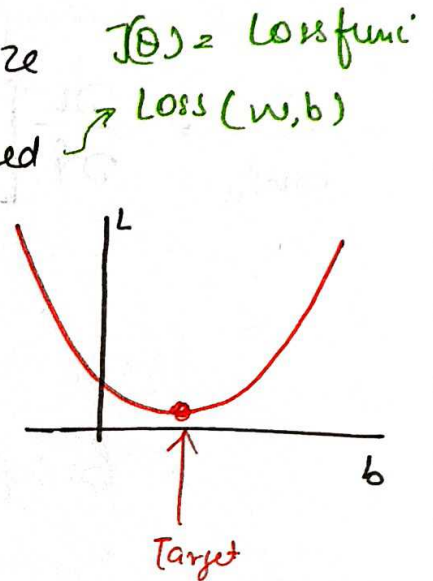


Gradient Descent

Gradient Descent: Gradient Descent is one of the most popular algorithms to perform optimization and by far the most common way to optimize neural networks.

Gradient descent is a way to minimize an objective function $J(\theta)$ parameterized by a model's parameters $\theta \in \mathbb{R}^d$ by updating the parameters in the opposite direction of the gradient of the objective function $\nabla_{\theta} J(\theta)$ w.r.t to the parameters. The learning rate α determines the size of the steps we take to reach a (local) minimum. In other words, we follow the direction of the slope of the surface created by the objective function downhill until we reach a valley.



Where we use gradient descent?

→ Backpropagation Algorithm

epochs = 5

for i in range(epochs):

for j in range(x.shape[0]):

→ select 1 row (random)

→ Predict (using forward prop)

→ Calculate loss (using loss function → mse)

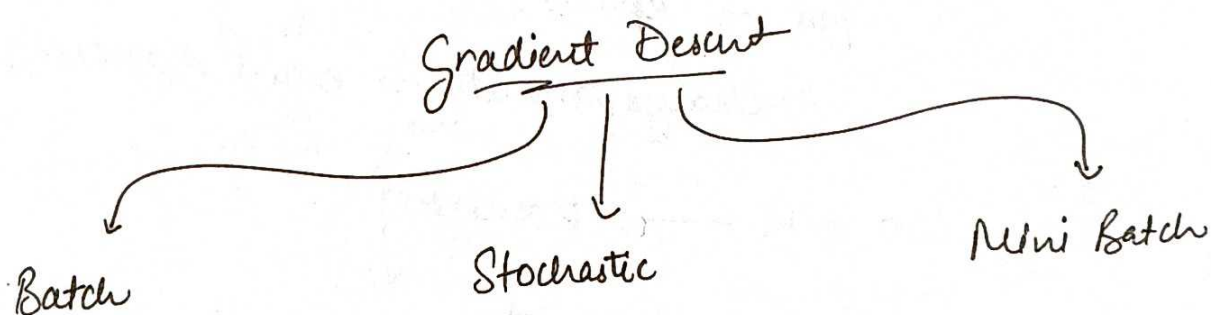
→ Update weight and bias using GP

$$w_n = w_0 - \eta \frac{\partial L}{\partial w}$$

→ calculate avg loss for the epoch

L_{avg}

gradient descent



There are three variants of gradient descent, which differ in how much data we use to compute the gradient of the objective function. Depending on the amount of data, we make a trade-off between the accuracy of the parameter update and the time it takes to perform an update.

(i) $\left\{ \frac{\partial L}{\partial w} \right\}$ → derivation → type depend on how much data

(ii) accuracy \leftrightarrow time
 \hookrightarrow trade off

Batch Gradient Descent

for i in range (no. epochs):

 params_grad = evaluate_gradient(loss_function, data, params)

 params = params - learning_rate * params_grad

In previous Gradient Descent code, we change the param for every row. for eg:- epochs \rightarrow 50, 100 \rightarrow data point

for \rightarrow epochs:

 for \rightarrow data: \rightarrow run 100 times and change parameters 100 times

But, In Batch Gradient we use entire data \rightarrow 100 in a first epoch and change parameter in single time using 100 data.

for \rightarrow epochs: ^{\rightarrow 50}

 change parameter \rightarrow entire data (100)

* Only 50 time weights and bias update

100 data \longrightarrow predict \downarrow

dot product

$$\hat{y}_{\text{hat}} = \text{np.dot}(X, w) + b$$

\hookrightarrow 50 points

$Y =$ 50 actual point

$$Y - \hat{y} = \text{loss}$$

$$\sum_{i=1}^{50}$$

* Same row ke liye ek sath predict krenge 50 points ka hoga fir 50 point ka ek sath loss find krenge

no of epoch =
no. of updates

Stochastic Gradient Descent

↳ $10 \times 50 = 500$ times
parameter change

epoch = 10

for i in range(10):

→ shuffle data

for j in range(x.shape[0]):

↳ 50

↳ 1 random point

↳ y-hat → forward

↳ loss

↳ w, b update → $w_n = w_0 - \eta \frac{\partial L}{\partial w}$

↳ Avg loss print → for the epoch

* frequency of weight
update is high

Where, in Batch Gradient Descent only 10 times change
parameter.
↳ epoch

Which is faster?

↳ Batch GD is faster

Code:

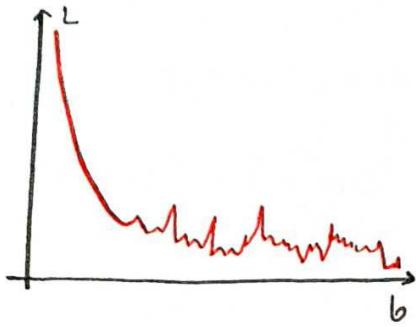
if batch-size is equal to x.shape[0] then
it is Batch GD.

If batch-size is equal to 1. Then it
is Stochastic GD.

Which is faster to reach answer?

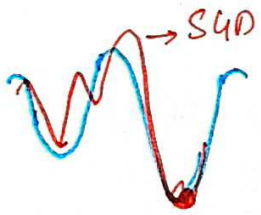
↳ Stochastic GD.

Stochastic GD



↳ Spiky SGD → Advantage

SGD → help the algo to move out of local minima cz of random spiky.



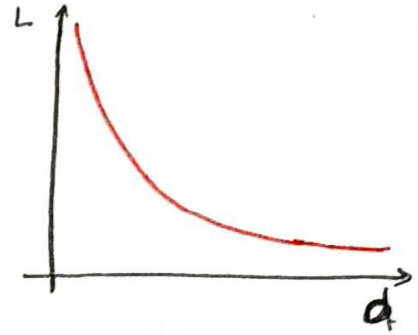
disadvantage

↳ Not find exact solution. Find approx solution cz of spiky nature.

Difference

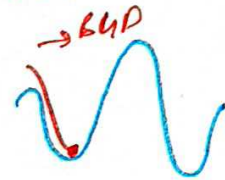


Batch GD



↳ Smooth BGD

BGD May be stuck in local minima



epoch = 0

for i in range(10):

$$\hat{y} = \text{np.dot}(x, w) + b$$

Vectorization technique

Batch GD use only epoch loop and not use another loop → and use dot product to find the prediction.

Smartest replacement to loop → faster than loop.

This is called vectorization

Vectorization technique disadvantage

(35)

↳ of datasets ~~etc~~ → 10 core

so we have to load dataset in RAM C2

BGD → find prediction → entire dataset

↳ higher memory usage.

Mini Batch Gradient Descent

↳ Best of both

BGD ↔ SGD

Example →

320 Rows → batch value → 32

↓
in every epoch
10 batches

update 10 times

for i in epochs: → shuffle data
for j in num of batch →
1 batch

↳ y-pred (vectorization)

↳ loss

↳ update

batch = x

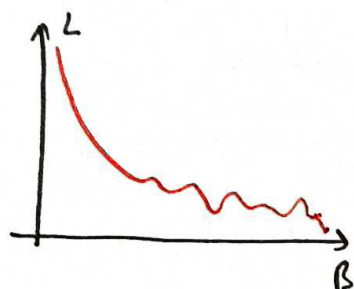
Total no. of row = n

batch-size = $\frac{n}{x}$

make batch
size

fact $\rightarrow BGD > \text{Mini-BGD} > SGD$

convergence $\rightarrow BGD < \text{Mini-BGD} < SGD$



\rightarrow why batch-size is provided in multiple of (2)?

2, 4, 8, 32, 64, ...

\hookrightarrow Bcz RAM-effective
 \hookrightarrow optimization

\rightarrow what if batch-size doesn't divide # rows properly

\hookrightarrow eg: no. of rows = 400
batch size = 150

$$\text{no. of batch} = \frac{400}{150} = 2.66$$

first batch $\rightarrow 150$

second batch $\rightarrow 150$


Third batch \rightarrow
left 100

Vanishing Gradient Problem

(36)

In machine learning, the vanishing gradient problem is encountered when training artificial neural networks with gradient-based learning methods and backpropagation. In such methods, during each iteration of training each of the neural network's weight receives an update proportional to the partial derivative of the error function with respect to the current weight. The problem is that in some cases, the gradient will be vanishingly small, effectively preventing the weight from changing its value. In the worst case, this may completely stop the neural network from further training.

2) $0.1 \times 0.1 \times 0.1 \times 0.1 = 0.0001$

2) Deep NN \rightarrow  \leftarrow Vanishing Gradient Problem in Deep NN

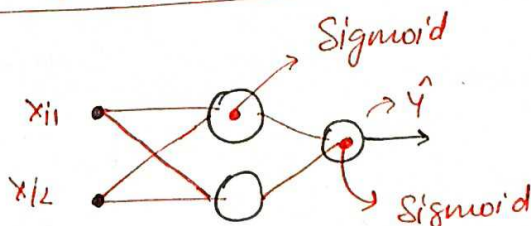
many hidden layers

3) Mostly vanishing gradient problem found in Sigmoid/tanh

derivative of sigmoid $\approx 0-0.5$

Problem

Backpropagation \rightarrow



$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} \times \frac{\partial z}{\partial a_1} \times \frac{\partial a_1}{\partial w'_{11}}$$

\downarrow \downarrow \downarrow \downarrow
 0-1 0-1 0-1 0-1

Very small \approx multiply of these no is very small ≈ 0.0001

Sometimes these no. lie between 0-1

$$w_{11} = w_0 - \eta \left[\frac{\partial L}{\partial w} \right]$$

\downarrow
derivative of L w.r.t weight

eg:-

$$W_{old} = 1$$

$$\eta = 0.01$$

$$\text{and } \frac{\partial L}{\partial w} = 0.0001$$

↳ calculate derivative

$$W_{new} = W_{old} - \eta \frac{\partial L}{\partial w}$$

$$W_{new} = 1 - 0.001 \times 0.0001$$

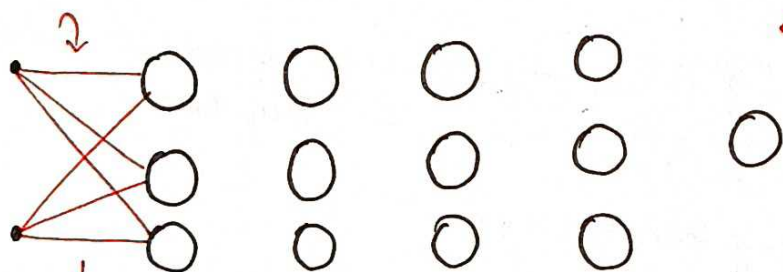
$$W_{new} = 0.99999$$

which means there is no change. old is 1. and new is 0.99999.

* If weights didn't change then Loss is same.

$$\hookrightarrow L = y - \hat{y}$$

eg: Previous example \rightarrow Only 1 hidden layer
what if 4 hidden layers.



Starting layer very small

$$\frac{\partial L}{\partial w_{ij}} = 0.000001$$

↳ derivative is very small

Backpropagation \rightarrow cannot converge

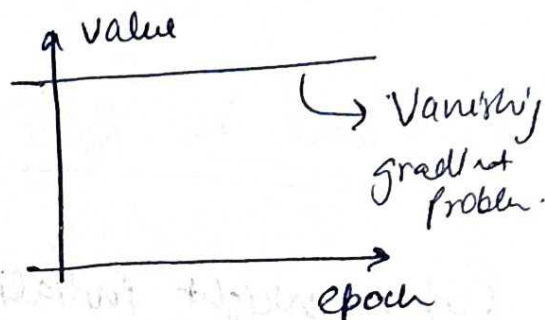
↳ model training X

Sigmoid func \rightarrow kitna bhi bada input ho ya kitna bhi chhota input ho use 0-1 ke beech mai leke ata hai. Isse Vanishing Problem hoti hai.

How to recognize Vanishing Gradient Problem?

1) Loss func \rightarrow epoch \rightarrow no changes in Loss
 \hookrightarrow \hookrightarrow Vanishing gradient Problem.

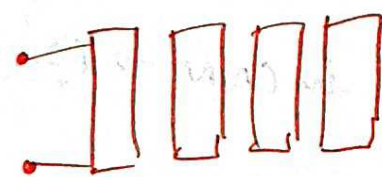
2) Weights \rightarrow Graph



How to reduce or handle Vanishing Gradient Problem \rightarrow

1) Reduce model complexity

but this is not applicable



$\hookrightarrow d_1 \times d_2 \times d_3 \times d_4$

$\hookrightarrow \frac{dL}{dw} = \text{small}$

reduce hidden layer



$\hookrightarrow d_1, d_2, d_3$

$\frac{\partial L}{\partial w} = \text{bigger}$

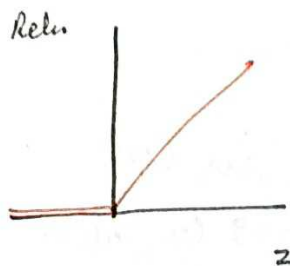
Sometimes this method is good but we increase hidden layer for complex pattern and find more accurate prediction. So, this method is not more useful.

2) Using ReLU Activation functions

$\hookrightarrow (0, \infty)$ any value. 1000, 2000

$m(0, \infty) \rightarrow$ if value is negative then change into 0
 if value is positive so value don't change

in sigmoid \rightarrow every value convert into (0-1)



→ derivative → if value is $-ve \neq 0$
if value is $pos \neq 1$

disadvantage → dying ReLU → if all derivative is $-ve \neq 0$
then w_1 is equal to w_0

↳ Discussed in ReLU Topic

3) Proper weight initialize → Glorot
→ Xavier
↳ future topic

4) Batch Normalization → future topic

5) Residual Network → Topic in CNN → RESNET
↳ building block

Exploding Gradient Problem → RNN Problem

↳ small introduction
derivative > 1
↳ (10), (10), (10), (10) → 10000

big number
 $\left(\frac{\partial L}{\partial w_{ij}} \right) = d_1 \times d_2 \times d_3 \times d_4$

* ~~loss~~ our model weight
increasing highly.

At some point, our model
behave randomly. And Loss \nrightarrow
not reduce

$w_n = w_0 - \eta \left(\frac{\partial L}{\partial w} \right) \rightarrow 1000$

$w_n = 99$

How to improve a Neural Network? → Kind of Roadmap (38)

1. Fine tuning NN hyperparameter

↳ eg:- 1. how many hidden layers?

2. No. of Neurons per layer

3. Learning Rate → LR

4. Optimizer

5. Batch Size

6. Activation

7. Epochs

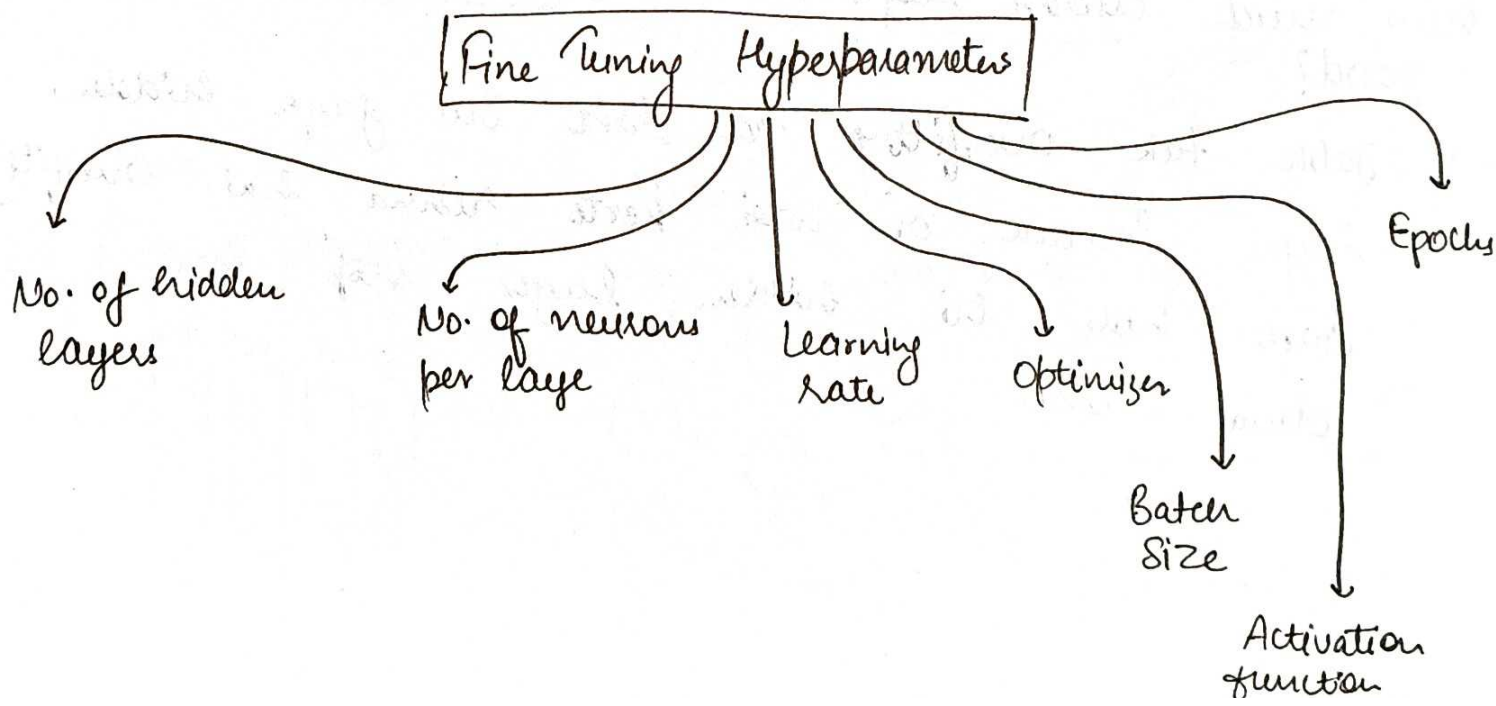
2. By Solving Problems:

→ Vanishing/Exploding gradient

→ Not Enough Data

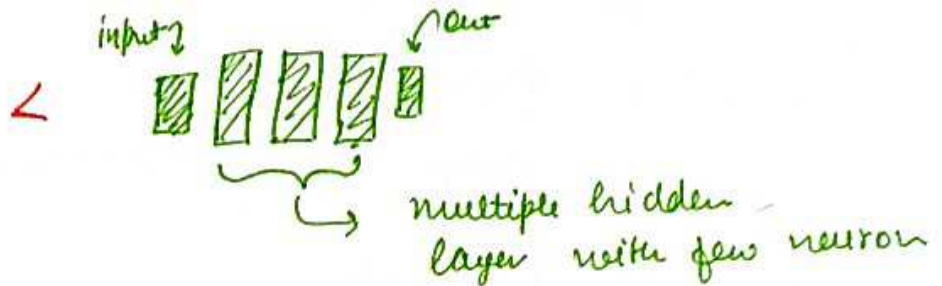
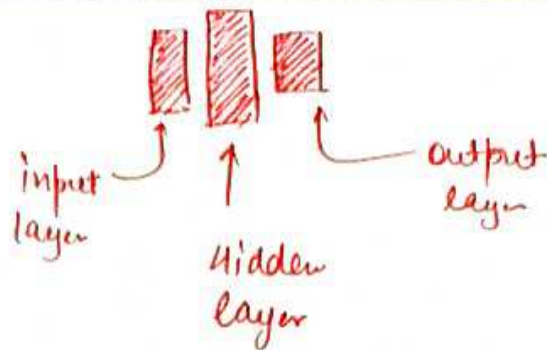
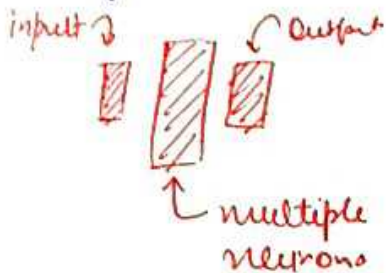
→ Slow Training

→ Overfitting



1) No. of hidden layers

1 hidden layer with 512 neurons \rightarrow complex less better than multiple layer with less neuron.

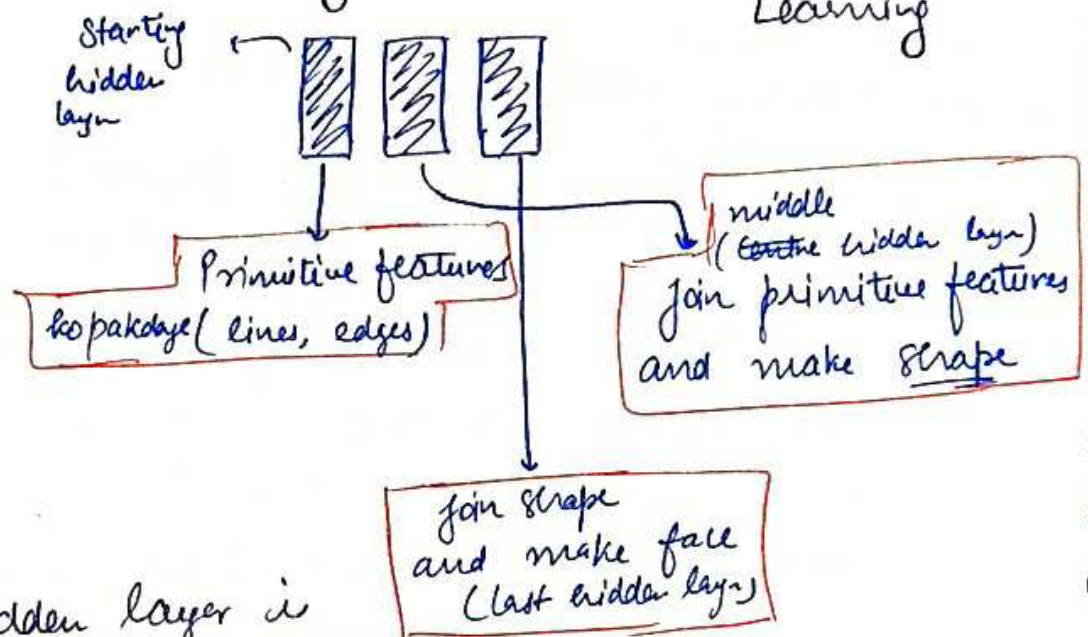


This is because Deep Learning uses \rightarrow Representation Learning

eg:-



face \uparrow



how much hidden layer is good?

\rightarrow Jabh tak overfitting na start ho jaye hidden layer increase or add karte rehna hai. Overfitting start hote hi hidden layer stop kar dena hai.



hidden layer

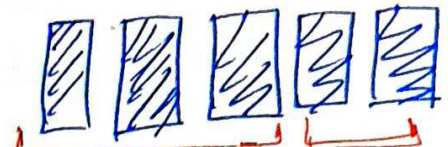
hidden layer mai

→ jaise aage jaahe hai ↑ complex pattern
Engh rahi hai. And this is a advantage

This advantage is called Transfer Learning.

Model → Human face → detection

→ reuse → Monkey face



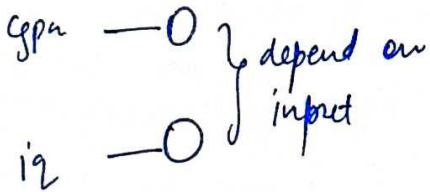
Line and Shape
Same

Took last layer
and Train on Monkey
face.

2 Neuron

Input layer ↓

eg:- Gpa/iq/pland



Output layer

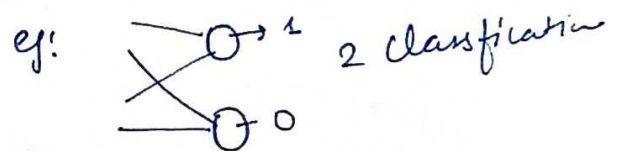
(i) Regression and Classification prob



only one output
neuron

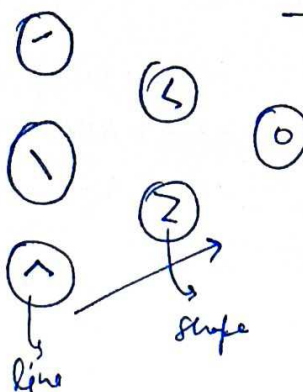
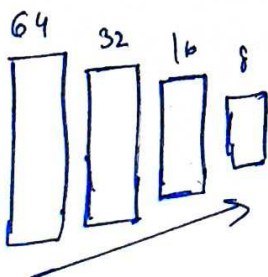
(ii) Multiclass classification

→ depend on output multiclass

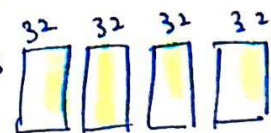


Hidden layer

(i) Pyramid Structure

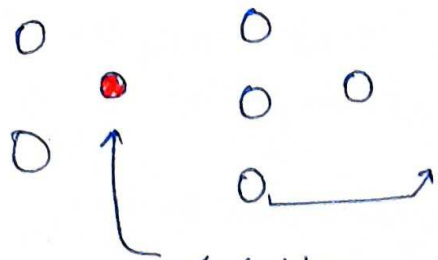


→ * If we don't
make pyramide
structure still our
prediction is good.



* No. of must be sufficient

eg:-



In this layer we can't recover the data

1 capture \rightarrow highly chance to lose the data

* No. of node should be more than what is require

\hookrightarrow if facing overfitting with more neuron than reduce some neuron.

3. Batch size \rightarrow Batch \rightarrow all row \Rightarrow weight update

\hookrightarrow Stochastic \rightarrow 1 row \Rightarrow weight update

Mini Batch \rightarrow (32) row \rightarrow 1 update
hyperparameters.

Slower

Smaller
(8 to 32)

New data

\rightarrow generalize result better

4. Epochs \rightarrow 100
 \rightarrow 500
 \rightarrow 1000

early stopping

Keras \rightarrow Stable ~~error~~ accuracy
 \downarrow
Stop

Keras callback feature

fast

Large

(8192)

We can also get better result

\rightarrow use small learning rate at starting epochs then increase the learning rate

depend
On GPU and RAM