⇒ Machine Learning ~~[crossed out]~~

→ Supervised : model training with labelled data
   eg: Regression (LR), classification (KNN)

→ Unsupervised : model training with unlabelled data
   eg: clustering (Kmeans)

→ Reinforcement : Model takes actions in the environment then received state updates and feedbacks

* if feature extraction (data preprocessing) is done by humans ⟶ it is Machine Learning

– Feature extraction + model = Deep Learning
                              ↳ prediction on clean data

DL does feature engineering + model training

⇒ Applications of DL:

• GAN → deep fake, medical domain, receptionist avatar
                    ↳ generating medical related images

• Vision Transformers → chat GPT

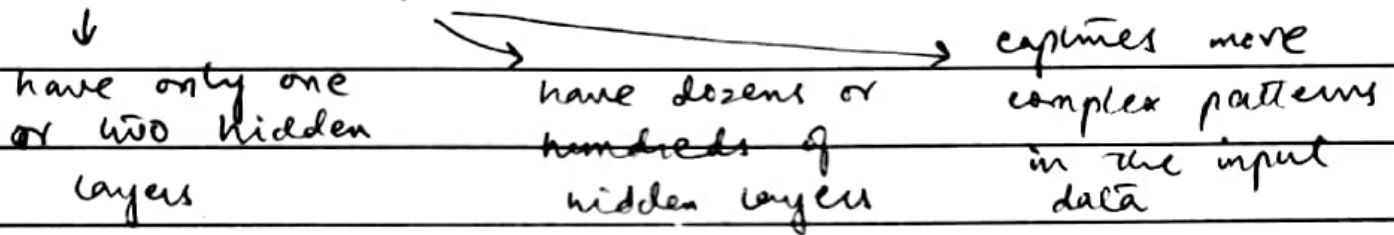• CNN → news

• Face detection vs Face Recognition
         ↳ either person or no      ↳ specific person

- Anamoly detection → surveillance videos
- RNN → predicts next word like in keyboard

→ Shallow vs Deep Neural Network :

↓

have only one
or two hidden
layers

have dozens or
hundreds of
hidden layers

captures more
complex patterns
in the input
data

→ Hyperparameters vs parameters :

↓

- allow the model to
learn the rules from
the data

- estimated during
training with historical
data

↓

- control how the
model is training

- values are set beforehand
- eg : w, b

- eg : # of layers, depth,
  iterations, epoch,
  neuron in each layer

→ Epoch :

Divide data into batches, each batch could be run
a certain no. of times

→ Cross validation : Divide data into ratios, no.) pairs we take might be biased

* K-fold cross validation → dataset is divided into k subsets/folds, model is evaluated and trained k time using a different fold as the validation set each time.

→ Overfitting vs Underfitting

↓            ↓

- good accuracy on training data but bad on testing

- bad accuracy on both training and testing data

- biasness ↑

- variance in data ↑

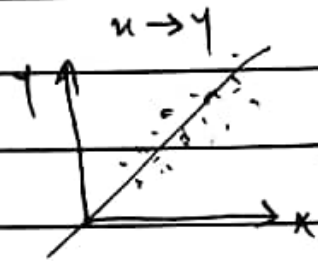| Train | Test | |
|-------|------|---|
| 0.9 | 0.4 | → overfitting |
| 0.3 | 0.3 | → underfitting → parameters are not properly tuned |

⇒ Linear Regression :

    $x \rightarrow$ independent variable

    $y \rightarrow$ dependent variable

    $y = mx + c$

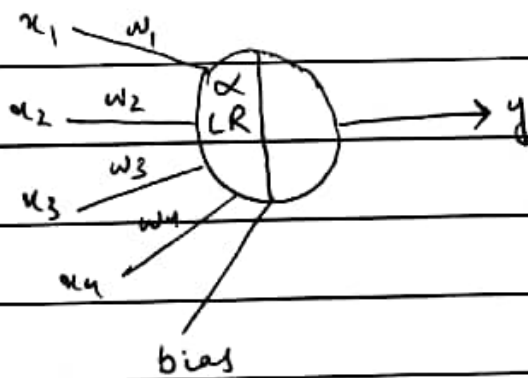    $y = b_1 x + b_0$  /  $y = w_1 x + w_0$

$x \rightarrow y$

Simple LR :   $\hat{y} = bx + a$

Multiple LR :   $\hat{y} = b_1 x_1 + b_2 x_2 + \cdots + b_k x_k + a$
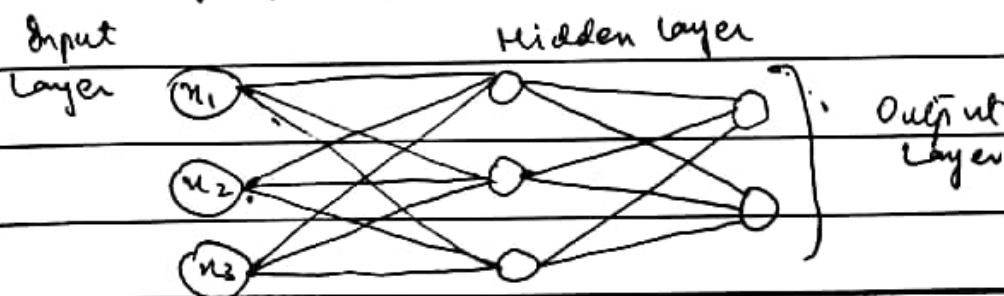
\* Objective is to find weights

⇒ The Perception : Forward Propagation

$x_1 \quad w_1$

$x_2 \quad w_2$   $\overset{\propto}{LR}$   $\rightarrow y$

$x_3 \quad w_3$

     $w_4$

$x_4$

bias

$$y = w_0 + \sum_{i=1}^{n} w_i x_i$$

$\rightarrow$ LR equation

\* range for $y, x, w \rightarrow -\infty$ to $+\infty$

Input
Layer         Hidden layer

$x_1$

$x_2$                      Output
Layer

$x_3$

\* No cyclic graph, moving forward

How many hidden layers ?

↳ could vary based on kispe result accha araha hai

\* weights tuning in back propagation — to minimize error

parameters are tuned in back — by optimising weights
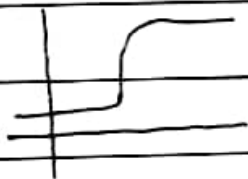
propagation

↳ weights, bias

→ ReLU, sigmoid, softmax, Tanh

→ Activation Functions — adds non-linearity, minimizes ranges

A neuron is LR + Activation Function ↳ $-\infty$ to $+\infty$

⇓

• sigmoid

0 to 1

sigmoid = $\dfrac{1}{1 + e^{-z}}$

$\dfrac{1}{1 + e^{-(-\infty)}} = \dfrac{1}{1 + e^{\infty}} = \dfrac{1}{\infty} = 0$

$\dfrac{1}{1 + e^{-(\infty)}} = \dfrac{1}{1 + e^{-\infty}} < \dfrac{1}{1 + \frac{1}{e^{\infty}}^{0}} = 1$

$g(z) = \dfrac{1}{1 + e^{-z}}$

$g'(z) = g(z)\left(1 - g(z)\right)$

- Hyperbolic Tangent

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$
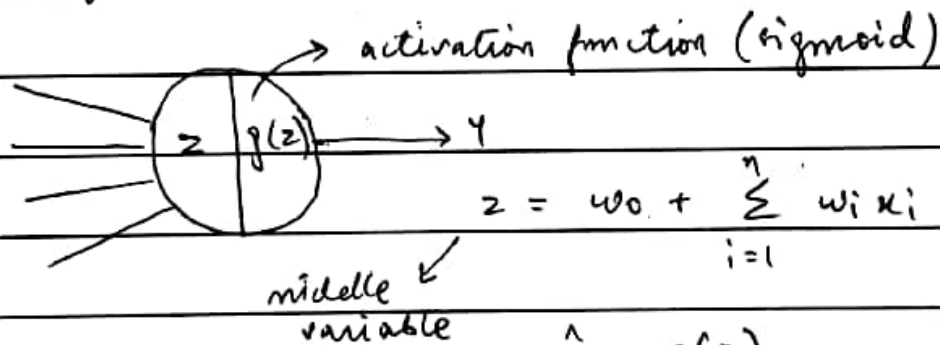
$$g'(z) = 1 - g(z)^2$$

- ReLU

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & u > 0 \\ 0, & \text{otherwise} \end{cases}$$

## LOGISTIC REGRESSION

- Single neuron is equivalent to Logistic Regression
- Classes can only be two
- Linear Regression + Activation Function = Logistic Regression



→ activation function (sigmoid)

$$z = w_0 + \sum_{i=1}^{n} w_i x_i \quad \Big| \quad wx + b$$

middle variable

$$\hat{y} = g(z)$$

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

- values in terms of logs
- probability of event occuring vs the ratio of not occuring

$$\ln\left(\frac{p}{1-p}\right)$$

↳ odds

- Instead of modeling the outcome, Y, directly, the method models the log odds (Y) using the logistic function

$$\text{LOGIT}(p) = \ln\left(\frac{p}{1-p}\right) = z \iff p = \frac{\exp(z)}{1+\exp(z)}$$

- Sigmoid $\Rightarrow$ maps $-\infty / +\infty$ to $0/1$
- Hyperbolic Tangent $\Rightarrow$ maps variables to $-1$ to $+1$
- ReLU $\Rightarrow$ negative values are clipped off
  positive values remain same
- activation function adds non-linearity (curvy line)

$$z = \left(\sum_{i=1}^{n} w_i x_i\right) + b$$

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

$$P(y=1) = \sigma(wx+b)$$
$$= \frac{1}{1+e^{-(wx+b)}}$$

$\left.\begin{array}{l}\end{array}\right\}$ probability of event occuring (class 1)

$$P(y=0) = 1 - \sigma(wx+b)$$
$$= 1 - \frac{1}{1+e^{-(wx+b)}}$$
$$= \frac{e^{-(wx+b)}}{1+e^{-(wx+b)}}$$

probability of event not occuring (class 2)

— mostly not used

⇒ Learning in Logistic Regression :

1) Loss Function / Cost Function : how close the current label ($\hat{y}$) is to the true label y.

2) Optimization algorithm : for iteratively updating weights

⇒ Loss / Cost Function

Cross entropy loss                    maximum likelihood

$$L(\hat{y}, y)$$                              ↳ maximize probability

predicted ↙        ↘ actual

- conditional max likelihood estimation prefers the correct clan labels of the training examples to be more likely.
- We choose w, b that maximize the log probability of the true y labels in the training data given the observations x.

→ Max Likelihood Estimation
- For samples labelled as 1 :   $\pi$          $\hat{y}$
  ↳ $\hat{y}$ is as close to 1          ≤ in $y_i = 1$
     as possible

- For samples labelled as $0$ : $\prod$   $(1-\hat{y})$

$\phantom{x}$ $s$ in $y_i = 0$

$\hookrightarrow$ $(1-\hat{y})$ should be

as close to $1$ as possible

* On combining both, $(w$ and $b)$ should be such that the product of both of these products in max over all the elements of the dataset.

$$L(\beta) = \prod_{s \text{ in } y_i=1} \hat{y} \quad * \quad \prod_{s \text{ in } y_i=0} (1-\hat{y})$$

This function is the one we need to optimise and in called the 'likelihood function'

$\rightarrow$ Loss Function :

if $y=1$ : $p(y|x) = \hat{y}$

if $y=0$ : $p(y|x) = 1-\hat{y}$

$$P(y|x) = \hat{y}^Y (1-\hat{y})^{(1-y)}$$

$\rightarrow$ Log likelihood :

this is a Bernoulli distribution since only two outcomes $(0,1)$

$$P(y|x) = \hat{y}^Y (1-\hat{y})^{1-y}$$

$$\log p(y|x) = \log \left[ \hat{y}^y (1-\hat{y})^{1-y} \right]$$
$$= y \log \hat{y} + (1-y) \log (1-\hat{y})$$

whatever values maximize a prob will also maximize the log of the probability

eg:

For class 1:
$$P(y|x) = \hat{y}^y (1-\hat{y})^{1-y}$$
$$= (0.7)^1 (1-0.7)^{1-1}$$

For class 0:
$$P(y|x) = (0.7)^0 (1-0.7)^{1-0}$$

- likelihood should be maximized
- loss should be minimized

→ Cross Entropy:

negative log likelihood loss = cross-entropy loss

$$L_{CE}(\hat{y}, y) = -\log p(y|x) = -\left[ y \log \hat{y} + (1-y) \log (1-\hat{y}) \right]$$

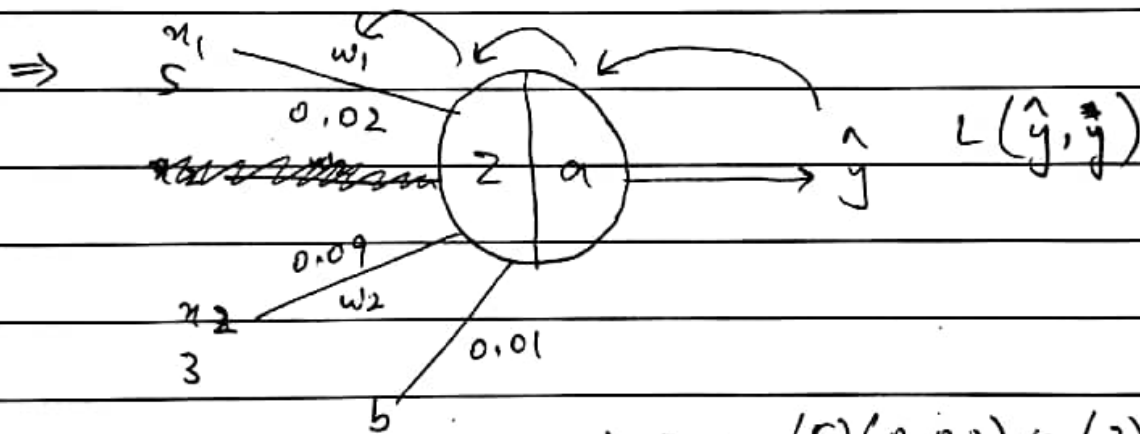$$L_{CE}(w,b) = -\left[ y \log \sigma(wx+b) + (1-y) \log (1- \sigma(wx+b)) \right]$$

→ Binary Cross Entropy Loss :

    ↳ working on two classes

    ↳ calculated as the average cross-entropy across all data examples

$$L = -\frac{1}{N}\left[ \sum_{j=1}^{N} \left[ y \log \hat{y} + (1-y) \log (1-\hat{y}) \right] \right]$$

⇒



$$x_1 \quad 5 \quad w_1$$
$$0.02$$
$$x_2 \quad 3 \quad w_2 \quad 0.09$$
$$b \quad 0.01$$
$$z \quad a \quad \rightarrow \hat{y} \quad L(\hat{y}, y)$$

→ $z = (5)(0.02) + (3)(0.09) + 0.01$

$= 0.38$

$\hat{y} = a = \dfrac{1}{1 + e^{-0.38}} = 0.5939$

→ Cross entropy loss :

$$L(\hat{y}, y) = -y \ln \left( \overbrace{6(wx + b)}^{a} \right)$$

$$= -y \ln (0.5939)$$

$$= 0.5210$$

\* $\log a = \dfrac{1}{a}$

→ Backpropogation :

$$\frac{dL}{dw_1} = \frac{dL}{da} \times \frac{da}{dz} \times \frac{dz}{dw_1}$$

$$L = \left[-y \log(\hat{y}) - (1-y)\log(1-\hat{y})\right]$$
$$= \left[-y \log a - (1-y)\log(1-a)\right]$$

$$\frac{dL}{da} = -y \cdot \frac{1}{a} - (1-y) \cdot \frac{1}{a} \cdot (-1)$$

$$= \frac{-y}{a} + \frac{(1-y)}{(1-a)}$$

$$a = \frac{1}{1+e^{-z}} \qquad \frac{da}{dz} = a(1-a)$$

$$z = w_1 x_1 + w_2 x_2 + b$$

$$\frac{dz}{dw_1} = x_1$$

$$\frac{dL}{dw_1} = \left[\frac{-y}{a} + \frac{(1-y)}{(1-a)}\right] \times \left[a(1-a)\right] \times \left[x_1\right]$$

$$= \frac{(1-a)(-y) + a(1-y)}{a(1-a)} \times a(1-a) \times x_1$$

$$= (-y + \cancel{ay} + a - \cancel{ay}) x_1$$
$$= (-y + a) x_1$$

$a \Rightarrow$ predicted
$y \Rightarrow$ actual

$$\frac{dL}{dw_1} = ax_1 - yx_1 \ / \ (a-y) x_1$$

$$\frac{dL}{dw_2} = (a-y)x_2$$

$$\frac{dL}{db} = (a-y)$$

$$\frac{dL}{dw_1} = (0.5939 - 1)(5) = \cancel{4} \ -2.03$$

$$\frac{dL}{dw_2} = (0.5939 - 1)(3) = -1.218$$

$$\frac{dL}{db} = (0.5939 - 1) = -0.4061$$

$\Rightarrow$ update weights

$\longrightarrow$ learning rate

$$w_1^{new} = w_1^{old} - \eta \frac{dL}{dw_1}$$

$\eta = 0.005$

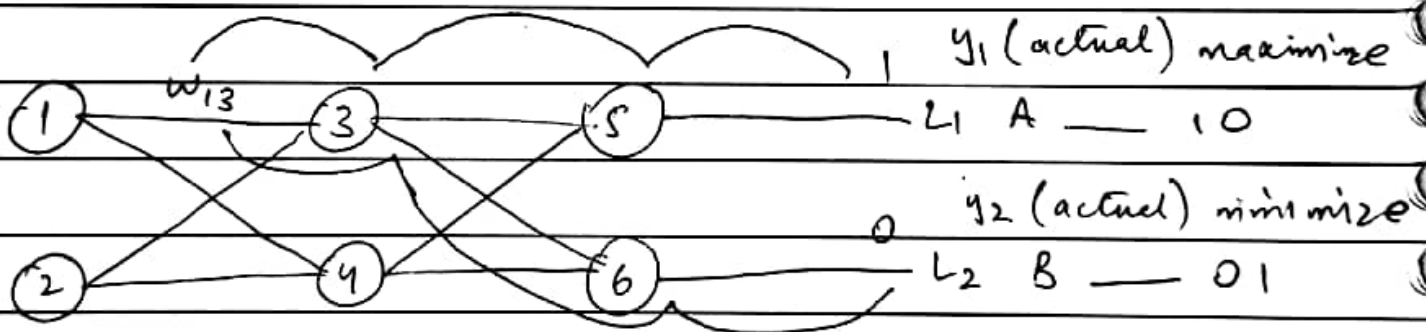$$w_1^{new} = 0.02 - 0.005(-2.03)$$
$$= 0.03$$

**Date:**

- for next iteration, use new $w_1$ as old $w_1$
- in each iteration, loss will be reduced

$$\frac{dL}{dw} = 0 \rightarrow \text{gradient descent } \cancel{\text{test}} \text{ slope zero}$$

## NEURAL NETWORKS

— from slides



$y_1$ (actual) maximize

L1  A ——— 1 0

$y_2$ (actual) minimize

L2  B ——— 0 1

$$\frac{\partial L}{\partial w_{13}} = \boxed{\frac{\partial L_1}{\partial a_5}} \times \frac{\partial a_5}{\partial z_5} \times \frac{\partial z_5}{\partial a_3} \times \frac{\partial a_3}{\partial z_3} \times \frac{\partial z_3}{\partial w_{13}}$$

Loss

$$+ \boxed{\frac{\partial L_2}{\partial a_6}} \times \frac{\partial a_6}{\partial z_6} \times \frac{\partial z_6}{\partial a_3} \times \frac{\partial a_3}{\partial z_3} \times \frac{\partial z_3}{\partial w_{13}}$$

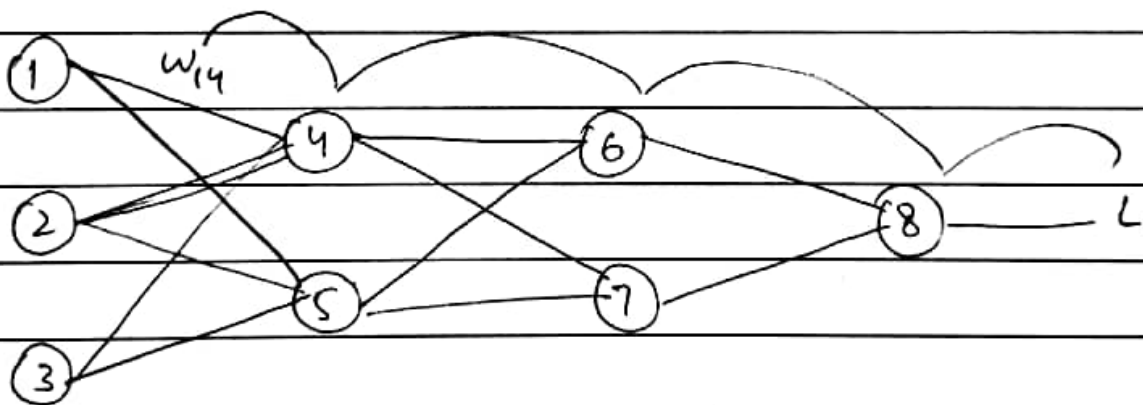One-hot encoding $\rightarrow$ 4, 7, 2 — A (1, 0)

5, 6, 3 — B (0, 1)

$$L = \frac{1}{2}(y-a)^2$$

$$= \frac{1}{2}(y_1 - a_5)^2 + \frac{1}{2}(y_2 - a_6)^2$$

$$= -(y_1 - a_5) \cdot a_5(1-a_5)a_3 \cdot a_3(1-a_3)x_1$$
$$- (y_2 - a_6) \cdot a_6(1-a_6)a_3 \cdot a_3(1-a_3)x_1$$
$$\hookrightarrow \frac{dL}{da_6}$$

$\Rightarrow$



$$\frac{dL}{dw_{14}} = \frac{dL}{da_8} \times \frac{da_8}{dz_8} \cdots\cdots \frac{da_4}{dz_4} \times \frac{dz_4}{dw_{14}}$$

# REGULARIZATION

overfitting — good accuracy on training, bad on testing
underfitting — bad accuracy on both

- NN are assumed to be overfitted models by default
- they are too complex
- training samples ko itne acchay se learn karlete hain k wo unseen data pe perform nahi karpata

$W_1 x_1 + W_2 x_2$                    there are so many variables
$0.1 x_1 + 0.3 x_2$                    that we cannot know
              ↙                        their importance
more dominantly
affecting

⇒ Regularization makes slight modifications to the
learning algorithm such that the model generalizes
better and in turn improves the model's
performance on the unseen data as well

- in ML, regularization penalizes the coefficients
- in DL, it penalizes the weight matrices of
the nodes

- assume that the regularization coefficient is so high that some of the weight matrices are nearly equal to zero
- this will result in a much simpler linear network and slight underfitting of the training data
- a large value of the regularization coefficient is not that useful — we need to optimize the value of the coefficient

Training, testing pe sahi      to obtain a well-fitted
kaam karay — appropiate      model
                  fitting

$\Rightarrow$ L1 L2 Regularization :

  L1 :

  Modified        = Loss Function + $\lambda \sum^{n}_{i=1} |w_i|$
  Loss Function

  L2 :

  Modified        = Loss Function + $\lambda \sum^{n}_{i=1} w_i^2$
  Loss Function


  $0.001 [ w_{14} + w_{15} + \dots ]$

  when overfitting, weights values are higher like 45, 30..
  $\hookrightarrow$ penalty to keep weights in range

⇒ Drop out :
- randomly skip some neurons
- random sampling of neurons
- har layer pe probability set karsaktay hain k
  iss layer pe kitne neurons chahiye
- testing time pe neurons are set as it is
ey: layer pe 50% ki prob set ki k iss layer
  k 50% neurons chahiye
- Back propogation is done on reduced network
- clipped off weights are not updated in B.P
- Different no. of weights are updated on each
  iteration
- Full network on testing
- Random selection of neurons on each iteration
- Overall accuracy is improved ___ overfitting pe
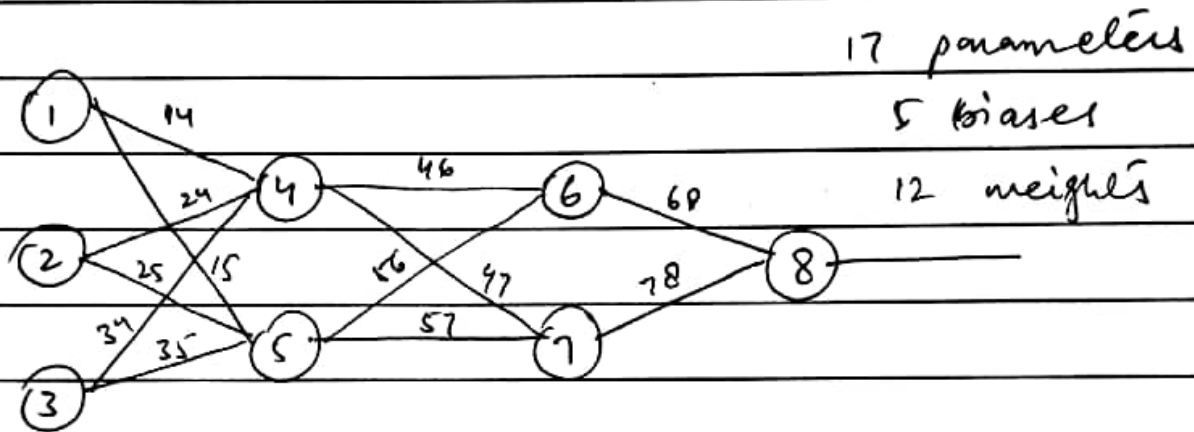                                        nahi jaata

→ Method 1 :
- Each unit is retained with a prob $p$. (train)
- weights multiplied with $p$ (test)

→ Method 2:
- weights multiplied with $\frac{1}{p}$ (train)
- no scaling (test)

— Dropout main poora neuron delete nota hai —
uske ingoing and outgoing connections bhi



17 parameters
5 biases
12 weights

remove ⑤ — new parameters =

— On which layer to apply dropout?
↳ could be different combinations
CNN
↳ sharp edges, corners
— ideally $p = 0.5$
— 100 samples pe dropout ho tou error ziada ayega
↳ less data — model couldn't learn
— merge methods like dropout + L2