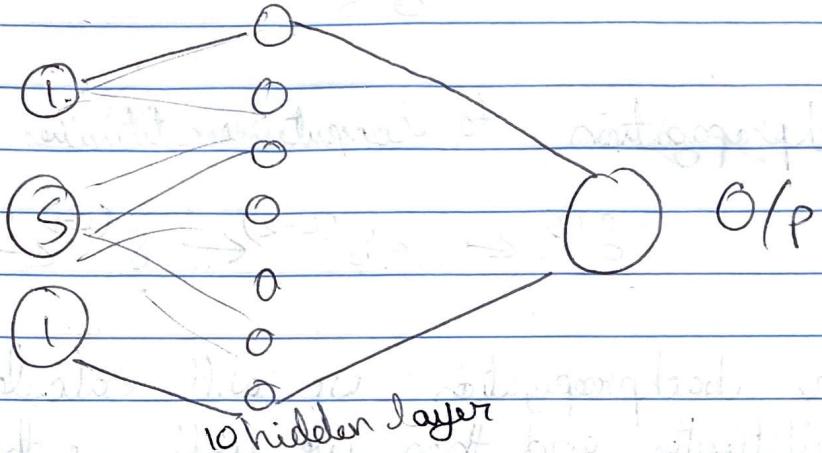
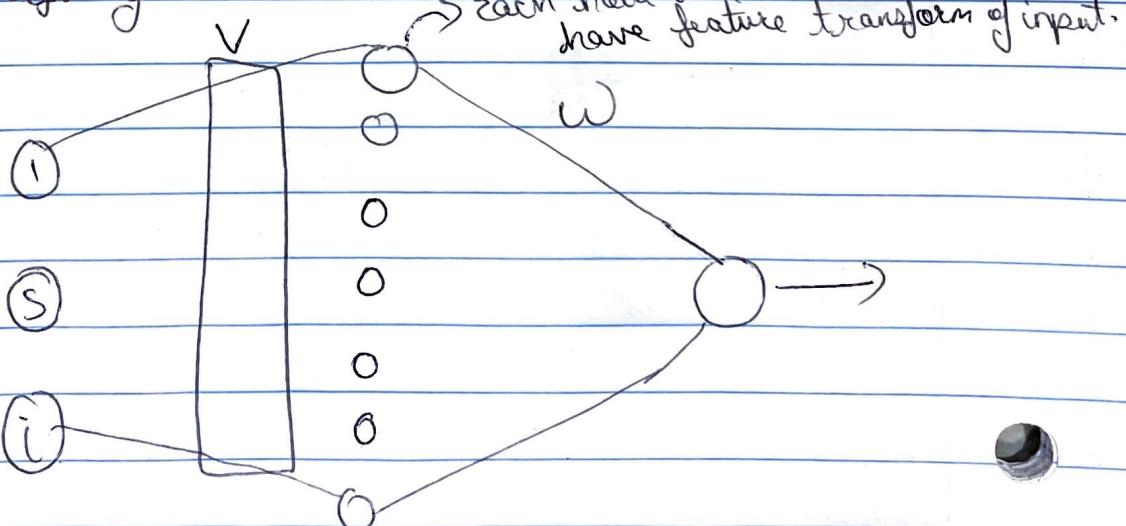


Lecture -21 etdCKb\* Gradient descent Algorithm:-

$$w(0) \quad t = 0 \quad (v - D_{avg})$$

\* Neural Networks Overfit:-

Neural network can create a classifier which will classify the points distinctly but <sup>may</sup> lead to overfitting.



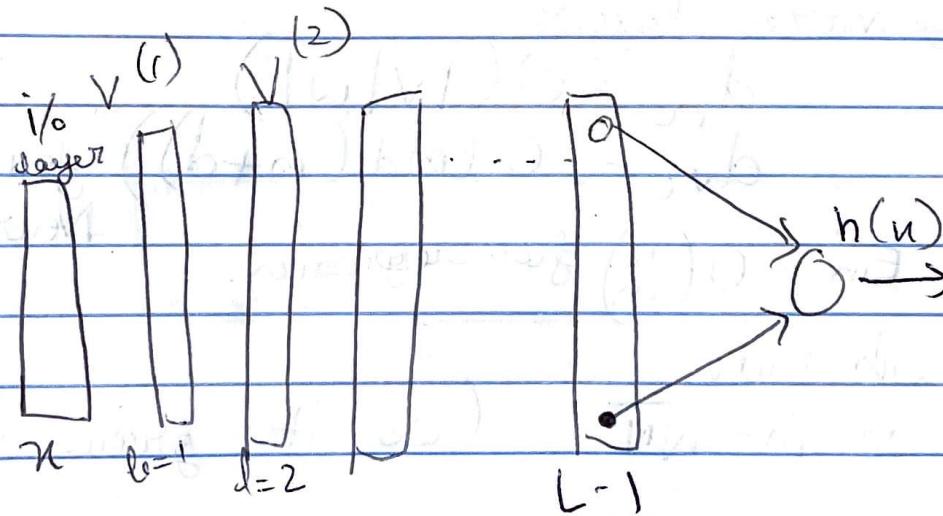
Neural Networks have a Tunable Transform.

- Neural Network

$$h(x) = \phi \left( w_0 + \sum_{j=1}^m w_j \phi(v_j^T x) \right)$$

- Non-linear Transform

$$h(x) = \phi \left( w_0 + \sum_{j=1}^m w_j \phi_j(x) \right)$$



- Approximation Capability of Neural Networks.

$$E_{in}(h) \leq \frac{(RC_f)^2}{m} = O\left(\frac{1}{m}\right).$$

\*

## Generalization

MLP (sign(.) activation func):

$$dvc = O(m d \log m)$$

where

d is # features

m is # hidden units

## Neural Networks

For  $\tanh h \rightarrow$  the dvc is more larger.

$$dvc = O(|v| |w|)$$

$$dvc = O(m d (m+d)) \text{ for 2 layer Neural Network}$$

$$Ein = O\left(\frac{1}{m}\right) \text{ for regression.}$$

\*

Thumb rule.

Set  $m = \sqrt{N}$  (Do it grows sublinearly)

$$Ein = O\left(\frac{1}{\sqrt{m}}\right) \text{ for classification.}$$

$$Ein \sim \frac{1}{N^{1/4}}, O\left(\sqrt{\frac{dvc}{N} \log N}\right) \sim \frac{1}{N^{1/4}}$$

As  $N \rightarrow \infty$ .

$$E_{\text{out}} \rightarrow E_{\text{in}} \quad \&$$

$$E_{\text{in}} \rightarrow E^*_{\text{out}}$$

\* Regularization - Weight Decay.

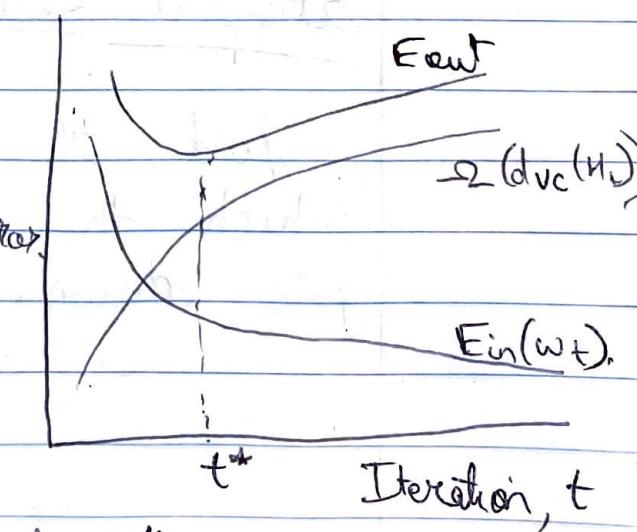
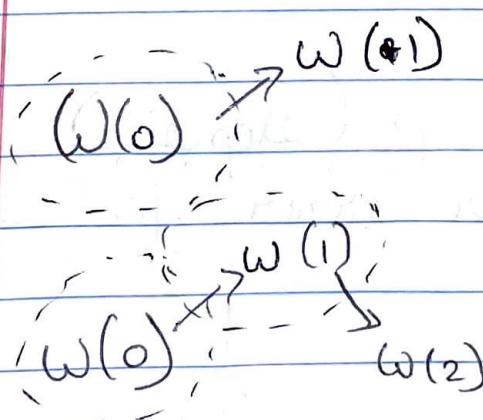
$\lambda = 0$   
(No regularization)  
overfitting

$\lambda = 0.01$   
(some regularization)  
less overfitting.

$$E_{\text{aug}}(\omega) = E_{\text{in}}(\omega) + \frac{\lambda}{N} (\text{square of wt})$$

\* Another way to regularize is to stop early.

when to stop?



Use validation set to find  $t^*$ .

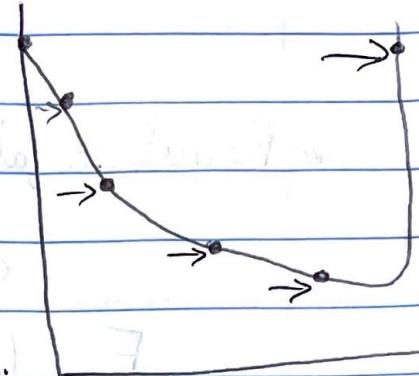
## \* Early stopping

## Variable learning Rate

$$\alpha > 1, \beta < 1$$

$$n = d n$$

learning rate will be  $\alpha$  times earlier value.



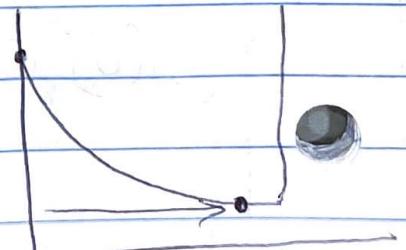
Initially we will have smaller  $n$  and slowly it will grow to a high value.

$\beta \rightarrow$  our angle.

We will go in the -ve direction of our gradient descent.

## Steepest Descent:-

find the largest  $n$  (step size) for which the value for error fun" still goes down.



In batch all of your data points ~~is~~  
In stochastic we pick random data points.

In stochastic steepest - We pick random data pt<sup>n</sup> & find the max steep we can go, for every random data point.

If all you comp are concerned about computational time, use Stochastic gradient descent rather than variable learning rate.