

* Linear Regression:

$$y = \theta_0 + \theta_1 x \quad \text{OR} \quad y = \theta_0 + \theta_1 x$$

slope ↑ intercept ↘

$$\theta_0 = b = \frac{n \sum (xy) - \sum x \sum y}{n \sum x^2 - (\sum x)^2}, \theta_1 = \bar{y} - b\bar{x}$$

$$\text{cost function: } J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [\hat{y}_i - y_i]^2$$

↑ predicted
actual
mean squared error

gradient descent:

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j}$$

$$\theta_j = \theta_j - \alpha \sum_{i=1}^m [(\hat{y}_i - y_i)x_i] \quad (\text{final formula})$$

* Logistic Regression:

when we use sigmoid function of linear regression, we get logistic regression.

$$(\text{Linear Reg}) \quad y = \theta_0 + \theta_1 x$$

$$(\text{sigmoid func}) \quad p = \frac{1}{1 + e^{-y}}$$

$$(\text{logistic reg}) \quad p = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$

OR

$$\ln\left(\frac{P}{1-P}\right) = \theta_0 + \theta_1 x$$

Date: _____

cost func: ~~_____~~

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m [-y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

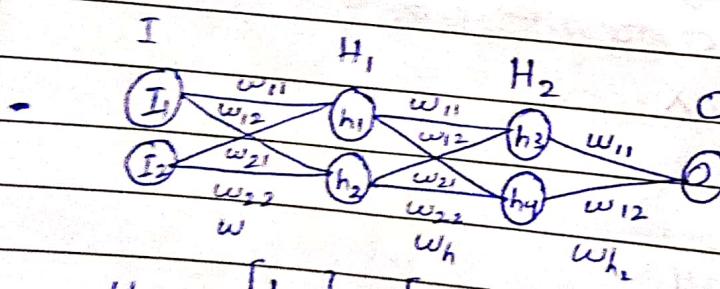
* we used mean squared error in linear regression,
but we use a different error calculation method
in logistic regression.

gradient descent:

$$\theta_j = \theta_j - \alpha \sum_{i=1}^m [\hat{y}_i - y_i] x_i$$

Neural Networks

Forward pass:



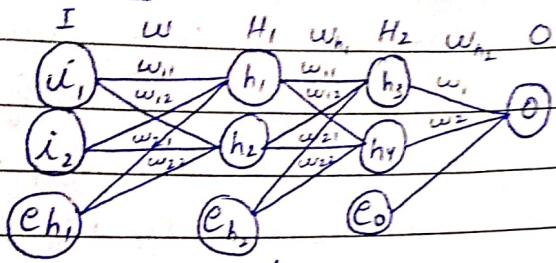
$$H_1 = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} i_1 w_{11} + i_2 w_{21} \\ i_1 w_{12} + i_2 w_{22} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}^T \begin{bmatrix} i_1 \\ i_2 \end{bmatrix}$$

$$H_1 = \begin{bmatrix} w_{11} & w_{21} \\ w_{12} & w_{22} \end{bmatrix} \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} = W^T I$$

$$H_2 = W_h^T H_1 , O = W_o^T H_2$$

Date: _____

Backward Pass:



$$e_o = \hat{o} - y \leftarrow \text{actual}$$

$$e_{h_2} = \begin{bmatrix} e_{h_21} \\ e_{h_22} \end{bmatrix} = \begin{bmatrix} w_{h_21} / (w_{h_21} + w_{h_22}) \times e_o \\ w_{h_22} / (w_{h_21} + w_{h_22}) \times e_o \end{bmatrix} = \underbrace{w_{h_2}}_{\text{sum of diag.}} e_o$$

$$e_{h_1} = \begin{bmatrix} e_{h_11} \\ e_{h_12} \end{bmatrix} = \begin{bmatrix} w_{11} / (w_{11} + w_{12}) \times e_{h_21} + w_{12} / (w_{11} + w_{12}) \times e_{h_22} \\ w_{21} / (w_{21} + w_{22}) \times e_{h_21} + w_{22} / (w_{21} + w_{22}) \times e_{h_22} \end{bmatrix}$$

$$e_{h_1} = \underbrace{w}_{\text{sof diag.}} e_{h_2}$$

$$\text{gradient descent: } w = w - \alpha \Delta w$$

Δw = error of next layer \times cost function \times input layer^T

$$\text{cost function} = y(1-y) \times \begin{matrix} \uparrow & \uparrow & \uparrow \\ \text{sigmoid of} & \text{output} & \text{transpose of} \\ \text{input layer.} & & \end{matrix}$$

$$\Delta w_{11} = e_{h_1} \times h_1 (1-h_1) \times \cancel{i_1} \quad (\cancel{E_{h_1}} \times \cancel{h_1} \times 1-h_1) \times \cancel{i_1}$$

$$\Delta w_{12} = e_{h_1} \times h_2 (1-h_2) \times \cancel{i_1} \quad \cancel{e_{h_1}} \cancel{h_2} \cancel{1-h_2} \cancel{i_1}$$

$$\Delta w_{21} = e_{h_1} \times h_1 (1-h_1) \times i_2 \quad \cancel{e_{h_1}} \cancel{h_1} \cancel{1-h_1} \cancel{i_2}$$

$$\Delta w_{22} = e_{h_1} \times h_2 (1-h_2) \times i_2 \quad \cancel{e_{h_1}} \cancel{h_2} \cancel{1-h_2} \cancel{i_2}$$

$$\Delta w = [E_{h_1} \times h_1 (1-h_1)] \times I^T$$

$$\begin{aligned} & [e_{h_1} \times h_1 + e_{h_1} \times h_2] (1-h_1) + [e_{h_1} \times h_1 + e_{h_1} \times h_2] (1-h_2) e_{h_1} \times h_1 + e_{h_1} \times h_2 \\ & [e_{h_2} \times h_1 + e_{h_2} \times h_2] (1-h_1) + [e_{h_2} \times h_1 + e_{h_2} \times h_2] (1-h_2) e_{h_2} \times h_1 + e_{h_2} \times h_2 \end{aligned}$$

ate: _____

* Logistic Regression

$$\log(\text{odds}) = \log\left(\frac{P}{1-P}\right)$$

↑ probability (p-value)
of success
↓ prob of failure

~~logit function~~

↓ bases of logistic regression

Concept:

odds are simply the ratio of the probability of success to probability of failure. $\log(\text{odds})$ makes the distribution symmetrical & easy to interpret, so we use that.

~~The typically used p-value is 0.05. Its basically a threshold to differ classes -~~

Date:

y
weight obese

50 0

100 1

20 0

30 0

110 1

90 1

80 0

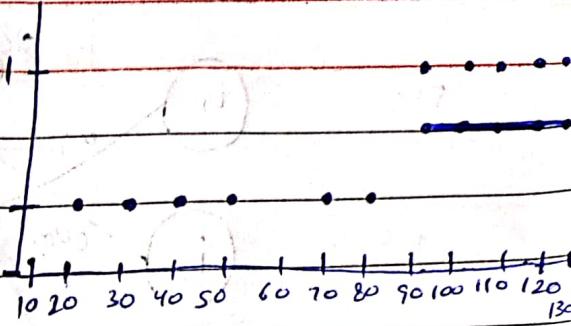
70 0

40 0

120 1

130 1

Logistic Regression



$$y = bx + a$$

$$b = \frac{n \sum xy - \sum x \sum y}{n \sum (x^2) - (\sum x)^2} \Rightarrow \frac{11(550) - 840(5)}{11(78200) - (840)^2}$$

$$n \sum (x^2) - (\sum x)^2 = 11(78200) - (840)^2$$

$$b = \frac{1850}{11} \Rightarrow 177.27 \Rightarrow 0.01196636481$$

$$a = \bar{y} - b\bar{x} = \frac{154600}{11} - \frac{37}{30.92} \Rightarrow 37$$

$$a = \bar{y} - b\bar{x} = 5 - \frac{37}{30.92} \left(\frac{840}{11} \right) \Rightarrow -0.4592445$$

$$= 0.45459 - 0.91379 \Rightarrow -0.4592445$$

$$\text{sigmoid} = \frac{e^{bx+a}}{1+e^{bx+a}} = \frac{e^{-0.4592445}}{1+e^{-0.4592445}} \Rightarrow 0.58287 \Rightarrow 0.582$$

$$\text{sig}\left(y = 0.01196x - 0.4592\right)$$

plug in a value to test: Let pval = 0.05

$$x = 150$$

$$\text{sig}(0.01196(150) - 0.4592) = 0.208$$

$$0.208 > 0.05 \text{ thus } y = 1$$

I ₁	1	I ₂	1	0
1	0	0	0	0
0	1	0	0	0
0	0	0	0	0

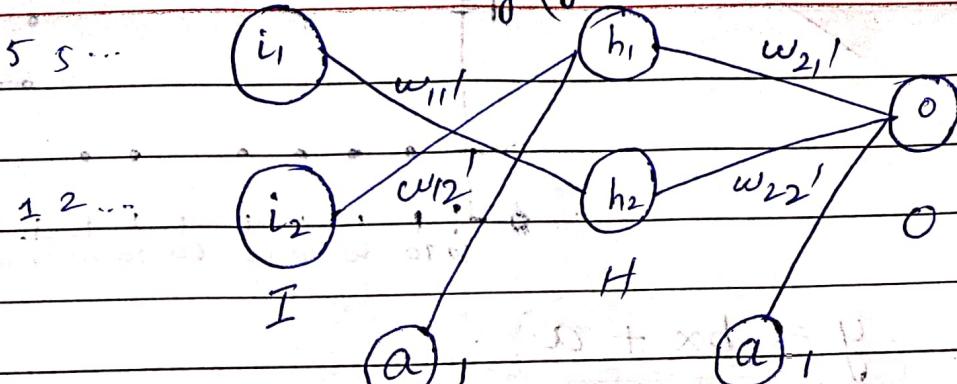
Date:

neural network

$$\text{sig } (y = a + bx)$$

$$y = a + wx$$

$$\frac{1}{1+e^{-y}}$$



forward pass for (1, 1):

output layer

$$o = a + w_1 h_1 + w_2 h_2$$

$$o = 1 + 1(0.88) + 1(0.88)$$

$$\text{sig}(o) = \frac{1}{1+e^{-(2.76)}} \Rightarrow 0.94$$

$$\text{sig}(2) = \frac{1}{1+e^{-2}} \Rightarrow 0.88$$

$$h_2 = 0.88 \cdot 0.88$$

forward pass for (1, 0):

hidden layer

$$h_1 = 1 + 1(0) \Rightarrow 1$$

$$h_2 = 1 + 1(1) \Rightarrow 2 \quad \text{sig}(1) = \frac{1}{1+e^{-1}} \Rightarrow 0.73 \quad \text{output}$$

$$\text{sig}(2) = \frac{1}{1+e^{-2}} \Rightarrow 0.88 \quad \text{sig}(0) = \frac{1}{1+e^{-(2.73)}} \Rightarrow 0.93$$

forward pass for (0, 1):

hidden

output

$$o = 1 + 1(0.88) + 1(0.73)$$

$$h_1 = 0.88, h_2 = 0.73 \quad \text{sig}(0) = 0.93$$

forward pass for (0, 0):

hidden

output

$$o = 1 + 1(0.73) + 1(0.73)$$

$$\text{sig}(0) = 0.92$$

Date: _____

weight updation using batch gradient descent

costs:

for q :

$$\text{Let } \alpha = 0.01, m = 4$$

$$\text{cost}_q = \alpha \frac{1}{2m} \sum (\hat{y} - y)^2$$

$$= 0.01 \left(\frac{1}{2(4)} (0.94 - 1)^2 + (0.93 - 0)^2 + (0.93 - 0)^2 \right)$$

$$= 0.01 (0.0036 + 0.8649 + 0.8649 + 0.8464)$$

$$= 0.01 (0.1896) \Rightarrow 0.00189$$

for w_1 :

$$\text{cost}_{w_1} = 0.00189$$

for w_2 :

$$\text{cost}_{w_2} = 0.00189$$

new values:

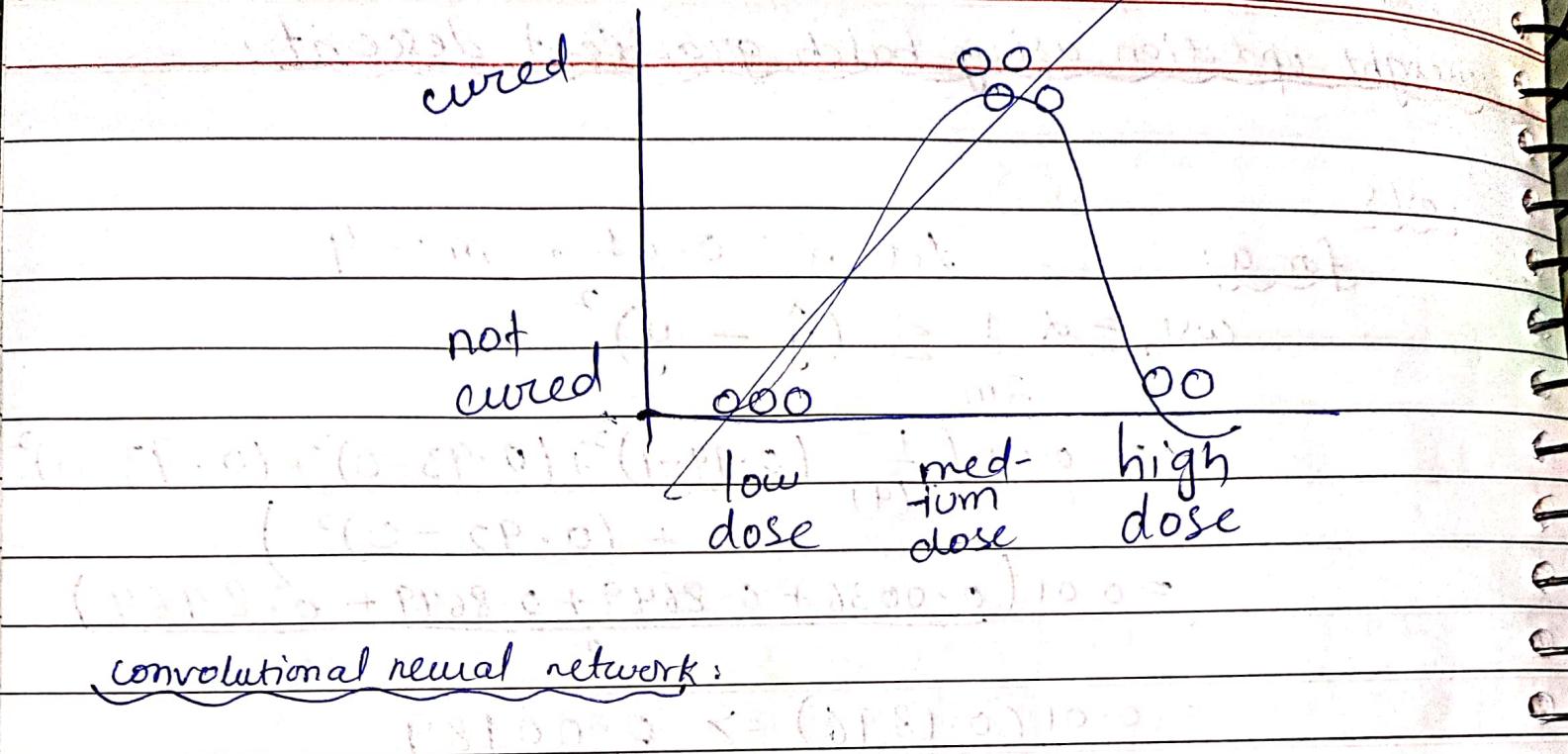
$$q = 1 - 0.00189 = 0.999$$

$$w_{11} = 1 - 0.00189 = 0.999$$

$$w_{21} = 1 - 0.00189 = 0.999$$

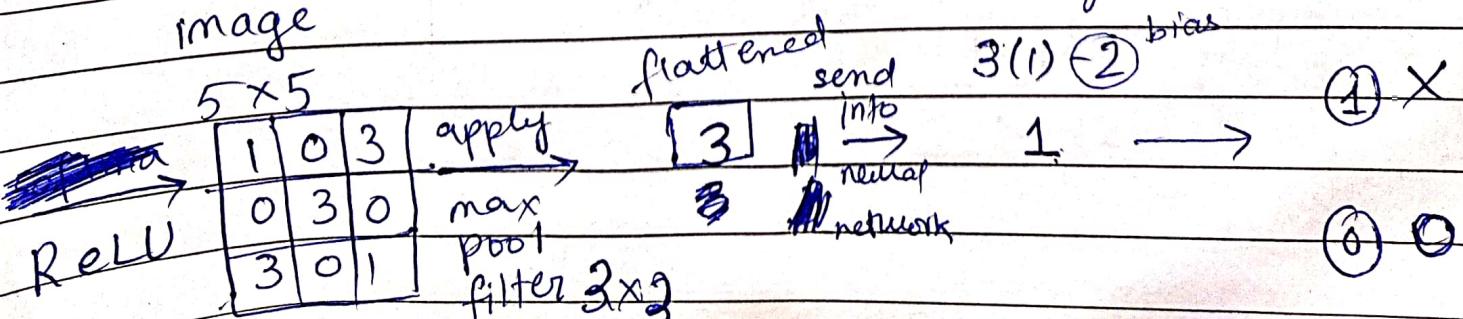
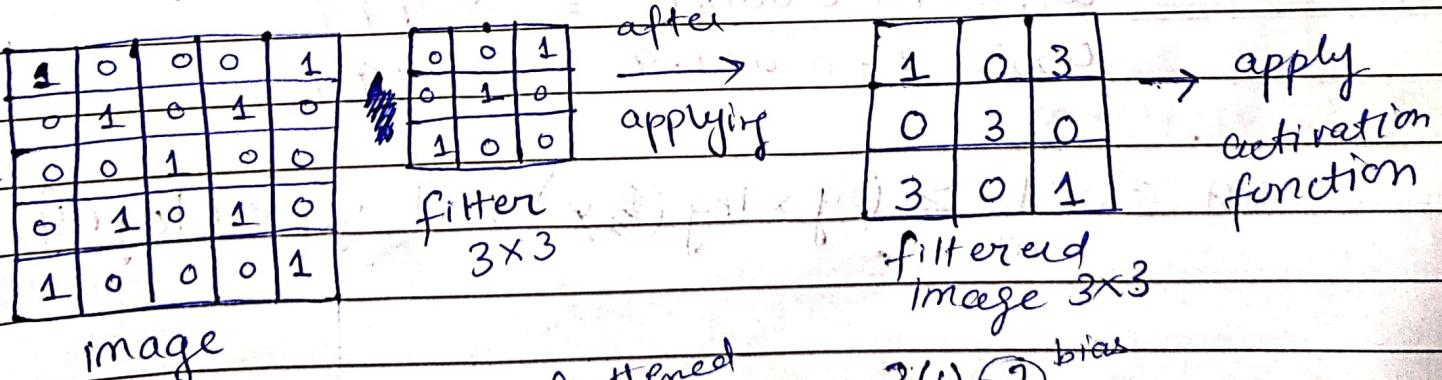
$$\text{new cost} = \frac{1}{m} \sum (y \times \log(\hat{y}) + (1-y) \times \log(\hat{y}))$$

Date: _____



convolutional neural network:

- ① choose a filter
- ② apply filter on image
- ③ apply suitable activation function (ReLU : convert negs to zero)
- ④ apply max pool filter
- ⑤ flatten the image
- ⑥ apply normal neural network.



Date: _____

* Recurrent Neural Network !

① It's like a normal neural network except there's a thing called feedback loop embedded in it. This loop makes considering previous data in the output of a new data point possible for the neural network.

② The feedback loop picks the output of hidden layers and feeds it right back to the first hidden layer.

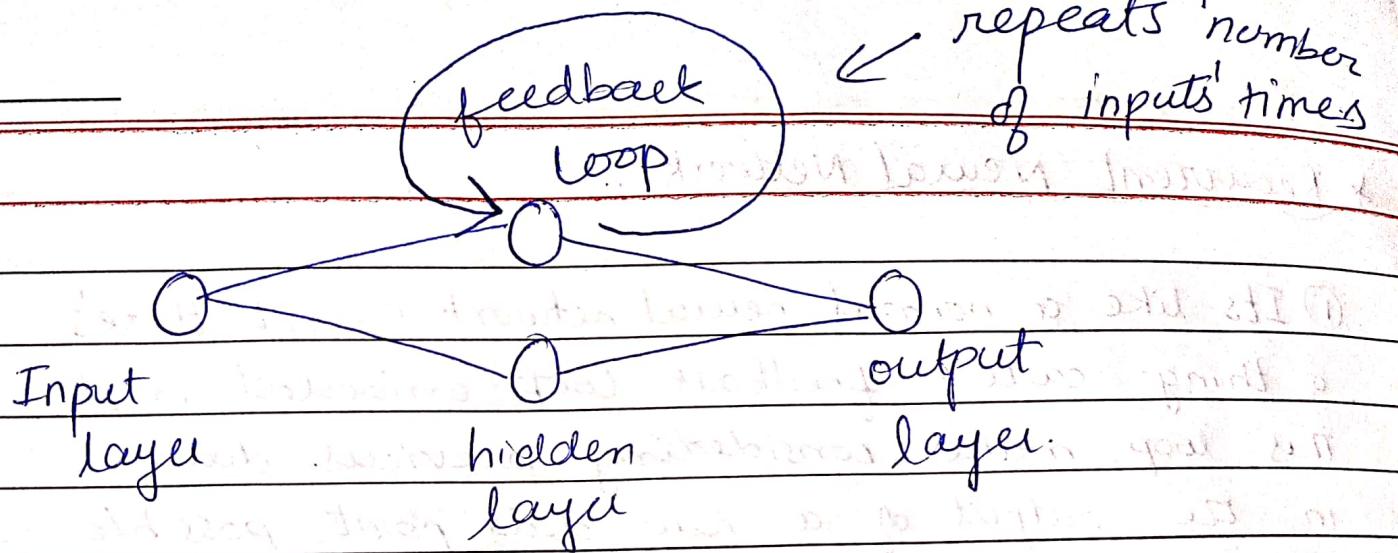
③ This could cause two types of problems;

- Exploding gradient problem
- Vanishing gradient problem

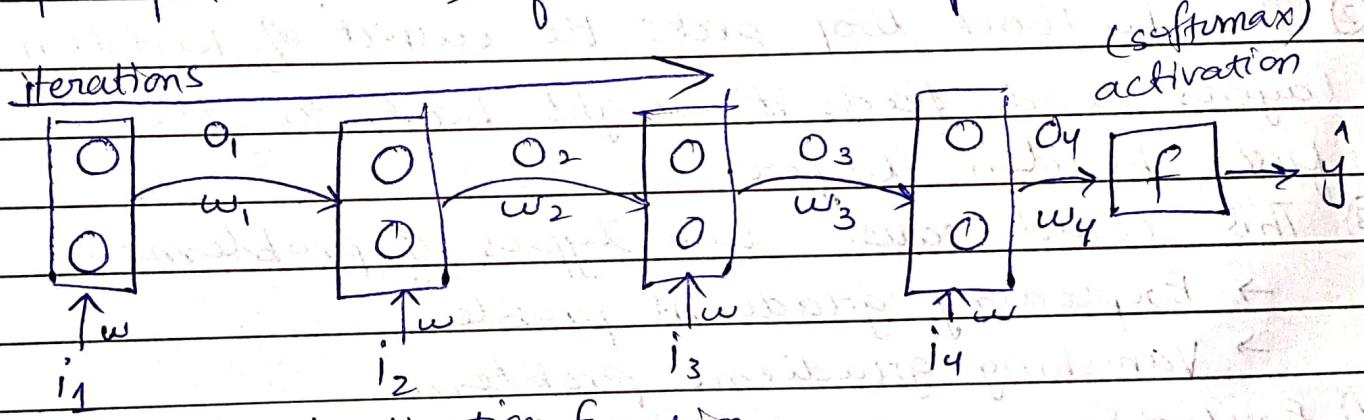
④ Exploding gradient problem occurs when the weight associated with the feedback loop grows exponentially. Its value becomes so large that it's nearly impossible to converge in the gradient descent.

⑤ Vanishing gradient problem occurs when the weight associated with the feedback loop grows exponentially 'small' (near to zero). Its value is so small that we fail to reach the gradient optimal as the convergence is too slow.

Date: _____



* if inputs are 4, feedback loop runs 4 times.



$$o_1 = f(i_1 w)$$

$$o_2 = f(i_2 w + o_1 w_1)$$

$$o_3 = f(i_3 w + o_2 w_2)$$

$$o_4 = f(i_4 w + o_3 w_3)$$

$$\hat{y} = \text{softmax}(o_4)$$