

Loss Functions

$$\cdot \text{MSE} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

$$\cdot \text{MAE} = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}$$

• Squared Error

$$= \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

• cross entropy. → logistic Regression

$$\boxed{\text{loss} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})} \quad \begin{matrix} \text{for} \\ \text{binary} \\ \text{cross entropy.} \end{matrix}$$

$$\begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

• Multi class cross entropy loss

$$l(x_i, y_i) = -\sum_{j=1}^c y_{ij} \log(\hat{y}_{ij})$$

categorical cross entropy.

c → no. of classes in o/p

O/P is multiclass e.g. [good, bad, neutral]

y_i	f_1	f_2	f_3	O/P	b_1	b_2	b_3	One hot encoding
-2	3	4		Good	1	0	0	
-4	5	7		Bad	0	1	0	
3	8	9		Neutral	0	0	1	

O/P y_i is one hot encoded vector.

$$y_i = \begin{bmatrix} y_{11} & y_{12} & y_{13} \\ y_{21} & y_{22} & y_{23} \\ y_{31} & y_{32} & y_{33} \end{bmatrix} \xrightarrow{?} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

where $y_{ij} = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ element is in class } j \\ 0 & \text{else} \end{cases}$

How to calculate \hat{y}_{ij}

use softmax Activation function $\sigma(z)$

$$\sigma(z) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

e.g. Against vector

$$[10, 20, 30, 40]$$

$$\text{for } 10 \quad \sigma(2) = \frac{e^{10}}{e^{10} + e^{20} + e^{30} + e^{40}}$$

$$\text{for } 20 \quad \sigma(2) = \frac{e^{20}}{e^{10} + e^{20} + e^{30} + e^{40}}$$

i.e. probability

Tanh

$$f(x) = \tanh(x) \cdot \frac{e^x}{1+e^{-x}} - 1 = \text{d sigmoid}(dx) - 1$$

$$f'(x) = 1 - f(x)^2$$

Sigmoid

$$f(x) = \frac{1}{1+e^{-x}}$$

$$f'(x) = \cancel{\frac{f(x)}{1-f(x)}}$$

$$f'(x) = f(x)(1-f(x))$$

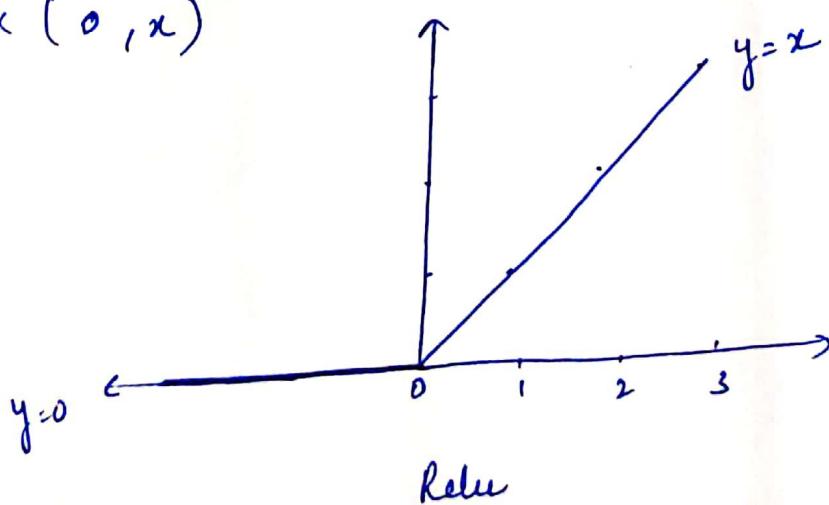
ReLU

Rectified Linear Unit

$$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

$$f'(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$

$$f(x) = \max(0, x)$$



Leaky ReLU

$$f(x) = \max(0.01 * x, x)$$

$$f(x) = \begin{cases} x & x > 0 \\ 0.01 * x & x \leq 0 \end{cases}$$

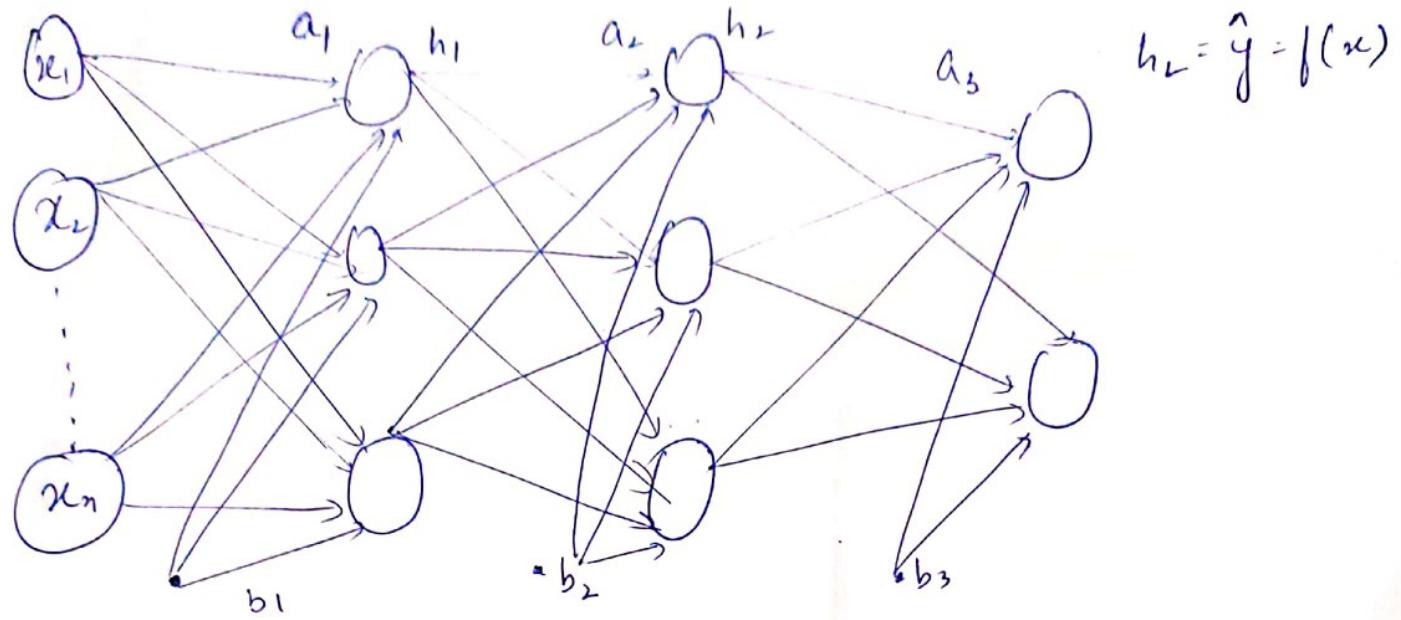
softmax

Probability distribution

$$f(x) = \frac{e^{x_i}}{\sum_{j=1}^k e^{x_j}}$$

category	car	bike	scooter
	0.7	0.2	0.1

FeedForward N/W



I/P layer

(L-1) hidden layers

here it is 2

Properties:

- ① n-dimensional vector as input to the n/w
- ② L-1 hidden layers having m neurons each.
- ③ One O/P Layer with K neurons
- ④ Each Neuron in hidden layer & O/P layer has 2 faces \rightarrow pre-activation (a_i) & activation (h_i)
- ⑤ I/P layer is 0th layer ; O/P layer is Lth layer
- ⑥ $w_i \in \mathbb{R}^{n \times n}$ and $b_i \in \mathbb{R}^n$ are the weights and bias b/w layers i-1 and i ($0 < i < L$)

$\rightarrow W_L \in \mathbb{R}^{n \times k}$ and $b_L \in \mathbb{R}^k$ are the weight and bias of the last hidden layer and the o/p layer ($L=3$ in this case)

① Pre-activation fxn at any layer i

$$a_i(x) = b_i + W_i h_{i-1}(x)$$

$$\begin{bmatrix} a_{11} \\ a_{12} \\ a_{13} \end{bmatrix} = \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} + \begin{bmatrix} w_{111} & w_{112} & w_{113} \\ w_{121} & w_{122} & w_{123} \\ w_{131} & w_{132} & w_{133} \end{bmatrix} \begin{bmatrix} h_{01} = x_1 \\ h_{02} = x_2 \\ h_{03} = x_3 \end{bmatrix}$$

$$a_1 = b_1 + W_1 h_0$$

$$(3 \times 1) \quad (3 \times 1) \quad (3 \times 3) \cdot (3 \times 1)$$

$$= \begin{bmatrix} \sum w_{11i} x_i + b_{11} \\ \sum w_{12i} x_i + b_{12} \\ \sum w_{13i} + b_{13} \end{bmatrix}$$

② Activation at layer i is

$$h_i(x) = g(a_i(x))$$

$$g(a_i) = \sigma(a_i) = \frac{1}{1 + e^{-a_i}}$$

Sigmoid
fxn

$$\begin{bmatrix} h_{11} \\ h_{12} \\ h_{13} \end{bmatrix} = g \begin{bmatrix} a_{11} \\ a_{12} \\ a_{13} \end{bmatrix} = \begin{bmatrix} g(a_{11}) \\ g(a_{12}) \\ g(a_{13}) \end{bmatrix}$$

$$g(a_{11}) = \frac{1}{e^{-a_{11}} + 1}$$

$g \rightarrow$ activation fxn.

(logistic, tanh, linear, ReLU-...)

③ Output at Lth layer

$$y = f(x) = h_L(x) = O(a_L(x))$$

O \rightarrow output activation fxn (softmax etc.)

④ Data = $\{x_i, y_i\}_{i=1}^N$

⑤ Model :

$$\hat{y} = f(x_i) = \Theta^T (\omega_0^3 g(\omega_1^2 g(\omega_0^1 x + b_1) + b_2) + b_3)$$

⑥ Parameters

$$\Theta = \omega_0, \omega_1, \dots, \omega_L \quad \text{size } \mathbb{R}^{n \times n}$$

$$b = b_0, b_1, \dots, b_L \quad \text{size } \mathbb{R}^{n \times 1}$$

⑦ Algorithm

Gradient Descent with BP

⑧ Loss / Error Fxn / Objective Fxn

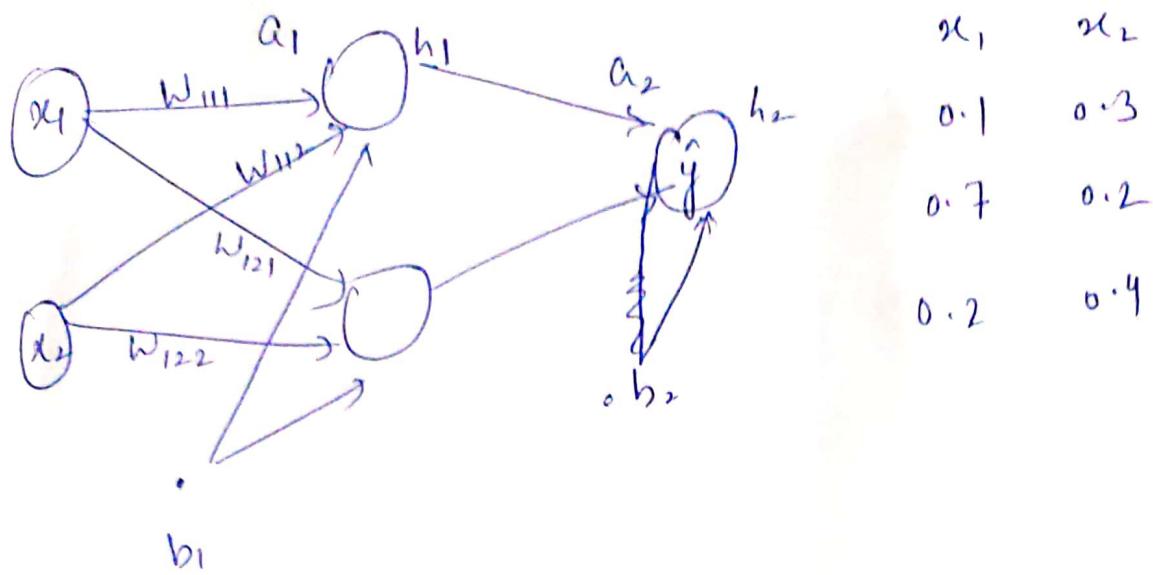
$$\min \frac{1}{N} \cdot \sum_{i=1}^N \sum_{j=1}^K (\hat{y}_{ij} - y_{ij})^2$$

[sum of squared f(x)]

$$\Leftarrow \min \mathcal{L}(\Theta)$$

5 101

Single hidden layer feed forward example.



$$a_1 = b_1 + w_1 h_0$$

$$\textcircled{1} \quad \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} = \begin{bmatrix} b_{11} \\ b_{12} \end{bmatrix} + \begin{bmatrix} w_{111} & w_{112} \\ w_{121} & w_{122} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$\uparrow \quad \uparrow \quad \uparrow$
 $a_1 \quad b_1 \quad w_1$

Randomly selected parameters

$$\left\{ \begin{array}{l} w_1 = \begin{bmatrix} 0.5 & 0.4 \\ 0.9 & 1.0 \end{bmatrix} \\ b_1 = \begin{bmatrix} -0.1 \\ -0.1 \end{bmatrix} \end{array} \right.$$

$$\begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} = \begin{bmatrix} -0.1 \\ -0.1 \end{bmatrix} + \begin{bmatrix} 0.5 & 0.4 \\ 0.9 & 1.0 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.3 \end{bmatrix} = \begin{bmatrix} 0.07 \\ 0.29 \end{bmatrix}$$

$$\textcircled{2} \quad h_1 = g(a_1) = g \begin{bmatrix} 0.07 \\ 0.29 \end{bmatrix} = \begin{bmatrix} \sigma(0.07) \\ \sigma(0.29) \end{bmatrix}$$

$$\sigma(0.07) = \frac{1}{1 + e^{-0.07}}$$

$$\sigma(0.29) = \frac{1}{1 + e^{-0.29}}$$

$$\textcircled{3} \quad \begin{bmatrix} a_{21} \\ a_{22} \end{bmatrix} = \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix} + \begin{bmatrix} w_{211} & w_{212} \\ w_{221} & w_{222} \end{bmatrix} \begin{bmatrix} \sigma(0.07) \\ \sigma(0.29) \end{bmatrix}$$

$$\textcircled{4} \quad y_{11} = h_2 = \theta(a_2) = \theta \begin{bmatrix} a_{21} \\ a_{22} \end{bmatrix}$$

$\theta \rightarrow \text{softmax}$.

h_2 is final predicted o/p denoted as \hat{y}

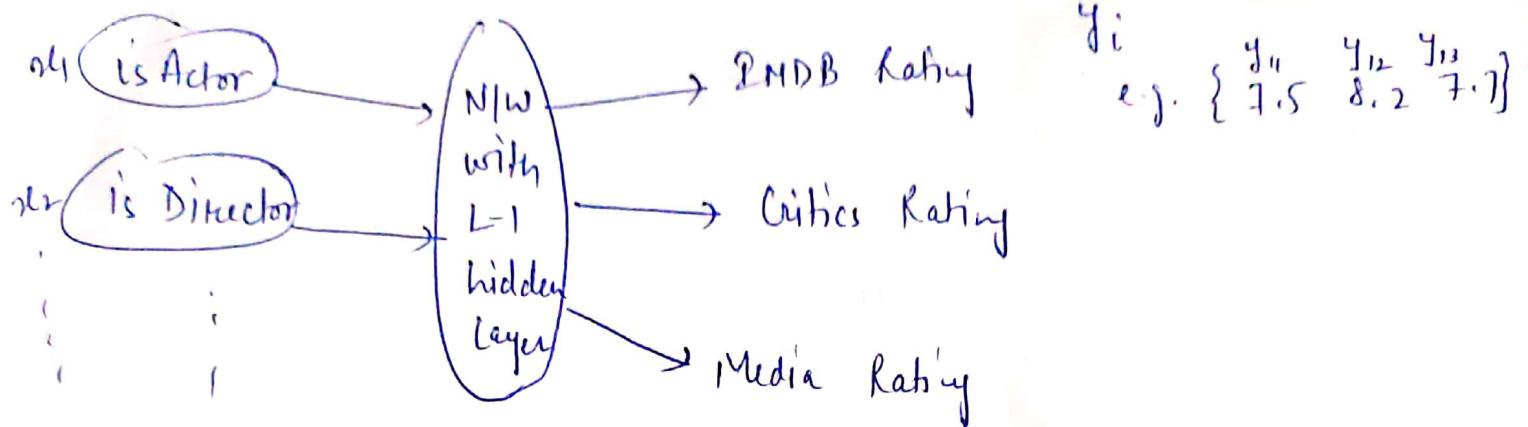
$$\textcircled{5} \quad \text{error} = \sum_{i=1}^n (y - \hat{y})^2$$

choice of loss function

4

→ depends on problem in hand.

e.g. ①



$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^3 (\hat{y}_{ij} - y_{ij})^2$$

O/P function should be such that it should not be bounded b/w 0 to 1.

∴ can't be logistic

Next option is O/P fxn to be linear

$$\begin{aligned} f(x) &= h_r = \theta(a_r) \\ &= w_0 a_r + b_0 \end{aligned}$$

∴ Which O/P function with suit.

Linear. As the problem is Regression
kind of Problem

Example 2

$$y_i = \begin{bmatrix} R & B & W \\ 0.35 & 0.25 & 0.4 \end{bmatrix}$$

True probability of drawing a ball from a knapsack with color Red/Blue/White.

suppose $\hat{y}_i = [0.25 \quad 0.45 \quad 0.3]$

$$L(\theta) = \sum_{i=1}^3 (y_i - \hat{y}_i)^2$$

OIP function should be such which should return bounded value b/w $0 \rightarrow 1$.

Ans can we use Sigmoid as it returns b/w $0 \rightarrow 1$

No. Because the probability distribution should sum to 1.

Softmax is the solution.

$$a_L = b_L + w_L \cdot h_{L-1}$$

$$\hat{y}_j = O(a_L)_j = \frac{e^{a_{L,j}}}{\sum_{i=1}^K e^{a_{L,i}}}$$

e.g. $a_L = \begin{bmatrix} a_{L1} & a_{L2} & a_{L3} \\ 10 & -20 & 30 \end{bmatrix}$

$$\hat{y} = [\hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3]$$

$$\hat{y}_i = \frac{e^{x_i}}{e^{x_1} + e^{x_2} + e^{x_3}}$$

[softmax fun]

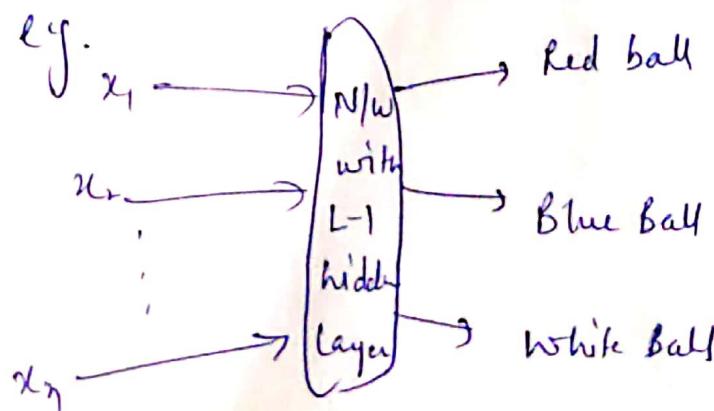
y, \hat{y} both are probability distributions

$\therefore L(\theta)$ preferred will be cross entropy
rather than squared error.

$$L(\theta) = - \sum_{i=1}^k y_i \log \hat{y}_i$$

$$y_i = 1 \quad \text{if } i = \text{true class (t)} \\ = 0 \quad \text{otherwise}$$

$$\therefore L(\theta) = - \log \hat{y}_t$$



$$y_t = \begin{bmatrix} 1 & 0 & 0 \\ R & B & W \end{bmatrix}$$

means ball drawn
is Red, the true
class here.

∴ for classification Problems where O/P
is one of the classes between 1 and K
objective Fxn is

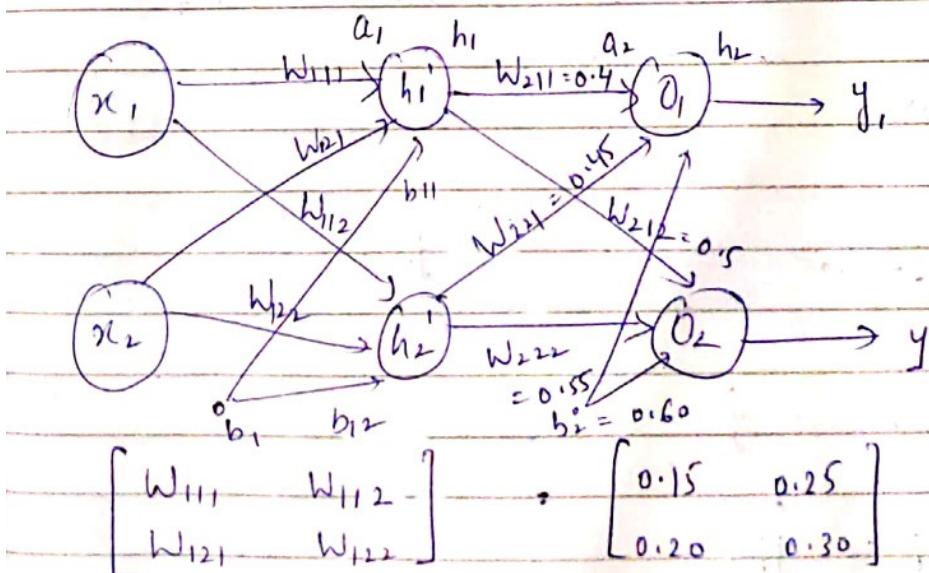
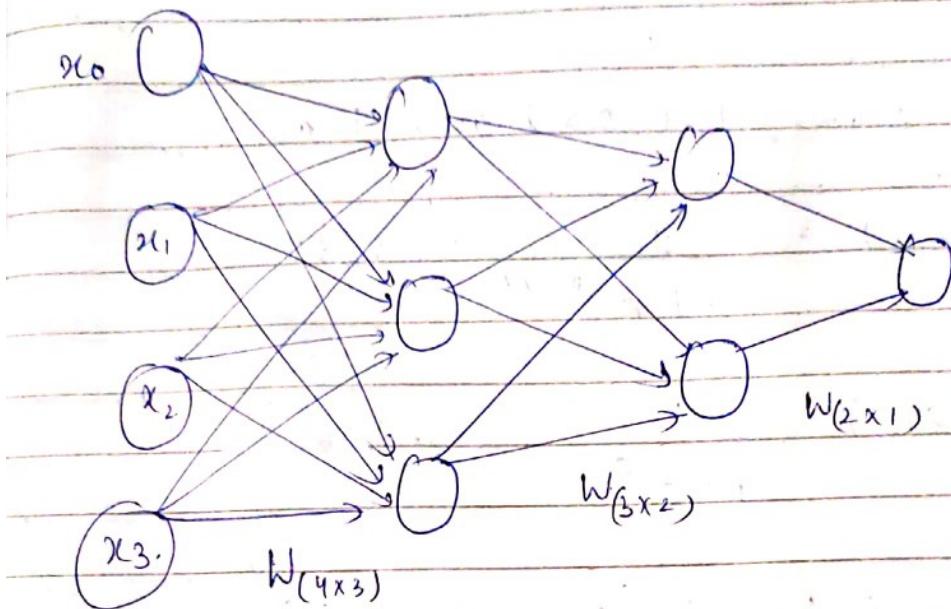
$$\text{minimize } \theta, L(\theta) = -\log \hat{y}_t$$

or

$$\text{maximize } \theta, -L(\theta) = +\log \hat{y}_t$$

$$\theta = [w_1, w_2, \dots, w_L, b_1, b_2, \dots, b_L]$$

Back Propagation: Example



For layer $i=1$

$$\begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} = \begin{bmatrix} b_{11} \\ b_{12} \end{bmatrix} + \begin{bmatrix} w_{111} & w_{112} \\ w_{121} & w_{122} \end{bmatrix}^T \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$a_{11} = b_{11} + (w_{111} \cdot x_1 + w_{121} \cdot x_2)$$

$$a_{12} = b_{12} + (w_{112} \cdot x_1 + w_{122} \cdot x_2)$$

$$a_{11} = ((0.15 * 0.05) + (0.2 * 0.1)) + 0.35 \\ = 0.377$$

$$a_{12} = ((0.25 * 0.05) + (0.3 * 0.1)) + 0.35$$

$$h_{11} = \sigma(a_{11}) = \frac{1}{1 + e^{-0.377}} = 0.5932$$

$$h_{12} = \sigma(a_{12}) = \frac{1}{1 + e^{-a_{12}}} = 0.5968$$

$$a_{21} = ((w_{211} * h_{11}) + (w_{221} * h_{22})) + b_{21}$$

$$= ((0.4 * 0.5932) + (0.45 * 0.5968)) + 0.6 \\ = 1.105$$

$$h_{21} = \frac{1}{1 + e^{-1.105}} = 0.7513 = \hat{y}_1$$

$$h_{22} = 0.7729 = \hat{y}_2$$

$$\text{error} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

$$= \frac{(0.01 - 0.7513)^2 + (0.99 - 0.7729)^2}{2}$$

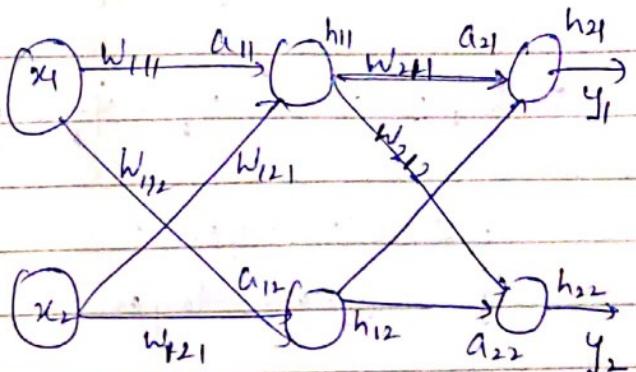
$$\text{error} = 0.2983$$

Minimize error

$$\frac{\partial \text{Error}}{\partial w_{21}} = \frac{\partial \text{Error}}{\partial y_1} \rightarrow \frac{\partial y_1}{\partial a_{21}} \times \frac{\partial a_{21}}{\partial w_{21}}$$

$$\frac{\partial \text{Error}}{\partial y_1} = \frac{\partial \text{Error}}{\partial h_{21}}$$

$$= \frac{\partial}{\partial h_{21}} \left(\frac{1}{2} (y_1 - h_{21})^2 \right)$$



[Apply chain rule]

$$\textcircled{1} \quad \frac{\partial \text{error}}{\partial h_{21}} = -(y_1 - h_{21})$$

OR

$$\boxed{\frac{\partial \text{error}}{\partial \hat{y}_k} = -(y_k - \hat{y}_k)}$$

$$\textcircled{2} \quad \frac{\partial y_1}{\partial a_{21}} = \frac{\partial}{\partial a_{21}} \frac{(1 + e^{-y_1})^{-1}}{2} = \frac{e^{-y_1}}{(1 + e^{-y_1})^2}$$
$$= \left(1 - \frac{1}{1 + e^{-y_1}} \right) \frac{1}{1 + e^{-y_1}}$$
$$= y_1 (1 - y_1)$$
$$\frac{\partial y_1}{\partial a_{21}} = h_{21} (1 - h_{21})$$

$$\frac{\partial A_{21}}{\partial w_{211}} = \frac{\partial ((w_{211} \cdot h_{11} + b_{2121}) + b_{11})}{\partial w_{211}}$$

$$= h_{11}$$

$$\nabla w_{211} = -(y_i - h_{21}) \cdot h_{21} (1-h_{21}) \cdot h_{11}$$

~~$$\nabla w_{kij} = (y_k - j_k) j_k (1-j_k) h^l$$~~

k \rightarrow output layer

j \rightarrow hidden layer

$$\nabla w_{lij} = -(y_{actuali} - h_{li}) \cdot h_{li} (1-h_{li}) \cdot h_{i+1}$$

l \rightarrow layer no.

i \rightarrow neuron no. I/P

j \rightarrow output neuron no.

$$\text{error} = 0.29837$$

$$\frac{\partial \text{error}}{\partial w_{211}} = \frac{\partial \text{error}}{\partial h_{21}} * \frac{\partial h_{21}}{\partial a_{21}} + \frac{\partial \text{error}}{\partial w_{211}}$$

$$(- (0.01 - 0.751365)) * (0.751365 * \\ (1 - 0.751365)) + 0.5932$$

$$\nabla w_{211} = 0.08216704$$

$$\text{New updated } w_{211} = w_{211} - \alpha \nabla w_{211}$$

$$= 0.4 - (0.60) * (0.08216704) \\ = 0.350699776$$

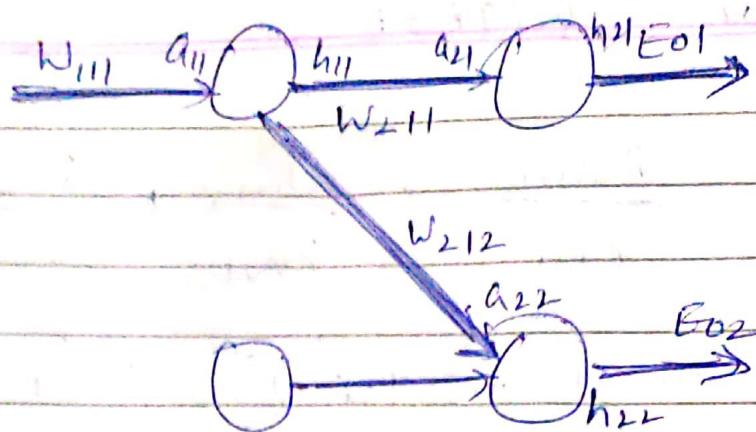
Likewise update:

$$w_{212}, w_{221}, w_{222}$$

~~B/w LIP & hidden layer~~

Now solving $w_{111}, w_{112}, w_{121}, w_{122}$

$$\frac{\partial \text{error}}{\partial w_{111}} = \frac{\partial \text{error}}{\partial h_{11}} * \frac{\partial h_{11}}{\partial a_{11}} + \frac{\partial \text{error}}{\partial w_{111}}$$



$$\frac{\partial \text{error}}{\partial w_{111}} = \left(\frac{\partial \text{error}}{\partial h_{11}} \right) + \frac{\partial h_{11}}{\partial a_{11}} + \frac{\partial a_{11}}{\partial w_{111}}$$

$$\frac{\partial \text{error}}{\partial h_{11}} = \left(\frac{\partial \text{error}}{\partial E_01} \right) + \left(\frac{\partial E_02}{\partial h_{11}} \right)$$

$$\frac{\partial E_01}{\partial h_{21}} + \frac{\partial E_01}{\partial a_{21}} + \frac{\partial a_{21}}{\partial h_{11}}$$

$$+ \frac{\partial E_02}{\partial h_{22}} + \frac{\partial h_{22}}{\partial a_{22}} + \frac{\partial a_{22}}{\partial h_{11}}$$

calculated
above

$$= -(y_{\text{actual}, 1} - h_{21}) \cdot h_{21} (1-h_{21}) + \frac{\partial ((w_{211} * h_{11}) + (w_{212} * h_{12}) + b_{11})}{\partial h_{11}}$$

$$= -(y_{\text{actual}, 1} - h_{21}) \cdot h_{21} (1-h_{21}) + w_{211}$$

$$\frac{\partial E_{02}}{\partial h_{22}} + \frac{\partial h_{22}}{\partial a_{22}} + \frac{\partial a_{22}}{\partial w_{11}}$$

$$-(y_{\text{actual}_2} - h_{22}) \cdot h_{22}(1-h_{22}) + w_{212}$$

(1)

put values in ①

$$(0.13849856 * 0.4) + (-0.0380982 + 0.5)$$

$$= 0.055399425 + (-0.01904919)$$

$$\boxed{\frac{\partial \text{error}}{\partial h_{11}} = 0.036350306.}$$

$$\boxed{\frac{\partial h_{11}}{\partial a_{11}} = a_{11}(1-a_{11}) \\ = 0.241300709}$$

$$\boxed{\frac{\partial a_{11}}{\partial w_{111}}} = \frac{\partial (w_{111} + x_1 + w_{121} + x_2 + b_{11})}{\partial w_{111}} \\ = x_1 = 0.05$$

$$\nabla w_{111} = w_{111} - \alpha \frac{\partial \text{error}}{\partial w_{111}}$$

$$= 0.15 - (0.6) * 0.00438568$$

$$= 0.1497368592$$