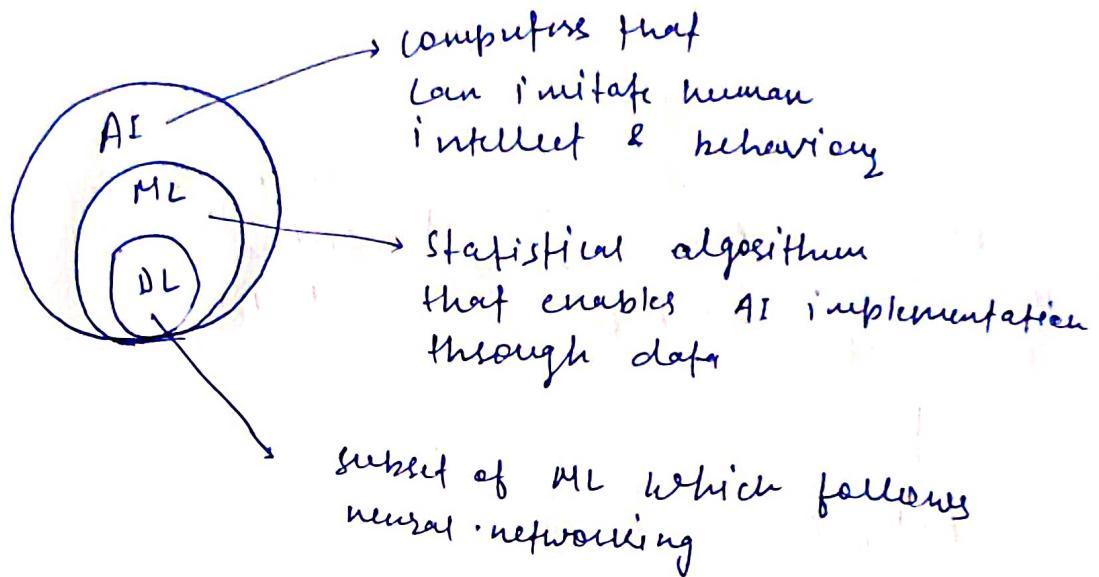


# AI

Smart knowledge

## # Basics

### \* AI vs ML vs DL



### \* ML

supervised

/  
Classification

- ↳ Logistic
- ↳ Decision tree
- ↳ Random forest
- ↳ KNNI
- ↳ SVM
- ↳ Naive Bayes

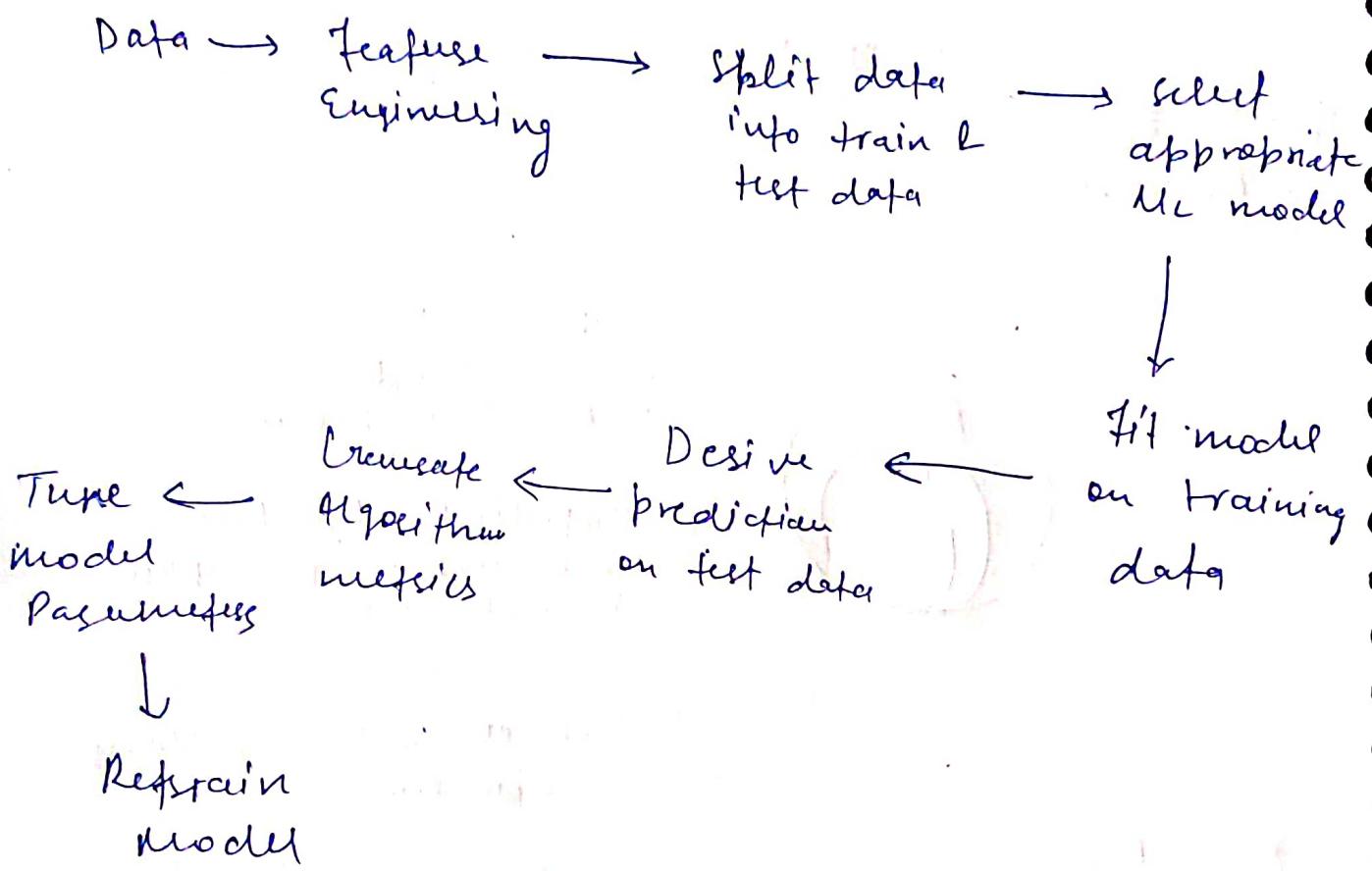
Regression

- ↳ simple linear reg.
- ↳ multiple linear reg
- ↳ Ridge reg
- ↳ Lasso reg

unsupervised

- ↳ K-means Clustering
- ↳ Hierarchical clustering

## \* Prediction pipeline



## \* Application of AI

- Image classification
- Image De-noising
- Object detection
- Chatbot
- Pose estimation
- Stock price prediction
- Image segmentation
- Text summarization
- Language translation

## \* Teachable machine

.. search infernet

## \* Working of a neuron

- A neuron accepts a signal and sends a message to the brain to act in a particular way.
- It accepts a signal, generates a particular response in correspondence from its nucleus and then triggers a response.
- The response may pass to another neuron and the cycle proceeds in this manner.

## + Neural Networks

- ① - Neurons
- ② - weights
- ③ - Bias
- ④ - Synapse
- ⑤ - Architecture of a neural network
- ⑥ - Layers in NN
  - a) Input layer
  - Hidden layers
  - Output layer
- ⑦ - General equation of a Neural network
- ⑧ - Activation functions
  - a) Linear
  - b) Sigmoid
  - c) ReLU
  - d) Leaky ReLU
  - e) Tanh
  - f) Softmax
- ⑨ - Optimizes & Optimization Techniques
  - a) Gradient descent
  - b) Stochastic Gradient descent (SGD)
  - c) Mini batch Gradient descent

- d.) SGD with momentum
- e.) Adagrad
- f.) Adadelta
- g.) RMS Prop
- h.) Adam

(10) - Loss function.

- a.) MSE
- b.) MAE
- c.) Binary Crossentropy
- d.) Categorical Crossentropy
- e.) Sparse-Categorical Crossentropy

(11) - forward Propagation

(12) Backward Propagation

(13) Epochs

(14) Classification Metrics - Accuracy, precision,  
Recall, F1-score

(15) Gradient descent from scratch (without tensorflow)

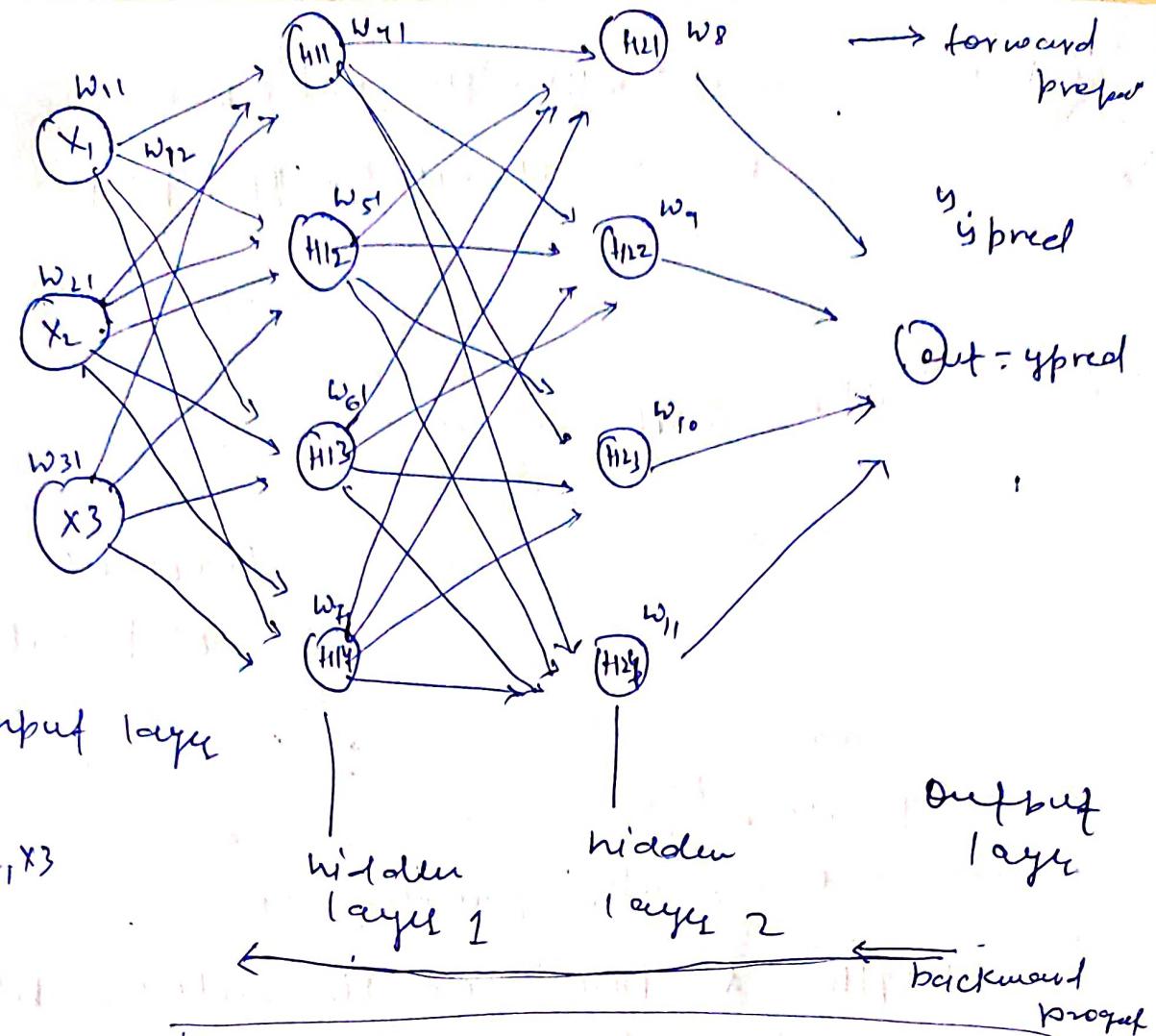
## Neural network working

- Neural networks are a set of algorithms, modeled loosely after human brain, that are designed to recognize patterns. They interpret data through a kind of machine perception and process input data to solve a certain problem. The patterns they recognize are numerical, contained in vectors, into which all-real world data i.e. text, images, sound can fit.
- Neural networks are a class of machine learning algorithms used to model complex patterns in datasets using multiple hidden layers & activation functions. A neural network takes an input, passes it through multiple layers hidden (mini-functions with unique coefficient that must be learned), and outputs a prediction representing the combined input of all the neurons.
- Neural networks are trained iteratively using optimization techniques like gradient descent. After each cycle of training, an error metric is calculated based on the difference between prediction and target. The derivatives of this error metric are calculated & propagated back-through the network using a technique called backpropagation. Each neuron's coefficient (weights) are then adjusted relative to how much they contributed to the total error. This process is repeated iteratively until the network error drops below an acceptable threshold.

- In a neural network, each neuron (except those in the input layer) is actually a sum of all its inputs; which are in fact the outputs from the previous layer multiplied by some weights. An additional term called bias is added to this sum. And a non-linear known as activation function is applied to the result.

## \* Layers in Neural Network

- ① - Input layer - Holds the data your model will train on. Each neuron in the input layer represents a unique attribute in your datasets (e.g. height, hair, color, etc.).
- ② - Hidden layer - sits b/w the input & output layers. and applies an activation function before passing on the result. There are often multiple hidden layers in a network. In general, hidden layers are typically full-connected layers - each neuron receives input from all the previous layer's neurons & sends its output to every neuron in the next layer.
- ③ - Output layer - The final layer in a ~~network~~ network. It receives inputs from the previous hidden layers, optionally applies an activation function, and returns an output representing your model's prediction.



$$Err = \frac{1}{n} \sum_{i=1}^n (y_i - y_{pred})^2$$

Total observation

$w_{ij}$  = weights

bias (constant term)

$$H11\_in = x_1 * w_{11} + x_2 * w_{21} + x_3 * w_{31} + b_{11}$$

$$H12\_in = x_1 * w_{12} + x_2 * w_{22} + x_3 * w_{32} + b_{12}$$

Input  $\uparrow$

$H11\_out = \text{Activation}(H11\_in)$

$H12\_out = \text{Activation}(H12\_in)$

output

$$H_{21\text{-in}} = H_{11\text{-out}} * w_{41} + H_{12\text{-out}} * w_{51} + H_{13\text{-out}} * w_{61} + H_{14\text{-out}} * w_{71} + b_{21}$$

- The update graph till now

→ forward propagation

propagation

till now

Iteration - 1

- ~~error~~ errors also generating here
- weights & bias are the entities which is going to update in every iteration to minimize errors by backward propagation.
- How weights & bias get updated optimization technique.

① - Gradient descent

$$w' = w - \alpha \frac{d(\text{err})}{d(w)}$$

- $w^8, w^9, w^{10}, w^{11}$  get updated first

- we add bias in case weights are 0, it will have no impact. So bias helps in shifting of Activation func. to left or right.
- Bias can be zero

\* Neuron - A neuron takes a group of weighted inputs applies an activation function, and returns an output. Inputs to a neuron can either be features from a training set or outputs from previous layer neurons.

Weight - Weights are applied to the inputs as they travel along synapses to reach the neuron & a bias term is added. The neuron then applies an activation function to the "sum of weighted inputs" added to the "bias" from each incoming synapse & passes the result on to all the neurons in the next layer.

Bias - Bias helps in controlling the value at which activation function will trigger. Bias value allows you to shift the activation func. to either left or right.

Activation function - ① - An Activation function is a neural network defines how the weighted sum of the input is transformed into an output from a node or nodes in a layer of the network.

② - An activation function is a function that is added into an artificial neural network in order to help the network learn complex patterns in the data.

③ - Activation function helps to normalize the output of any input in the range b/w -1 to +1 or 0 to 1

④ decides what is to be fired (activated) to the next neuron.

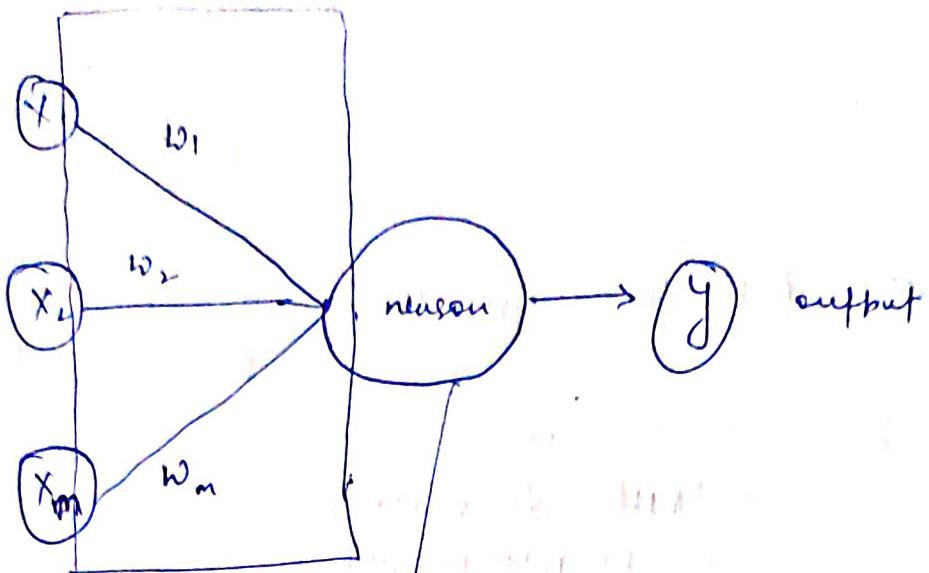
Synapse - Synapses are like roads in a neural network. They connect inputs to neurons to neurons, & neurons to outputs.

- Optimizers -
- ① - It's very important to <sup>→ change</sup> track the weights of the model during the training process, to make our prediction as correct & optimized as possible.
  - ② - This is where optimizers come into picture. They tie together the loss function & model parameters by updating the model in response to the output of the loss function.

### Loss function / cost function

- ① - It is a parameter in model's predict function that tells us "how good" the model is at making predictions for a given set of parameters.
- ② - The loss function has its own curve & its own derivatives. The slope of this curve tells us how to change our parameters to make the model more accurate. we use the cost function to update our parameters.

## \* General Equation of Neural Network



$$\sum_{i=1}^m w_i x_i$$

$$y_{\text{pred}} = \sigma(z) = \sigma(\sum_{i=1}^m w_i x_i + b)$$

where

$z = w \cdot x + b$

$w$  = weight

$x$  = input

$b$  = bias

$\sigma$  = Activation function

### \* Activation function

- properties
  - ① - Non-linear
  - ② - Differentiable

### \* Types

- Linear
- sigmoid
- Tanh
- ReLU

- Leaky ReLU
- ELU
- softmax

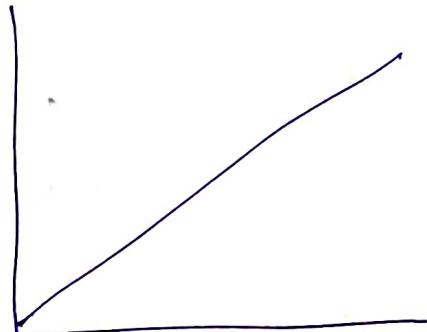
→ Linear Activation function

- Range  $-\infty$  to  $+\infty$
- It has constant derivative, so gradient has no relation with input

- Used in ANNs for regression

↓  
Artificial  
neural  
network

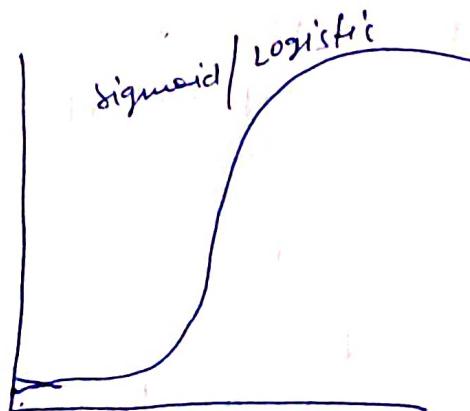
- $f(n) = n$



→ Sigmoid A. f.

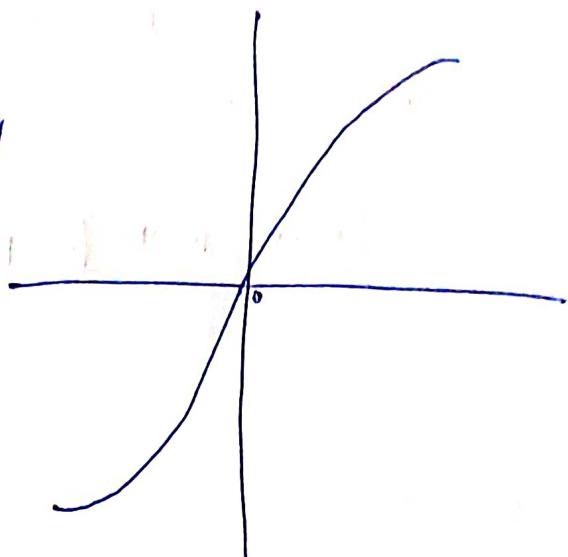
- Range 0 to 1
- Used for binary classification models
- It suffers from vanishing Gradient (small slopes)

- 
- $$f(n) = \frac{1}{1 + e^{-n}}$$



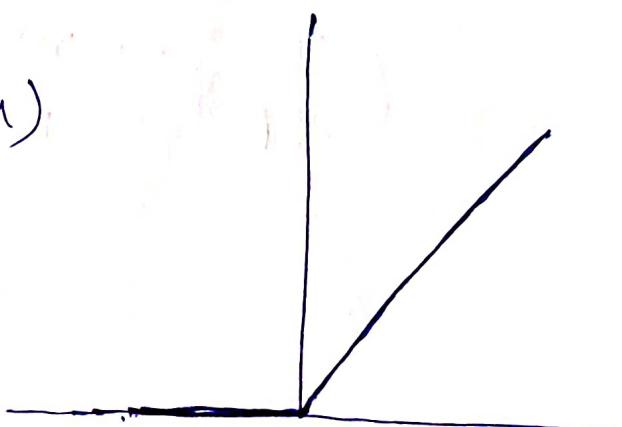
## → Tanh

- range -1 to 1
- used in RNNs
- It suffers from vanishing gradient.
- It is zero centered
- $f(u) = \frac{(e^u - e^{-u})}{(e^u + e^{-u})}$



## → ReLU

- ReLU means rectified linear unit.
- Range 0 to +infinity
- used in hidden layers
- It suffers from dying ReLU problem
- $f(x) = \max(0, u)$



→  $HII\_out = \text{Activation}(HII\_in)$

$$HII\_in = -vc$$

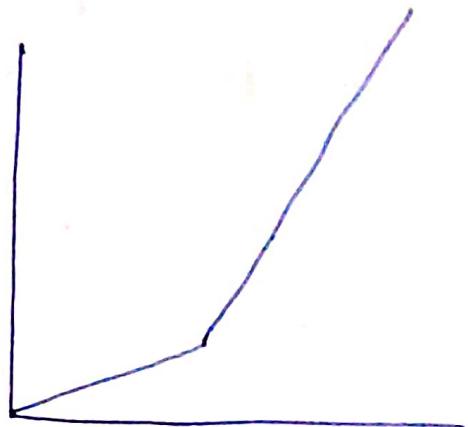
$$\text{Activation} = \text{ReLU}$$

$HII\_out = 0$   
→ so no combination of this neuron; so it's not activated

∴ suffers dying ReLU problem

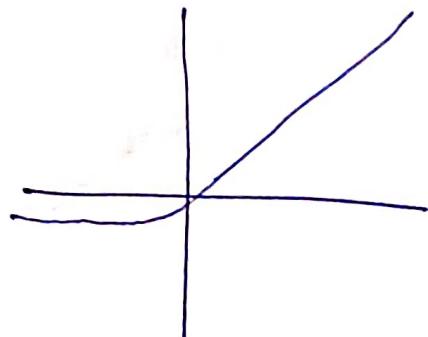
## → Leaky RELU

- Leaky RELU is another variant of RELU
- used in hidden layers
- it overcomes the dying RELU problem
- $f(u) = \max(0.1u, u)$



## → ELU

- it stands for Exponential Linear Unit
- used in hidden layers
- it overcomes the Dying RELU problem
- Suffers from exponential gradient.
- $$\begin{cases} u & \text{for } u \geq 0 \\ \alpha(e^u - 1) & \text{for } u < 0 \end{cases}$$



- → Activation function used in output layers is linear

## \* Vanishing gradient

- like the sigmoid function certain activation functions squash an ample input space into a small output space b/w '0 & 1'
- Therefore, a large change in the input of the sigmoid function will cause a small change in the output. Hence, the derivative becomes small. for shallow networks with only a few layers that use these activations, this isn't a big problem.
- However, when more layers are used, it can cause the gradient to be too small for training to work effectively.

## \* Exploding gradient

- Exploding gradient are problems where significant error gradients accumulate & result in very large updates to neural network model weights during training
- An unstable network can result when there are exploding gradients, & the learning cannot be completed.
- The values of the weights can also become so large as to overflow & result in something called NaN values.

- \* Optimizers
- point ① & ② already discussed previously
  - Optimizers shape & mold your model into its most accurate possible form by maneuvering with the weights. The loss function is the guide to the terrain, telling the optimizer when it's moving in the right or wrong direction.

## ③ Optimizers functions

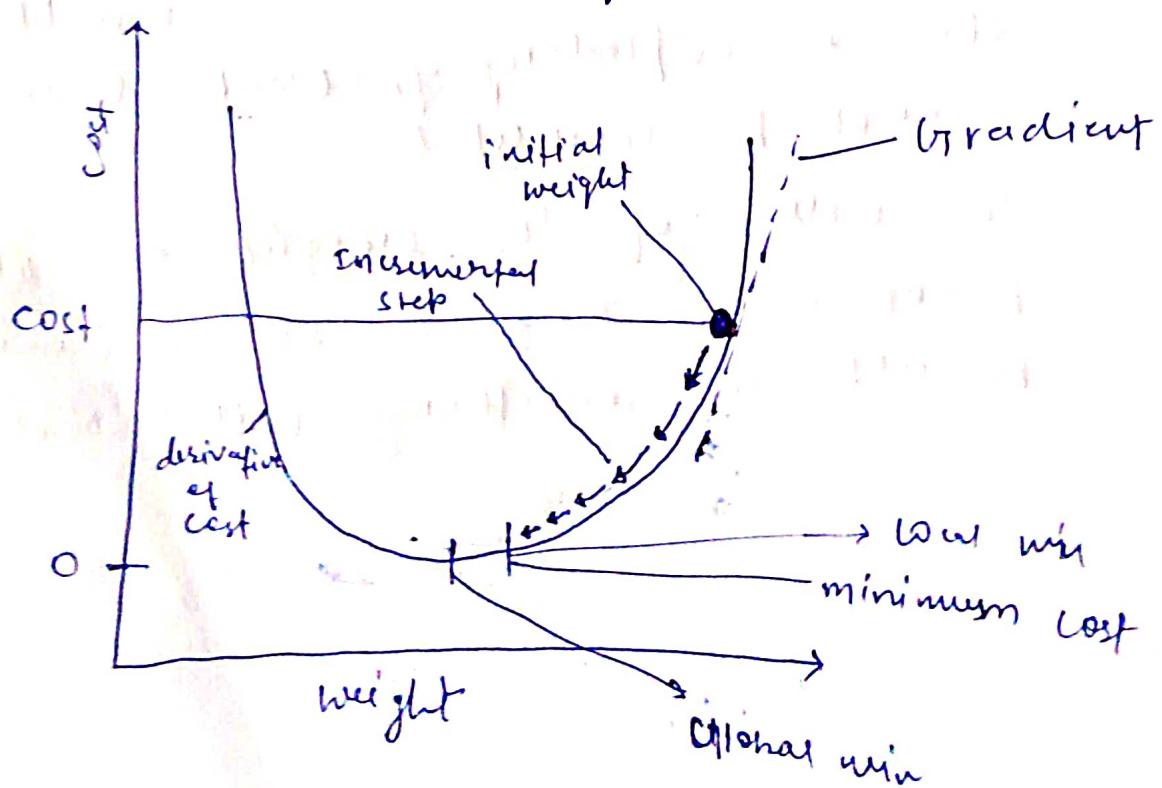
### → Gradient Descent

- It minimizes the given function to its local minimum.
- Gradient Descent iteratively reduces a loss function by moving in the direction opposite to that of steepest ascent.

$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{\delta(\text{Loss})}{\delta(w_{\text{old}})}$$

$\downarrow$   
Learning Rate

$\boxed{\text{Loss} = \text{Cost} = \text{Error function}}$

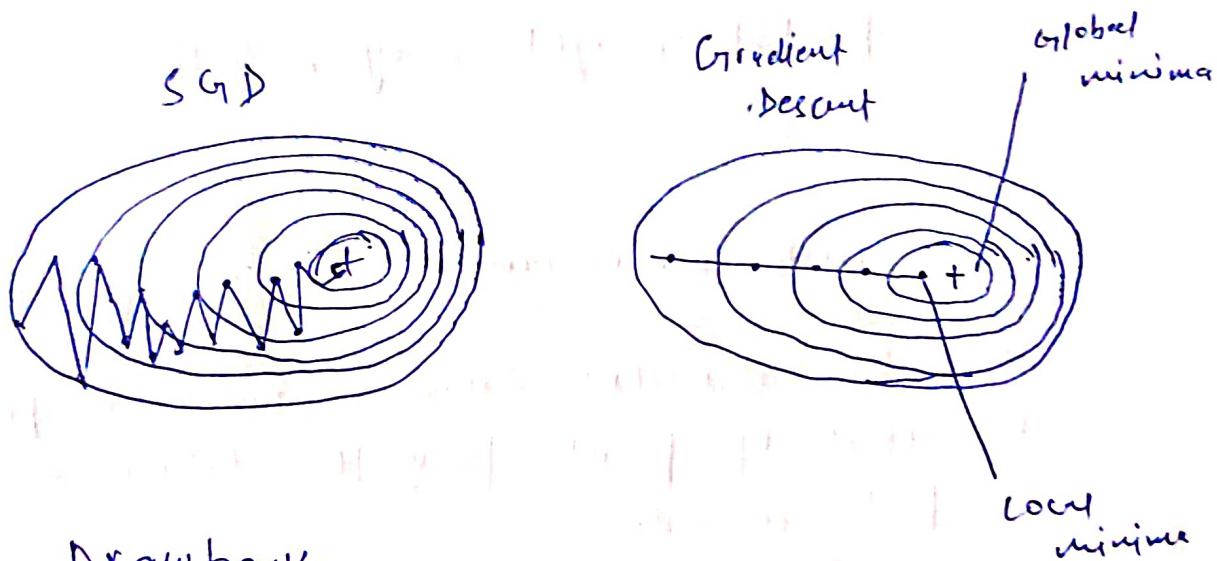


## Drawbacks

- slow in computation
- Requires large memory
- computationally expensive

## → Stochastic Gradient Descent

- In the SGD algorithm derivative is computed taking one point at a time.
- Memory requirement is less than that of GD



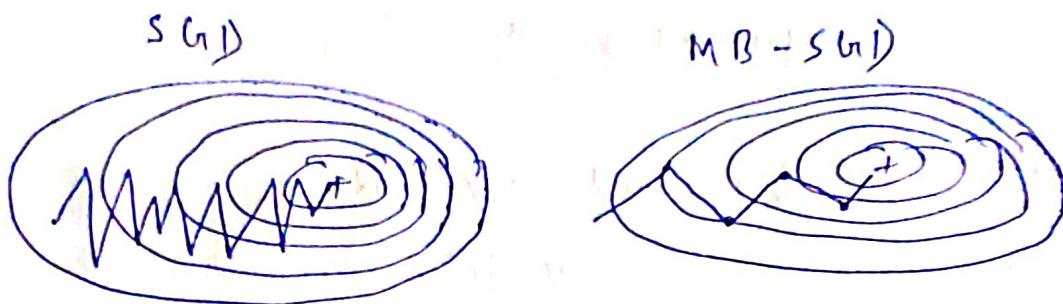
## Drawbacks

- Takes a long time to converge
- May stuck at local minima

## → Mini-batch Gradient Descent

- MB-SGD algorithm takes a batch of points or subset of points from the dataset to compute derivative.
- MB-SGD divides the dataset into various batches & every batch, the parameters are updated.

- less time complexity as compared to SGD & more stab.



### Drawbacks

- may stuck at local minima
- updated weights may be too noisy

→ SGD with momentum

- MB-SGD algorithm takes a batch of points as subset of points from the dataset to compute derivate.
- MB-SGD divides the dataset into various batches & after every batch, the parameters are updated.
- The idea is to denoise derivate using exponential weighting average that is to give more weightage to recent update compared to the previous update.

$$V_{dw} = \beta \cdot V_{dw} + (1-\beta) \cdot d_w$$

$$V_{db} = \beta \cdot V_{db} + (1-\beta) \cdot d_b$$

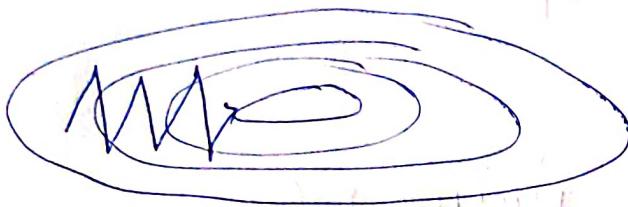
$$w = w - \alpha \cdot V_{dw}$$

$$b = b - \alpha \cdot V_{db}$$

- The problem with SGD is that while it tries to reach minima because of the high oscillation we can't increase the learning rate so it takes time to converge

### Drawbacks

- may stuck at local minima
- updated weights may be too noisy



SGD without momentum



SGD with momentum

### → Adam (adaptive Momentum Estimation)

- It is a method that computes adaptive learning rates for each parameter. It stores both the decaying average of the past gradients, similar to momentum & also the decaying average of the past squared gradients, similar to RMS-Prop & Adadelta

$$\begin{cases} V_{dw_t} = \beta V_{dw_{t-1}} + (1-\beta) \frac{\partial L}{\partial w_{t-1}} \\ V_{db_t} = \beta V_{db_{t-1}} + (1-\beta) \frac{\partial L}{\partial b_{t-1}} \end{cases}$$

exponential  
weighted avg for  
past gradients

$$S_{dw_t} = \gamma S_{dw_{t-1}} + (1-\gamma) \left( \frac{\delta L}{\delta w_{t-1}} \right)^2$$

exponential avg for past squared gradient.

$$S_{db_t} = \gamma S_{db_{t-1}} + (1-\gamma) \left( \frac{\delta L}{\delta b_{t-1}} \right)^2$$

$$w_t = w_{t-1} - \frac{n}{S_{dw_t} + \epsilon} \times V_{dw_t}$$

$$b_t = b_{t-1} - \frac{n}{S_{db_t} + \epsilon} \times V_{db_t}$$

Adam  
optimizing

alpha = neta = learning rate

$$w' = w - \text{alpha} \times \left( \frac{\partial L}{\partial w} \right)$$

## \* Loss functions

loss or cost or error function

- In a neural network; the weights get multiplied with the inputs & then activation function is applied to the element before going to the next layer. Finally we get the predicted value ( $\hat{y}$ ) through the output layer. But prediction is always closer to the actual ( $y$ ), which we define the loss/cost functions to capture

the errors & try to optimize it through backpropagation.

- The error function that compares actual values with its corresponding predicted value is referred to as loss function or cost function

- Types

- MSE (mean square error)

- It is used in regression

- MSE is the 'preferred' loss function if the distribution of the target variable is Gaussian. Mean square error is defined as the average of the squared differences b/w the predicted & actual values

$$L = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

→ predicted value

n is batch size.

- MAE (Mean absolute error)

- It is used in Regression

- It is defined as the average of the absolute difference b/w the actual & predicted values.

$$L = \frac{1}{n} \sum_{i=1}^n |y - \hat{y}|$$

## → Binary Cross Entropy

- It is used in Binary classification.
- It measures the difference b/w two probability distributions. If the cross entropy is small, it suggests that two distributions are similar to each other.
- In case of a binary classification the predicted probability is compared to the target/actual (0 or 1). Binary cross entropy calculates a score that provides the negative average difference b/w the actual & predicted probabilities for predicting the class. 1. This score penalizes the probability based on the distance from the expected value.

$$L = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

$$= \begin{cases} -y * \log(\hat{y}), & \text{if } y=1 \\ -(1-y) * \log(1-\hat{y}), & \text{if } y=0 \end{cases}$$

## → Categorical Class Entropy

- It is used in Multi-class classification.
- In case of a multi-class classification

the predicted probability is compared to the target/actual, where each class is assigned a unique integer value ( $0, 1, 2, \dots, t$ ), assuming data has  $t$  unique classes. It calculates a score that provides the negative average difference b/w the actual & the predicted probabilities for all classes

- for multi-class Cross Entropy ; actual targets ( $y$ ) are one-hot encoded. For a 3-class classification  $[0,0;1], [1,0;0] [0;1,0]$

$$L = - \sum_{j=1}^t y_{ij} * \log(\hat{y}_{ij})$$

$i$  = no. of rows

$j$  = no. of categories

$y_i$  = is one hot encoded target vector

### → Categorical Cross Entropy

- it is used in multiclass classification
- in case of a multi-class classification the predicted probability is compared to the target/actual, where class  $i$  is assigned a unique integer value ( $0, 1, 2, \dots, t$ ), assuming data has  $t$  unique classes. It calculates a score that provides the negative average difference b/w

the actual & predicted probabilities for all classes.

- for multi-class cross entropy, actual targets ( $y$ ) are one-hot encoded. For a 3-class classification

$$[[0,0,1], [1,0,0], [0,1,0]]$$

$$L = - \sum_{j=1}^t y_{ij} * \log(\hat{y}_{ij})$$

where  $i$  = no. of news

$j$  = no. of categories

$y_{ij}$  = is one hot encoded target vector.

## # CNN (Convolutional Neural Network)

- CNN image classifications take an input image, process it & classify it under certain categories (E.g., Dog, Cat, Tiger, Lion).
- Computers see an input image as array of pixels & it depends on the image resolution. Based on the image resolution, it will see  $h \times w \times d$  ( $h$ =height,  $w$ =width,  $d$ =dimensions). E.g. An image of  $6 \times 6 \times 3$  array of matrix of RGB (3 refers to RGB values) & an image of  $4 \times 4 \times 1$  array of matrix of grayscale image.
- When Deep learning CNN model's are applied to train & test, each input image will pass through a series of convolution layers with filters (Kernels), Pooling, flatten layer, fully connected layers (FC). ~~pass through~~ and apply softmax function to classify an object with probabilistic values b/w 0 & 1. For multiclass classification, we can use <sup>loss as</sup> categorical cross entropy, sparse categorical cross entropy.

|  
label encoded  
data

→ used for  
? not encoded  
data

## \* Convolution

- Let us suppose this in the input matrix of  $5 \times 5$  & a filter of matrix  $3 \times 3$ , for those who don't know what a filter is a set of weight applied on an image or a matrix to obtain the required feature.

- In convolutional layer, the image pixels (image matrix) are multiplied with filter matrix to generate a feature matrix.
- Consider the handwritten digit 2 & its pixel representation.



0	0	0	0	0
0	24	20	0	0
0	195	0	0	0
0	99	0	0	0
0	145	27	184	140
0	0	0	0	0

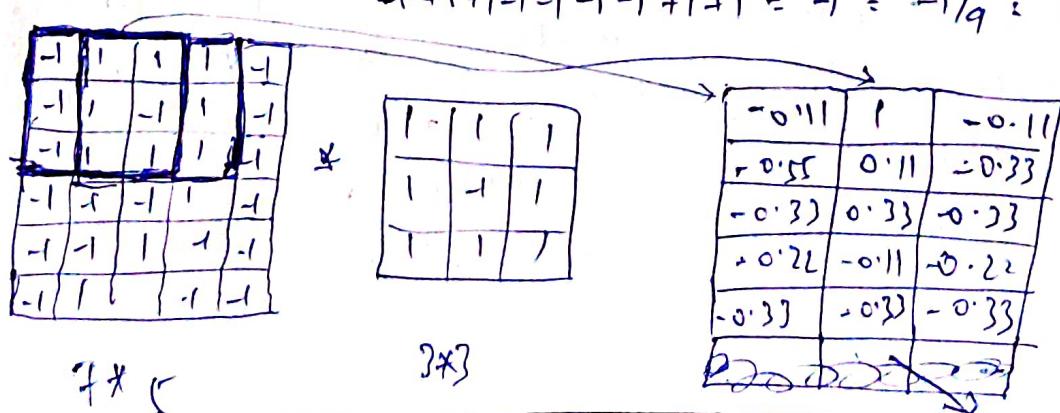
- For sake of simplicity of mathematical operations. consider image matrix of '9' as follows

-1	1	1	-1	-1
-1	1	1	-1	-1
-1	1	1	-1	-1
-1	1	1	-1	-1
-1	1	1	-1	-1

-1	1	1	-1	-1
-1	1	1	-1	-1
-1	1	1	-1	-1
-1	1	1	-1	-1
-1	1	1	-1	-1

- Convolution operation is applied using a filter. Let's assume that a  $3 \times 3$  filter is used for this operation. In convolution layer, the image pixels (image matrix) are multiplied with filter (kernels) matrix to generate feature map.

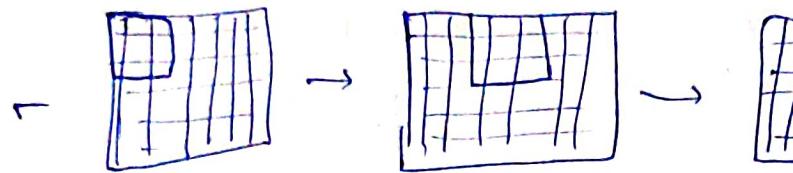
$$-1 + 1 + 1 - 1 - 1 - 1 + 1 + 1 = -1 = -1/9 = -0.11$$



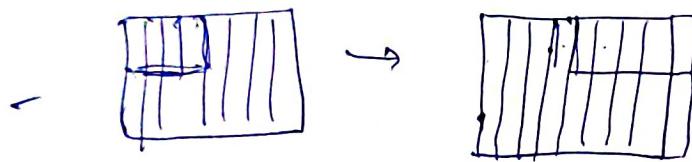
5x3

$$\boxed{\begin{aligned} \text{img} &:= (\text{ih}, \text{iw}) \\ \text{K} &:= (\text{kh}, \text{kw}) \\ \text{fm} &:= (\text{ih} - \text{kh} + 1, \text{iw} - \text{kw} + 1) \end{aligned}}$$

feature map



$\text{Stride} = 1$  (1 cell movement)



$\text{Stride} = 2$   
(2-cell movement)

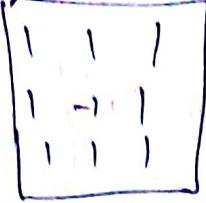
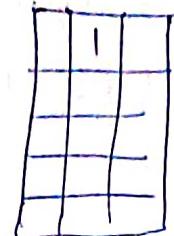
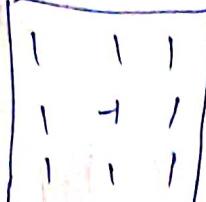
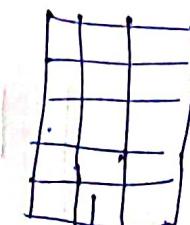
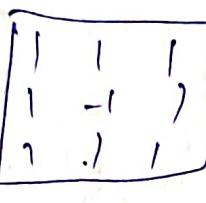
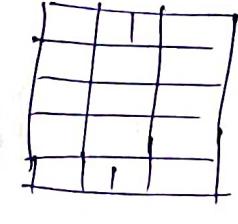
for  $\text{Stride} = 2$

$$\boxed{\begin{aligned} \text{img} &:= \text{ih}, \text{iw} \\ \text{K} &:= \text{kh}, \text{kw} \\ S &:= s \\ \text{fm} &:= [(\text{ih} - \text{kh})/2 + 1, (\text{iw} - \text{kw})/2 + 1] \end{aligned}}$$

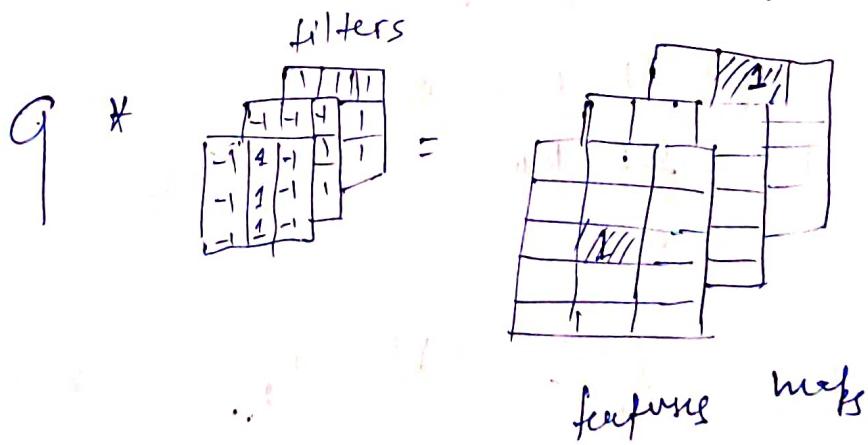
- This is how result of convolution process detects different patterns in the images using filters.

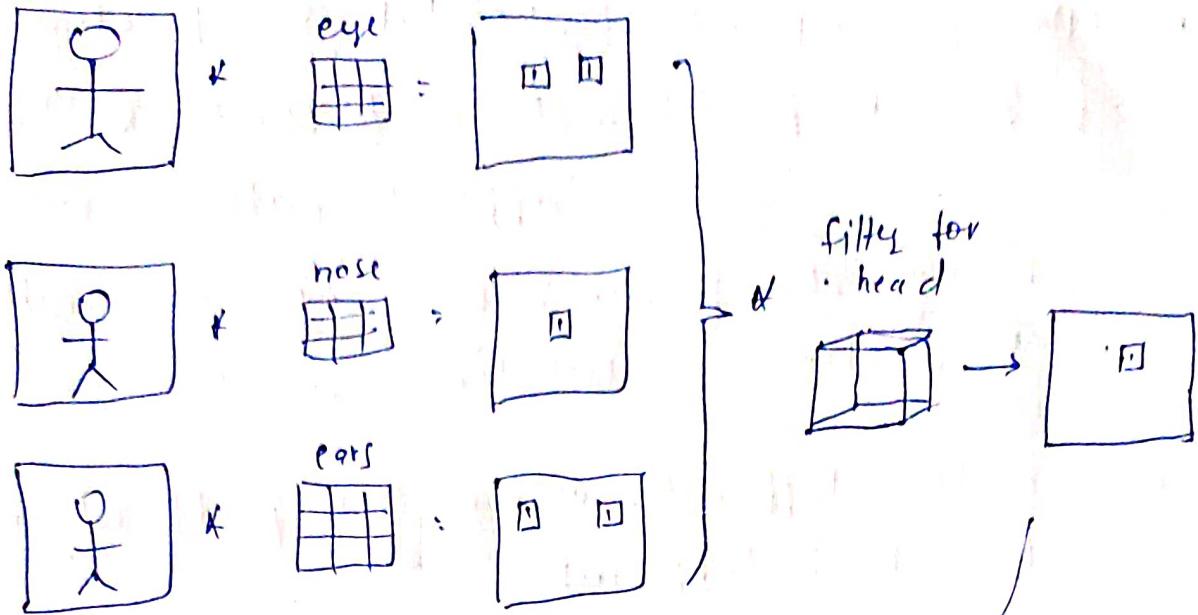


## Loopy pattern detector

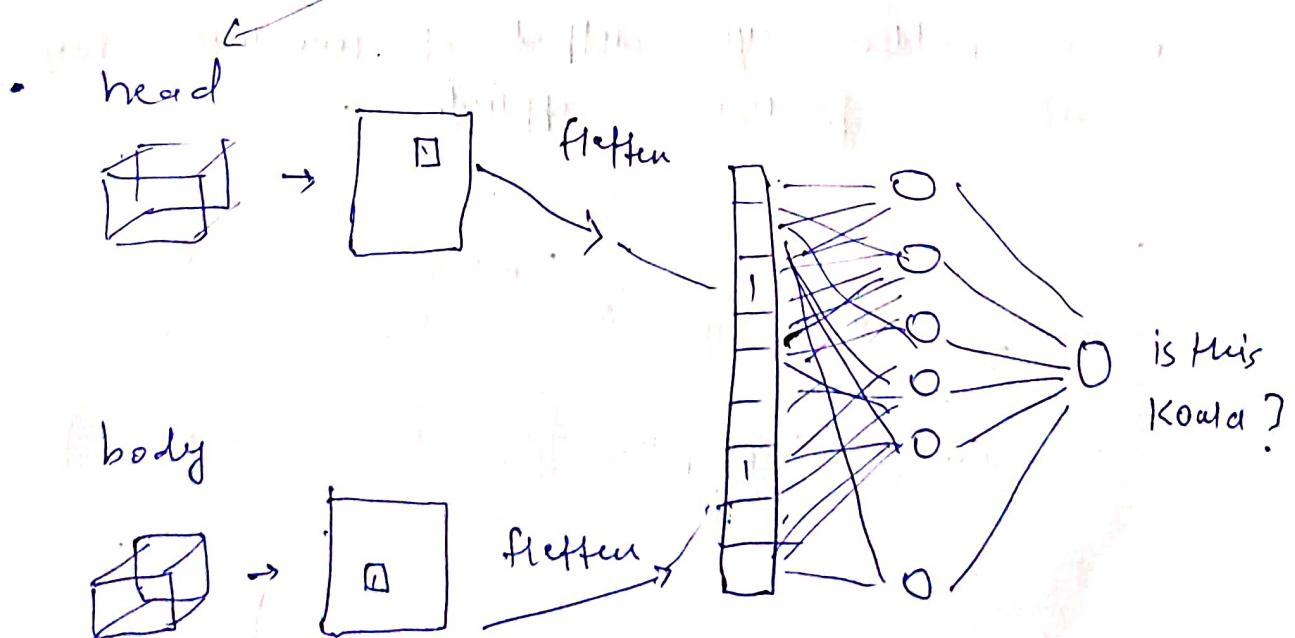
- $9 \times$   = 
- $6 \times$   = 
- $3 \times$   = 

- $f$  filters are used in each convolution layer to detect different features





- In case of above example, filters might end up detecting eyes, nose, ears, Then feature map matrix obtained is again subjected to another convolution possibly to detect head of animal.



- After few repeated convolutional operation, the feature map obtained is flattened & then a Dense layer is applied with sigmoid (binary classification) or softmax (multi classification) activation function.
- With different images, the flattened matrix will be different.
- Convolutional layer is responsible for feature extraction
- Dense output layer is responsible for classification

- \* Activation function is applied to the feature map.
- on hidden layers (output of convolution layer) activation function is applied.

-1	1	1	1	-1
-1	1	1	1	-1
-1	1	1	1	-1
-1	1	1	1	-1
-1	1	1	1	-1

\*

1	1	1
1	1	1
1	1	1

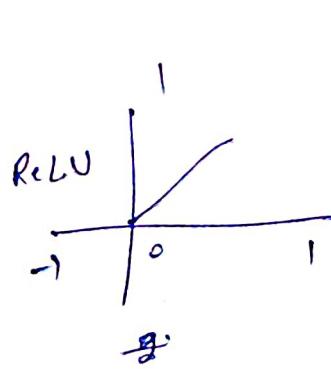
→

0.11	1	-0.11
-0.55	0.11	-0.33
-0.33	0.33	-0.33
-0.11	-0.11	-0.22
-0.33	-0.33	-0.22

Roopy pattern  
filter

0	1	0
0	0.11	0
0	0.33	0
0	0	0
0	0	0

←



## \* Peeling

- peeling layer is used to reduce the dimension on the feature map obtained after convolutional layers, MaxPooling layer is applied to select the maximum value from that window of given stride.

max pooling

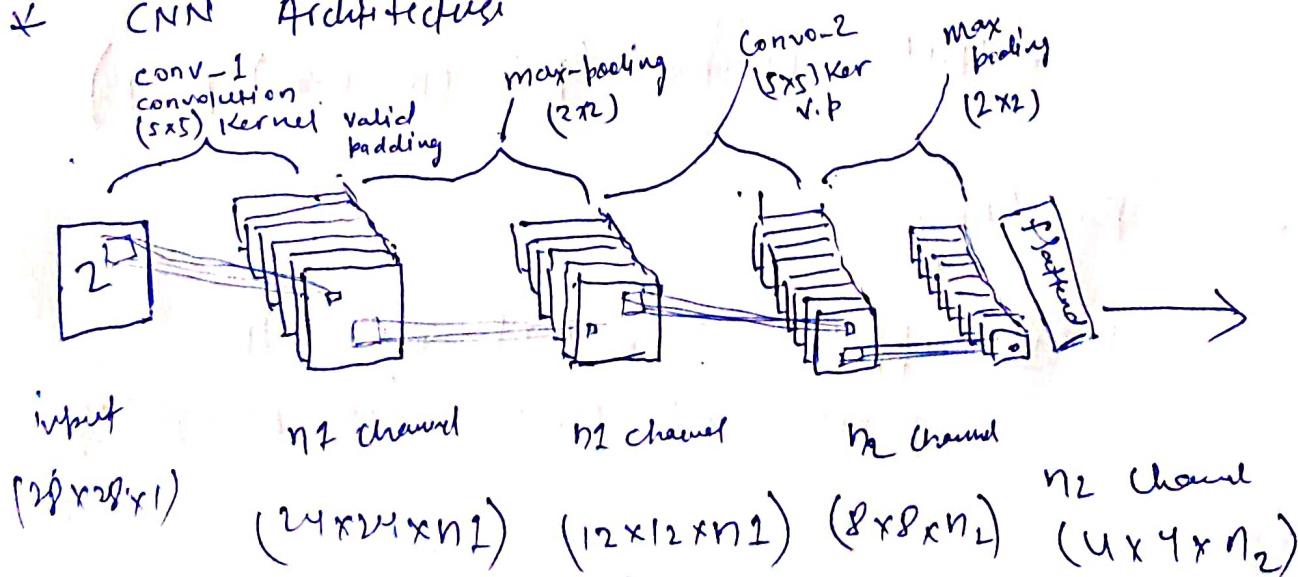
1	1
0.33	0.33
0.33	0.33
0	0

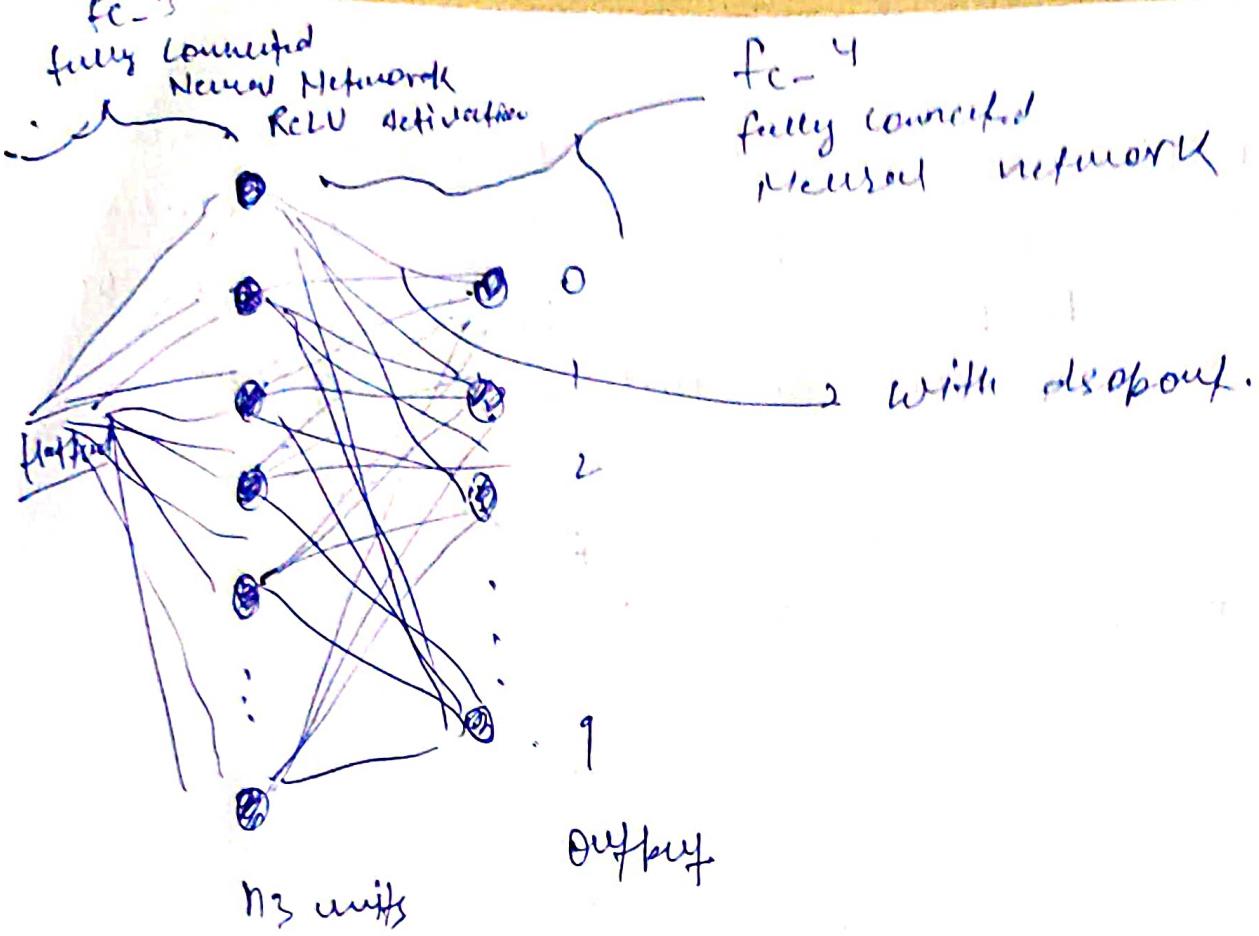
pooling size = (2,2)

stride = 1

- changing the position of 9 to ~~back~~ back to 1 cell, this type of positional movement is known as translational.
- Max pooling along with convolution helps detect position invariant feature detection.
- Pooling also prevents overfitting as there are less parameters.

## \* CNN Architecture





#### \* Stride

- Stride is number of pixels shifts over the input matrix. When the stride is 1 then we move filters to 1 pixel at a time. When the stride is 2 then we move the filters to 2 pixel at a time & so on.

Input  $(m, n)$ , Filter  $(p, q)$  Stride =  $p$

$$\text{Output} = \left( \frac{m-p+1}{p}, \frac{n-q+1}{q} \right)$$

#### \* Padding

It is applied so that even the corner pixels get their due weightage.