

Agenda (ANN):

- 1/2/s {
- ① Why DL / what is DL, what is ANN.
 - ② What is a Neuron, why ANN / DL
 - ③ What is MLP
 - ④ Forward / Backward Propagation.
 - ⑤ Training steps of ANN.

4/5/s {

implement
NN in code.

Tensorflow | Keras

Dropout, Regularization, Batch Normalization, losses,
Callbacks, etc., optimizers

7/8/9 {

Code implementation, hyperparameter tuning,

Optimization in NN.

{ Explainability / Autoencoding (Aeoirsing).

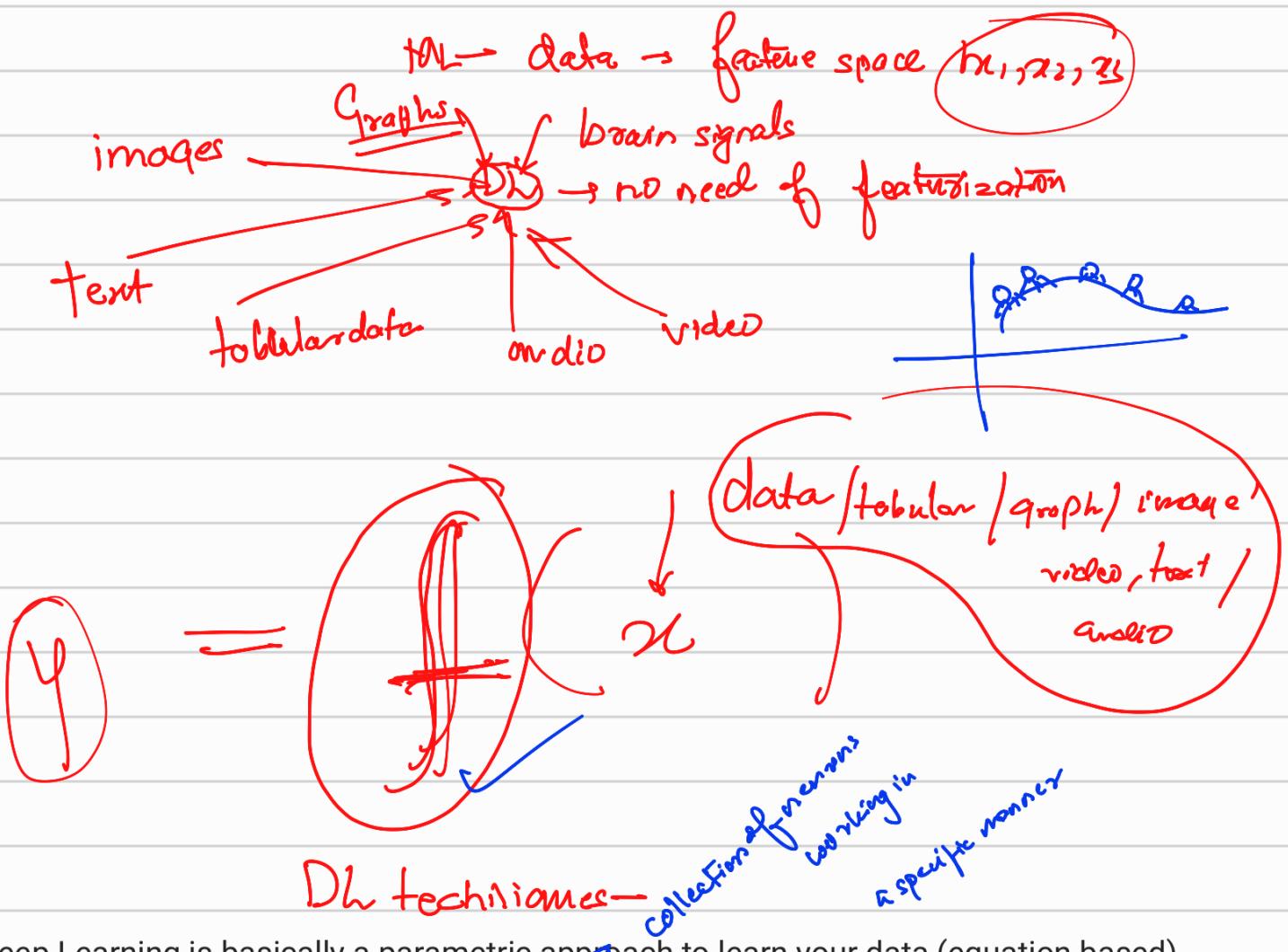
What is Deep learning, Why DL, v/s ML.

DL
E

DL

Deep learning is basically a field of study where we use something known as neurons for training the data.

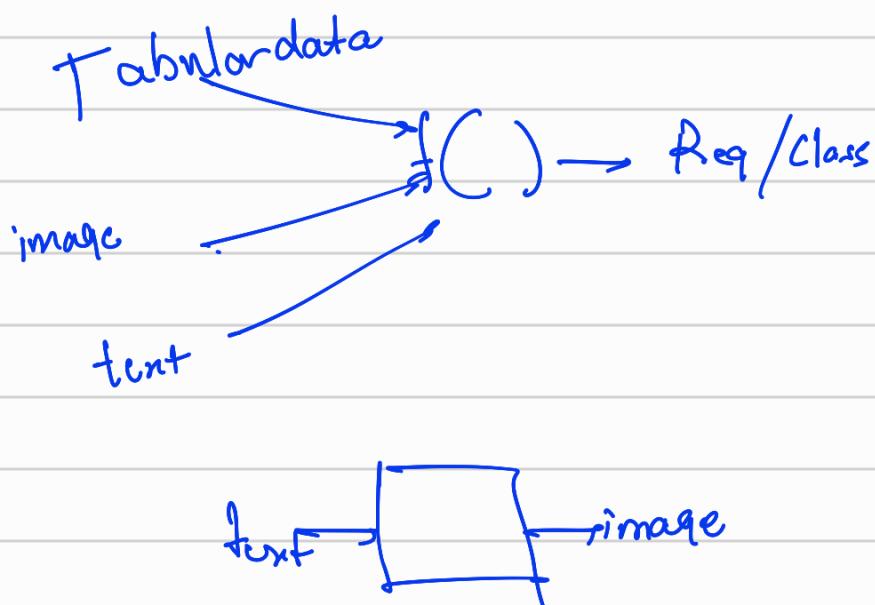
Automated feature engineering.

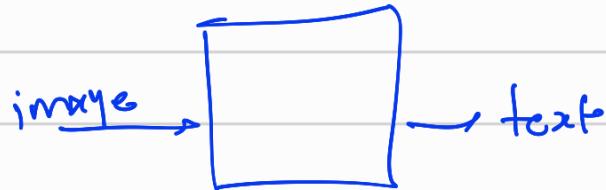


Deep Learning is basically a parametric approach to learn your data (equation based)

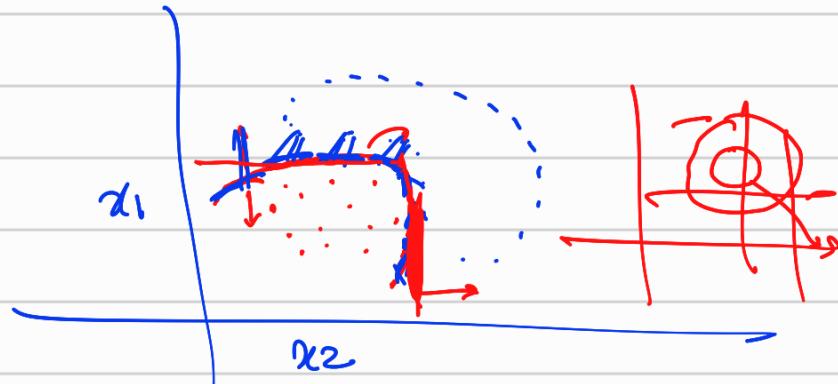
- Artificial Neural Networks - NN is designed to extract information from tabular data
- Convolution Neural Networks - NN is designed to extract information from image data
- Recurrent Neural Networks (LSTM/GRU? Transformers) - NN is designed to extract information from sequence data (like text, audio, videos)
- Graph Neural Network for Graphs (molecule)

Multi Modality:

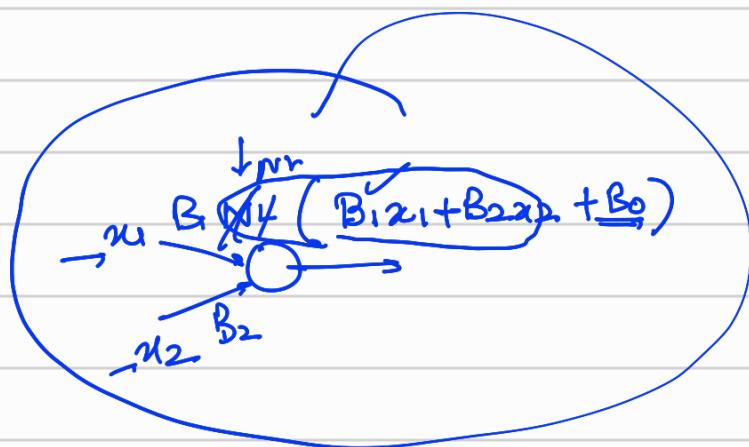




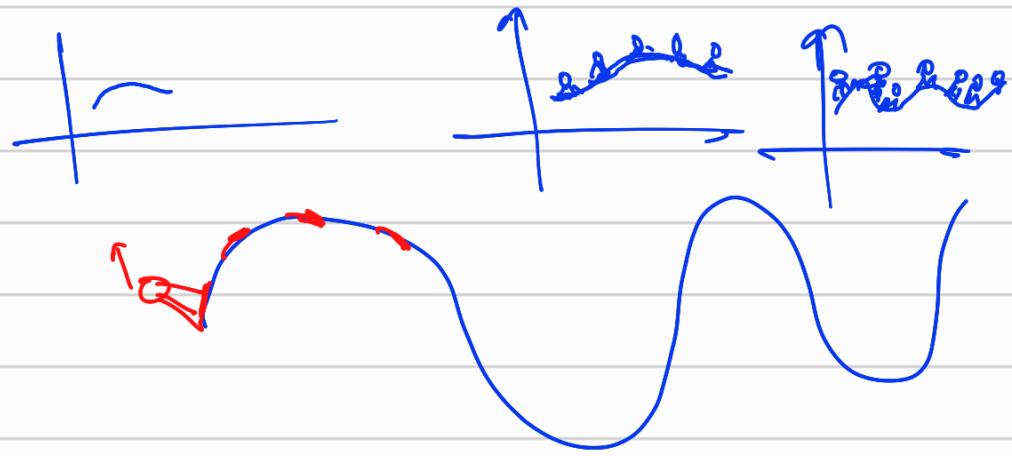
what is Neural network:



Can i fit a decision boundary along the red curve which can be thought of multiple connected curves, and each curve has the independence to move up or down and has the liberty to change o the amount of curvature



$\frac{d}{dx} f(x) = \text{some value}$



22:55 pm Break

15-20 mins
is important

Proving why non-linearity

- ② Why the non-linearity is absolutely vital
- ① Gradient descent → underlying algo for training them.

Input layer

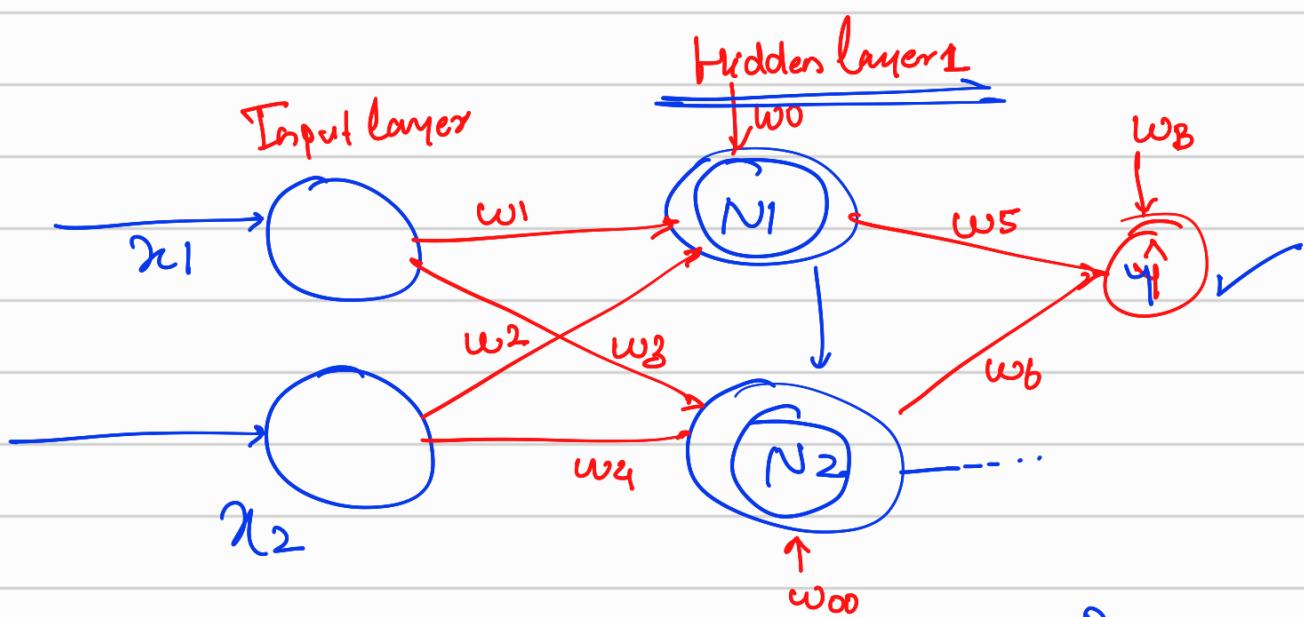


Diagram illustrating the activation function (tanh) and bias terms:

$$N_1 = \tanh(w_1x_1 + w_2x_2 + w_{01})$$

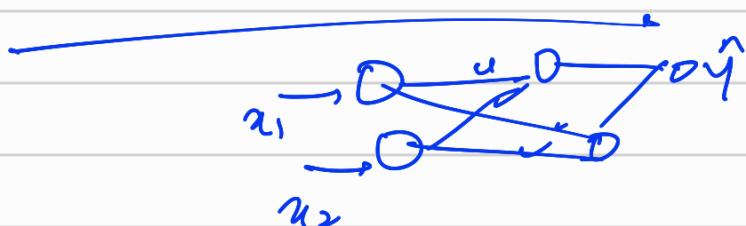
$$N_2 = \tanh(w_3x_1 + w_4x_2 + w_{02})$$

$$\hat{y} = w_5(N_1) + w_6(N_2) + w_B$$

$$\hat{y} = \frac{w_5 w_1 x_1 + w_5 w_2 x_2 + w_6 w_3 x_1 + w_6 w_4 x_2}{w_B w_{00} + w_B}$$

$$\hat{y} = x_1(w_5 w_1 + w_6 w_2) + x_2(w_5 w_2 + w_6 w_4) + w_B + w_B w_{00}$$

$\hat{y} = x_1 A + x_2 B + C$



glorot uniform, glorot normal

How do we train such a network of neurons.

1. Initialise the architecture (define the number of hidden layers, neurons etc.)
2. Randomly initialise the weights
3. Given the weights calculate the loss. → **Forward Propogation**
4. Given the loss, update the weights - Gradient Descent Algorithm + Chain Rule → **BP**
5. using the updated weights go to step 3 and keep on doing it till you stabilise your training loss to the min possible (min Val loss)

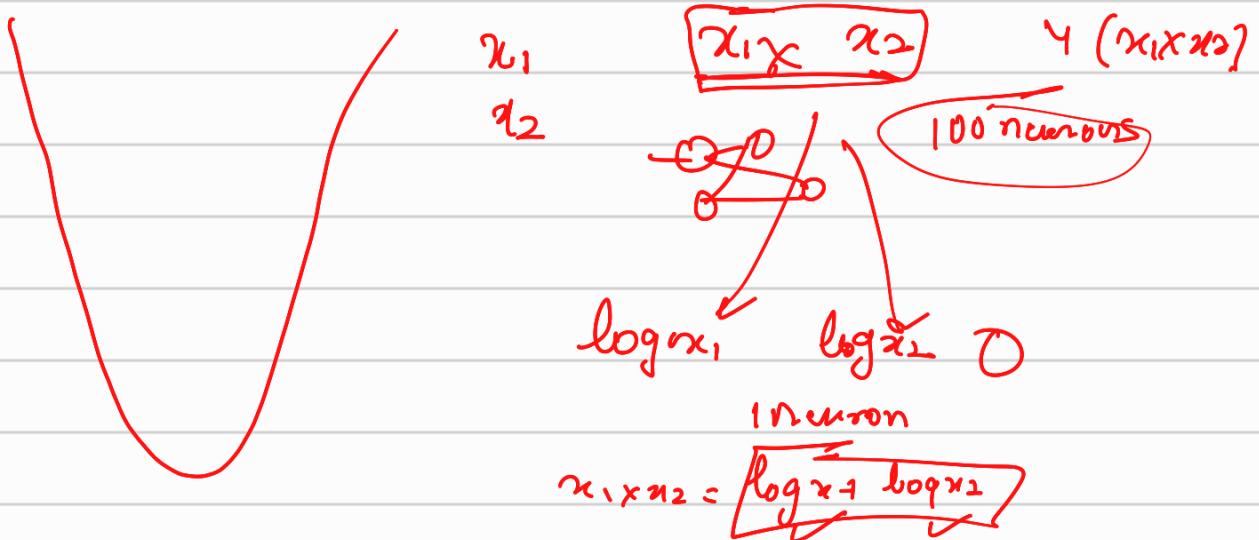
$$\text{loss} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{loss} = f(\text{? ? ?})$$

(bias & weights), label, input(s) ↓

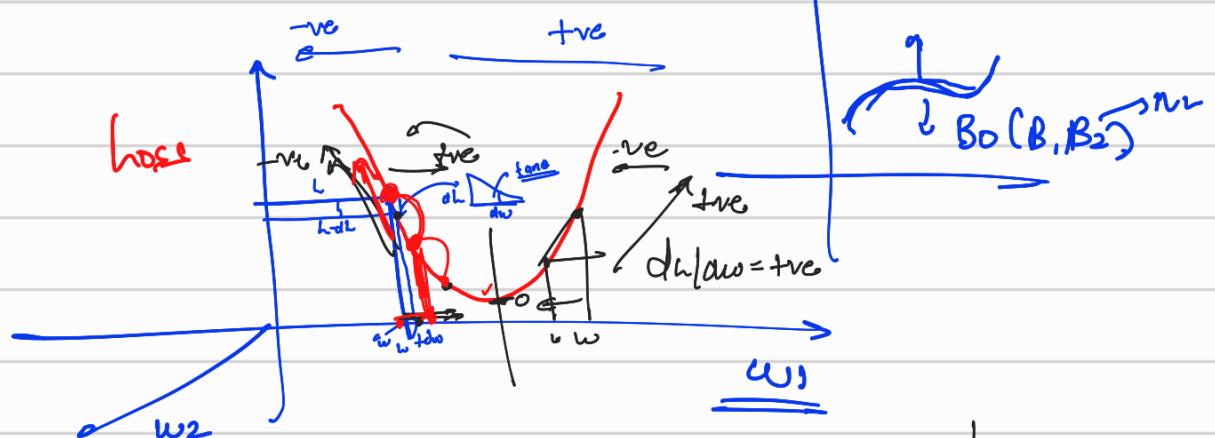
minimise the loss

Gradient Descent Algorithm and Chain Rule:



✓ 13/12/2022

Gradient Descent Algorithm:



$$\underline{dh/dw_1 = -ve}$$

$$y = f(x)$$

$$\hat{y} = f(x) + \underline{\underline{\epsilon}}_2$$

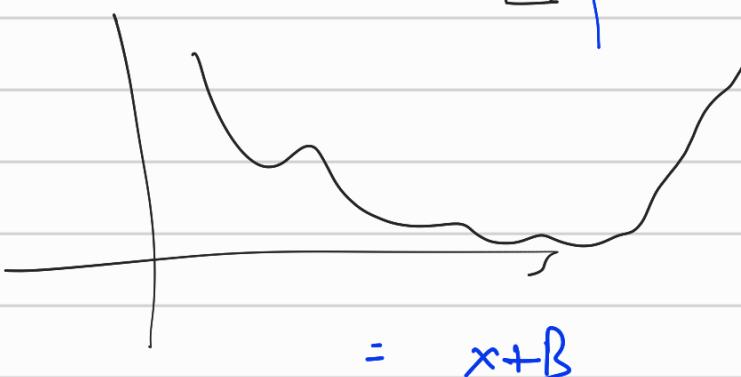
irreducible error

update rule for mg GD for weight w_1

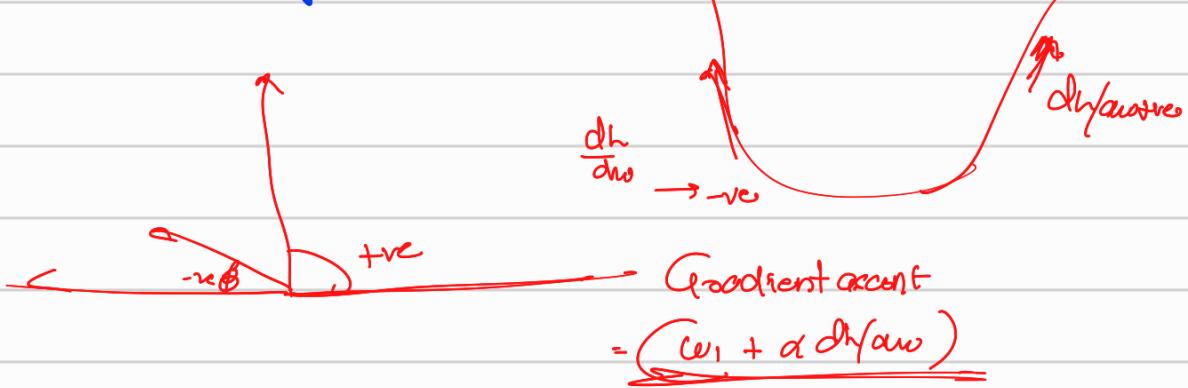
for all w_i is independently and bis

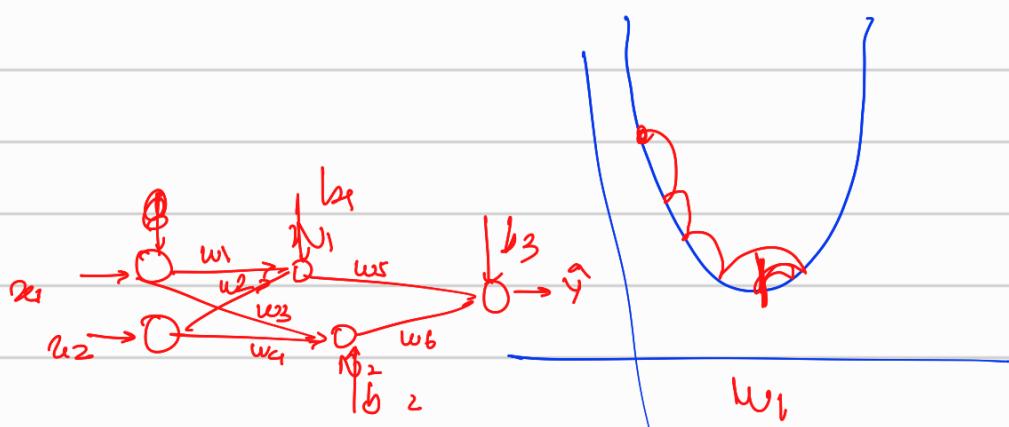
$$\Rightarrow \underline{\underline{w_1(t)}} = \underline{\underline{w_1(t-1)}} - \alpha \left[\frac{dh}{dw_1(t-1)} \right]$$

learning rate
partial derivative



$$= x + \underline{\underline{B}}$$





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

$$\frac{\partial L}{\partial w_1}$$

$$L = \sum_{i=1}^n \frac{1}{n} (y_i - (w_5 u_1 + w_6 u_2 + b_3))^2$$

$$\frac{\partial L}{\partial w_2}$$

$$\frac{\partial L}{\partial w_3}, \dots, \frac{\partial L}{\partial w_6}$$

$$L = \sum_{i=1}^n \frac{1}{n} (y_i - \tanh(w_5 u_1 + w_6 u_2 + b_3))$$

$w_5 + w_6 + b_3$

Chain Rule

$$\begin{aligned} y &= 5z \\ z &= 5x \\ x &= 5m \end{aligned}$$

$$\frac{dy}{dz} = 5 \quad y = f(z)$$

$$\frac{dz}{dx} = 5 \quad z = f(x)$$

$$\frac{df}{dm} = 5 \quad x = f(m)$$

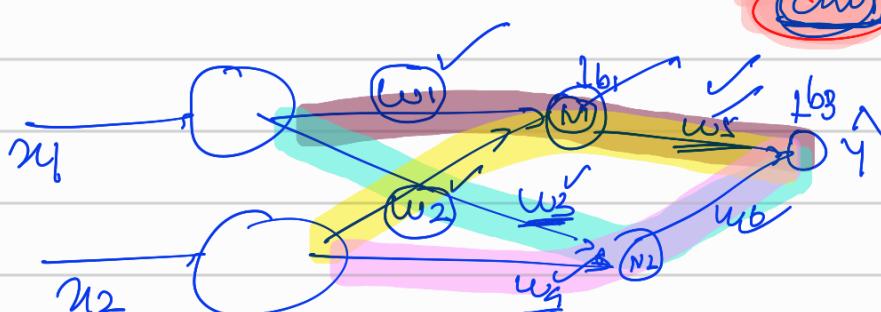
$$\frac{dy}{dx} = \frac{dy}{dz} \times \frac{dz}{dx} \rightarrow \text{Chain Rule}$$

$y \approx \text{implicit } f(x)$

$$\begin{aligned} y &= 5 \times 5x \\ y &= 25x \\ \frac{dy}{dx} &= 25 \end{aligned}$$

Replacement

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial w} \times \frac{\partial w}{\partial w}$$



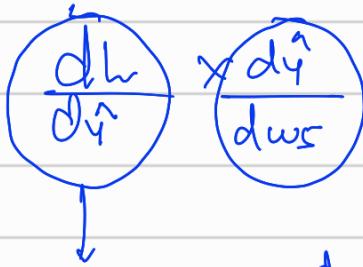
$$\frac{dy}{dx} = \frac{\partial y}{\partial w_1} \times \frac{\partial w_1}{\partial u_1} + \frac{\partial y}{\partial w_2} \times \frac{\partial w_2}{\partial u_1} + \dots$$

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

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial w_1}$$

$$\frac{dh}{dw_1} \dots \frac{dh}{dw_5} \quad \frac{dh}{dw_6}$$

$$\hat{y} = \underline{\underline{w_5 N_1 + w_6 N_2 + b_3}}$$



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

$$\frac{dh}{dy} = \sum_{i=1}^n \frac{2}{n} (y_i - \hat{y}) \times (-1)$$

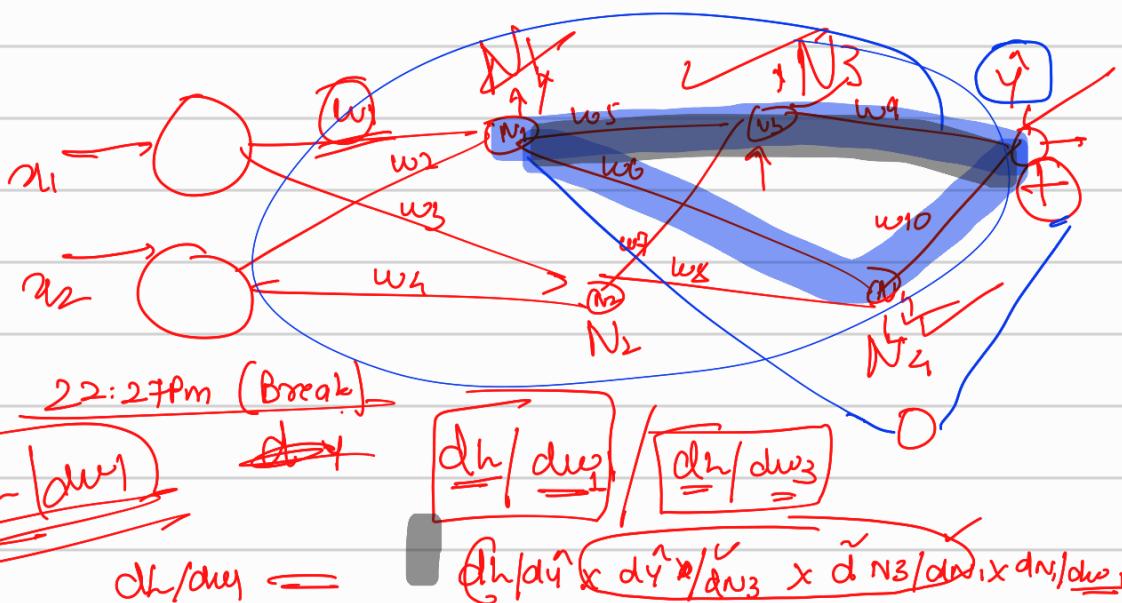
$$\hat{y} = \underline{\underline{w_5 N_1 + w_6 N_2 + b_3}}$$

$$\frac{dy}{dw_5} = \underline{\underline{N_1 \times}}$$

$$\frac{dh}{dw_5} = \sum_{i=1}^n \frac{2}{n} (y_i - \hat{y}) \times N_1 y_i$$

$$\frac{dh}{dw_6} = \sum_{i=1}^n \frac{2}{n} (y_i - \hat{y}) \times N_2 y_i$$

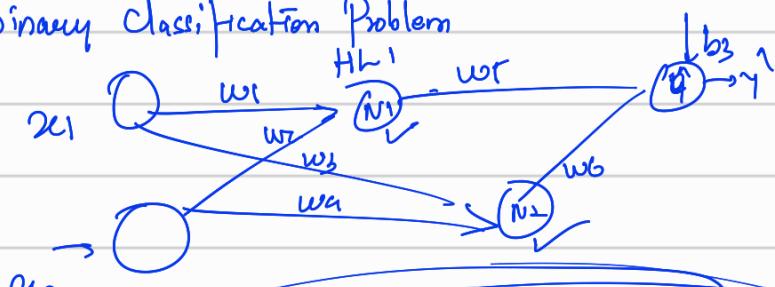
$$\frac{dh}{dw_1} =$$



$$+ \frac{dy/d\hat{u}}{d\hat{u}/du_1} \times \frac{d\hat{u}}{du_1} \times \frac{dN_1}{du_1} \times \frac{dN_1}{du_1}$$

$$\text{d}m/dv =$$

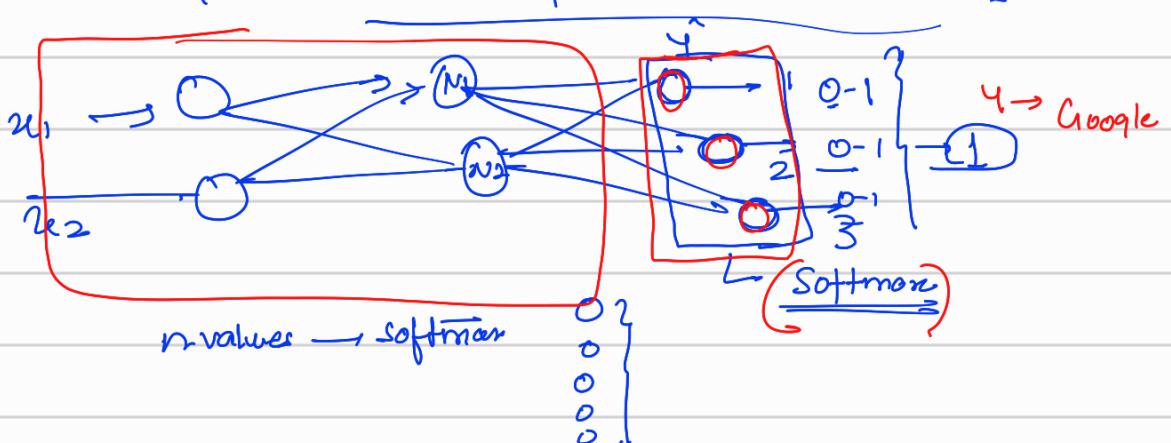
Binary Classification Problem



$$\hat{y} = \sigma(w_5 n_1 + w_6 n_2 + b_3)$$

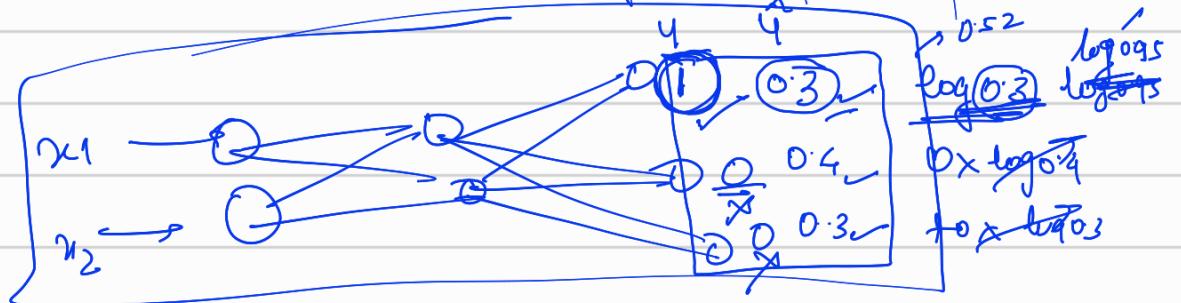
$$\text{loss} = - \sum_{i=1}^n [y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)]$$

Classification → Multi Problems (Three Classes)



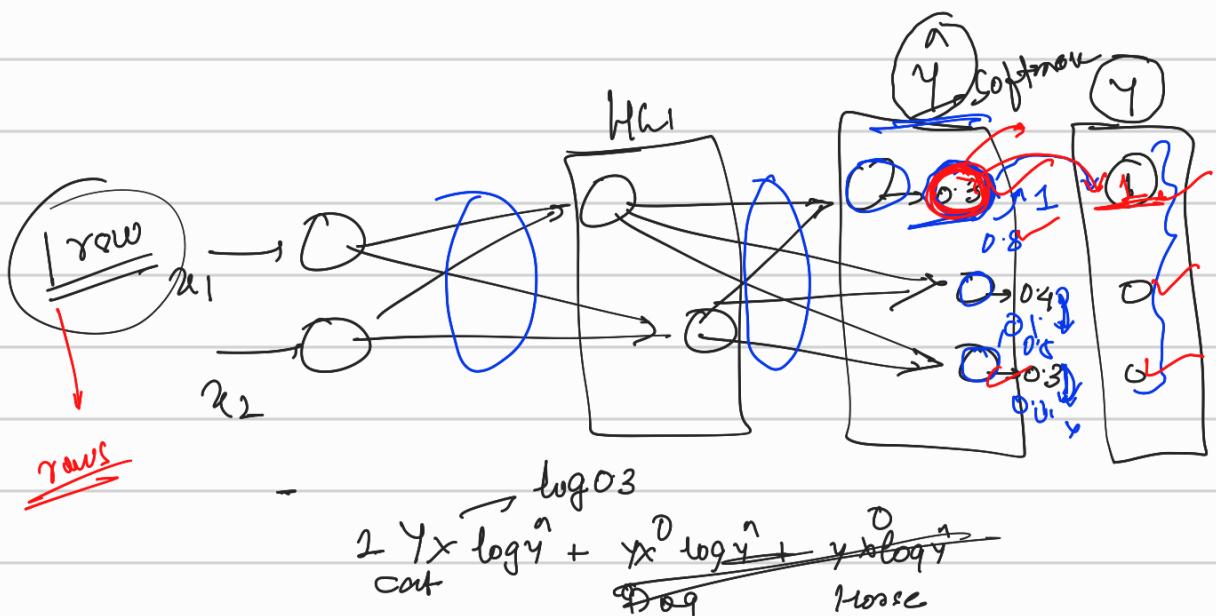
$$\text{loss} = \text{log loss} = \sum_{i=1}^n \frac{1}{n} \sum_{j=1}^K y_j \log g_i^n$$

	x_1	x_2	Google	Micro	Walmart
	✓	✓	✓	✓	✓
	✓	✓	0	1	0
	✓	0	0	0	1
	0	0	1	1	1



$$\sum_{i=1}^n \frac{1}{n} - \sum_{k=1}^3 y_k \times \log y = -\log 0.3 = 0.52 = 0.022$$

	hen	Height	Cat	Dog	Horse	
1	✓	✓	0	0	0	1 = 1000
2	✓	✓	0	1	0	2 = 5000
3	✓	✓	0	0	1	3 = 3000



$$2 Y \times \log 0.3 + Y^0 \times \log 0.1 + Y^0 \times \log 0.6$$

$$-\log 0.3 = 0.52$$

Rey

HL

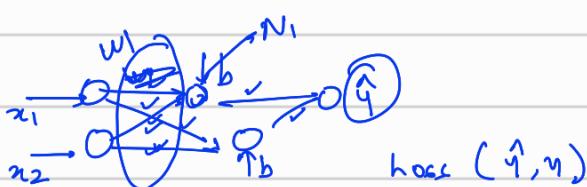
Output layer (\rightarrow activation) loss = mSE

BC

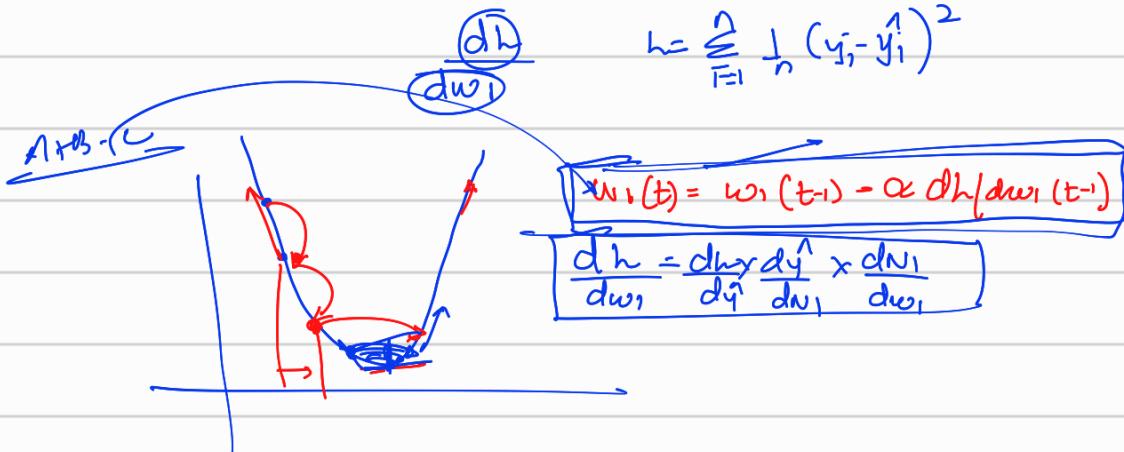
Output layer (sigmoid) loss = binom log loss.

MC output layer (softmax) - loss (softmax loss) $Y = [1 \ 2 \ 0]$

Summary of previous class



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

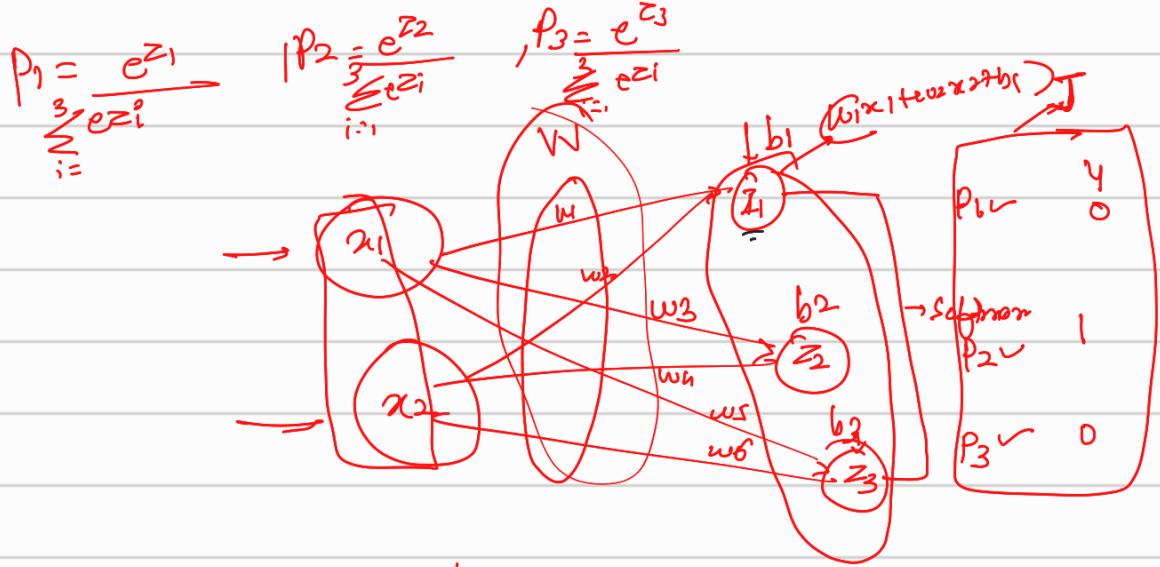


- ① How do we create 1 layer of neural network by basic python functions.

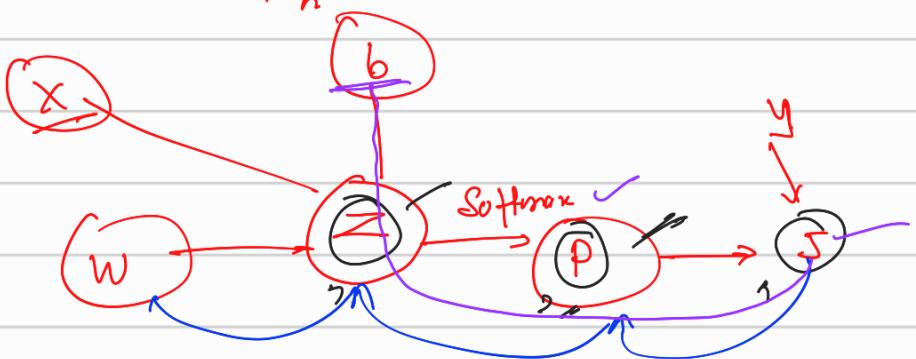


FP, BP

- ② How do we code our first nn from scratch.
- ③ Activation functions



$$\text{Loss } (J) = - \sum_{i=1}^n y_i \log P_i$$



Using Chain Rule

$$\frac{dJ}{dw} = \frac{dJ}{dp} \times \frac{dp}{dz} \times \frac{dz}{dw}$$

denoting $\frac{dJ}{dw} = \frac{dJ}{dw}$

denoting $\frac{dJ}{dz} = dz$

$$\frac{dz}{dw} = \frac{z - w_x + b}{dz/dw} = x$$

$$dz = \frac{dJ}{dz} = [P_1 \rightarrow P_2, P_3] - [0, 1, 0]$$

$z = wx + b$ $\frac{dJ}{dB} = 1$

$$dw = (p - y) \times x$$

$w = w - \alpha dw$

$$\frac{dJ}{dB} = \frac{dJ}{dp} \times \frac{dp}{dz} \times \frac{dz}{dB}$$

$$\frac{dJ}{dB} = dz \times 1$$

Probabilities:

$$[P_1, P_2, P_3] - [0, 1, 0] = 1$$

0
1
2

$$= P_{1=0} P_{2=1} P_{3=0}$$

$P_{\text{prob}}[y] = \text{prob}[y] - 1$

$\text{d}z$ Pools

$$\text{d}z[\text{range}(m), y] = \text{d}z[\text{range}(m), y] - 1$$

$$\text{d}z = \begin{bmatrix} 0.2 & 0.3 & 0.5 \end{bmatrix} - \begin{bmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

$y=2$

$y=1$

$y=0$

$\begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$

$\begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$

$\underline{\text{Probs}}[2] = 1$

Pools

$$\begin{bmatrix} 0.2 & 0.3 & 0.5 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 & 0.3 & 0.5 \end{bmatrix}$$

$$= -0.8 \quad 0.3 \quad 0.5$$

$$\begin{bmatrix} 0.2 & 0.3 & 0.5 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 & 0.3 & 0.5 \end{bmatrix}$$

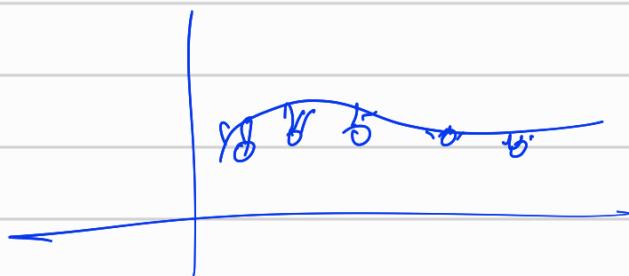
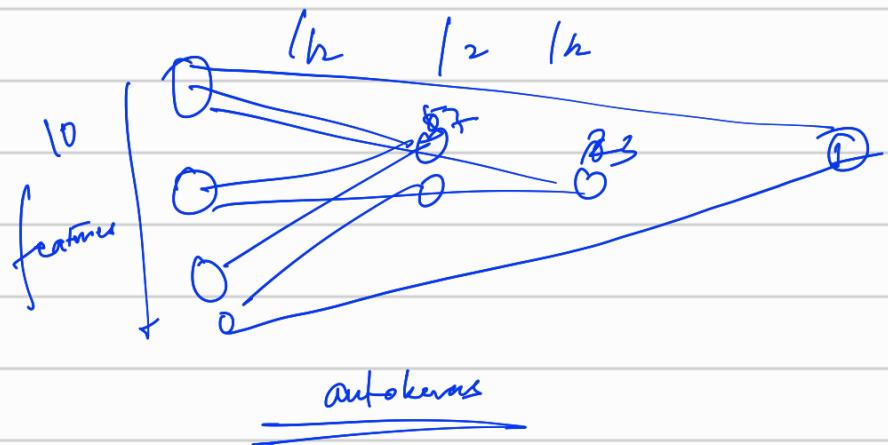
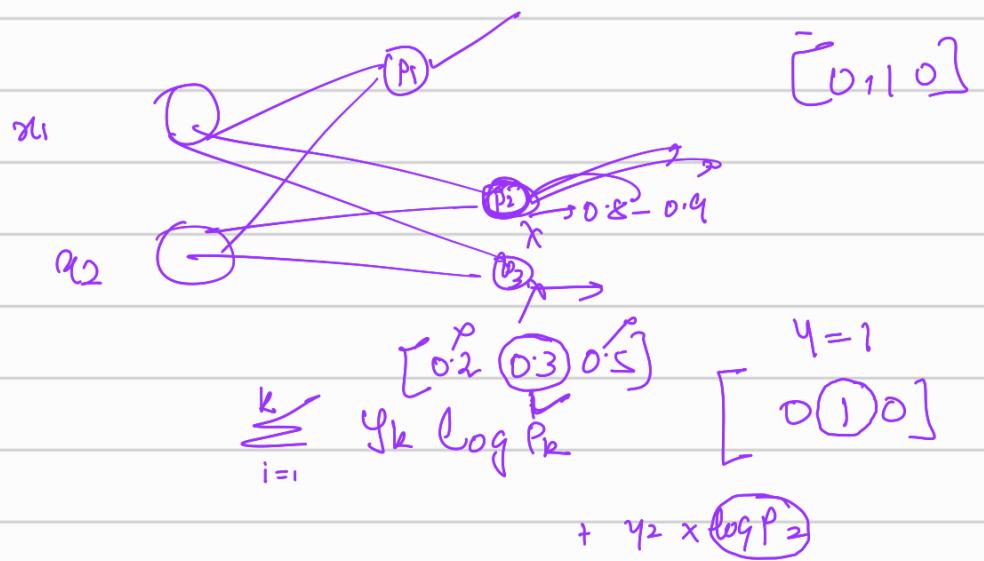
$$\begin{bmatrix} 0.2 & 0.3 & 0.5 \end{bmatrix}$$

$$\text{d}z = \begin{bmatrix} 0.2 & 0.3 & 0.5 \end{bmatrix} - \begin{bmatrix} 0 & 1 & 2 \\ 0 & 0 & 0 \end{bmatrix}$$

22:22 → 22:32 pm

$Bc \rightarrow O_2$

$Mg \rightarrow$ caten Dog, Horse



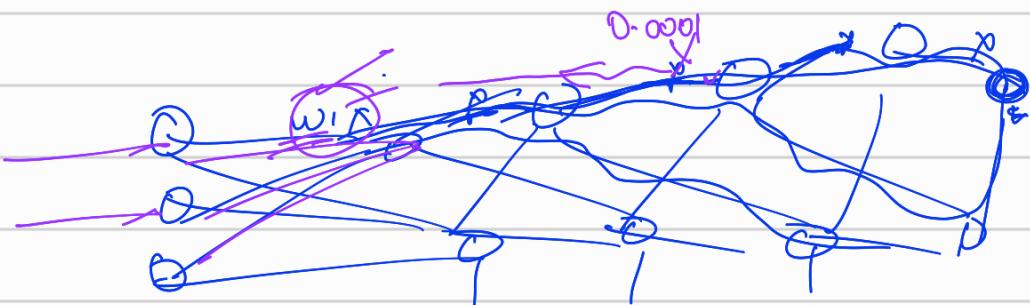
Activation function: what would be a choice of a good activation function.



Or \tilde{o}

Chain rule $x \times x \times$

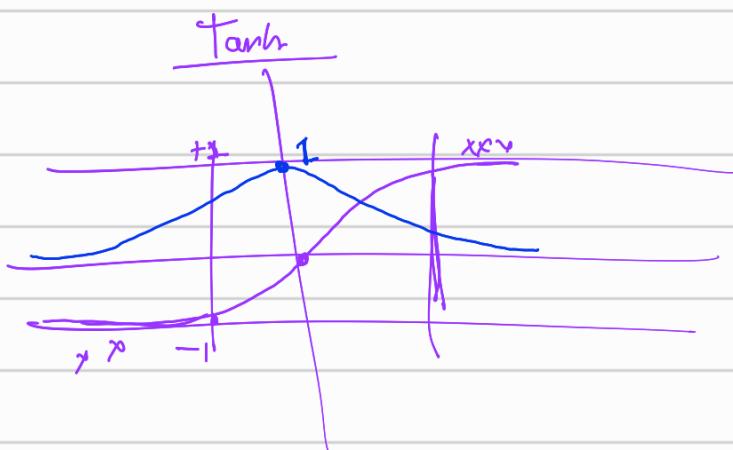
Tanh



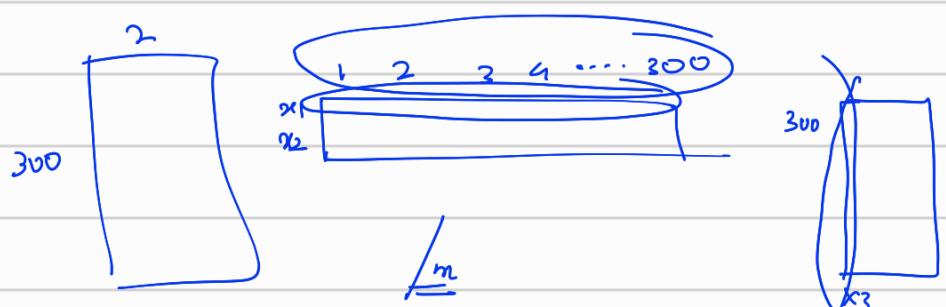
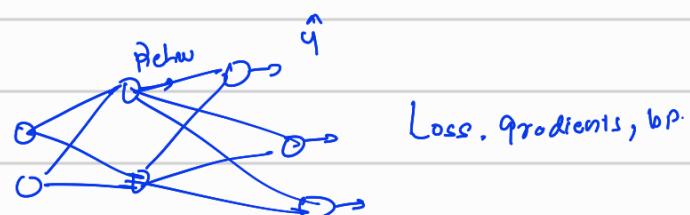
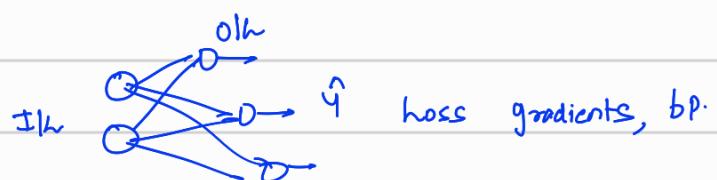
$$\frac{dJ}{dw_1} = \text{Value}$$

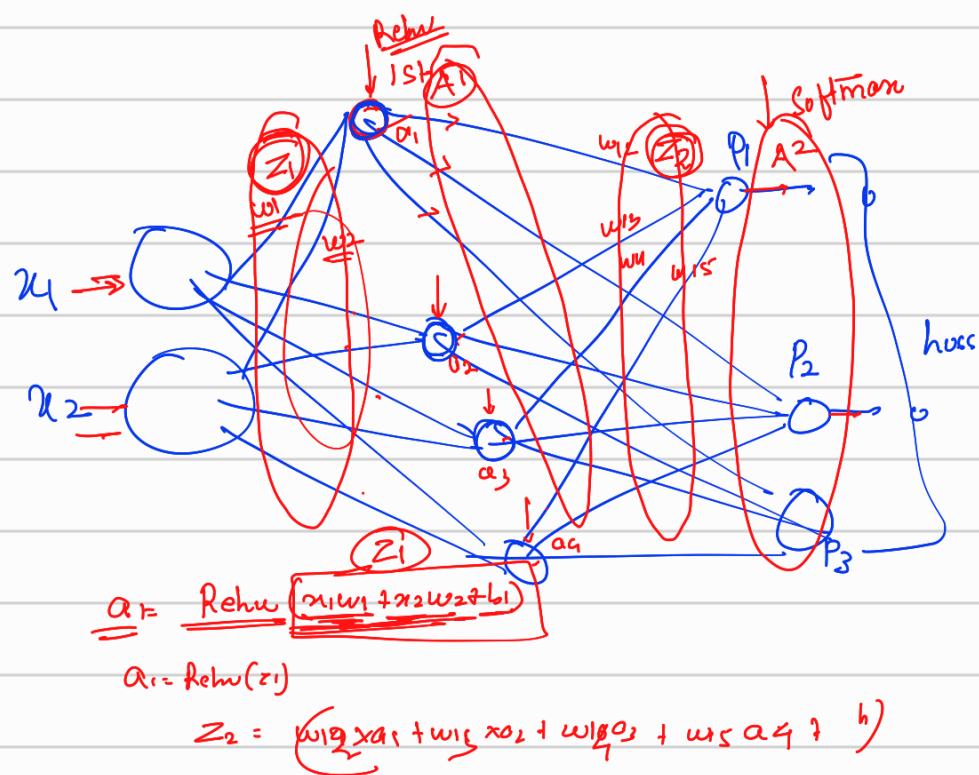
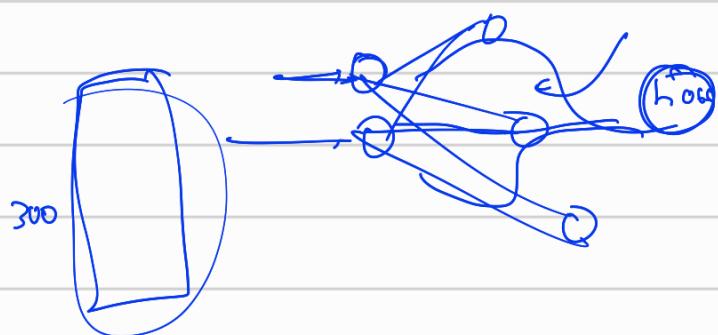
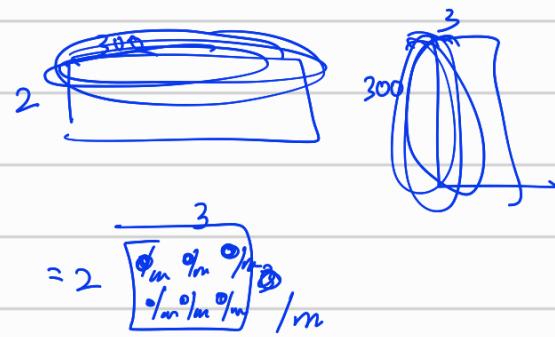
Exploding gradient

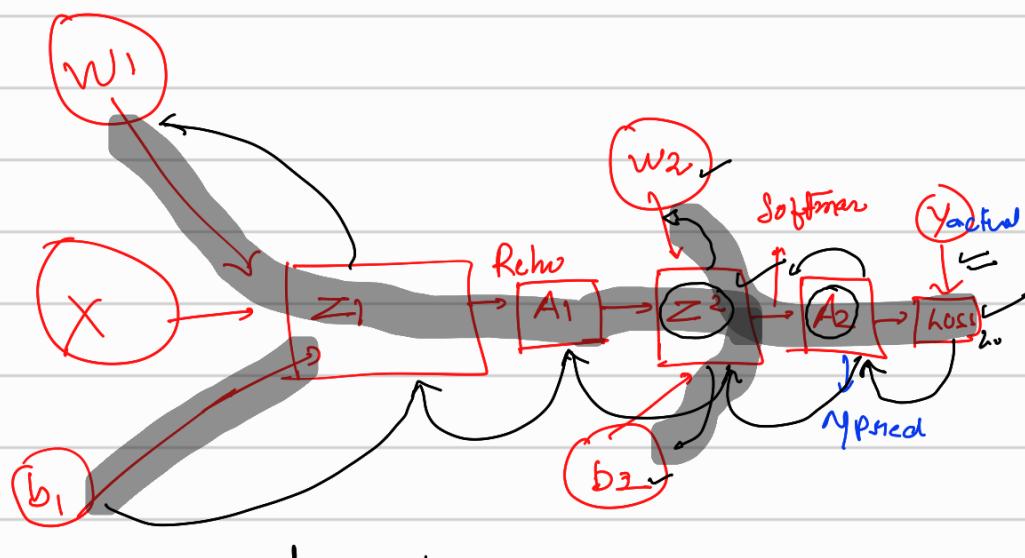
$$w_1 = w_1 - \alpha \frac{dJ}{dw_1}$$



0
0

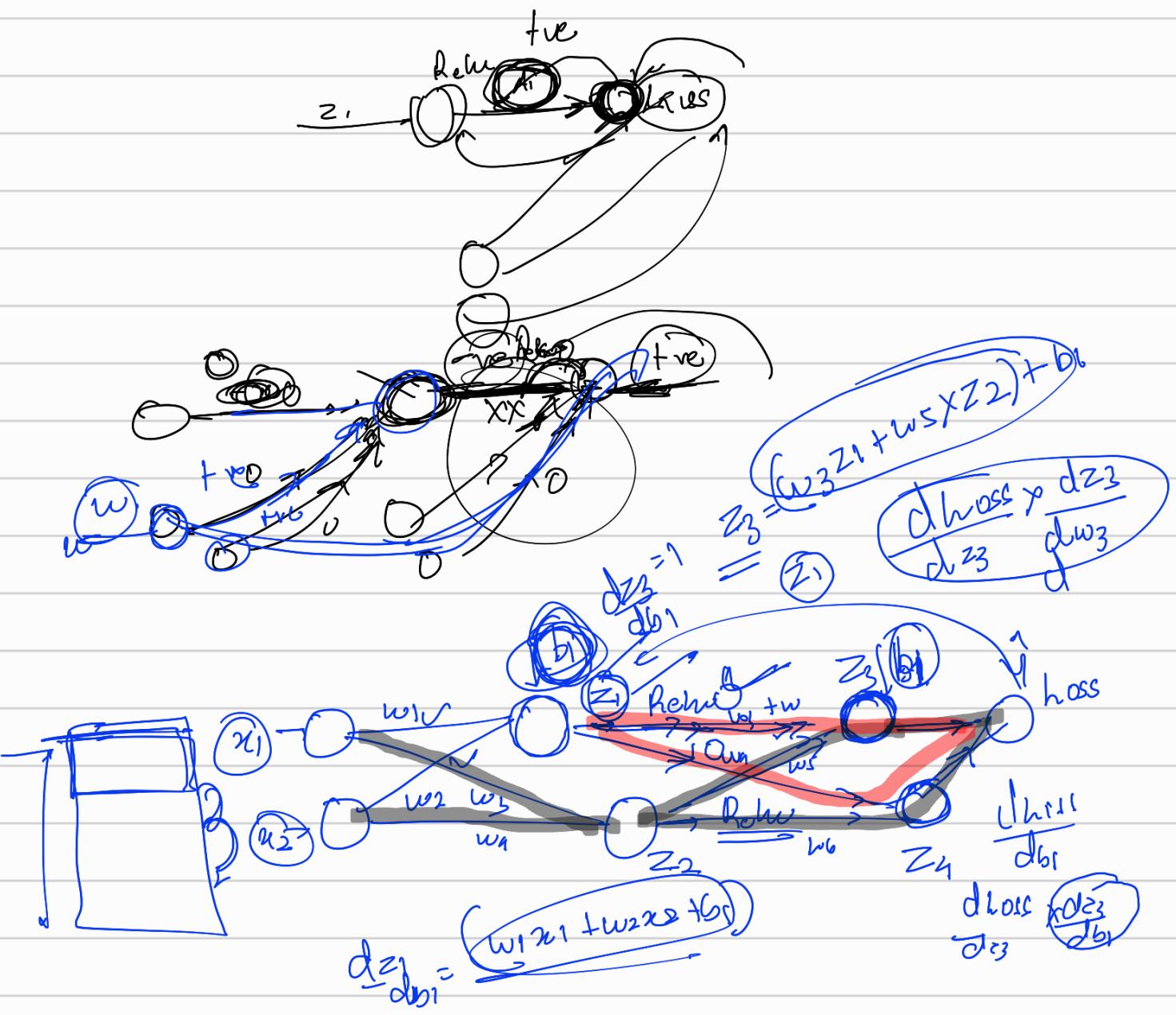


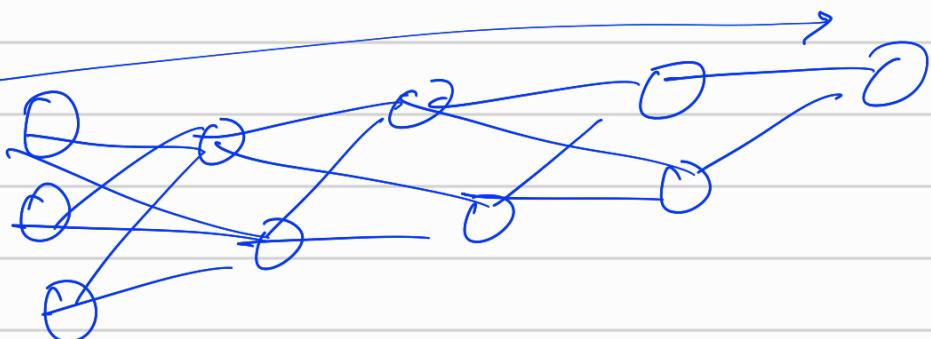
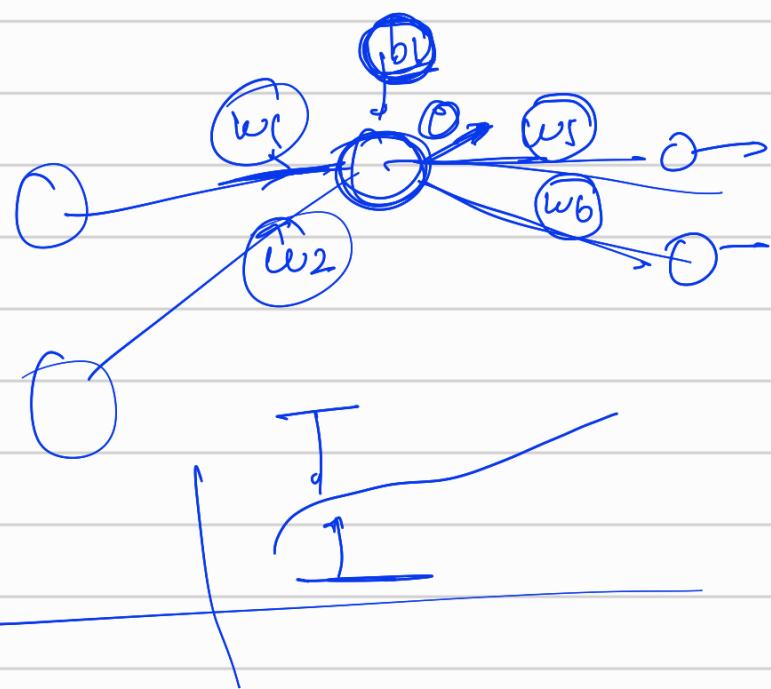
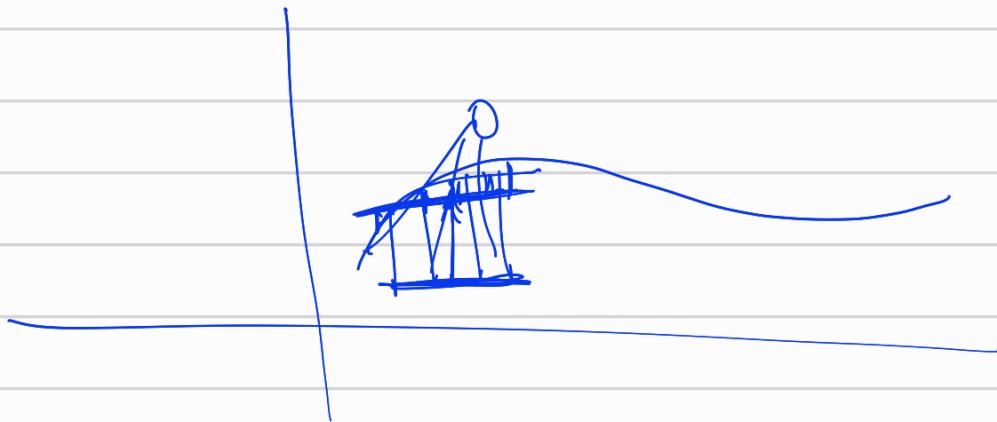




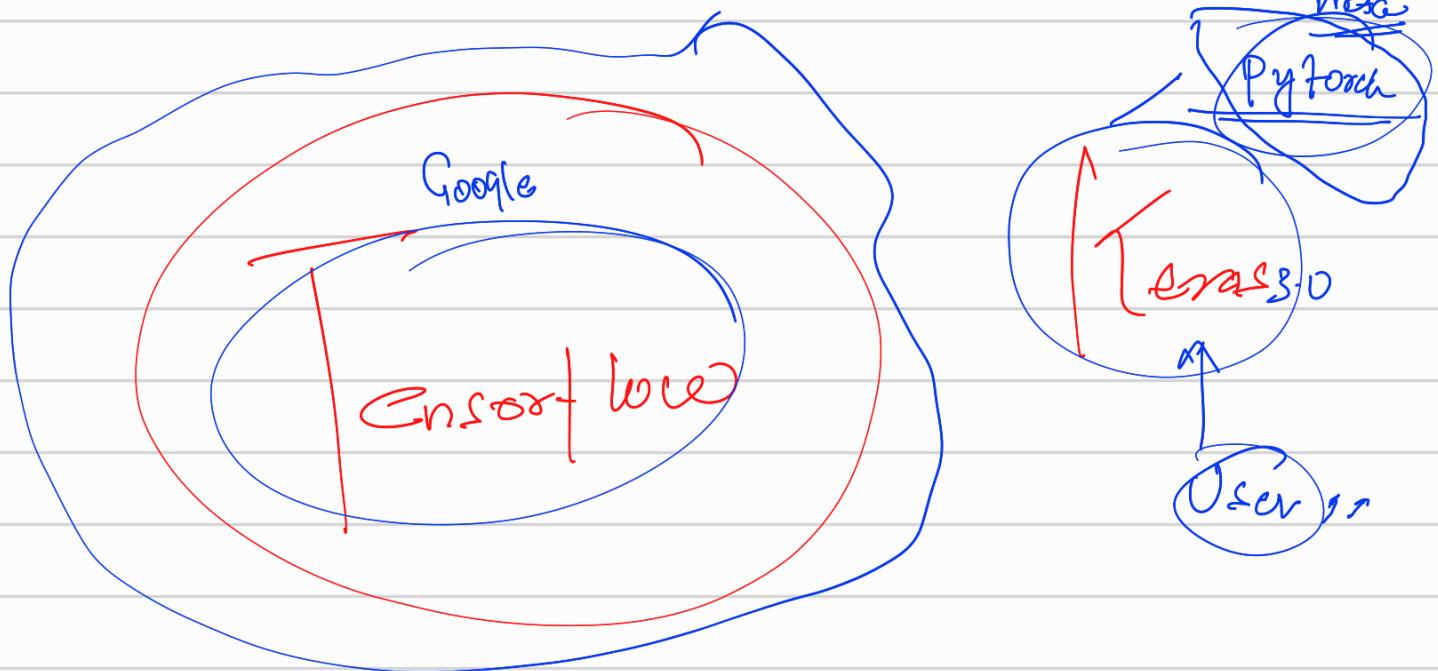
$$\frac{dL}{dA_2}$$

$$\frac{dL}{dz_2} = dz_2 = (p - y)$$





- ① Tensorflow & Keras
- ② Dropout, Batch normalization, Regularization, Optimizers, Callbacks
Hyperparameter Tuning.



Class : $\{0: 5, 1: \underline{\underline{70}}\}$

$$\sum_{i=1}^n \left(y_i \log y_i + (1-y_i) \log (1-y_i) \right) \times 5$$

\downarrow

$$f(x) = 0.2, f(x) = 0.05, f(x) = 0.7, f(x) = 0.3, f(x) = 0.9$$

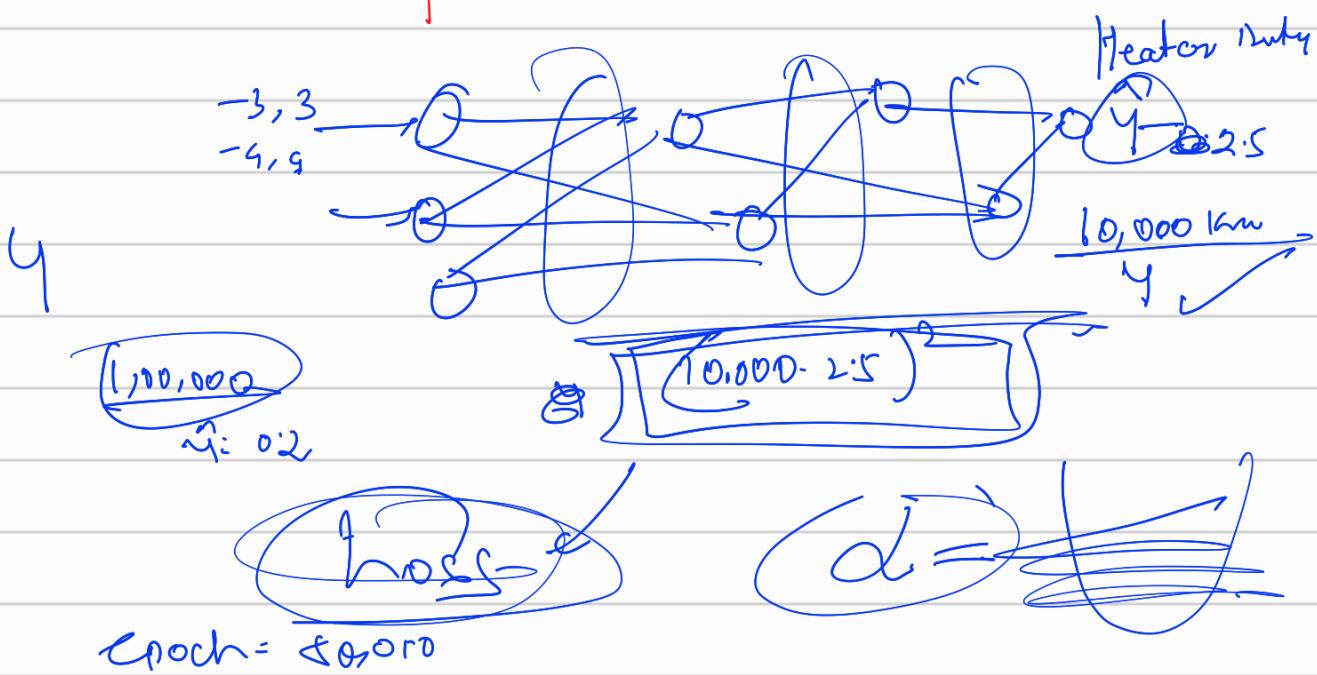
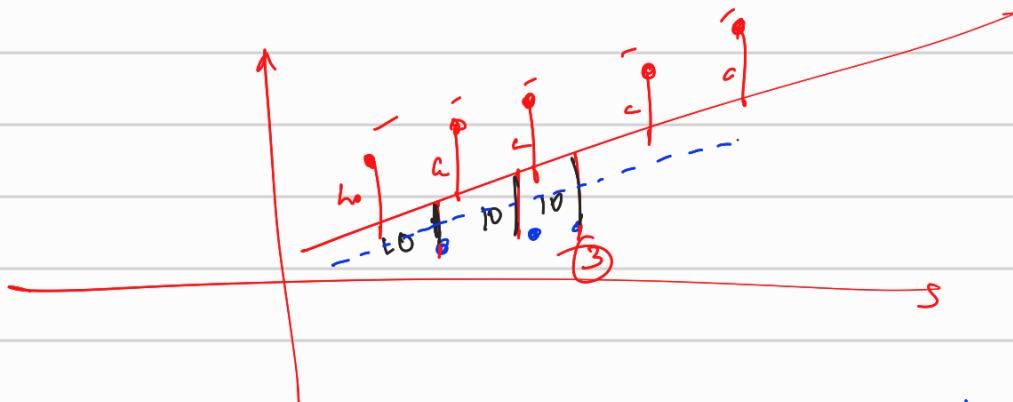
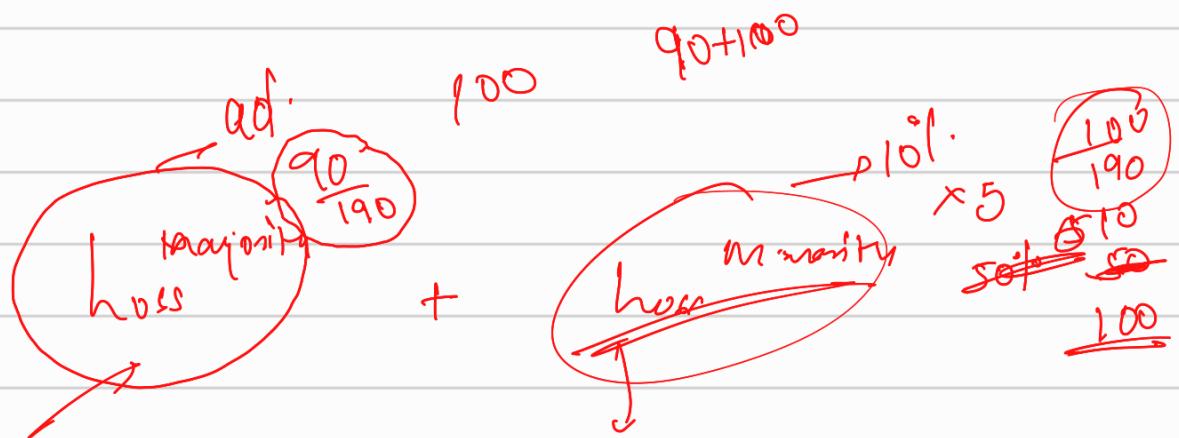
$$\frac{1}{n} \sum_{i=1}^n$$

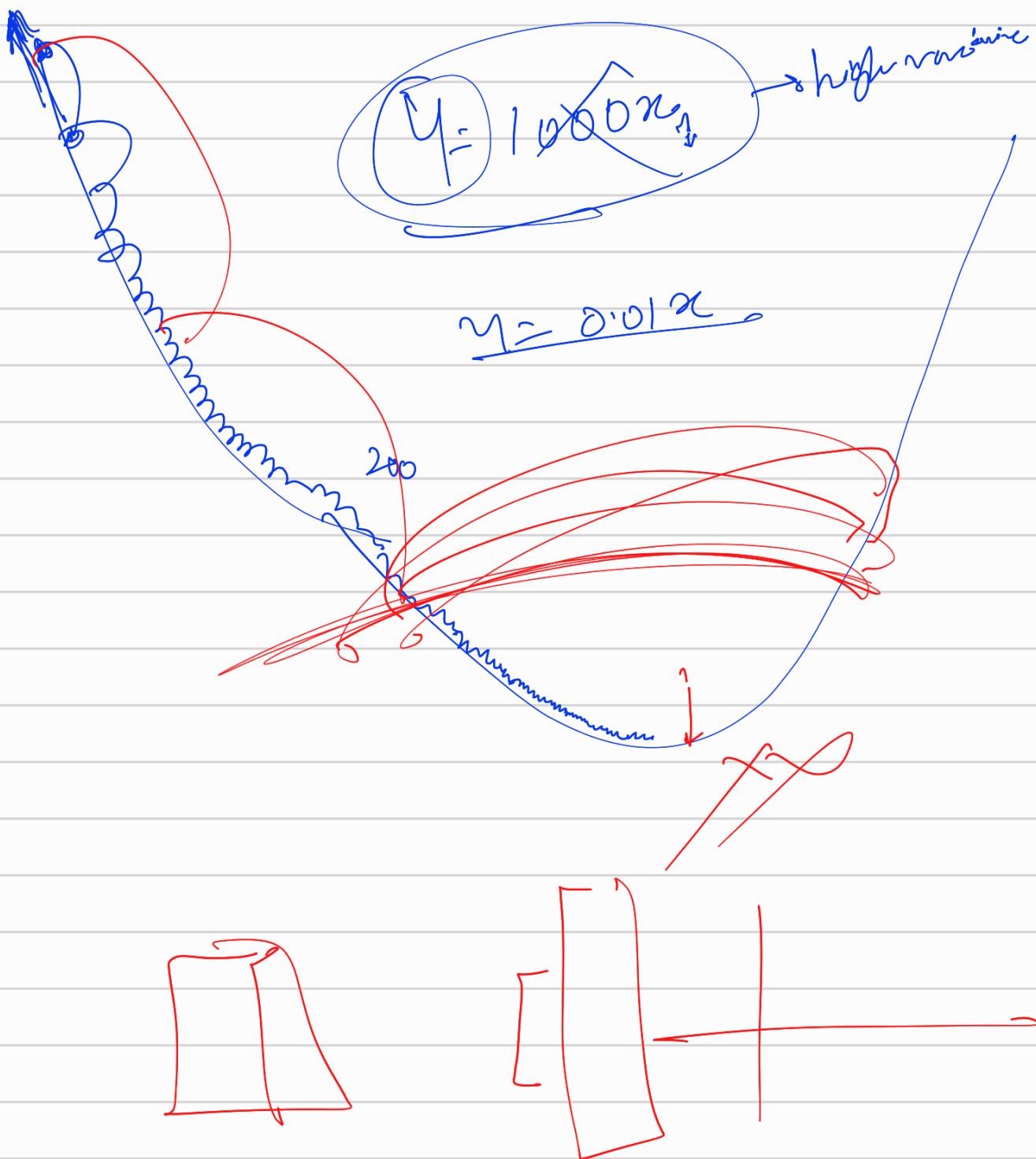
$$\sum_{k=1}^p y_k \log y_k$$

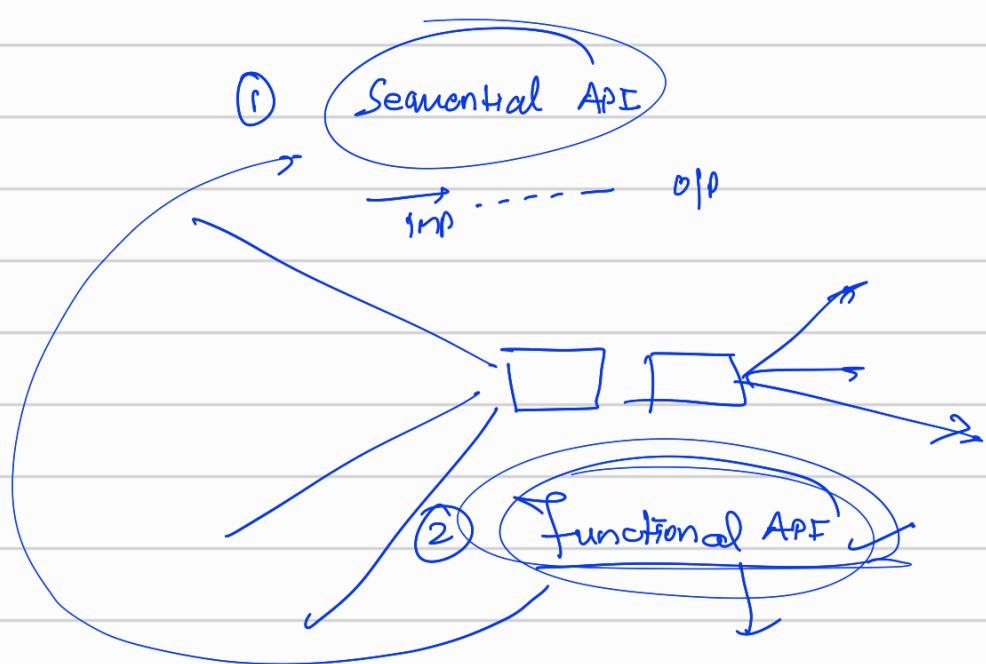
$$y_{i1} \log y_{i1} + y_{i2} \log y_{i2}$$

Breal -

22:27 PM







- Dropout
- Regularization
- Batch normalization
- Optimizers → Variants of GD.
- Callbacks

Functional API

`inp1 = Input (shape=(11,))`

`inp2 = Input (shape=(5,))`

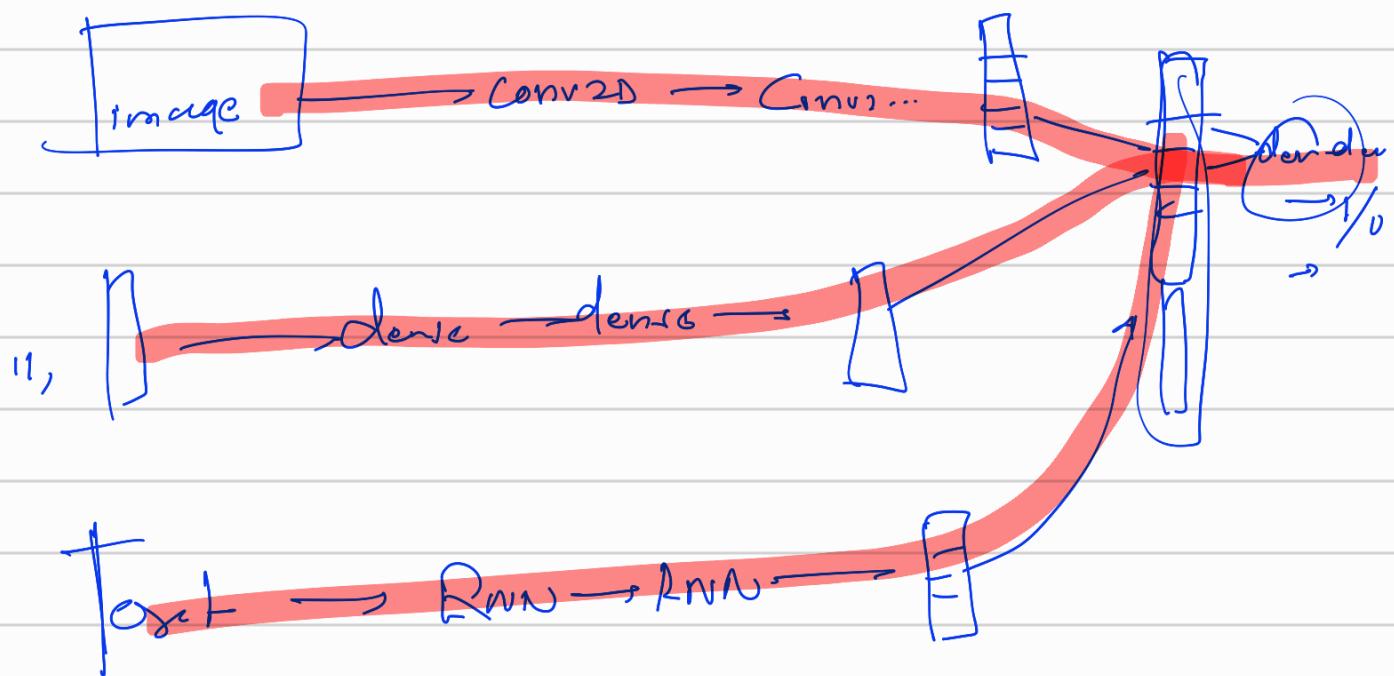
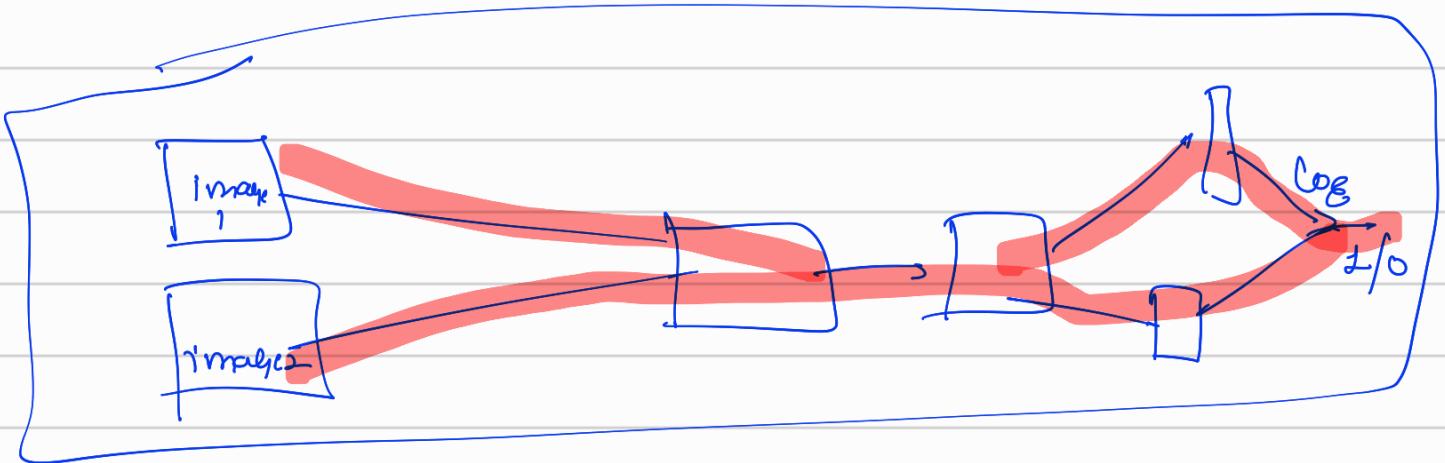
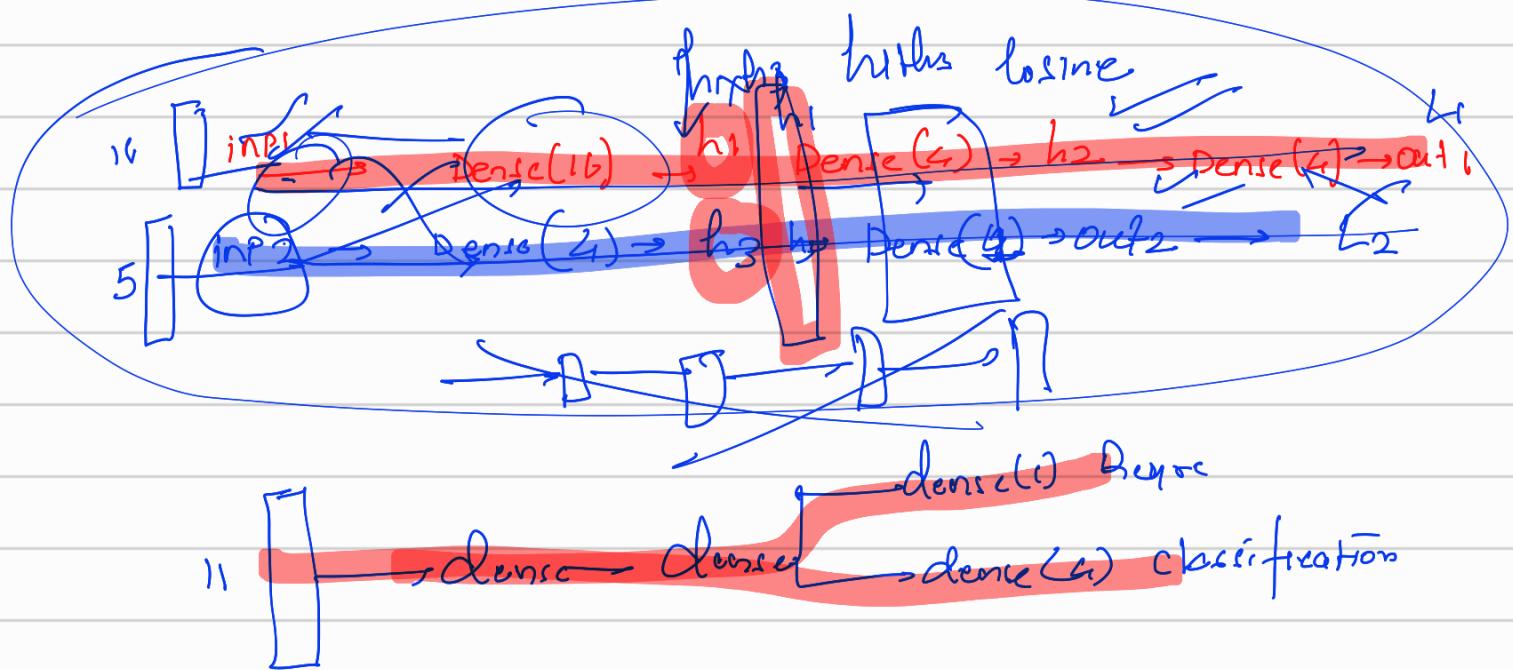
`h1 = Dense (16, activation='relu', name=' ') (inp1)`

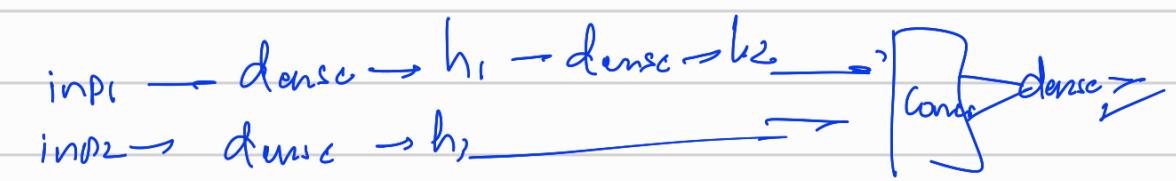
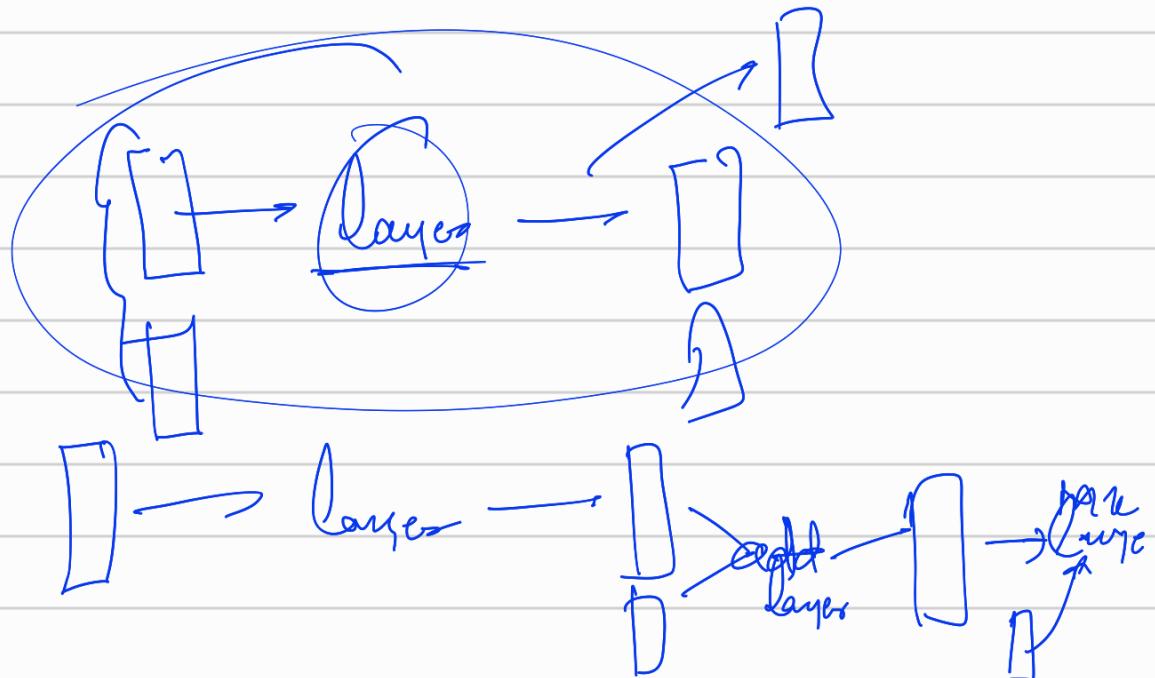
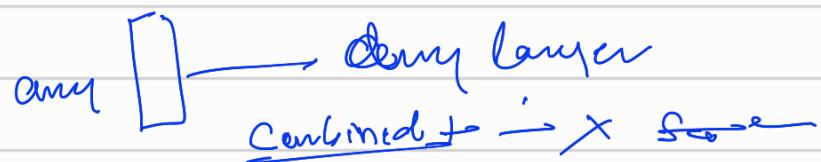
`h2 = Dense (4, activation='...', name='...') (h1)`

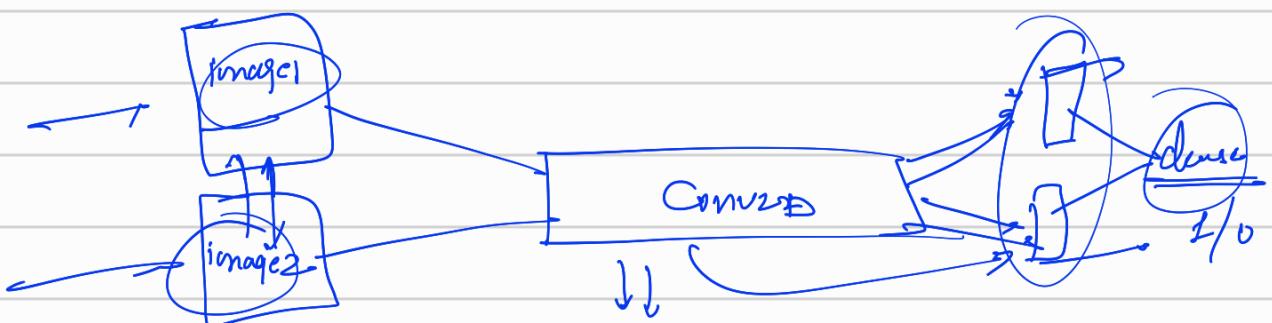
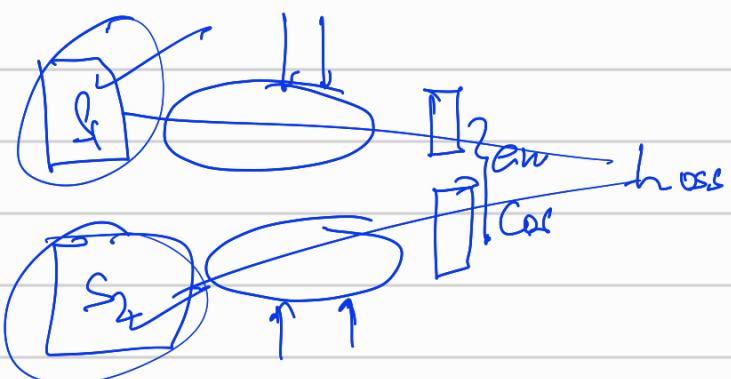
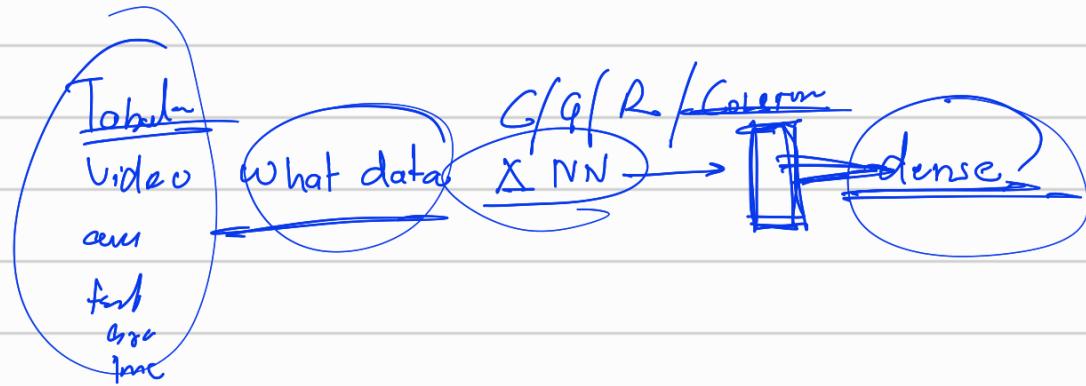
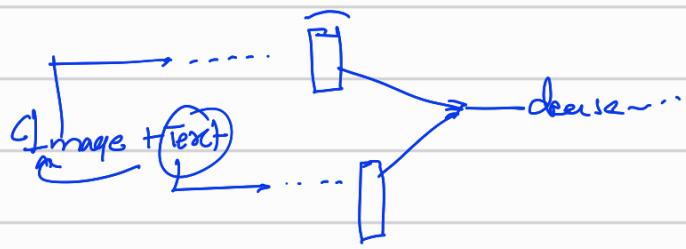
`h3 = Dense (4, activation='...') (inp2)`

`out1 = Dense (4, activation='softmax') (h2)`

`out2 = Dense (1,) (h3)`





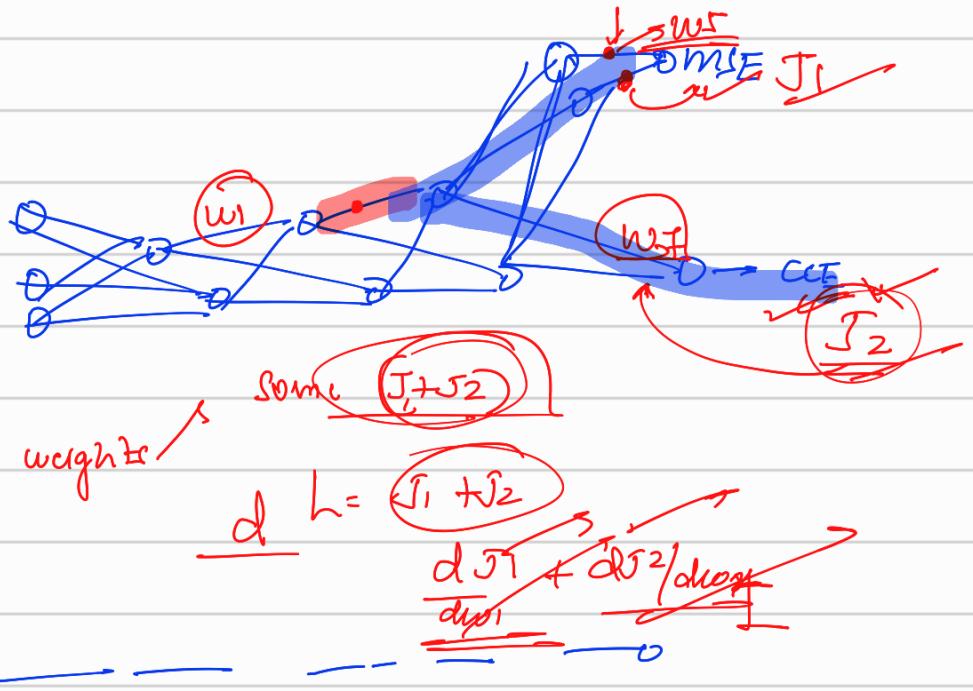


22:36 Pros.

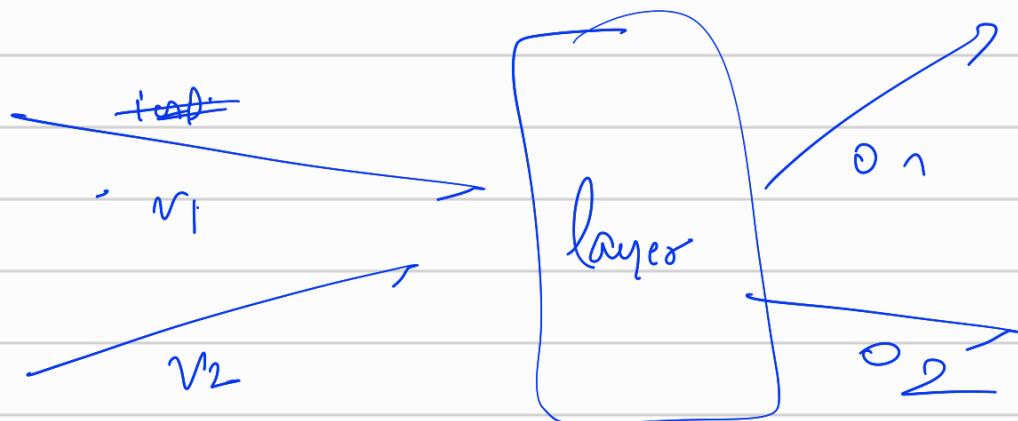
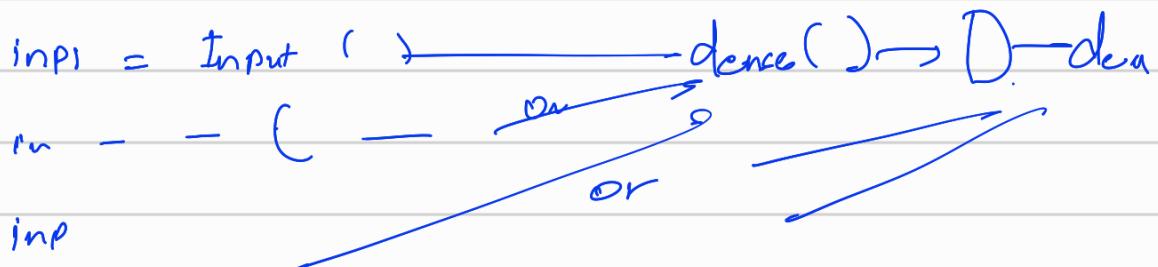
$$\begin{array}{c} J_1 \\ \rightarrow \\ J_2 \end{array}$$

Dropout:

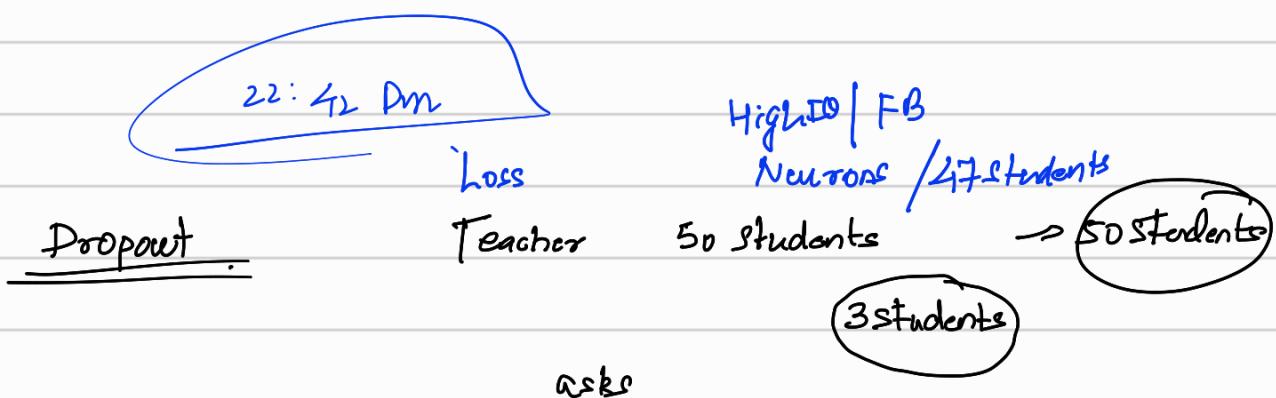
$$L = (J_1 + J_2)$$



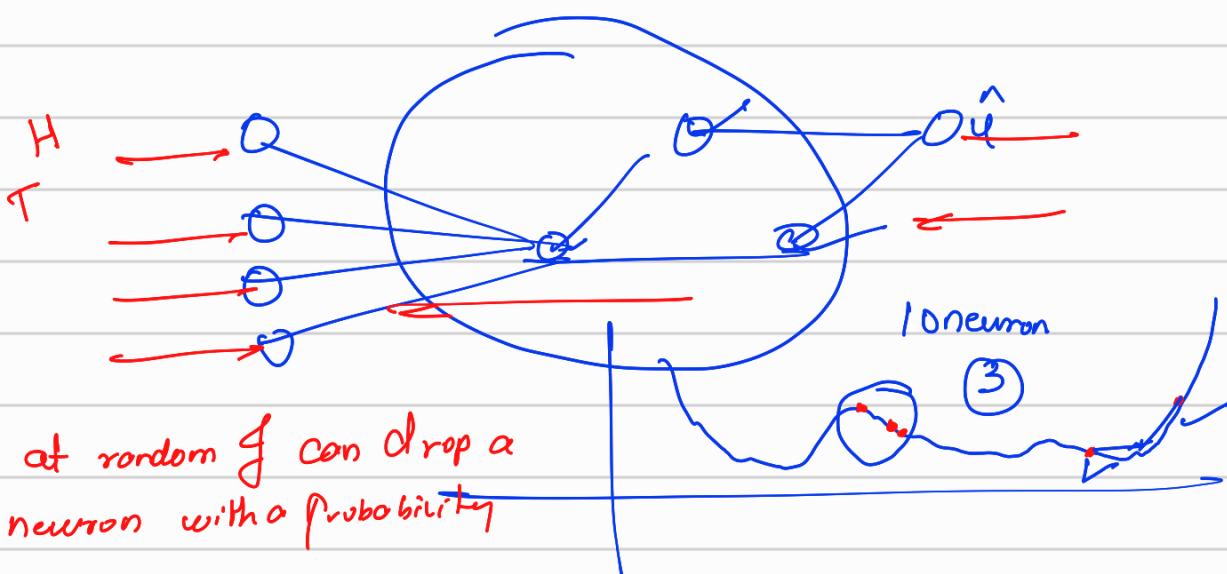
$$L = J_1 + J_2$$
$$\frac{dL}{dw_1} = \underline{\frac{dJ_1}{dw_1}} + \underline{\frac{dJ_2}{dw_1}}$$



10:36 pm



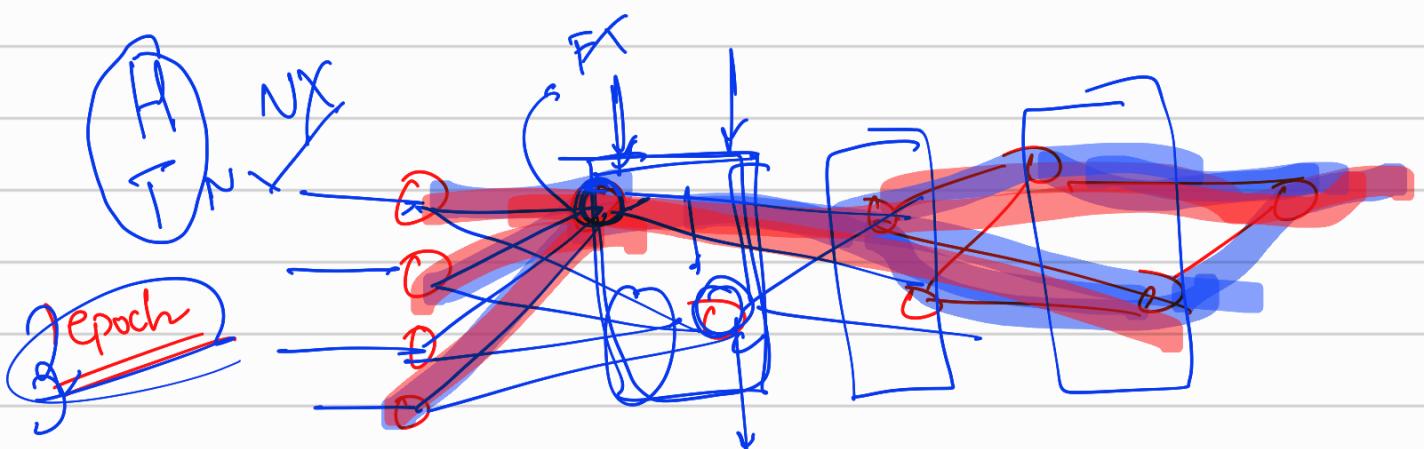
Teacher \rightarrow cold calling Dropout



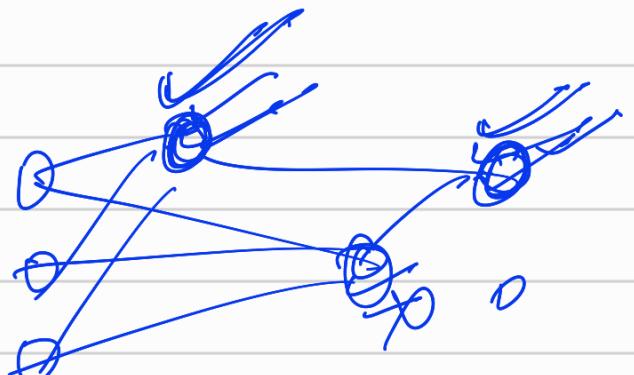
User / loss

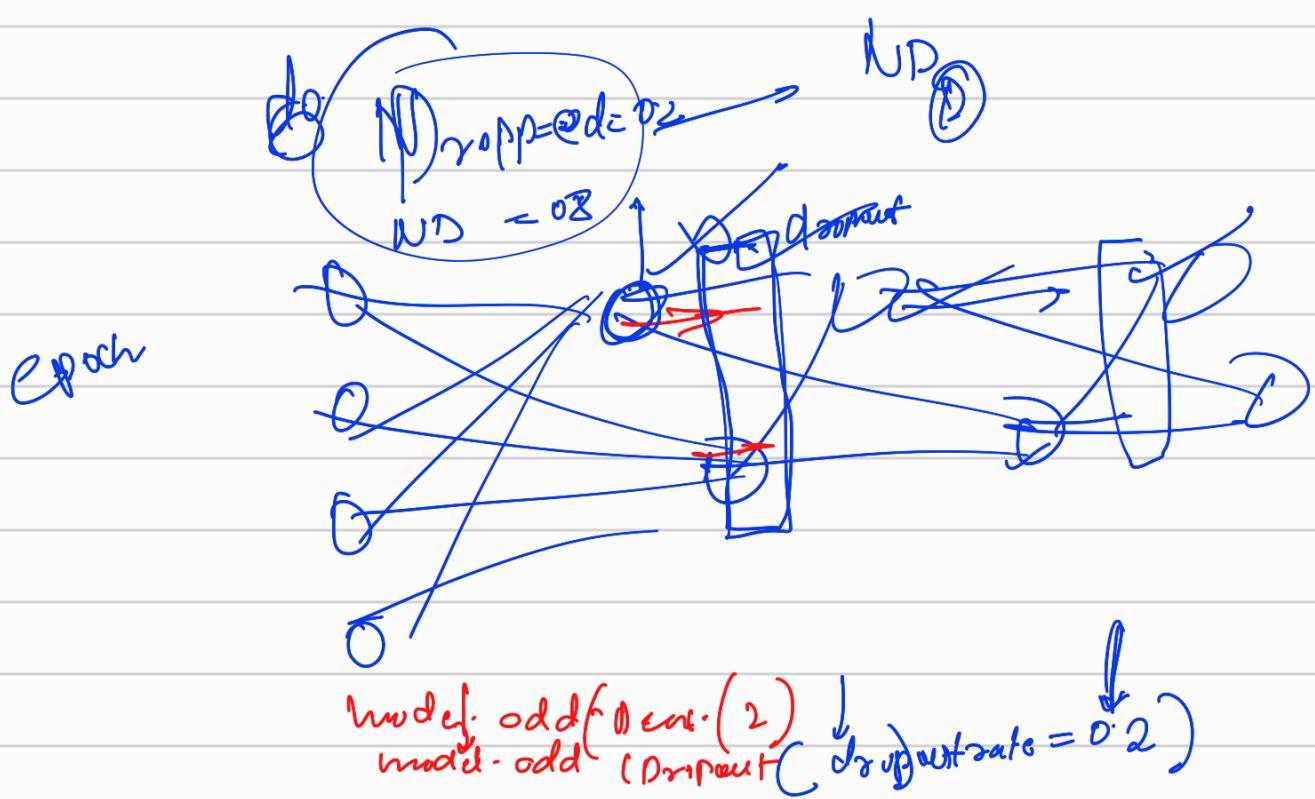


→ during that iteration the neuron will be dropped with a coin toss decision, if he is dropped then he will not(all weights in coming and outgoing) participate in the forward and backward propagation



$H \rightarrow T$ drop out rate = 0.3
 $100 \rightarrow L_0$
~~30~~





model. add Dense(2)

model.add(Dropout

2/ 3 dense 1 Dropout

