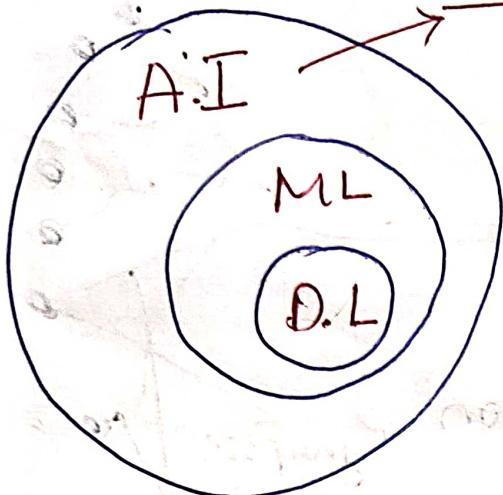


# DEEP LEARNING

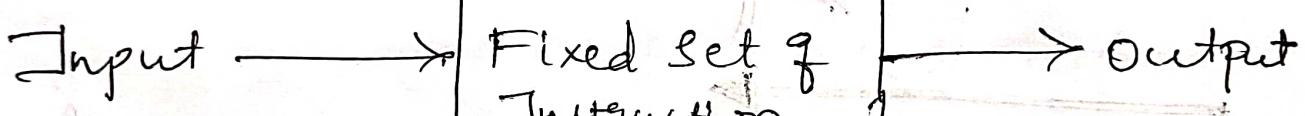
## Neural Networks

RPA



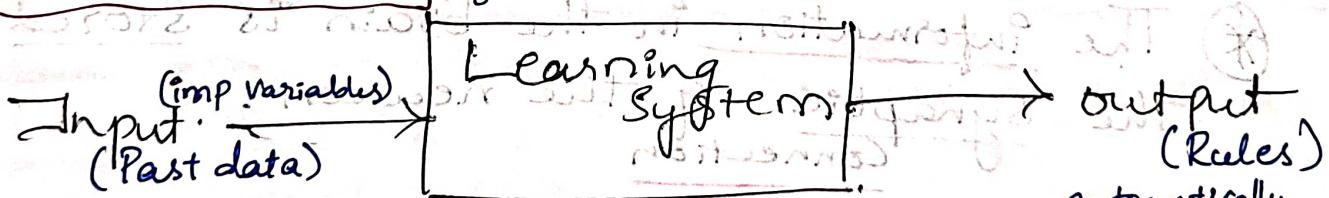
\* Deep Learning is the subset of M. Learning which is a subset of Artificial Intelligence.

### A.I. Intelligence :-



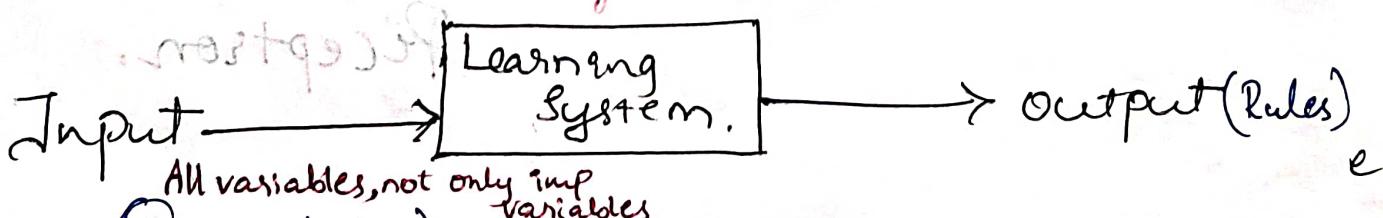
\* Rules based decision Making (Robotic Process Automation)

### Machine Learning :-



\* Machine learning involves figuring out the rules automatically from the past data.

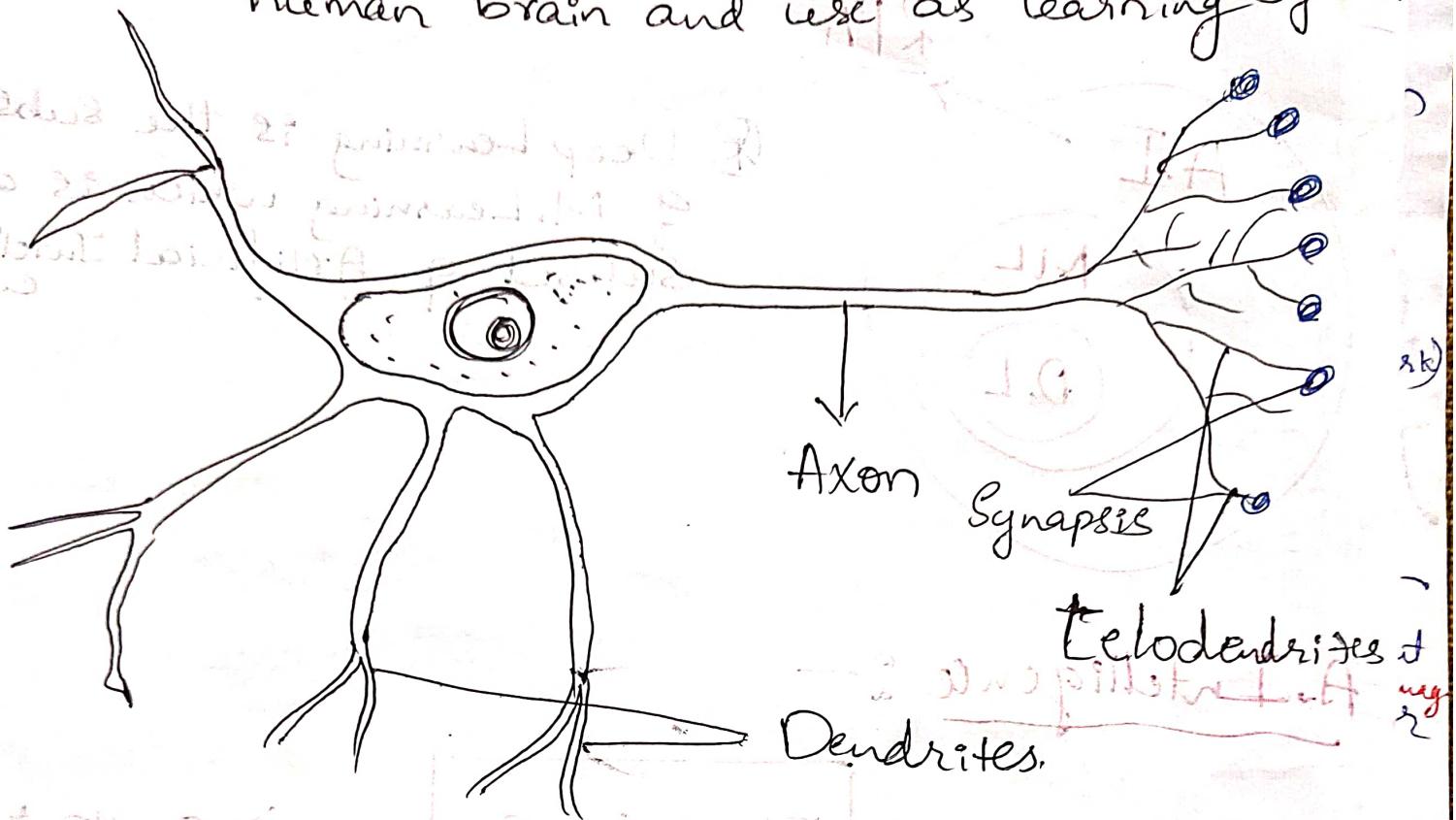
### Deep Learning :-



\* The algorithm automatically figures out the imp variables and create the rules.

# \* The Best Learning System is Human Brain

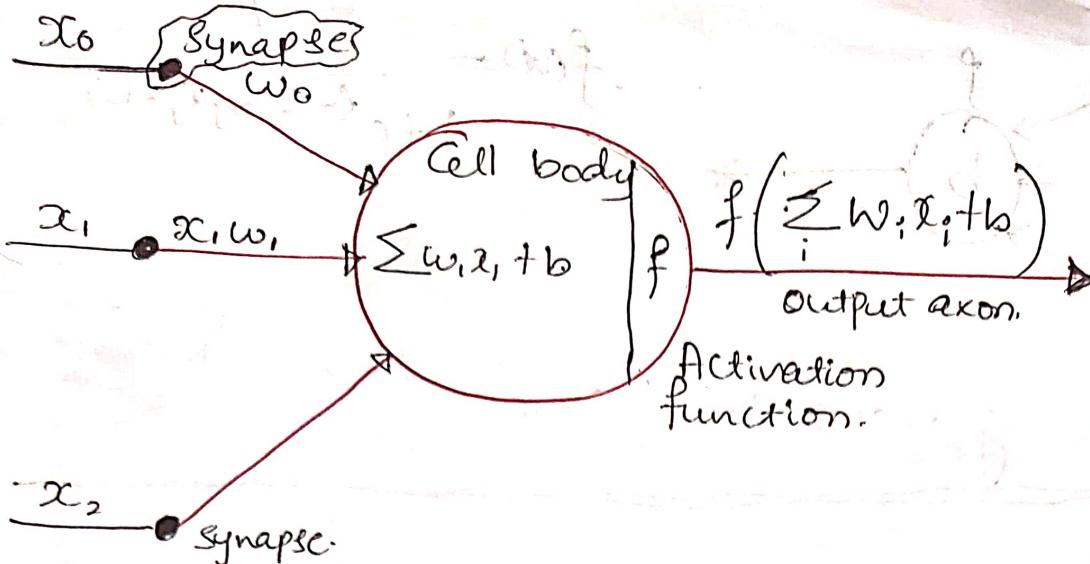
↳ Scientist were in the urge to replicate human brain and use as learning system



The direction of Electro-Chemical Signal from Dendrites to telodendrites by the diff concentration  $\text{Na}^+$  &  $\text{K}^+$  ions

\* The information in the brain is stored in the Synaptic connection.

\* In 1949, Donald Hebb tried to create the Mathematical model of Neuron called Perceptron.



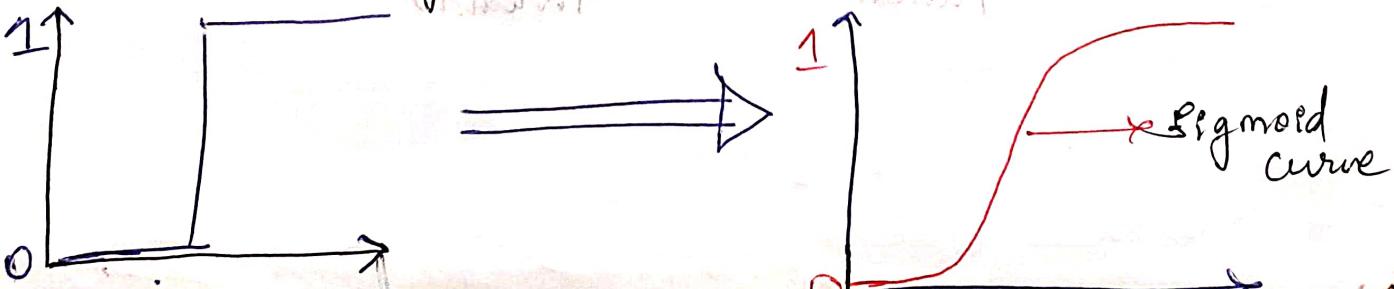
- \* The input from the Synapse is Attenuated (adds weight) and sent into the body.
- \* The output will be 1 if the Signal coming in meets the threshold, otherwise zero(0).

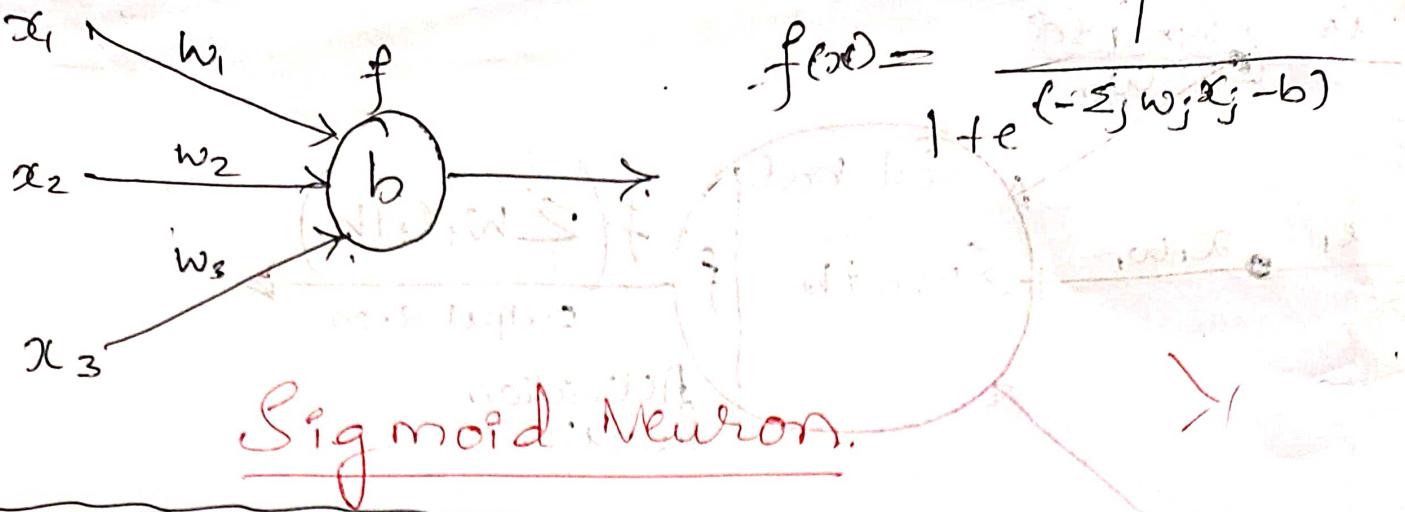
$$\text{Output} = \begin{cases} 0 & \text{if } \sum w_i x_i \leq \text{threshold} \\ 1 & \text{if } \sum w_i x_i > \text{threshold.} \end{cases}$$

- \* For each problem neurons will have different set of weights, for figuring out the right set of weights, tuning of the neuron must be done.

This tuning process is difficult if we have digital output (0, 1).

∴ we modify this neuron to a Sigmoid neuron which has output from '0' to '1' in gradually increasing manner.

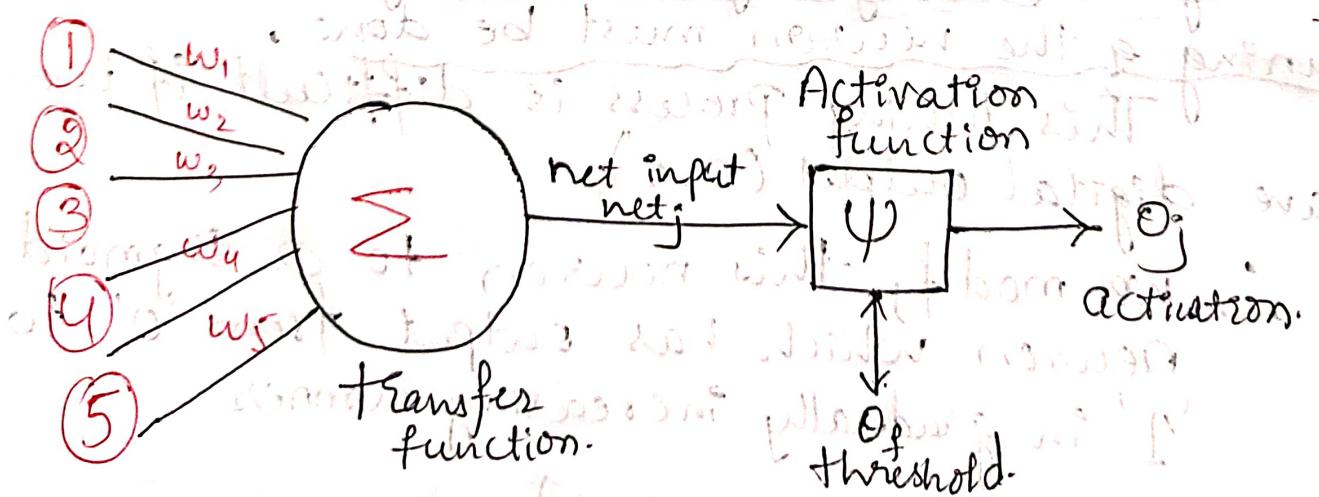




Example:-

Independent Variable	weight coeff
1 Salary	$w_1$
2 Credit rating	$w_2$
3 Age	$w_3$
4 Education	$w_4$
5 Account balance	$w_5$

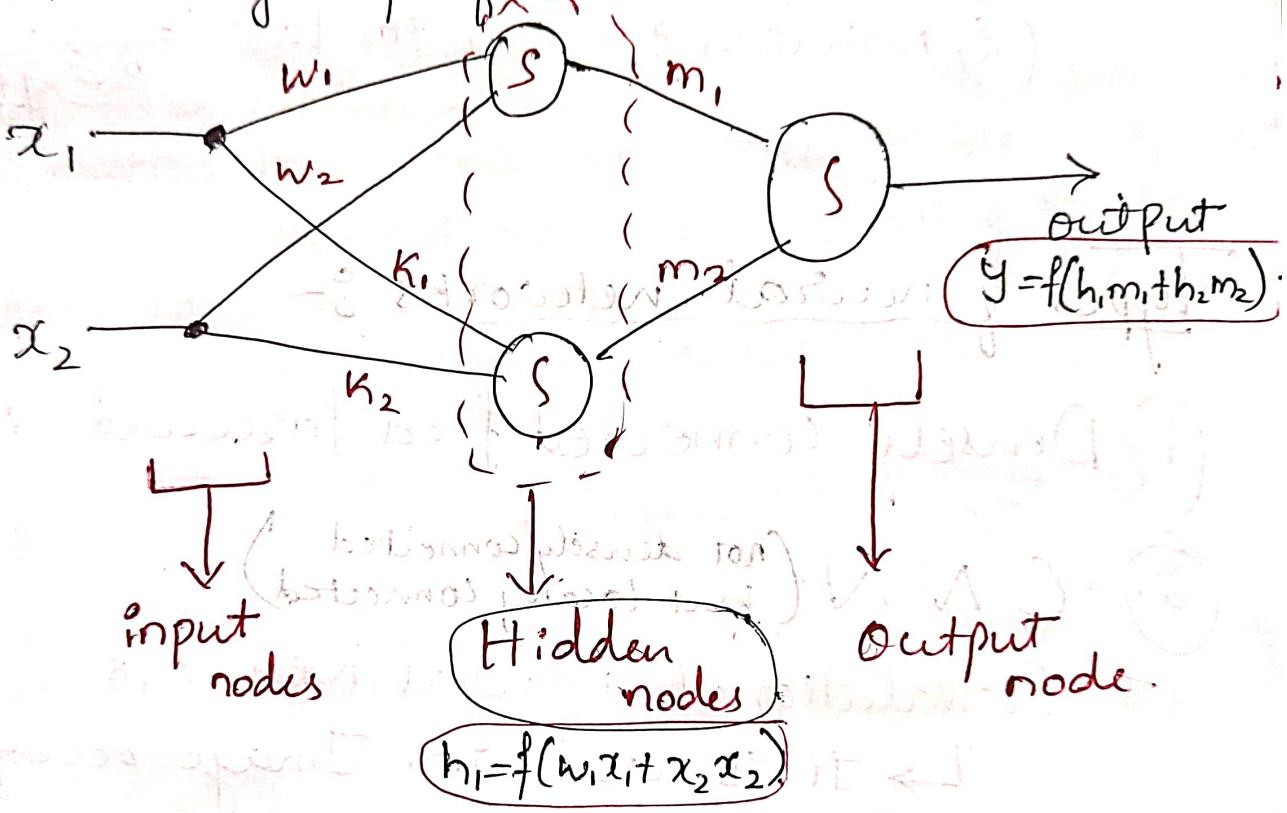
To determine loan approval



\* This artificial neuron does the same job of Log. regression

∴ Sigmoid neuron failed to be different  
 → it also failed in XOR outputs

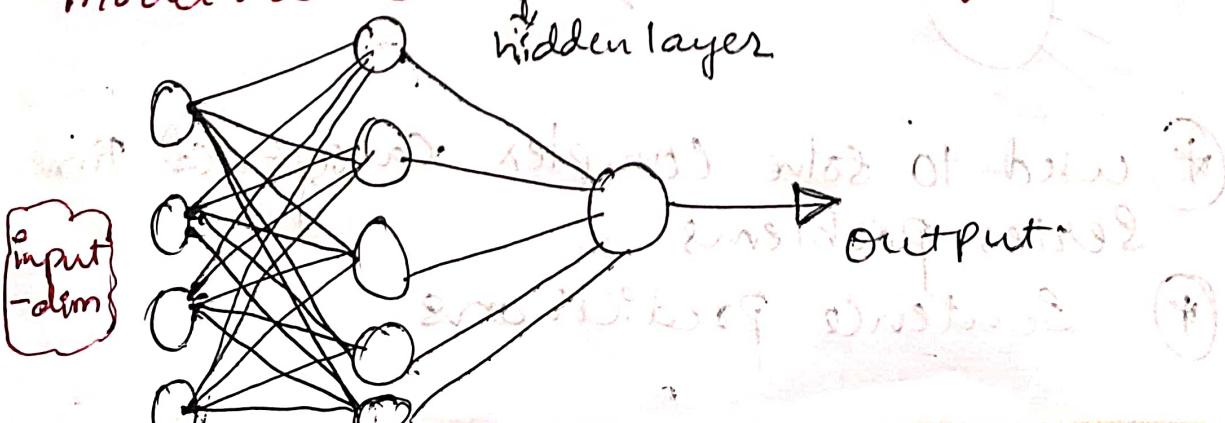
\* Due to the failure of perceptron, scientists thought of trying with the group of neurons



model = Sequential()

model.add(Dense(5, input\_dim=4, activation='sigmoid'))

model.add(Dense(1, activation='sigmoid'))



## Output Layer

- \*) Activation function in output layer
  - ① Linear  $\rightarrow$  for regressions.
  - ② Sigmoid  $\rightarrow$  classification problems.

## (\*) Multiclass Classification problems

- ① Multiple output neurons
- ② Activation function for last layer  $\rightarrow$  Softmax.

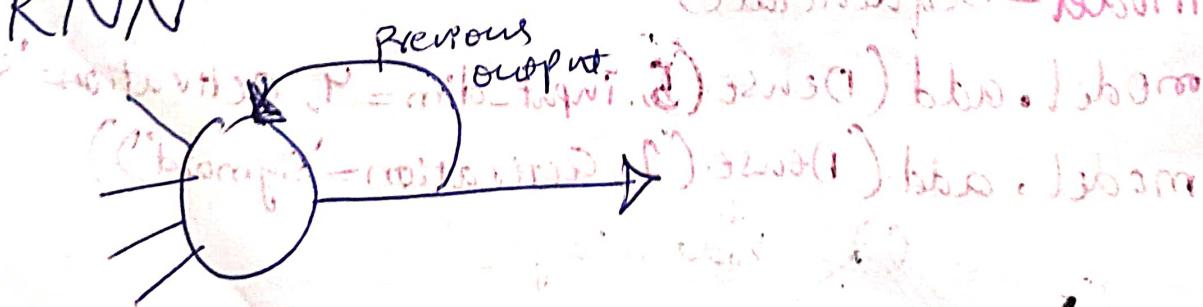
## Types of neural networks :-

- ① Densely connected feed forward N.N
- ② C N N (not densely connected but locally connected)
  - height width
  - number of neurons

Convolutional Neural Network

$\hookrightarrow$  If it is used in Image recognition

- ③ Recurrent Neural network RNN



- \*) used to solve complex categorical time series problems
- \*) Sentence predictions

# Hyperparameters in NN

- ① Network structure.
- ② Activation function.
- ③ Number of hidden layers.
- ④ Number of neurons  $\rightarrow$  through cross validation.
- ⑤ Cost function.

\* There is a mathematical theorem which guarantees no matter how complicated the problem is, we can achieve desired accuracy by using single hidden layer neural network (by setting the hyper parameter  $\rightarrow$  no. of neurons in hid. layer)

But why need more layers  $\rightarrow$

- $\rightarrow$  Learning will become slower in one layer
- $\rightarrow$  One layer learn one pattern and other will learn other patterns so learning becomes easy.
- $\rightarrow$  The power of backprop will begin in many layers.

# How to Assign the weights to the neuron

→ we use Gradient decent method to assign the correct set of weights

$$w_1^0 = w_1^0 - \delta \frac{\partial SSE}{\partial w_1}$$

→ If this value is positive,

then that means if weight increases the error will increase

and vice versa

$$w_1^0 =$$

→ Firstly we assume the random set of weights and get the error. After this we try to increase or decrease the weights accordingly by using above formula till we reach the least error and consider the corresponding weights.

$$G_1 = G_1 - \delta \nabla SSE$$

ANN → Artificial Neural network

multiple method of updating weights :-

Full batch → Regular Training as done for all other ML Algos. Entire input is used to train. Update weights using Grad Descent.

Online → Show one input row at a time... adjust weights, show another input and adjust weights... Once the input is all over, start with row 1 if needed.

Mini-Batch → Pick a Small Random Sample of inputs perform batch update most preferred. Then pick another set of samples randomly... Eg. perform update.

```
import Keras  
from Keras import regularizers, optimizers  
from Keras.layers import Dense  
from Keras.callbacks import EarlyStopping,  
ReduceLROnPlateau
```

```
model-Sig = Sequential() → 1st hidden  
model-Sig.add(Dense(10, input_shape=(21,),  
activation='sigmoid'))  
model-Sig.add(Dense(18, activation='sigmoid'))  
model-Sig.add(Dense(1, activation='sigmoid',  
kernel_regularizer=regularizers.  
l2(0.04)))  
2nd hidden  
output layer
```

```
from Keras.optimizers import SGD
```

```
Opt = SGD() → Adam or RMSprop → Adam-initM  
model-Sig.compile(loss='binary_crossentropy',  
optimizer=opt, metrics=['accuracy'])  
model-Sig.fit(x_train, y_train, epochs=30,  
batch_size=100)
```

# \* What is Deep Learning?

1 Deep learning is a Neural Network.

which has more than 2 hidden layers.

$N > 2$  is Deep

Less data  
computational power

→ Since we know that from math. theorem a single layer is enough with 'n' neurons to achieve desired Accuracy, but in this case the learning in the neurons is like memorising, but in deep learning we need shallow learning to understand the data set

Therefore we use greater than 2 hidden layers which enables neurons to learn general patterns and understand the dataset better

## → Problems in Deep learning

### ① Vanishing gradients :-

As we add more hidden layers, back-propagation becomes less and less effective in passing information to the gradient. Lower layers are scaled down exponentially.

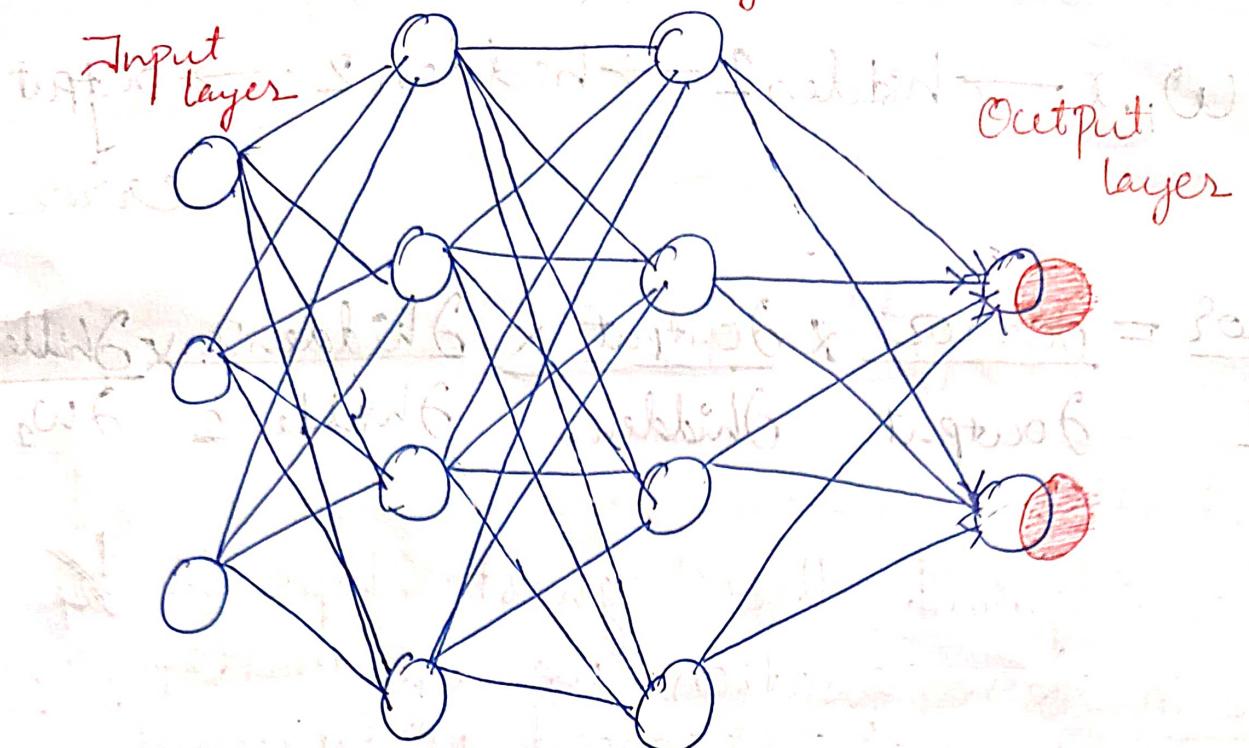
→ ReLU Activation Function will solve this problem.

### ② Difficulty in optimization :-

In deeper networks, the parameter space is large, and the surface being optimized has complex structures. Gradient descent becomes very inefficient.

### ③ Overfitting :-

- training accuracy ~~is~~ becomes high
- test accuracy ~~is~~ becomes low



④ Compute the error in the output layer

Idle  
Point  
ep grad  
allow grad

$\Sigma \text{sq}(y_i - \hat{y}_i)$

outwards diff.

process for layer

bottom

forward pass

backward pass

compute error

# Understanding the Vanishing Gradient

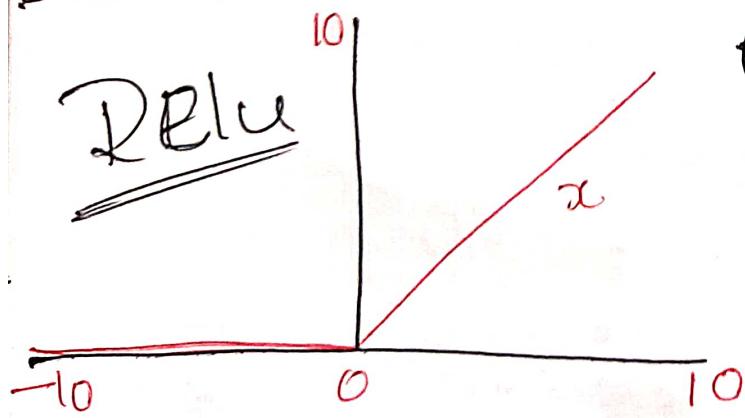
$W_1 \cdot x \rightarrow \text{hidden 1} \rightarrow \text{hidden 2} \rightarrow \text{output}$   
↓  
error

$$\frac{\partial \text{error}}{\partial W_1} = \frac{\partial \text{error}}{\partial \text{output}} \times \frac{\partial \text{output}}{\partial \text{hidden 2}} \times \frac{\partial \text{hidden 2}}{\partial \text{hidden 1}} \times \frac{\partial \text{hidden 1}}{\partial W_1}$$

For Sigmoid, the highest slope is  $\frac{1}{4}$

∴ in above equation we are multiplying one derivative of Sigmoid with another derivative of Sigmoid, therefore the product goes on decreasing since the slope is  $\frac{1}{4}$ th.

∴ If we add more & more layers the derivative values are getting smaller & smaller ∴ termed as Vanishing Gradient.



The derivative of the ReLU slope ( $x$ ) is 1

∴ the derivative will not become smaller.

∴ The hidden layer must be built using ReLU Activation function.

~~Instance name = class name~~

~~S<sub>i</sub> = Student (i)~~

## Optimization Difficulty :-

→ The surface over which we need to find minimum can be very complicated.

① Saddle points → The point which is min for one and max for another layer

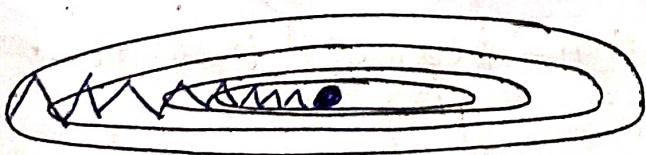
② Steep gradients

③ Extremely shallow gradients

→ Since in deep neural network we will be searching for minima in million dimensions this simple Gradient descent will not solve the problem of optimization.

In normal G.D.

→ Below shown contour plot to find the minima



④ We take tangent from each circle and measure gradient. After measuring we move in opp direction of gradient to find the minima.

⑤ Since heights of the slope varies, we move in bit diff direction and finally reach the minima which takes much longer duration.

∴ Idea of momentum came into consideration

④ Vector addition is done using parallelogram & triangular law of vector addition.



⑤  $\therefore$  Instead of taking the latest gradient, we take the average of past few gradients.



Gradient descent with Momentum

```
from keras.optimizers import SGD, RMSprop
model = Sequential()
model.add(Dense(10, input_dim=5, activation='relu',
               kernel_initializer='normal'))
model.add(Dense(2, activation='sigmoid', kernel_initializer='normal'))
Opt = SGD(momentum=0.9)
model.compile(loss='binary_crossentropy',
              optimizer=Opt, metrics=['accuracy'])
model.fit(X_train, y_train)
```

\* G. helton thought of taking the step depending on the gradient i.e., at steeper gradient, taking smaller step size and at larger gradient, taking larger step size and therefore combining both to get resultant step size direction vectors. This idea lead <sup>to</sup> get appropriately correct direction and named it as RMSprop.

→ RMSprop work well at Saddle point

Opt = RMSprop(learning\_rate=0.1)

model.compile(loss='binary\_crossentropy',  
optimizer=opt, metrics=['accuracy'])

model.fit(X\_train, y\_train, epochs=30, batch\_size=100)

\* The combined momentum (i.e both momentum & RMS prop) is called Adam (Adaptive moment)

→ It has two hyper parameters.

i) momentum ( $\beta_1$ )

ii) Squared momentum ( $\beta_2$ )

Adam with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,

learning\_rate =  $1e^{-3}$  or  $5e^{-4}$  is great starting point for many models.

Opt = Adam(learning\_rate=0.001, beta\_1=0.9,  
 $\beta_2=0.999$ , epsilon= $1e^{-07}$ , usegrad=False)

It is good to use Adam to overcome difficulty in optimization.

Weight initialization can do better instead  
of choosing random set of weights  
↑ techniques

→ It is the active area of research.

## → Xavier/Glorot initialization

In this weight initialization method,  
if the neuron is getting ~~many~~ inputs  
from many <sup>input nodes</sup> of them, then its weights  
should be smaller and if inputs are  
less, then the weights must be larger.

Model = Sequential()  
model.add(Dense(10, input\_dim=5, activation='relu'))

model.add(Dense(10, input\_dim=5, activation='relu'))  
Kernel\_initializer='glorot\_normal')

model.add(Dense(2, activation='Sigmoid',  
Kernel\_initializer='glorot\_normal'))

Model output and fit  
(train) and (test)

(loss) function base (2)

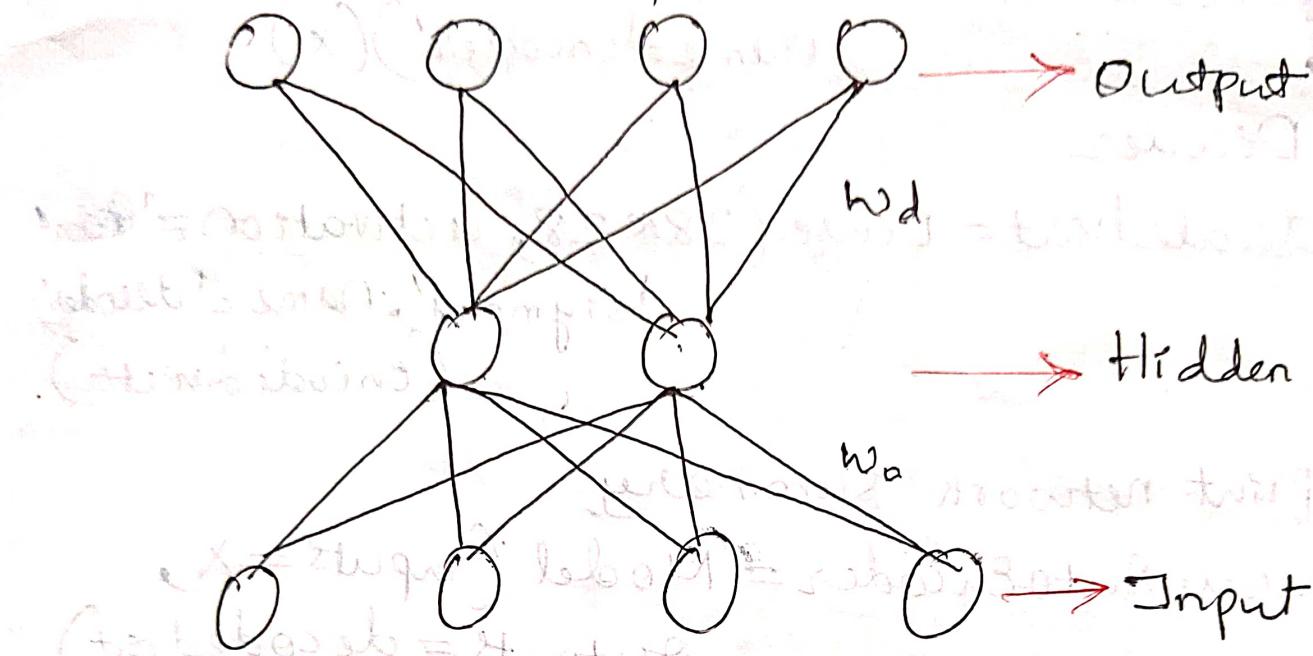
Model output and fit  
(train) and (test)

Model output and fit  
(train) and (test)

# AutoEncoders :-

Restricted Boltzmann  
Machines (RBM)

- AutoEncoder is a neural network which is used to solve unsupervised data. has same number of input & output nodes.



- In Auto Encoder the Output is same as Input.  
In hidden layer, reduced Outputs are sent to Output neurons which are same as input  
⇒ It says that the number of truly required inputs are two (in above case) since we reconstructed back from those.  
→ It has done the dimensionality reduction (non-linear)

- Therefore if input variables are correlated non-linearly, Auto Encoders can be used for dimensionality reduction

# Input

$X = \text{Input}(\text{name}='inputs', \text{shape}=[28 \times 28],$   
28x28,  $\text{dtype}=\text{float32})$

# Encoder

$\text{encodedvect} = \text{Dense}(32, \text{activation}=\text{'relu'},$   
 $\text{name}='encoder')(x)$

# Decoder

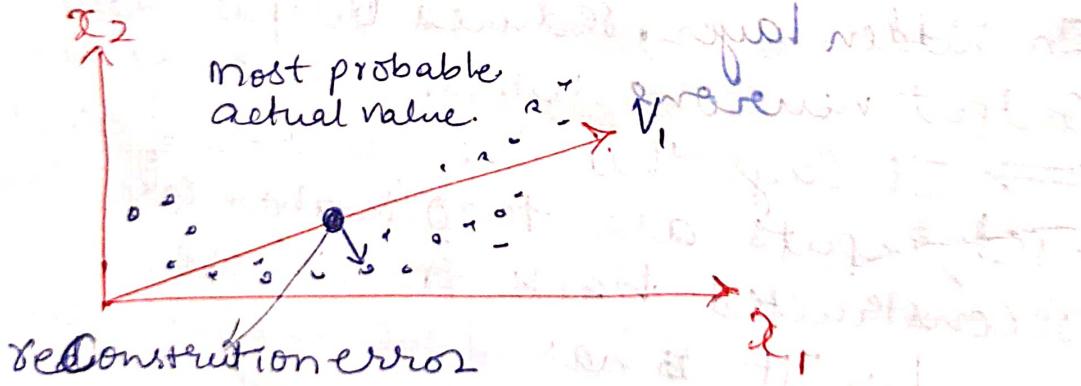
$\text{decodedvect} = \text{Dense}(28 \times 28, \text{activation}=\text{'tanh'}$   
 $\text{'sigmoid'}, \text{name}='decoder')$   
( $\text{encodedvect}$ )

# Print network summary

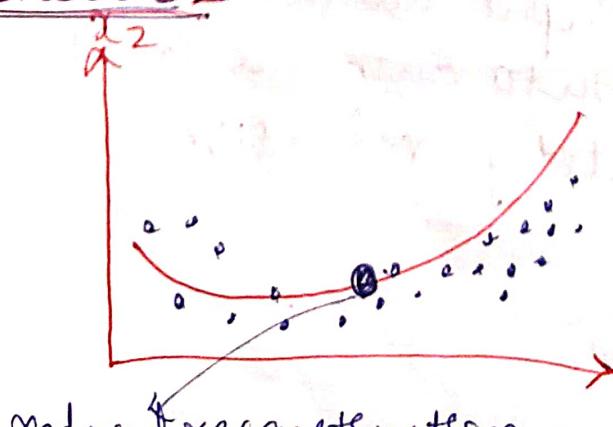
$\text{MNISTAutoEncoder} = \text{Model}(\text{inputs}=X,$   
 $\text{outputs}=\text{decodedvect})$

MNISTAutoEncoder.summary()

PCA



Auto Encoders



1, A normal auto-encoder with multiple Sigmoid/  
ReLU layers will map the input data into a  
non-linear subspace and have lower  
reconstruction error

④ Auto Encoder can be used as replacement of PCA.

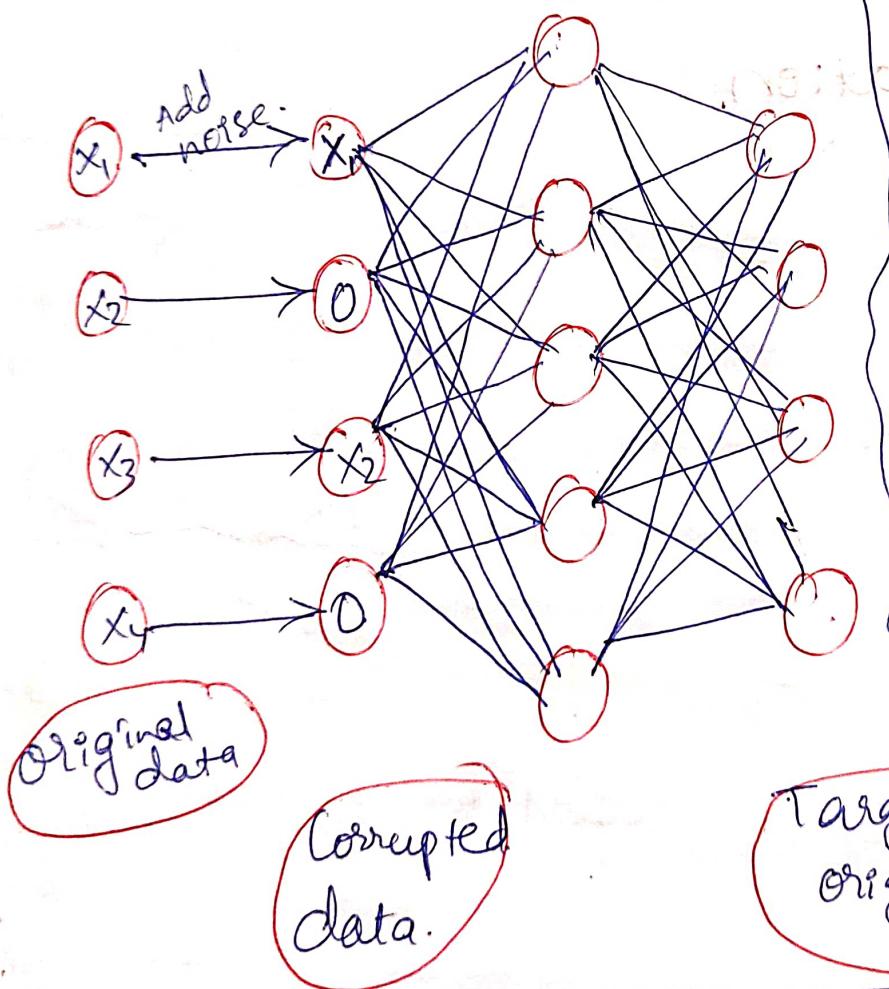
⑤ Auto Encoder is used to clean the noise in the image

↳ In each pixel the intensity are stored ranging from 0 - 255 (255 is most dark)

↳ Noise is the incorrect pixel intensity stored

⑥

Hidden layer



Denoising AutoEncoder

# Autoencoder application:-

## ① Image processing application

→ De-noising

→ Auto filling missing value.

## ② Anomaly detection

## ③ Feature generation

## ④ Learning generative model

## ⑤ Text translation

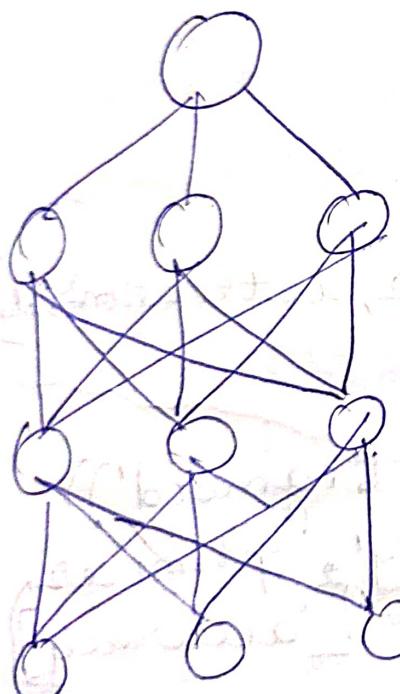
Autoencoder is trained with same input & output tuning. If any data containing other variable as input, then the output will be diff. Thus used in Anomaly detection



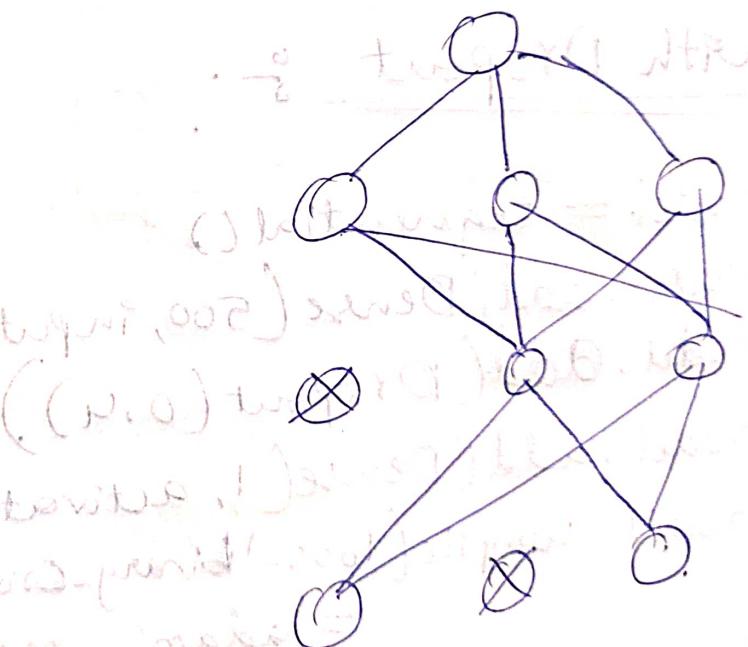
# Dropout

(Controlling overfit.)

- ↳ In this method, start with the complex neural network and after which cutting out the neurons randomly is done which helps neurons not to memorise but learn.
- ↳ This step is done using dropout technique during training the model.



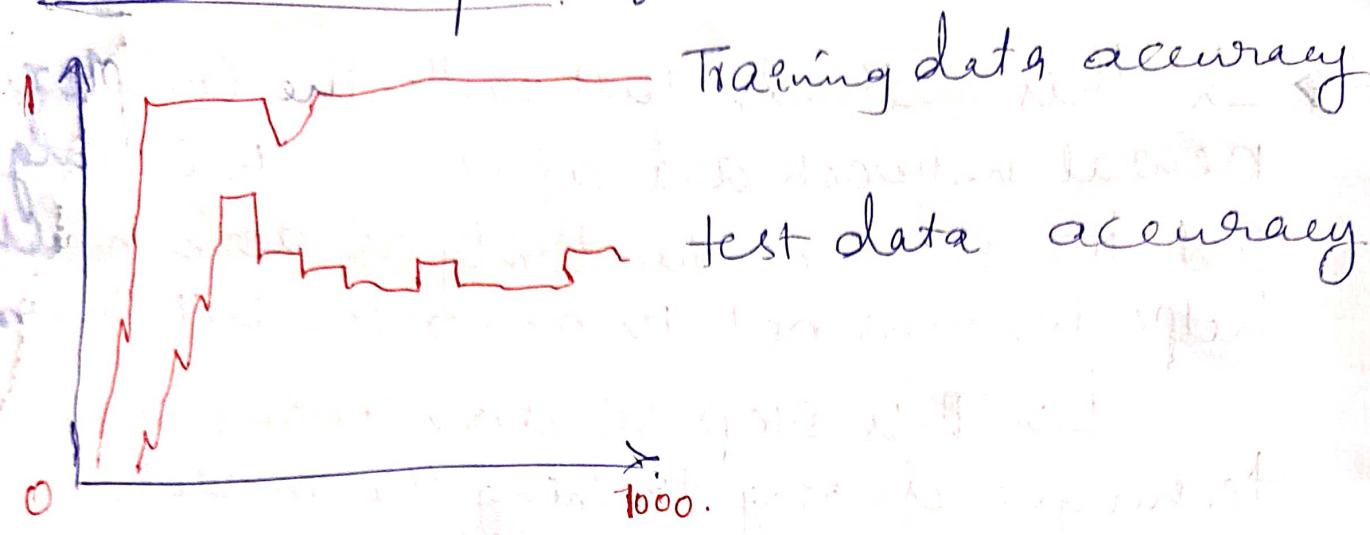
(a) Standard Neural Net



(b) After applying dropout

# dropout implementation.

without dropout :-



with Dropout :-

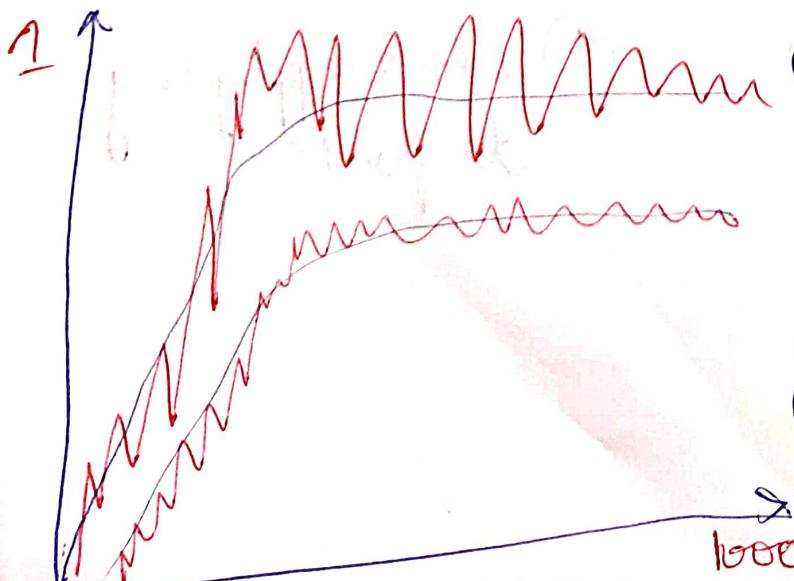
```
model = Sequential()
```

```
model.add(Dense(500, input_dim=2, activation='relu'))
```

```
model.add(Dropout(0.4))
```

```
model.add(Dense(1, activation = 'sigmoid'))
```

```
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
```



① Training will be lot more noisy becaz diff set of neurons are disappearing

∴ On an average there will be less diff b/w train & test accuracy.  
∴ more variance.

# Normalization :-

(Batch normalization)

- ④ usually scaling of the data is done to train the model faster although scaling is optional.
- ④ likewise in neural networks Batch normalization is done after every layers to train the model faster and for better accuracy.

model=Sequential().

Dense(64, input\_shape=(4,), activation='relu')

BatchNormalization()

Dense(128, activation='relu')

BatchNormalization()

Dense(128, activation='relu')

BatchNormalization()

Dense(64, activation='relu')

BatchNormalization()

Dense(1, activation='relu')

④ Can be applied to a state-of-the-art image classification model.

④ Achieves same accuracy with 14 times fewer training steps

④ Acts as a regularizer too.