$h_{11}$

$a_{11}$

$h_{12}$

$a_{12}$

$h_{13}$

$a_{13}$

$W_{111}$

$W_{121}$

$W_{131}$

$W_{112}$

$W_{122}$

$W_{132}$

21

22

23

k = Layer

i = current layer — neuron no.
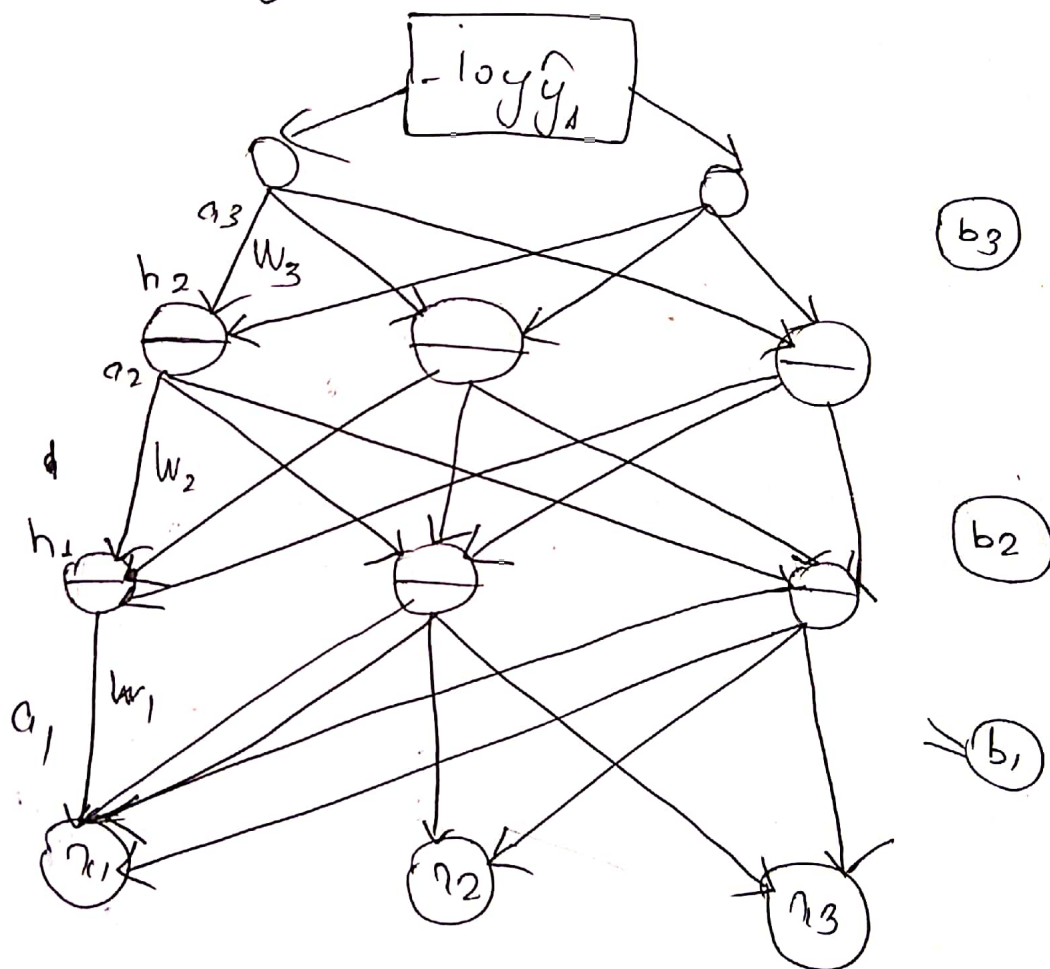
j = previous (Input layer — neuron no.

$f(x) = x2$

$\dfrac{df(x)}{dx}$   $f(x) = x2 + y2$

$\dfrac{\partial}{\partial x}$ , $\dfrac{\partial}{\partial y}$ } partial derivative

$\nabla_\theta = \left[ \dfrac{\partial}{\partial x} \quad \dfrac{\partial}{\partial y} \right]$

$-\log \hat{y}_A$

$a_3$  $W_3$  $h_2$  $b_3$

$a_2$

$W_2$

$h_1$  $b_2$

$w_1$

$a_1$  $b_1$

$x_1$  $x_2$  $x_3$

$$\dfrac{\partial L(\theta)}{\partial W_{111}} = \dfrac{\partial L(\theta)}{\partial \hat{y}} \dfrac{\partial \hat{y}}{\partial a_3} \quad \dfrac{\partial a_3}{\partial h_2} \dfrac{\partial h_2}{\partial a_2} \dfrac{\partial a_2}{\partial h_1} \dfrac{\partial h_1}{\partial a_1} \dfrac{\partial a_1}{\partial W_{111}}$$

$\underbrace{\qquad\qquad}$   $\underbrace{\qquad\qquad\qquad\qquad}$   $\underbrace{\qquad}$

der. w.r. to output layer

der. w.r.t to hidden layer

wts

First

we have, $L(\theta) = \sum_i \log \hat{y}_i$

Then,

$$\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{\partial(-\log \hat{y}_i)}{\partial a_{Li}}$$

$$= \frac{\partial(-\log \hat{y}_i)}{\partial \hat{y}_i} \times \frac{\partial \hat{y}_i}{\partial a_{Li}} \qquad - (1)$$

where, $L =$ layer no.

$\qquad i =$ neuron no. ($1$ to $k$)

$\qquad l =$ index of correct output.

The first part of derivative in $ea(1)$ is straight forward $\frac{\partial(-\log \hat{y}_l)}{\partial \hat{y}_l} = -\frac{1}{\hat{y}_l}$

So,

$$\frac{\partial L(\theta)}{\partial a_{Li}} = \frac{\partial(-\log \hat{y}_l)}{\partial \hat{y}_l} \times \frac{\partial \hat{y}_l}{\partial a_{Li}}$$

$$= -\frac{1}{\hat{y}_l} \frac{\partial}{\partial a_{Li}} softmax(a_L)_l$$

$$= -\frac{1}{\hat{y}_l} \frac{\partial}{\partial a_{Li}} \frac{\exp(a_L)_l}{\sum_{i} \exp(a_L)_l}$$

$$\frac{\partial\left(\frac{g(x)}{h(x)}\right)}{\partial x} = \frac{\partial g(x)}{\partial x}\frac{1}{h(x)} - \frac{g(x)}{h(x)^2}\frac{\partial h(x)}{\partial x}$$

Using division rule,

$$g(x) = exp(o_L)_l$$

$$= \frac{-1}{\hat{y}_L}\left(\frac{\frac{\partial}{\partial o_{Li}}exp(o_L)_l}{\sum_{i'}exp(o_L)_{i'}} - \frac{exp(o_L)_l\left(\frac{\partial}{\partial o_{Li}}\sum_{i'}exp(o_L)_{i'}\right)}{\sum_{i'}\left(exp(o_L)_{i'}\right)^2}\right)$$

conside,

$$\left(\frac{\partial}{\partial o_{Li}}exp(o_L)_l\right)$$ This value is 0 for all i: 0 token except for l

so, an indicator can be used $1(l=i)$

$exp(o_L)_i$ to denote all the values except $i=l$, resolve to 0

Now, it is simply derivative of exponent.

$$\frac{\partial L(\theta)}{\partial o_{Li}} = \frac{-1}{\hat{y}_l}\left(\frac{1(l-i)exp(o_L)_l}{\sum exp(o_L)_{i'}} - \frac{exp(o_L)_l}{\sum_{i'}exp(o_L)_{i'}}\frac{exp(o_L)_i}{\sum_{i'}exp(o_L)_{i'}}\right)$$

writing interms of softmax

$$\frac{\partial L(\theta)}{\partial o_{Li}} = \frac{-1}{\hat{y}_l}\left(1_{(l=i)}softmax(o_L)_i - softmax(o_L)_i \, softmax(o_L)_i\right)$$

softmax is $\hat{y}$

$$\frac{\partial L(\theta)}{\partial o_{Li}} = \frac{-1}{\hat{y}_i}\left(1_{(l=i)}\hat{y}_l - \hat{y}_l\,\hat{y}_i\right)$$

after concellation,

$$\frac{\partial L(\theta)}{\partial o_{Li}} = -\left(1_{(l=i)} - \hat{y}_i\right)$$

So for we have,

$$\frac{\partial L(\theta)}{\partial c_{Li}} = -\left( \mathbb{1}_{(l=i)} - \hat{y}_i \right)$$

Now, gradient w.r. to vector $a_L$.

$$a_L = [a_{L1}, a_{L2} \ldots a_{LK}]$$

by indictor variable, it resolves to 0, for all values of $i$ except $i = l$.

Let $k = 4$, $l = 2$

$$\frac{\partial L(\theta)}{\partial a_{L1}} = -\left( 0 - \hat{y}_i \right)$$

$$\frac{\partial L(\theta)}{\partial a_{L2}} = -\left( 1 - \hat{y}_i \right)$$

$$\frac{\partial L(\theta)}{\partial a_{L3}} = -\left( 0 - \hat{y}_i \right)$$

$$\frac{\partial L(\theta)}{\partial L4} = -\left( 0 - \hat{y}_i \right)$$

The gradient w.r. to $a_L$ is $\nabla_{a_L} =$

$$\nabla_{a_L} = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{L1}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{LK}} \end{bmatrix}$$

$$\nabla_{a_L} = \begin{bmatrix} \mathbb{1}_{(l=i)} - \hat{y}_i \\ \vdots \\ \mathbb{1}_{(l=k)} - \hat{y}_i \end{bmatrix}$$

It is simply the diff. bet" $[0 \, 0 \, \text{---} \, 0 \ldots 0k]$ and $\hat{y}$

In reality it is the diff bet" true dis$(y)$ and pred dis $\hat{y}$

$$\nabla_{a_L} L(\theta) = -(y - \hat{y}_i)$$

$$\frac{\partial L(\theta)}{\partial W_{111}} = \frac{\partial L(\theta)}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial a_3} \qquad \frac{\partial a_3}{\partial h_2} \times \frac{\partial h_2}{\partial a_2} \qquad \frac{\partial a_2}{\partial h_1} \times \frac{\partial h_1}{\partial a_1} \qquad \frac{\partial a_3}{\partial W_{111}}$$

$$\underbrace{\qquad\qquad}$$
$$-(y - \hat{y_i})$$

⑤

So, Now,

→ deriv. w.r. to hidden layers;

$$\frac{\partial L(\theta)}{\partial h_{ij}} = \sum_{m=1}^{k} \frac{\partial L(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}}$$

$$\underbrace{\qquad\qquad}$$

already
computed
earlier.
$$\sum_{m=1}^{k} \frac{\partial L(\theta)}{\partial a_{i+1}} W_{i+1,m,j}$$

How $W_{i+1,m,j}$ is coming?

Supp. $\dfrac{\partial a_{31}}{\partial h_{22}}$

$$\frac{\partial \left( W_{311} h_{21} + W_{312} h_{22} + W_{313} h_{23} \right)}{\partial h_{22}}$$

$$= W_{312}$$

$$\sum_{m=1}^{k} \frac{\partial L(\theta)}{\partial a_{i+1}} W_{i+1,m,j}$$

$$\nabla_{a_{i+1}} L(\theta) = \begin{bmatrix} \dfrac{\partial L(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \dfrac{\partial L(\theta)}{\partial a_{i+1,k}} \end{bmatrix} \qquad W_{i+1 \cdot j} \begin{bmatrix} W_{i+1,1,j} \\ | \\ \vdots \\ | \\ W_{i+1,k,j} \end{bmatrix}$$

↓
grad. of loss func. w.r.t all
output neurons. from $a_{i+1,1}$ to $a_{i+1,k}$.

↓
It refers all rows of
$j^{th}$ column.

dot product of $\left( W_{i+1 \cdot j} \right)^T \nabla_{a_{i+1}}^{a} L(\theta) = \sum_{m=1}^{k} \dfrac{\partial L(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}$

Here,

derivative of Loss function w.r.t hidden Layer is only the dot product

but" gradient of Loss w.r.t output layer and corresponding wts.

So,

$$\frac{\partial L(\theta)}{\partial h_{ij}} = \left(W_{i+1 \cdot j}\right)^T \nabla_{o_{j+1}} L(\theta)$$

$$\nabla_{h_i} L(\theta) \quad \begin{bmatrix} \left(W_{i+1,j}\right)^T \nabla_{o_{i+1}} L(\theta) \\ \\ \vdots \\ \\ \vdots \\ \\ \left(W_{i+1,n}\right)^T \nabla_{o_{i+1}} L(\theta) \end{bmatrix}$$

also, $\left(W_{i+1}\right)^T \nabla_{o_{i+1}} L(\theta)$

It is special case for Loss & hidden Layers.

Now, lets mak it more generic for all hidden Layer.

1. consider for next layer & $a_i$

$$\frac{\partial L(\theta)}{\partial a_{ij}} = \frac{\partial L(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}} \quad \left\{ h_{ij} \text{ is application of} \atop \text{activation} \right.$$

already
computed ← So, $\frac{\partial L(\theta)}{\partial a_{ij}} = \frac{\partial L(\theta)}{\partial h_{ij}} g'(a_{ij})$ ↗ activation

$$\nabla_{a_i} L(\theta) \quad \begin{bmatrix} \frac{\partial L(\theta)}{\partial h_{i1}} g'(a_{i1}) \\ \\ \vdots \\ \\ \frac{\partial L(\theta)}{\partial h_{in}} g'(a_{in}) \end{bmatrix}$$

$$\nabla_{a_i} L(\theta) = \nabla_{h_i} L(\theta) \odot \begin{bmatrix} \cdots g'(a_{ik}) \cdots \end{bmatrix}$$

↑
element wisve vector Mult.

Till now

$$\frac{\partial L(\theta)}{\partial w_{111}} = \frac{\partial L(\theta)}{\partial \hat{y}} \frac{\partial a_3}{\partial a_3} \frac{\partial a_3}{\partial h_2} \frac{\partial h}{\partial a_2} \frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a} \frac{\partial a_3}{\partial w_{111}}$$

Uptil now,

$$\nabla h_i L(\theta) = (W_{i+1})^T \nabla a_{i+1} L(\theta)$$

$$\nabla a_i L(\theta) = \nabla h_i L(\theta) \odot [\cdots g'(a_{ik}) \cdots]$$

Now, deri of loss wrt wt and bias

$$a_k = b_k + W_k h_{k-1}$$

$$\frac{\partial a_{ki}}{\partial W_{kij}} = h_{k-1, j}$$

- $k$ : layer
  $i$ = current layer neron no
  $j$ = prev / input layer neuron no.

$$\frac{\partial L(\theta)}{\partial W_{kij}} = \frac{\partial L(\theta)}{\partial a_{ki}} h_{k-1, j}$$

$$W_{kij} = W_{kij} - \eta \frac{\partial L(\theta)}{\partial W_{kij}}$$

$\nabla_{W_k} L(\theta)$

$$\begin{bmatrix} \dfrac{\partial L(\theta)}{\partial W_{k11}} & \dfrac{\partial L(\theta)}{\partial W_{k12}} & \dfrac{\partial L(\theta)}{\partial W_{k13}} \\ \dfrac{\partial L(\theta)}{\partial W_{k21}} & \dfrac{\partial L(\theta)}{\partial W_{k22}} & \dfrac{\partial L(\theta)}{\partial W_{k23}} \\ \dfrac{\partial L(\theta)}{\partial W_{k31}} & \dfrac{\partial L(\theta)}{\partial W_{k32}} & \dfrac{\partial L(\theta)}{\partial W_{k33}} \end{bmatrix}$$

$\nabla W_{kij}$

$$L(\theta) = \frac{\partial L(\theta)}{\partial a_{ki}}$$

$\nabla_{W_k} L(\theta)$

$$\begin{bmatrix} \dfrac{\partial L(\theta)}{\partial a_{k1}} h_{k-} & \dfrac{\partial L(\theta)}{\partial a_{k1}} h_{k-1, 2} & \dfrac{\partial L(\theta)}{\partial a_{k1}} h_{k-1,3} \\ \dfrac{\partial L(\theta)}{\partial a_{k2}} h_{k-} & \dfrac{\partial L(\theta)}{\partial a_{k2}} h_{k-1,2} & \dfrac{\partial L(\theta)}{\partial k_2} h_{k-1,3} \\ \dfrac{\partial L(\theta)}{\partial k_3} h_{k-} & \dfrac{\partial L(\theta)}{\partial a_{k3}} h_{k-1,2} & \dfrac{\partial L(\theta)}{\partial a_{k3}} k-1,3 \end{bmatrix}$$