \\ w'_{10} & w'_{11} \\ w'_{20} & w'_{21} \end{bmatrix} \cdot \begin{bmatrix} a^0_0 \\ a^0_1 \end{bmatrix} + \begin{bmatrix} b'_0 \\ b'_1 \\ b'_2 \end{bmatrix}$$

↓
det product

$(w^1)_{(mat)}$ (3×2) (3×1) (2×1) $(3 \times 1)_{(mat)}$

$$z' = w^1 a^0 + b'$$

$$a' = \text{activation}(z')$$

Same thing happens in the each and every layer in the neural network, so that we can represent the process of the forward propagation in the 2 layer ANN as follows,

Forward Prop.

$$z' = w^1 a^0 + b^1$$

$\xrightarrow{\quad}$

$$a^1 = \text{act}^1(z')$$

$\xrightarrow{\quad}$

$$z^2 = w^2 a^1 + b^2$$

$\xrightarrow{\quad}$

$$a^2 = \text{act}^2(z^2)$$

$$a^0 \quad w^1 \quad b^1 \quad a^1 \quad w^2 \quad b^2 \quad a^2 \leftarrow \text{matrices.}$$

$(2 \times 1) \quad (3 \times 2) \quad (3 \times 1), (3 \times 1) \quad A(2 \times 3) \quad (2 \times 1) \quad (2 \times 1)$

Backward prop.

- In the backward pass, the gradient or the loss with respect to each weight in the networks in the networks is computed
- This is done by applying the chain rule of the calculus to calculate the partial derivatives of the loss with respect to the weights and biases of the networks.
- The gradients indicate how much the loss would increase or decrease if the corresponding weights and biases were adjusted.
- Since the chain rule is important, let's practice it.
- Example for chain rule in partial derivative. (Ex01)

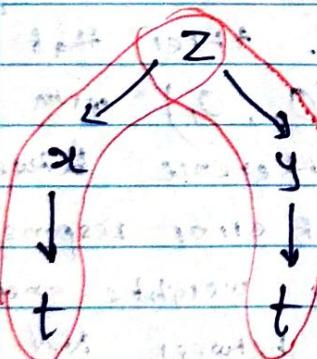
$$z = x^3y + 3x^2 + y^2 \quad | \quad x = 4 + t^2 \quad | \quad y = st^3$$

Find the

$$\frac{dz}{dt}$$

d - derivative

∂ - partial derivative.

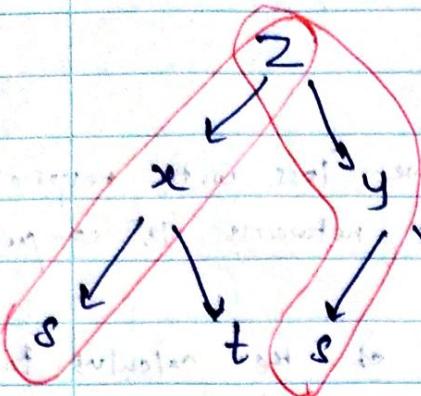


$$\frac{dz}{dt} = \left(\frac{\partial z}{\partial x} \cdot \frac{dx}{dt} \right) + \left(\frac{\partial z}{\partial y} \cdot \frac{dy}{dt} \right)$$

- Example for chain rule in partial derivative (Ex02).

$$z = 4x^3y^5 \quad | \quad x = 3s^2t \quad | \quad y = 2s^3t^2$$

Find, $\frac{\partial z}{\partial s}$

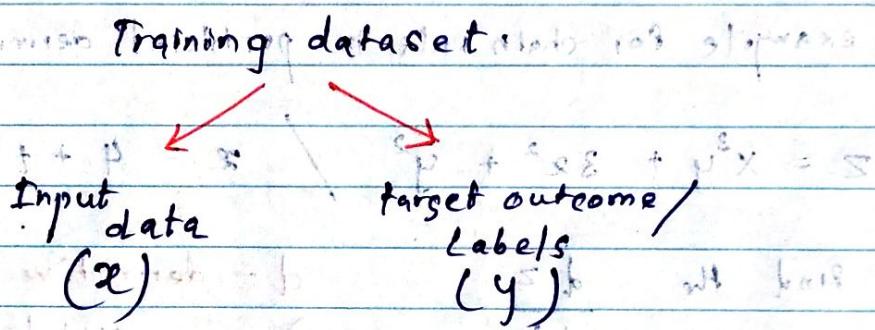


$$\frac{\partial z}{\partial s} = \left(\frac{\partial z}{\partial x} \cdot \frac{\partial x}{\partial s} \right) \left(\frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial s} \right)$$

The chain rule of differentiation is being used in the backward pass of the NN training process.

Again to the backward pass process...

In a training dataset, there are two parts,



In the forward pass, we let the NN to predict the output by passing the input data through the NN. After that we compares the output of the NN (O_0, O_1) with the target outcome/labels and find the difference between them. After that we calculate the value of error responsible for each weigh and bias and update the weights and bias in order to minimize the difference between NN output and target outcomes.

using "Cost Function" we find the overall error made by the NN. Cost Function give a quantitative measure of the error made by the Neural network.

Cost function:

cost function give a quantitative measures of all the error made by the neural networks. There are different types of cost functions used in different ANN problems,

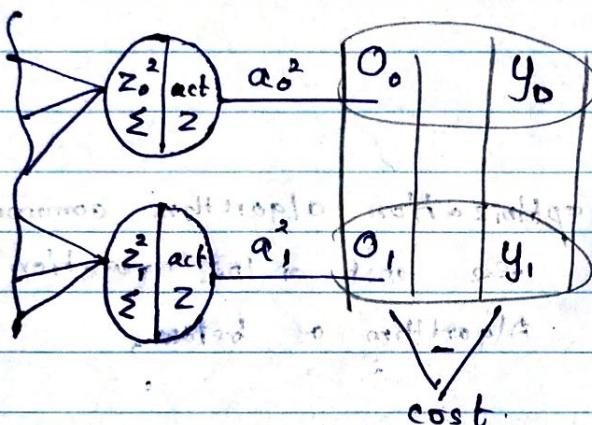
- MSE - mean square error. - used in this example.

at 3: Binary crossentropy

with hinge loss.

$$\text{MSE} = \frac{1}{N} \sum_{i=0}^{N-1} (O_i - y_i)^2$$

Layer 02.



O - output of ANN

y - target output

target / labels

→ For our example NN, the cost function $C(x)$ will be,

$$C(x) = (O_0 - y_0)^2 + (O_1 - y_1)^2 \quad (\text{since } y \text{ is not changing it, is a constant.})$$

since the "y" is constant, the $C(x)$ is a function of O.

$$C(O_i)$$

and also, the O_i is dependent on the values of parameters (weight and bias). Cost is function of all weight and bias.

$$C(w_{00}, w_{01}, w_{10}, \dots, w_{12}, b_1) \rightarrow C(p_i)$$

parameters of NN (p_i) ← any parameter.

[we use p_i to represent all parameters in ANN such as w_{00} , w_{01} , b_1^2 , ...]

Now we know that, cost is the function of NN's parameter such as all weights and biases.

$$C(W_{00}, W_{01} \dots W_{12}, b_i) \xrightarrow{Pi} C(P_i)$$

(→) the main purpose of ANN training process is to minimize the cost or loss function during the training the model. To minimize the cost, we use the Gradient Descent algorithm.

Gradient Descent.

Gradient Descent is a optimization algorithm commonly used in deep learning to minimize the cost or loss function during the training the model. Algorithm as below,

$$\text{updated parameter } P_i = \frac{\text{new}}{\text{old}} - \eta \frac{\partial C}{\partial P_i}$$

learning rate.

η - learning rate.

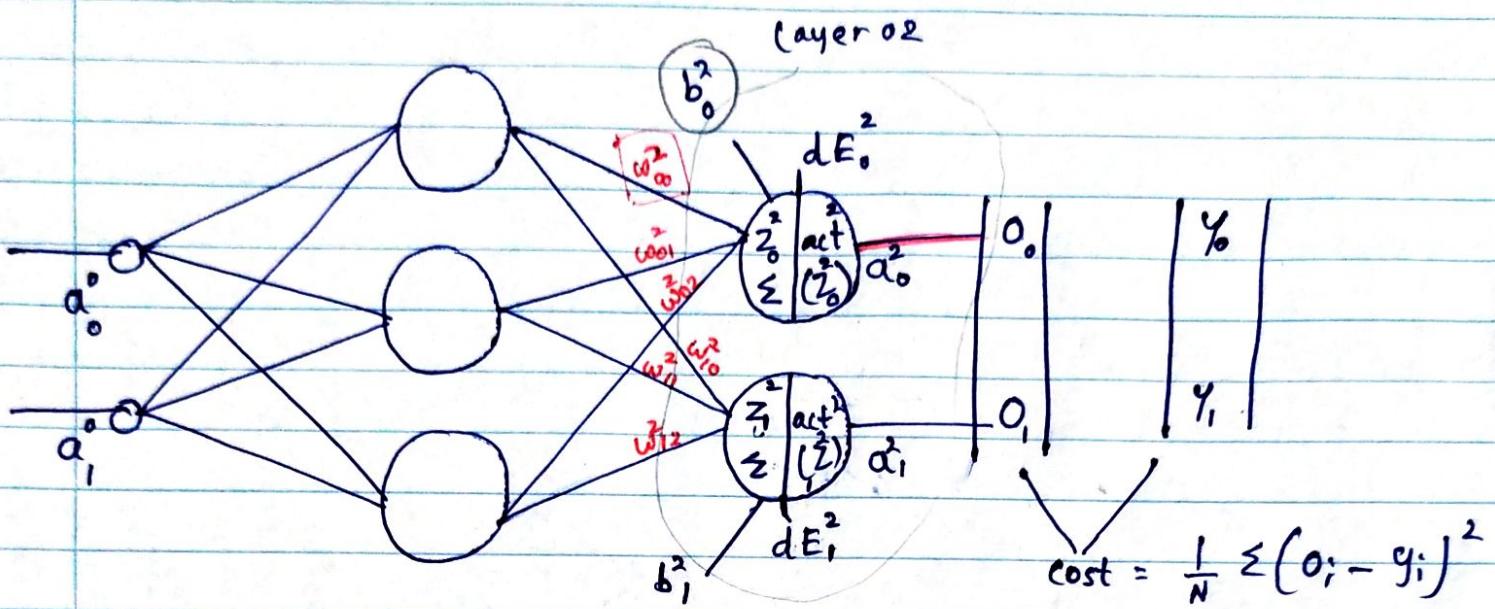
$\frac{\partial C}{\partial P_i}$ } = the slope of cost function
with respect particular parameter. (P_i)
(partial derivative)

P_i is each and every parameter of the ANN.

for example for W_{00} of first layer of the ANN.

$$W_{00}^{\text{(new)}} = W_{00}^{\text{(old)}} - \eta \frac{\partial C}{\partial W_{00}}$$

→ Let's assume the layer 02 of the neural network.



→ Let's assume the b_0^2 bias of the layer 02, first neuron.

$$b_0^2 \text{ (new)} = b_0^2 \text{ (old)} - \eta \frac{\partial C}{\partial b_0^2}$$

← we cannot find this straightforward.
we have to use the partial derivative chain rule to find this.

$$\frac{\partial C}{\partial b_0^2} = \left| \frac{\partial C}{\partial o_0} \cdot \frac{\partial o_0}{\partial z_0^2} \cdot \frac{\partial z_0^2}{\partial b_0^2} \right| \quad (\text{A'' equation.})$$

term $\left(\frac{\partial C}{\partial o_0} \cdot \frac{\partial o_0}{\partial z_0^2} \right)$ is common to the

all the parameters of the particular neuron, so that it can be represented as " dE^2 " error factor.

discuss more in the parts that coming up.

Cost → Out → a → z → b^2
chain rule

each neuron of the ANN has separate " dE^2 " value.

for example,

this is common.

$$\frac{\partial C}{\partial b_0^2} = dE_0^2 \cdot \frac{\partial z_0^2}{\partial b_0^2}$$

$$\frac{\partial C}{\partial w_{00}^2} = dE_0^2 \cdot \frac{\partial z_0^2}{\partial w_{00}^2}$$

$$\frac{\partial C}{\partial b_0^2} = \frac{dE^2}{\partial b_0^2} \cdot \frac{\partial z_0^2}{\partial b_0^2}$$

(eq01) $dE_0^2 = \frac{\partial C}{\partial o_0} \cdot \frac{\partial o_0^2}{\partial z_0^2}$ ← This is common for particular neuron.

$$dE_0^2 = (o_0 - y_0)^2$$

$$\frac{\partial C}{\partial o_0} = \frac{\partial \sum (o_i - y_i)^2}{\partial o_0} \rightarrow \frac{\partial [(o_0 - y_0)^2 + (o_1 - y_1)^2]}{\partial o_0} \times \frac{1}{2}$$

(eq02) $\frac{\partial C}{\partial o_0} = (o_0 - y_0)$ ← (2)

(eq03) $\frac{\partial o_0}{\partial z_0} = \text{act}'_2(z_0^2)$ ← (3) derivative of the activation function used in the particular layer. act_2 — activation function of layer 2.]

using eq01, eq02 and eq03,

$$dE_0^2 = (o_0 - y_0) \cdot \text{act}'_2(z_0^2)$$

derivative of the activation function.

above given the error vector for the layer 02, 0th node/neuron.

For the Layer 02, 1st node,

$$dE_1^2 = (o_1 - y_1) \cdot \text{act}'_2(z_1^2)$$

And also, error vector for the layer 2 can be represent using a matrix as below,

$$dE^2 = \begin{bmatrix} dE_0^2 \\ dE_1^2 \end{bmatrix}$$

Since we know the (δE_0^2) and (δE_1^2) now, we can find the loss for $\frac{\partial C}{\partial b_o^2}$ and all the parameters (w and b) of layer 02.

As an example let's find the b_o^2 and w_{00}^2 .

$$\rightarrow \frac{\partial C}{\partial b_o^2} / \frac{\partial C}{\partial b_o^2} = \delta E_0^2 \cdot \frac{\partial z_o^2}{\partial b_o^2} \leftarrow (\text{using equation given in the 7 page})$$

$$\delta E_0^2 = (O_o - Y_o) \cdot \text{act}'_2(z_o^2)$$

$$\frac{\partial z_o^2}{\partial b_o^2} = \frac{\partial (w_{00}^2 a_o^1 + w_{01}^2 a_1^1 + w_{02}^2 a_2^1 + b_o^2)}{\partial b_o^2} \leftarrow (\text{partial derivative.})$$

$$\frac{\partial z_o^2}{\partial b_o^2} = 1$$

so that,

$$\boxed{\frac{\partial C}{\partial b_o^2} = (O_o - Y_o) \cdot \text{act}'_2(z_o^2) / \delta E_0^2}$$

$$\rightarrow \frac{\partial C}{\partial w_{00}^2} / \frac{\partial C}{\partial w_{00}^2} = \delta E_0^2 \cdot \frac{\partial z_o^2}{\partial w_{00}^2}$$

$$\frac{\partial z_o^2}{\partial w_{00}^2} = \frac{\partial (w_{00}^2 a_o^1 + w_{01}^2 a_1^1 + w_{02}^2 a_2^1 + b_o^2)}{\partial w_{00}^2}$$

$$\frac{\partial z_o^2}{\partial w_{00}^2} = a_o^1$$

so that,

$$\text{so that, } \frac{\partial C}{\partial w_{00}^2} = \frac{\partial C}{\partial w_{00}^2} = \frac{\partial C}{\partial w_{00}^2}$$

$$\boxed{\frac{\partial C}{\partial w_{00}^2} = (O_o - Y_o) \cdot \text{act}'_2(z_o^2) \cdot a_o^1 / \delta E_0^2}$$

In a similar way, we can calculate the each and every parameter errors in the layer 2. (1st layer of the backward pass / Last layer).

→ below I have given the error calculations for each and every weight parameter.

$$\frac{\partial C}{\partial w_{00}^2} = \delta E_0^2 \cdot a_0'$$

$$\frac{\partial C}{\partial w_{01}^2} = \delta E_0^2 \cdot a_1'$$

$$\frac{\partial C}{\partial w_{02}^2} = \delta E_0^2 \cdot a_2'$$

$$\frac{\partial C}{\partial w_{10}^2} = \delta E_1^2 \cdot a_0'$$

$$\frac{\partial C}{\partial w_{11}^2} = \delta E_1^2 \cdot a_1'$$

$$\frac{\partial C}{\partial w_{12}^2} = \delta E_1^2 \cdot a_2'$$

→ As you can remember, our ultimate goal is to, below, (for weights)

$$P_i^{\text{new}} = P_i^{\text{old}} - \eta \frac{\partial C}{\partial P_i}$$

for weights

$$w_i^{\text{new}} = w_i^{\text{old}} - \eta \frac{\partial C}{\partial w_i}$$

As you can remember, we the weights on the last layer with a 2×3 matrix. Likewise, we have represent build the 2×3 matrix for $(\eta \frac{\partial C}{\partial w})$ this as well. Then only we can subtract.

$$\delta w^2 = \begin{bmatrix} \delta E_0^2 \cdot a_0' & \delta E_0^2 \cdot a_1' & \delta E_0^2 \cdot a_2' \\ \delta E_1^2 \cdot a_0' & \delta E_1^2 \cdot a_1' & \delta E_1^2 \cdot a_2' \end{bmatrix}$$

$$\delta w^2 = \begin{bmatrix} \delta E_0^2 \\ \delta E_1^2 \end{bmatrix} \cdot \begin{bmatrix} a_0' & a_1' & a_2' \end{bmatrix}$$

$$\boxed{\delta w^2 = \delta E^2 \cdot (a')^T} \quad \xleftarrow{2 \times 3 \quad 3 \times 1 \quad 1 \times 2} \text{for weights.}$$

layer 02

$$\boxed{\delta b^2 = \delta E^2} \quad \xleftarrow{\text{for bias of layer 02}}$$

so that we can run the Gradient Decent algorithm to update the weights and bias for layer 02 (last layer).

δw^2

$$\boxed{w^2_{\text{new}} = w^2_{\text{old}} - \eta \delta w^2_{\text{old}}}$$

$$\boxed{b^2_{\text{new}} = b^2_{\text{old}} - \eta \delta b^2_{\text{old}}}$$

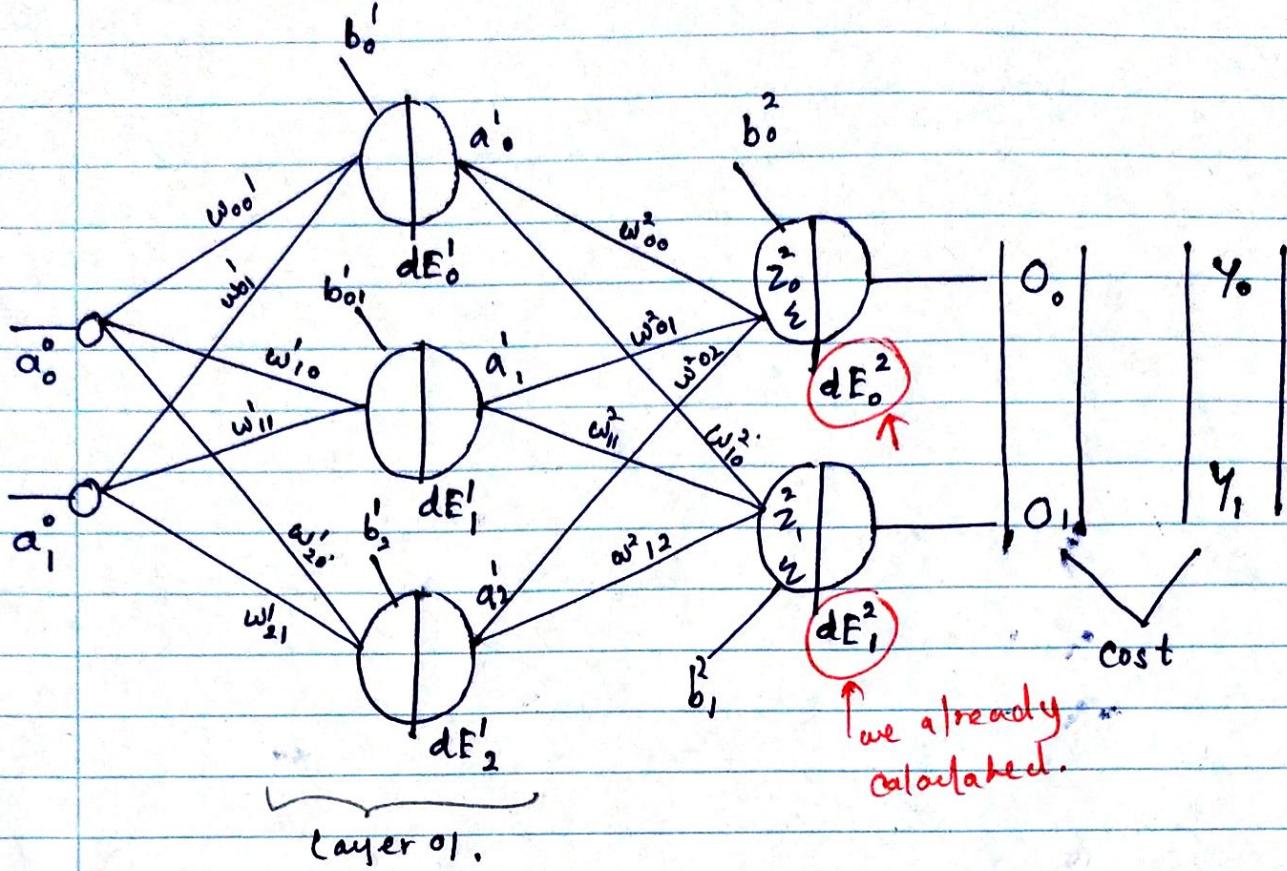
now we are done with the layer 02. Now we have to do the similar thing to the layer 01 as well.

$$\boxed{w_1 = w_1 - \eta \delta w_1}$$

values from left to right are 36, 36, 36, 36, 36, 36, 36, 36
label categories from left to right are 8x8, 8x8, 8x8, 8x8, 8x8, 8x8, 8x8, 8x8
so price east. Max is 36. (36) is 8x8 value 8x8 and
is 8x8 max

$$\boxed{w_1 = w_1 - \eta \delta w_1}$$

$$\boxed{w_1 = w_1 - \eta \delta w_1}$$



→ we know that we are going through the backward pass of the neural network, so that we already have known the dE^2_0 and dE^2_1 .

→ let's assume the " b_0^1 " bias of the (layer 0), first neuron.

$$\frac{b'_0}{(\text{new})} = \frac{b'_0}{(\text{old})} - \eta \frac{\partial C}{\partial b'_0} \quad \leftarrow \text{like we did before, we have to find this using the chain rule or partial derivative.}$$

$$b'_0 \rightarrow z'_0 \rightarrow a'_0 \rightarrow z'_1 \rightarrow a'_1 \rightarrow \text{cost.}$$

chain rule.

$$\frac{\partial C}{\partial b'_0} = \left[\underbrace{\frac{\partial C}{\partial a'_0} \cdot \frac{\partial a'_0}{\partial z'_0}}_{dE^2_0} \frac{\partial z'_0}{\partial a'_0} + \underbrace{\frac{\partial C}{\partial a'_1} \cdot \frac{\partial a'_1}{\partial z'_1}}_{dE^2_1} \frac{\partial z'_1}{\partial a'_0} \right] \cdot \frac{\partial a'_0}{\partial z'_0} \cdot \frac{\partial z'_1}{\partial b'_0}$$

are derived these two in the Layer 02 process.

$$\frac{\partial C}{\partial a'_0} = \left[dE_0^2 \cdot \frac{\partial z_0^2}{\partial a'_0} + dE_1^2 \cdot \frac{\partial z_1^2}{\partial a'_0} \right] \frac{\partial a'_0}{\partial z'_0}, \frac{\partial z'_0}{\partial b'_0}$$

$\frac{\partial C}{\partial a'_0}$

First we should find the $\frac{\partial C}{\partial a'_0}$

$$\frac{\partial C}{\partial a'_0} = \left[dE_0^2 \cdot \frac{\partial z_0^2}{\partial a'_0} + dE_1^2 \cdot \frac{\partial z_1^2}{\partial a'_0} \right]$$

dE'

error vector.

$$\frac{\partial z_0^2}{\partial a'_0} = \frac{\partial (w_{00}^2 a'_0 + w_{01}^2 a'_1 + w_{10}^2 a'_2 + b_0^2)}{\partial a'_0}$$

$$\frac{\partial z_0^2}{\partial a'_0} = \frac{\partial a'_0}{\partial a'_0}$$

$$\frac{\partial z_0^2}{\partial a'_0} = w_{00}^2$$

Likewise,

$$\frac{\partial z_1^2}{\partial a'_0} = w_{10}^2$$

$$\therefore \frac{\partial C}{\partial a'_0} = dE_0^2 \cdot w_{00}^2 + dE_1^2 \cdot w_{10}^2$$

Likewise we can find the $\left(\frac{\partial C}{\partial a'_1} \right)$, $\left(\frac{\partial C}{\partial a'_2} \right)$ for each node layer a'_1 .

$$\frac{\partial C}{\partial a'_0} = \left| \begin{array}{l} dE_0^2 \cdot w_{00}^2 + dE_1^2 \cdot w_{10}^2 \\ \vdots \end{array} \right|$$

← matrix.

$$\frac{\partial C}{\partial a'_1} = \left| \begin{array}{l} dE_0^2 \cdot w_{01}^2 + dE_1^2 \cdot w_{11}^2 \\ \vdots \end{array} \right|$$

$$\frac{\partial C}{\partial a'_2} = \left| \begin{array}{l} dE_0^2 \cdot w_{02}^2 + dE_1^2 \cdot w_{12}^2 \\ \vdots \end{array} \right|$$

$$\begin{bmatrix} \omega_{00}^2 \delta E_0^2 + (\omega_{10}^2 \delta E_1)^2 \\ \omega_{01}^2 \delta E_0^2 + \omega_{11}^2 \delta E_1^2 \\ \omega_{02}^2 \delta E_0^2 + \omega_{12}^2 \delta E_1^2 \end{bmatrix} = \begin{bmatrix} \omega_{00}^2 & \omega_{10}^2 \\ \omega_{01}^2 & \omega_{11}^2 \\ \omega_{02}^2 & \omega_{12}^2 \end{bmatrix} \cdot \begin{bmatrix} \delta E_0^2 \\ \delta E_1^2 \end{bmatrix}$$

$[\omega^2]^T$

$$\boxed{\delta a' = [\omega^2]^T \delta E^2}$$

As the equation we build in the previous page,

$$\frac{\partial c}{\partial b'_0} = \left[\delta E_0^2 \cdot \frac{\partial z_0^2}{\partial a'_0} + \delta E_1^2 \cdot \frac{\partial z_1^2}{\partial a'_0} \right] \cdot \frac{\partial a'_0}{\partial z'_0} \cdot \frac{\partial z'_0}{\partial b'_0},$$

$$\frac{\partial c}{\partial b'_0} = \left[\frac{\partial c}{\partial a'_0} \right] \cdot \frac{\partial a'_0}{\partial z'_0} \cdot \frac{\partial z'_0}{\partial b'_0}$$

derivation of function
activation function
used in layer 0

$$\boxed{\frac{\partial c}{\partial b'_0} \neq \delta E_0^2 = \delta a'_0 \cdot \text{act}'(z'_0)} \quad \leftarrow \text{for bias of layer 0}$$

As such, we can represent the error for the bias using below equation,

$$\boxed{\frac{\partial b'}{\partial b} = \delta E^I}$$

error matrix
for layer 0

$$\delta E^I = \delta a' \cdot \text{act}' \cdot \begin{bmatrix} z'_0 \\ z'_1 \\ z'_2 \end{bmatrix}$$

$$\boxed{\delta E' = \delta a' \cdot \text{act}_i(z')}$$

for bias in layer 0

$$\frac{\delta c}{\delta b'_0} = \delta E'_0 \cdot \underbrace{\frac{\delta z'_0}{\delta b'_0}}_1 = \delta E'_0$$

matrix.

$$\frac{\delta c}{\delta i} = \delta E'_i$$

$$\boxed{\delta b' = \delta E'}$$

$$\frac{\delta c}{\delta z'_2} = \delta E'_2$$

for weight in layer 0.2

$$\delta w' = \begin{bmatrix} \frac{\delta c}{\delta w'_{00}} & \frac{\delta c}{\delta w'_{01}} \\ \frac{\delta c}{\delta w'_{10}} & \frac{\delta c}{\delta w'_{11}} \\ \frac{\delta c}{\delta w'_{20}} & \frac{\delta c}{\delta w'_{21}} \end{bmatrix} = \begin{bmatrix} \delta E'_0 \cdot \frac{\delta z'_0}{\delta w'_{00}} & \delta E'_0 \cdot \frac{\delta z'_0}{\delta w'_{01}} \\ \delta E'_1 \cdot \frac{\delta z'_1}{\delta w'_{10}} & \delta E'_1 \cdot \frac{\delta z'_1}{\delta w'_{11}} \\ \delta E'_2 \cdot \frac{\delta z'_2}{\delta w'_{20}} & \delta E'_2 \cdot \frac{\delta z'_2}{\delta w'_{21}} \end{bmatrix}$$

3×2

$$\frac{\delta z'_0}{\delta w'_{00}} = \frac{\partial (w'_{00} a'_0 + w'_{01} \cdot a'_1 + b'_0)}{\partial w'_{00}} = a'_0$$

$$\frac{\delta z'_0}{\delta w'_{00}} = a'_0$$

$$\frac{\partial z'_0}{\partial w'_{00}} = a_0^\circ \quad \frac{\partial z'_0}{\partial w'_{01}} = a_1^\circ \quad \text{and then use backprop}$$

$$\frac{\partial z'_1}{\partial w'_{10}} = a_0^\circ \quad \frac{\partial z'_1}{\partial w'_{11}} = a_1^\circ$$

$$\frac{\partial z'_2}{\partial w'_{20}} = a_0^\circ \quad \frac{\partial z'_2}{\partial w'_{21}} = a_1^\circ$$

$$\therefore \delta w' = \begin{bmatrix} \delta E'_0 a_0^\circ & \delta E'_0 a_1^\circ \\ \delta E'_1 a_0^\circ & \delta E'_1 a_1^\circ \\ \delta E'_2 a_0^\circ & \delta E'_2 a_1^\circ \end{bmatrix}$$

$$= \begin{bmatrix} \delta E'_0 \\ \delta E'_1 \\ \delta E'_2 \end{bmatrix} \cdot \begin{bmatrix} a_0^\circ & a_1^\circ \end{bmatrix}$$

matrix

$$\boxed{\delta w' = \delta E' \cdot [a^\circ]^T} \quad \leftarrow \text{for weight at layer 01}$$

$$\boxed{\delta b' = \delta E'} \quad \leftarrow \text{for bias of layer 01 (previously found)}$$

so that we can run the gradient Descent algorithm to update the weights and bias for layer 01 (First layer).

$$\boxed{w'_{\text{new}} = \frac{w'_{\text{old}}}{\eta} - \eta \delta w'}$$

$$\boxed{b'_{\text{new}} = \frac{b'_{\text{old}}}{\eta} - \eta \delta b'}$$

now we are done with the layer 01, now we have to do the similar thing for the layer 02 as well.

Now we will summarize the entire process.

To update w1
To repeat

$$[0.0] \cdot 1.76 = 1.76$$

To repeat for w2
(Forward pass)

$$1.76 \cdot 1.86 = 1.86$$

No. _____

Summarize the entire training process.

Date: _____

→ Forward Pass

$$z' = w^1 a^0 + b^1 \quad \left. \begin{array}{l} \text{calculated.} \\ \text{page 02} \end{array} \right\}$$

$$a' = \text{act}^2(z') \quad \left. \begin{array}{l} \text{calculated.} \\ \text{page 02} \end{array} \right\}$$

$$z^2 = w^2 a' + b^2 \quad \left. \begin{array}{l} \text{calculated.} \\ \text{page 02} \end{array} \right\}$$

$$a^2 = \text{act}^2(z^2) \quad \left. \begin{array}{l} \text{calculated.} \\ \text{page 02} \end{array} \right\}$$

→ Backward Propagation

- 1) compute gradients. (sequence is given below)

$$\delta E^2 = (o - y) \text{act}'_2(z^2) \quad \text{page 08}$$

$$\rightarrow \delta b^2 = \delta E^2 \quad \text{page 11}$$

$$\rightarrow \delta w^2 = \delta E^2 \cdot [a']^T \quad \text{page 11}$$

$$\delta a' = [w^2]^T \cdot \delta E^2 \quad \text{page 14}$$

$$\delta E'_{out} = \delta a' \cdot \text{act}'(z') \quad \text{page 15}$$

$$\rightarrow \delta b' = \delta E' \quad \text{page 15}$$

$$\rightarrow \delta w' = \delta E' \cdot [a^0]^T \quad \text{page 16}$$

2) update weights and bias.

$$w' = w' - \eta \times \delta w'$$

$$b' = b' - \eta \times \delta b'$$

$$w^2 = w^2 - \eta \times \delta w^2$$

$$b^2 = b^2 - \eta \times \delta b^2$$

η is learning rate that should be defined by the programmer. Usually it is 0.01 or 0.1 (constant).

Flow diagram

