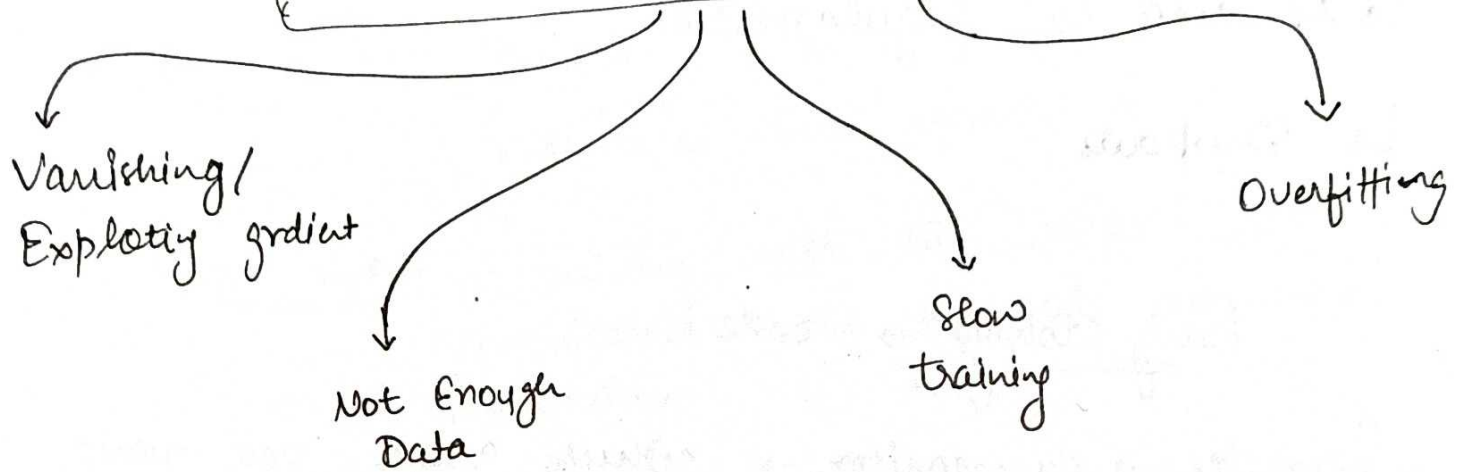


# Problems with Neural Networks



## 1) Vanishing / Exploding Gradient

- Weights initialize
- Activation function
- Batch Norm
- Gradient Clipping for Exploding Gradient only

## 2) Not enough Data

- ↳ Transfer learning
- ↳ Unsupervised pretraining

## 3) Slow Training

- ↳ optimizers
- ↳ Learning rate scheduler

↳ epoch change  $\leftrightarrow$  Learning Rate change

## 45 Overfitting

↳  $l_1$  and  $l_2$  regularization

↳ Dropouts

### Early Stopping → code

Parameter → (i) monitor → which value you want to monitor to check whether changing in loss or not.  
Mostly "val-loss" use.  
↳ uniform data

(ii) min-delta  $\Rightarrow$  how much ~~co~~ value consider as a change.

eg:- 0.00001 → If val-loss value less change than 0.00001 then we stop epochs after patience value.

(iii) patience  $\Rightarrow$  wait after checking the loss.

for eg: val-loss change less than 0.00001. So, we can't stop the epoch. we wait and run 3 more epochs if val-loss still less than 0.00001 then we stop epochs. Otherwise continue the epochs.

patience value is 3 ↙

(iv) Verbose  $\Rightarrow$  when verbose = 0 (output don't show in written "epoch 161: early stopping").

If verbose = 1 then show early stopply message.



(v) Mode  $\Rightarrow$  one of {"auto", "min", "max"}.

$\rightarrow$  "Min" mode: Training will stop when the quantity monitored has stopped decreasing.

$\rightarrow$  "Max" mode: Training will stop when the quantity monitored has stopped ~~decreasing~~ increasing.

$\rightarrow$  "Auto" mode: The direction automatically inferred from the name of the monitored quantity.

mostly this mode use.  
Automatically figured out

(vi) Baseline  $\Rightarrow$  Baseline value for the monitored quantity. Training will stop if the model doesn't show improvement over the baseline. Generally, not give any baseline cz nobody knows actual loss or accuracy during training.

(vii) Restore\_best\_weights: whether to restore model weights from the epoch with the best value of the monitored quantity. If false, the model weights obtained at the last step of training are used.

Restore-best-weights: false  $\rightarrow$  Training  $\rightarrow$  Patience epoch value  
 $\rightarrow$  True  $\rightarrow$  Training  $\rightarrow$  when the Patience epoch start value  
 0.7782  $\rightarrow$  0.77821  $\rightarrow$  0.77822  
 True take this value. Patience epoch  $\uparrow$  2 epoch  $\rightarrow$  False  $\rightarrow$  take this value

# CHEAT SHEET

## Improving Neural Network Performance

- 1) Vanishing Gradients
  - Activation Functions
  - Weight Initialization
- 2) Overfitting
  - Reduce complexity / increase Data
  - Dropout layers
  - Regularization ( $L1$  &  $L2$ )
  - Early stopping
- 3) Normalization
  - Normalizing inputs
  - Batch Normalization
  - Normalizing Activations
- 4) Gradient checking and Clipping
- 5) Optimizers
  - Momentum
  - Adagrad
  - RMSprop
  - Adam
- 6) Learning rate scheduling



## ⇒ Hyperparameter Tuning

- No. of hidden layers
- nodes / layer
- Batch size

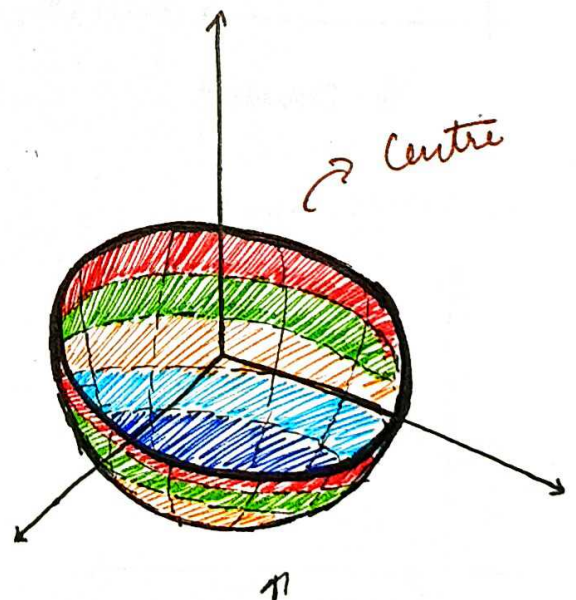
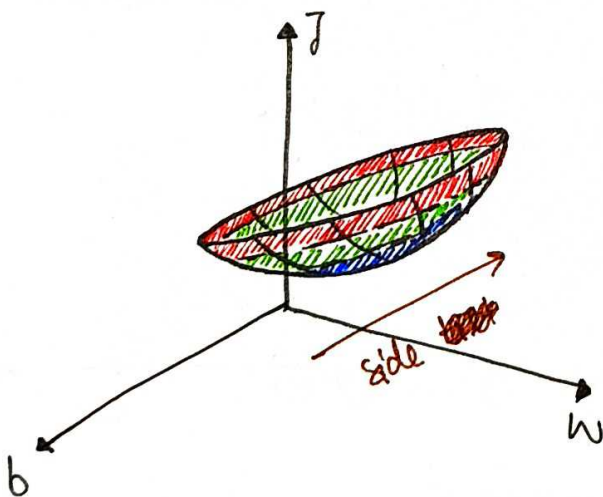
## Normalizing Input

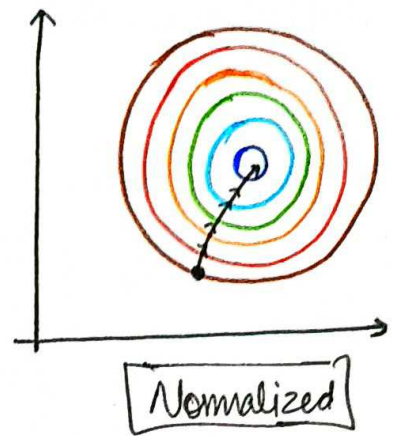
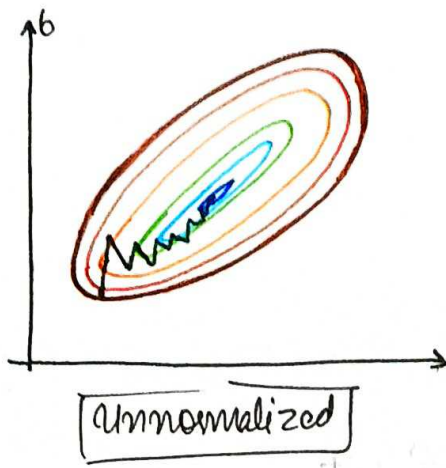
### Feature Scaling

$x_1$  (age)                       $x_2$  (salary)  
       ↓                                      ↓  
 weight ( $w_1$ )                      ( $w_2$ )

$$w_2 = w_2 - \eta \frac{\partial L}{\partial w_2}$$

↳ When model update the weight. Model give more focus on  $w_2$  than  $w_1$ , because salary in lakhs and easily see the changes and age is not big value. So, model may can ignore. Model cannot predict accurately. Not train accurately.





Solution

Standardization

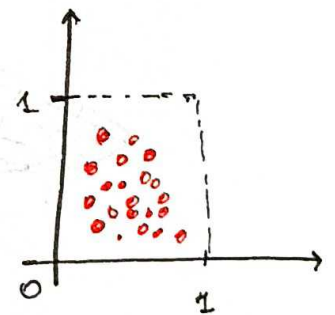
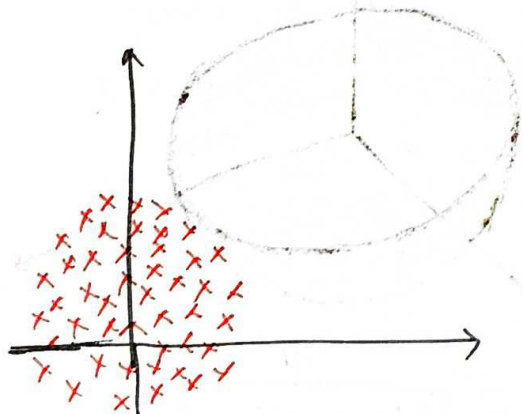
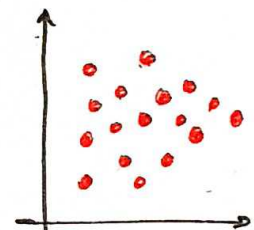
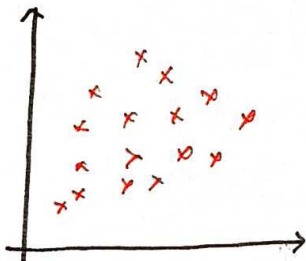
Normalization

$$\frac{x_i - \mu}{\sigma}$$

$$\frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

-1 to 1

0 to 1



When we use Standardization and Normalization? (93)

Standard → When we <sup>don't</sup> know min and max value of the column. like  $cgpa \rightarrow 0 \text{ to } 10$  (Normal)

Salary → we don't know how much salary is min and max use Standardization

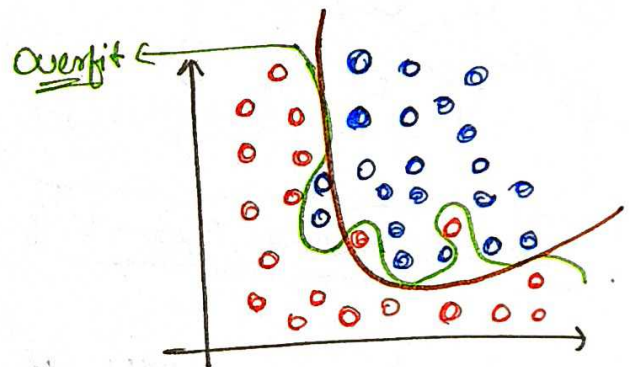
Normal → When we know min and max value of the column → use Normal  
 $cgpa = 0 \text{ to } 10 \rightarrow$  Normal

data → pdf → Normally distributed or close to normally distributed → use Standardization

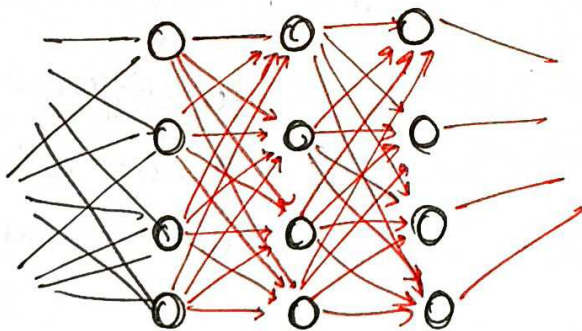
## Overfitting

### The problem of Overfitting

Overfit → Train data → good  
↳ Test data → bad



ANN → overfitting → complex architecture



run on many epochs

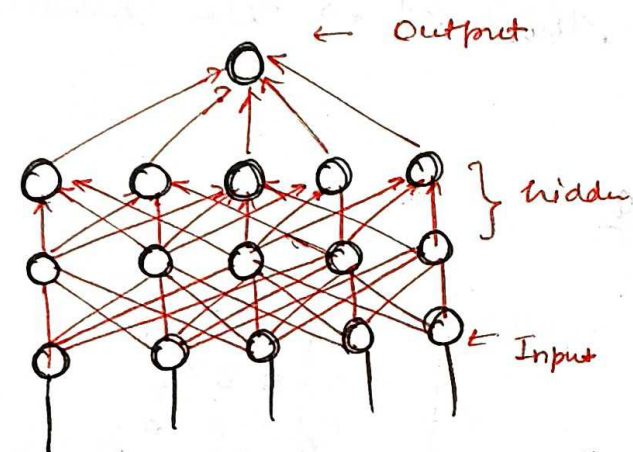
Overfit many times



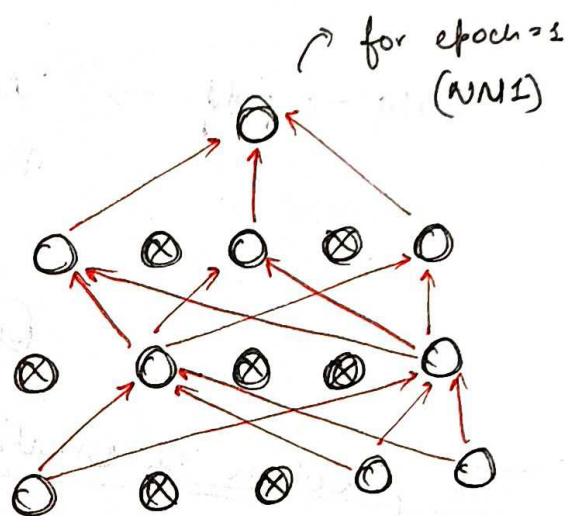
## Possible Solution

- 1) Add more data
- 2) Reduce complexity
- 3) Early stopping
- 4) Regularization  $\begin{cases} \rightarrow L1 \\ \rightarrow L2 \end{cases}$
- 5) Dropout

## The concept of Dropouts



Standard Neural Network

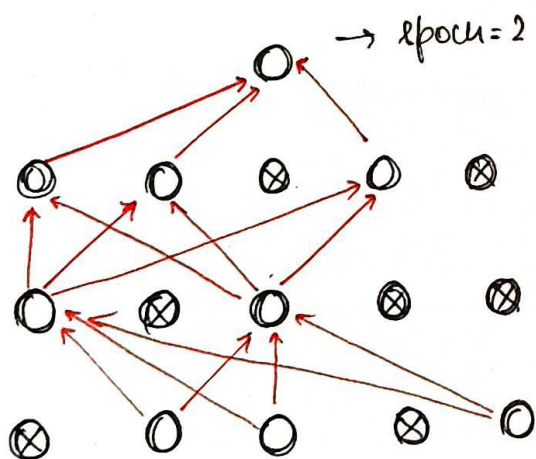


After Applying Dropout

Drop random nodes of input and hidden layer in every epoch.

eg:  $\rightarrow$  10 epochs

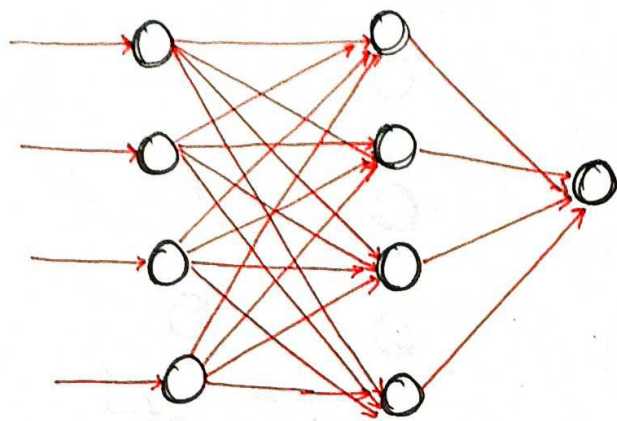
$\downarrow$   
10 different  
Neural Networks  
with same data



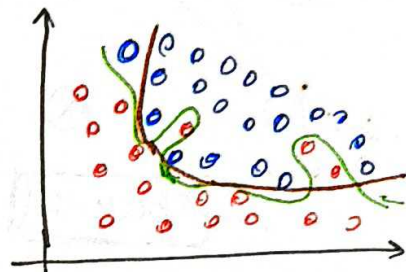
different NN2 from NN1



Why overfitting happens?



because nodes of neural network connected to each other and capture all small patterns.

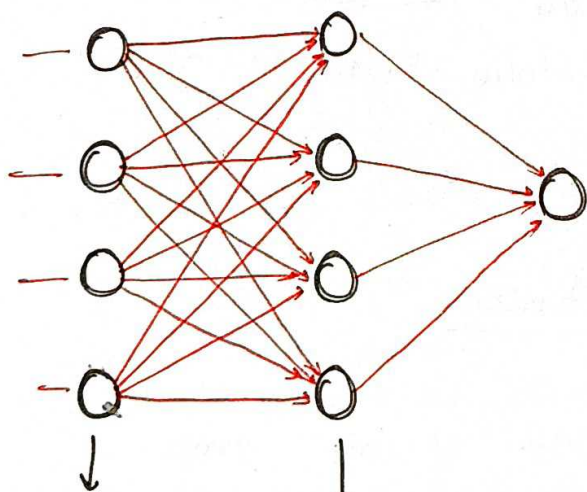


### Solution

↳ Reduce the Nodes

↳ Changes in Nodes. So, Model cannot focus on single thing or single pattern. Model should be focus on every thing and every pattern.

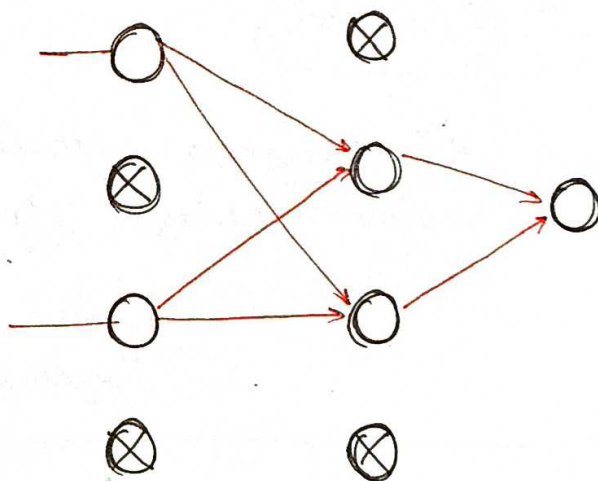
Why this works? → Reduce the Nodes



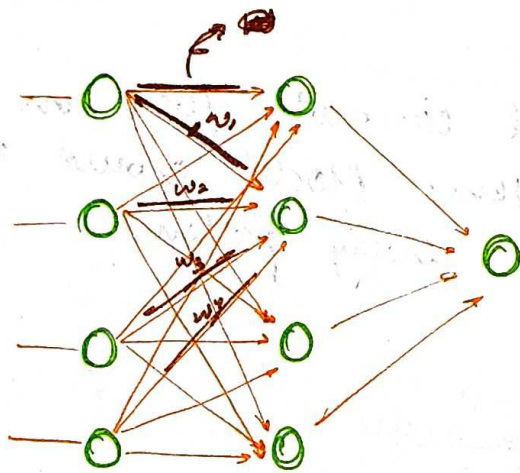
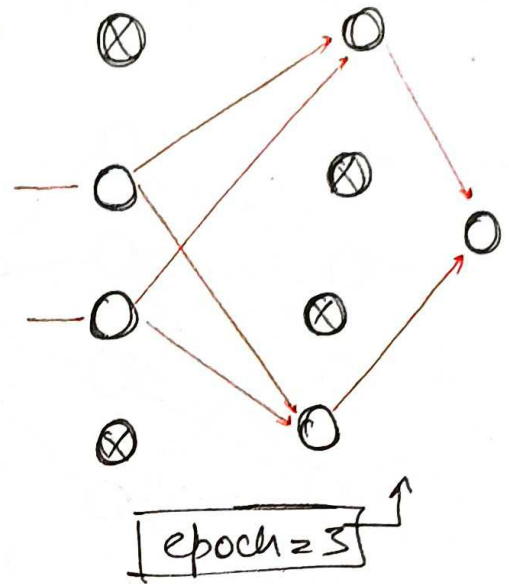
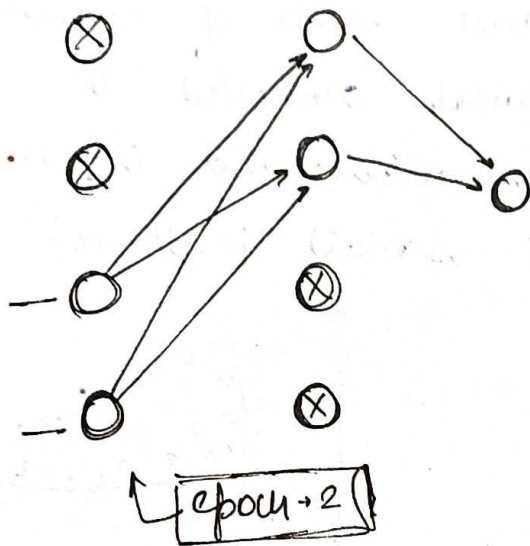
$p=0.5$

$p=0.5$

↳ Remove 50% of node in this layer



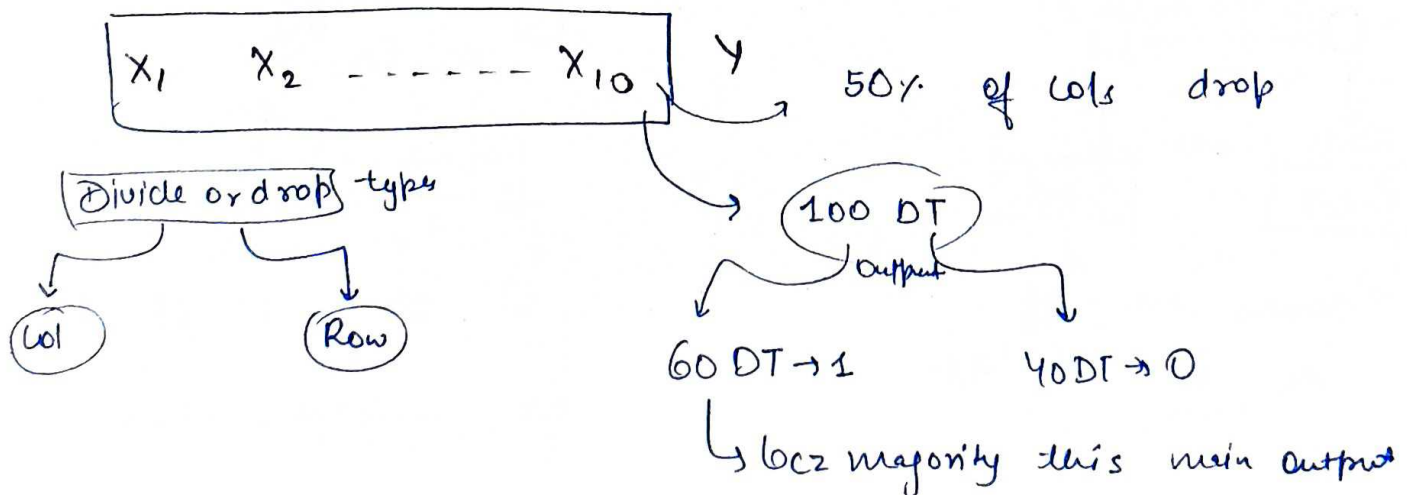
$[epoch = 1]$



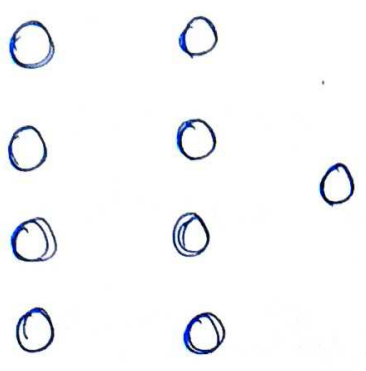
→ let assume,  $w_1$  weight is more than other weights. So, in first epoch model focus on  $w_1$  weight and first node of input layer. In epoch, 1st node is canceled so, our model can focus on other node.

Using this method model can focus on every node and weight. Reduce overfitting.

### Random Forest Analogy







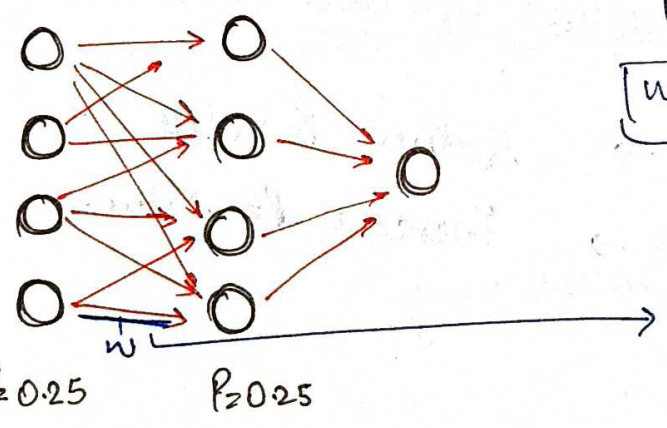
→ No. of combination of Neural Network is  $2^8$  → no. of nodes  
→  $2^8$  why 8? not 9? because output node never cancel.

$P_z 0.5$

epoch = 100

Training a united NN

How prediction works?



100 epochs → Train model

Weight =  $W(1 - P_z)$

$W \times 0.75$

why using this?

\* When we write  $P_z 0.25$  which means Probability of node not present in 100 epochs is 0.25.  
OR  
\*  $\uparrow$  0.75 time  $\wedge$  hi the probability  $\wedge$  Training ke time. 100 epochs  
0.25 time drop the training ke time.

↓  
Next page Summary



## Summarize

Training time  $\rightarrow$  ~~not~~ nodes drop kar sktte hai.

Testing time  $\rightarrow$  all nodes are available.

But we cannot assign whole node bcz every all nodes are not available. To fix time available the risk probability se multiply kar diya. That's why we multiply  $0.75$   $\rightarrow$  Probability

$P \downarrow \rightarrow$  overfitting

$P \uparrow \rightarrow$  underfitting

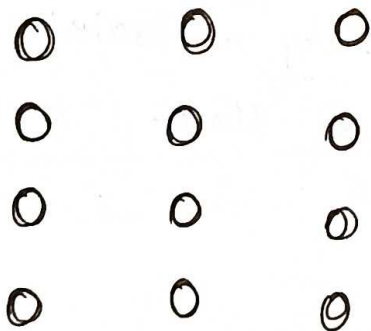
Best value  $p$  is between  
 $0 - 0.5$

## Practical Tips and Tricks

$\hookrightarrow$  Overfitting in model  $\rightarrow$  Reduce  $p$ -value  
underfitting in model  $\rightarrow$  Increase  $p$ -value.

2) last layer  $\rightarrow$  dropout

$\hookrightarrow$  ~~For~~ first try dropout on last layer then use every layer.



$\leftarrow$  First try dropout on last layer

- 3) CNN  $\rightarrow$  40-50% of value get good result (46)  
ANN  $\rightarrow$  10-50% of P value get good result

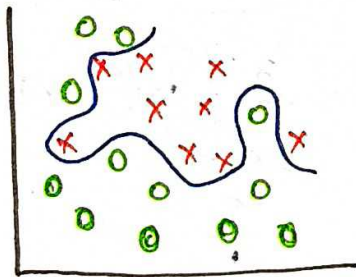
## Drawbacks

- 1) Convergence  $\rightarrow$  delay / slow training  
 $\rightarrow$  reach at accurate weights and bias
- 2) Loss function changes  $\rightarrow$  Face difficulties then debugging

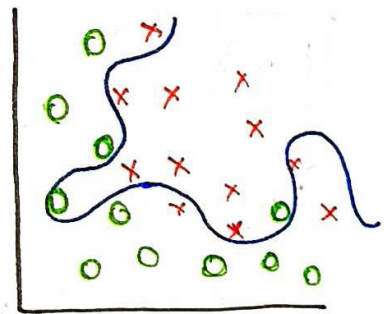
When you debugging your model you can face difficulties because many nodes are dropout and difficult to understand weight and bias.

## Regularization

### Overfitting



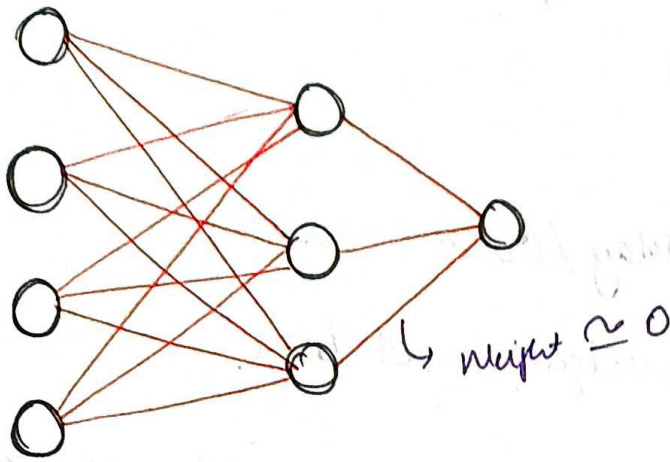
Train Dataset



Test Dataset

Why neural Network overfit?

Because complex <sup>Architecture</sup> pattern  
caused small pattern



## Regularization

ANN  $\rightarrow$  weight / bias

$\rightarrow$  min lossfunction  $\rightarrow$

$L = \text{mse}$

$\rightarrow$  binary

$L_2$

$L_1$

$$C = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{y}_i) + \text{Penalty Term}$$

$$C = L + \left[ \frac{\lambda}{2n} \sum_{i=1}^k \|w_i\|^2 \right]$$

$$\frac{\lambda}{2n} [w_1^2 + w_2^2 + \dots + w_{10}^2]$$

let 10 weight  
 $w_1 \rightarrow w_{10}$

$\lambda \rightarrow$  hyperparameter

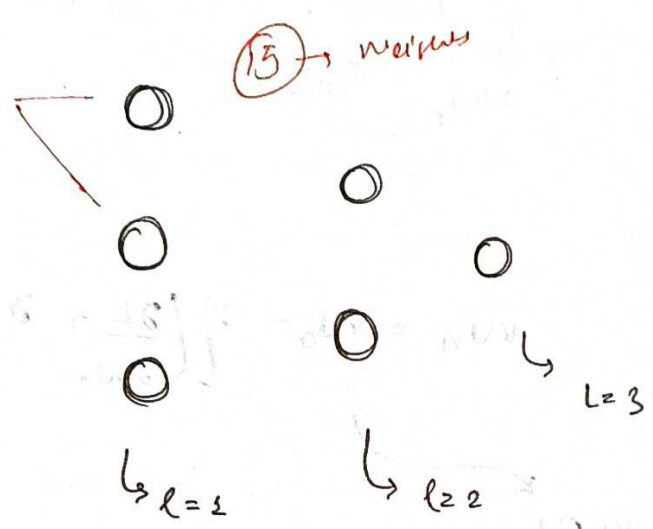
if  $\lambda \rightarrow$  higher  
lower  $\rightarrow$  underfitting  
 $\rightarrow$  still overfitting

When we apply regularization,  
value of weight going to  
0.  $w \approx 0$



for  $L_1 \rightarrow C = L + \frac{\lambda}{2n} \sum_{i=1}^n ||w_i||^2$

Most of the book  $\rightarrow C = \sum_{i=1}^n L(y_i, \hat{y}_i) + \sum_{l=1}^L \sum_{i=1}^n \sum_{j=1}^n ||w_{ij}^l||^2$



entering layer by layer

layer  $\rightarrow L$

row  $\rightarrow i$

column  $\rightarrow j$

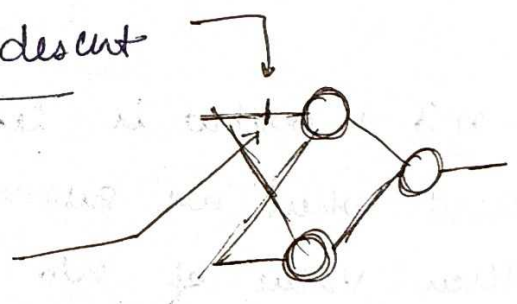
$w_1^2 + w_2^2 + \dots$

$\hookrightarrow$  only weights add  
not bias

# Intuition behind Regularization

In Backpropagation  $\rightarrow$  use  $\rightarrow$  Gradient descent

$w_n = w_0 - \eta \frac{\partial L}{\partial w_0}$   $\leftarrow$  change the particular weights

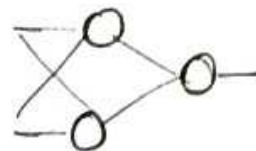


Loss function  $\rightarrow$

$L' = L + \frac{\lambda}{2} \sum ||w_i||^2$

$\hookrightarrow$  not dividing by  $n$  because we are solving only for one.

$$L' = L + \frac{\lambda}{2} \sum \|w_i\|^2$$



$$\frac{\partial L'}{\partial w_0} = \frac{\partial L}{\partial w_0} + \frac{\lambda}{2} \frac{\partial w_0^2}{\partial w_0}$$

$$w_i^2 + (w_1^2 + w_2^2 + \dots)$$

differentiate w.r.t.

so will be 0.

output

$$= \frac{\partial L}{\partial w_0} + \lambda w_0$$

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

$$w_n = w_0 - \eta \left( \frac{\partial L}{\partial w_0} + \lambda w_0 \right)$$

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

$$w_n = (1 - \eta \lambda) w_0 - \eta \frac{\partial L}{\partial w_0} \rightarrow \text{almost same}$$

difference

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

$1 - \eta \lambda \rightarrow$  output is less than 1

and when we subtract with  $w_0$

then value of  $w_0$  will decrease.

also called weight decay.

$$w_0 \gg (1 - \eta \lambda) w_0$$

$w_0 \downarrow \downarrow$

old  $w_0$

new  $w_0$

Jitne No. of epochs chalega utni bar  $w_0$  ki value decrease hogi. Zero ke paas ja skta hai but completely zero nhi hoga.

L1 is a sparse Model

↳ where many node will be 0  
And Make Simple Model or Neural Network.

But in L2 → Nodes are near ~~two~~ zero but  
not completely zero. So, our model not become  
very simple neural Network.

L2 is useful but <sup>try</sup> with both L1 and L2