

Deep learning

UNIT - I

- 1) Biological neuron
- 2) Artificial neuron
- 3) Difference b/w them.
- 4) Different hyper Parameters
- 5) Limits of Traditional Computer Program.
- 6) feed-forward neural network
- 7) Different Activation functions : Sigmoid, Tanh, ReLU
- 8) Perception \rightarrow Working. \rightarrow Example / logic example.
- 9) Linear Perception as Neurons. & limitations.
- 10) Back Propagation algorithm, in Reducing Error.
- 11) Gradient descent with Sigmoid Neuron
 \hookrightarrow Derivation
- 12) delta rule for training linear neurons.
- 13) Preventing Overfitting, in DL methods & Demonstration.
- 14) Stochastic & minibatch Gradient.
 \downarrow

Companion
- 15) Mechanics of ML.

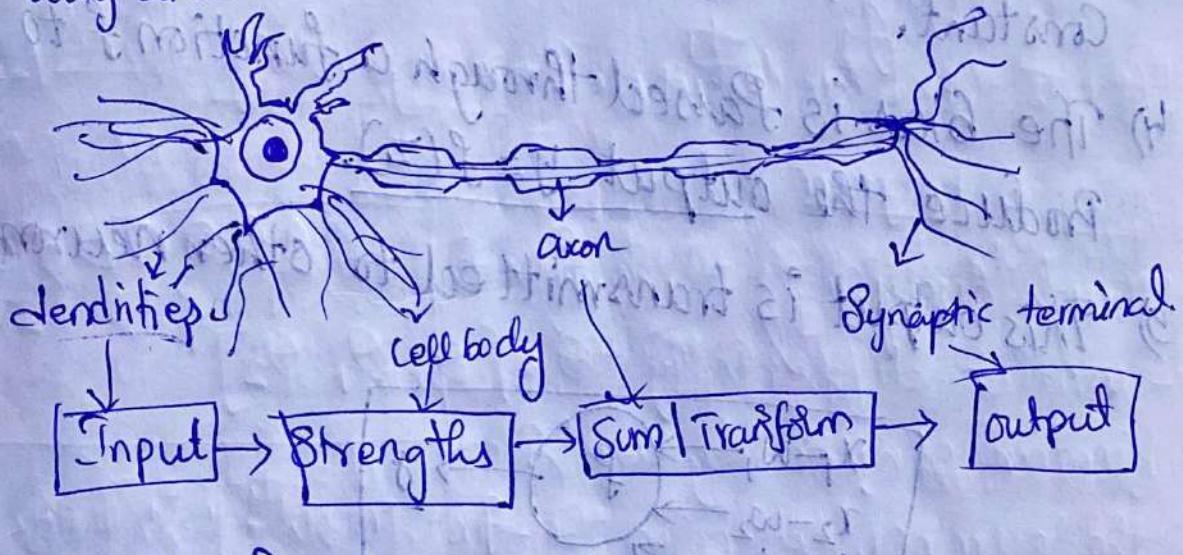
UNIT 1 Chapter 1

(Q1, 2, 4, 5, 6, 7, 8,

9, 15).

1) Biological Neurons

- 1) The foundational unit of human brain is the neuron.
- * It is the massive biological network that enables us to experience outside world.
- * As its core, the neuron is optimized to receive info from other neuron / process the info in a unique way and send its results to other cell.



→ Explain Process :-

1) The neuron receives input along dendrites.

2) Each of incoming connection is dynamically strengthened or weakened based on how it is used.

3) After being weighted by the strength of their respective connection, the input goes to cell body.

4) This sum is then transformed into a new signal that's propagated along the cell axon's & sent off to other neuron.

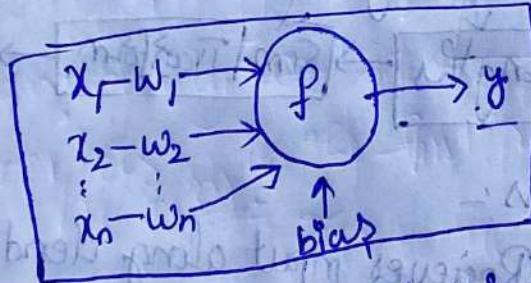
~~Artificial Neuron~~

~~Artificial Neural Network~~

② Artificial neuron:

4703A $\frac{2}{3}$

- 1) Such as Biological neuron the (a.N) takes in some no of inputs x_1, x_2, \dots, x_n each of which is multiplied by weight, $w_1, w_2, w_3, \dots, w_n$.
- 2) These weighted inputs are, as before, summed together to produce the logit of the neuron, $z = \sum_{i=0}^n w_i x_i$.
- 3) The logit is also called as bias, which is constant.
- 4) The bias is passed through a function f to produce the output $y = f(z)$.
- 5) This output is transmitted to other neurons.



We can express this in Vector form

input vector $\mathbf{x} = [x_1, x_2, \dots, x_n]$

weight vector $\mathbf{w} = [w_1, w_2, \dots, w_n]$

Expression: $y = f(\mathbf{x} \cdot \mathbf{w} + b)$ where b is bias term

3)

Artificial Neuron

- 1) It is a mathematical model which is mainly inspired by the biological neuron system in the human brain.
- 2) Its Processing way Sequential and centralized.
- 3) Small in size
- 4) It process the info at a faster speed
- 5) Cannot perform Complex Pattern Recognition
- 6) It doesn't provide any feedback
- 7) No fault tolerance
- 8) Operating Environment is well defined and well-constrained
- 9) Reliability : It is very Vulnerable
- 10) Response time is measured in milliseconds.
- 11) It has very accurate structure

Biological Neuron

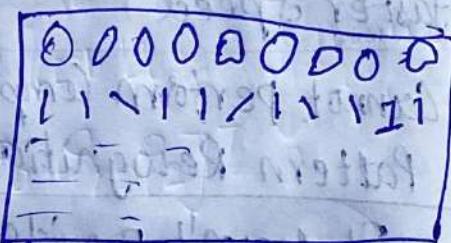
- It is also composed of several Processing pieces known as Neurons.
- It Process the info in a Parallel and distributive manner.
- large in size.
- It process the Info at a slow speed.
- Perform Complicated tasks.
- It Provides feedback.
- fault tolerance is their
- Operating Environment is poorly defined and unconstrained.
- It is robust.
- It's Response time is measured in nanoseconds
- They are tolerant & ambiguity.

③ Limits of Traditional Computer Programs.

- * Traditional computer programs are designed to be very good at two things
 - 1) Performing arithmetic really fast.
 - 2) Explicitly following a list of instructions;
- * The T.C.P are good at heavy financial number crunching.

Let us take an example: Recognizing hand written digits Program

although every digit in the figure is written differently the human can easily identify the number as 0 in the first row, 1 in second row etc.,



But if we are said to write program to automatically read someone's handwriting, it will be very difficult.

If we might state that we have a "zero," if our image only has a single, closed loop.

but this isn't really a sufficient condition what if someone write a messy zero, it will be a very difficult task to identify it.

like the fig 2 describes



It can't always identify it as 0 or 6 in this case if we even keep the distance between them in the loop.

④ 5) Mechanics of Machine Learning

- To tackle the above Handwritten problem we have to use different kind of approach.
- from childhood to adulthood we often learn things and we learn what is wrong and we will learn from it and train from it.
- we learn how to multiply numbers, solve Equations, take derivates by set of instructions.
- we learned to recognize a dog by being shown multiple examples and being corrected when we made the wrong guess. if our parents say we are wrong we will modify our model to incorporate this new information.
- Deep Learning is a subset of a more general field of AI called ML. which is predicted on this idea of learning.
- we give model to train and learn instead of set of large instructions.

Let's define our model to be a function $h(x, \theta)$.

The Input (x) is an example Expressed in Vector form.

$\theta \rightarrow$ vector of the parameter that model use.

Let's take an example: Sleep Predict Exam Performance on ~~Sleep~~ no of hours of Sleep.

We collected a lot of data for each dataset.

Point $x = [x_1 \ x_2]^T$, $x_1 \rightarrow$ no of hours of Sleep.

$x_2 \rightarrow$ no of hours we spent studying.

Our goal, the model $h(x, \theta)$ with

Parameters Vector $\theta = [\theta_0 \ \theta_1 \ \theta_2]^T$

Such that:

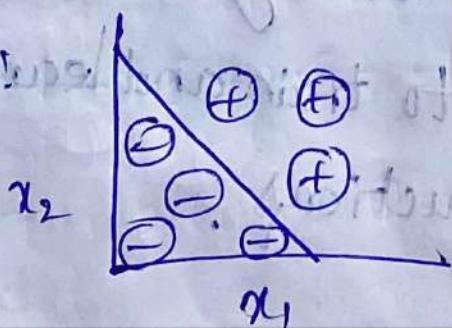
$$h(x, \theta) = \begin{cases} -1 & \text{if } x^T \cdot \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} + \theta_0 < 0 \\ 1 & \text{if } x^T \cdot \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} + \theta_0 \geq 0. \end{cases}$$

-1: below average

1: Otherwise.

This model is called linear Perceptron?

Sample Data:

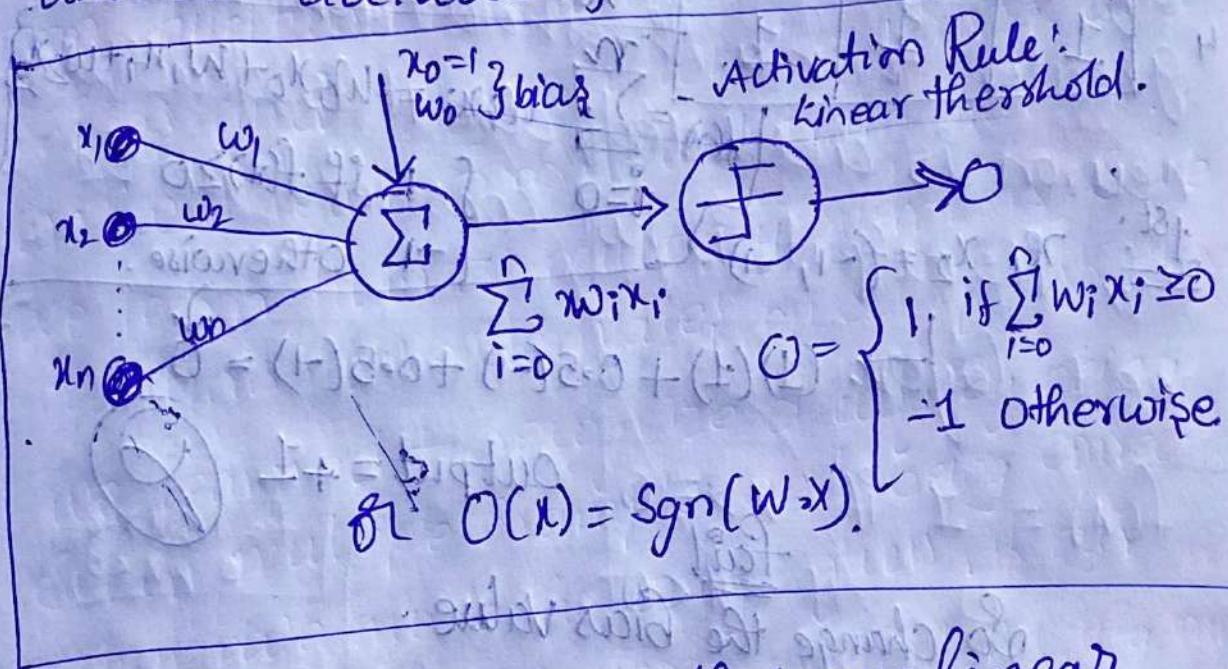


② Linear Perception as Neurons &

limitations:

A linear perception or perceptron is a ML algorithm for supervised learning of various binary classification tasks.

The single layer neural network has four parameters: input values, weights, bias, and an activation function.



It is very easy to show that our linear perception & the neuronal model are perfectly equivalent. And it is quite simple to show that singular neurons are strictly more expressive as linear perception.

$$O = \text{tanh}(w \cdot x) \quad O = 0 = 2 \cdot 0 + 2 \cdot 0 + 1 \leftarrow (1, 1)$$

Trying the Linear Perceptron by using
and gate:

	x_1	x_2	y
1	-1	-1	-1
2	1	-1	-1
3	-1	1	-1
4	1	1	1

Let us take

$$w_0 = 1 \text{ (bias)} \quad x_0 = 1$$

$$w_1, w_2 = 0.5 \text{ (assum)}$$

We know formula

$$f(x) = \sum_{i=0}^n w_i x_i = w_0 x_0 + w_1 x_1 + w_2 x_2$$

1st: $x_1, x_2 = (-1, -1)$

$$\hookrightarrow (1)(1) + 0.5(-1) + 0.5(-1) = 0$$

Output = +1



So change the bias value.

$$w_0 = -1$$

$$w_1, w_2 = 0.5$$

$$x_0 = 1$$

$$P(-1, -1) \rightarrow -1 + 0.5 - 0.5 = -1 < 0 \quad \boxed{\text{Output} = -1}$$

$$P(1, -1) \rightarrow -1 + 0.5 - 0.5 = -1 < 0 \quad \boxed{\text{Output} = -1}$$

$$P(-1, 1) \rightarrow -1 + 0.5 + 0.5 = -1 < 0 \quad \boxed{\text{Output} = -1}$$

$$P(1, 1) \rightarrow -1 + 0.5 + 0.5 = 0 = 0 \quad \boxed{\text{Output} = 1}$$

Limitations :-

- 1) Lack of complexity: They are incapable of handling complex patterns.
- 2) Limited Expressiveness:
- 3) No feature hierarchy: due to linear nature they cannot form such hierarchies.
(Simple to Complex)
- 4) Inadequate for Complex tasks:
- 5) Overfitting to linear Assumptions:

Activation functions:

The Activation function is one that outputs a smaller value for tiny inputs and a higher value if its inputs are greater than their threshold.

There are two types of activation functions:

- 1) linear Activation function.
- 2) Non-linear Activation function.

Linear Activation Function:

The function is linear.

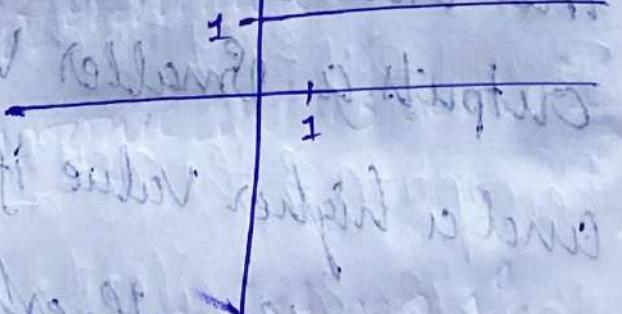
* Binary Step Function:

If the input to the activation function is greater than a threshold then the neuron is activated, else it is deactivated. i.e., its output is not considered for the next layer.

$$f(x) = 1 \text{ if } x >= 0$$

$$0 \text{ if } x < 0.$$

$$\begin{cases} 1 & , x \geq 0 \\ 0 & , x < 0. \end{cases}$$



drawbacks:

- 1) It can be only used while creating a binary classifier.

* This function is not useful when there are multiple classes in the target variable.

* Non Linear Activation functions:

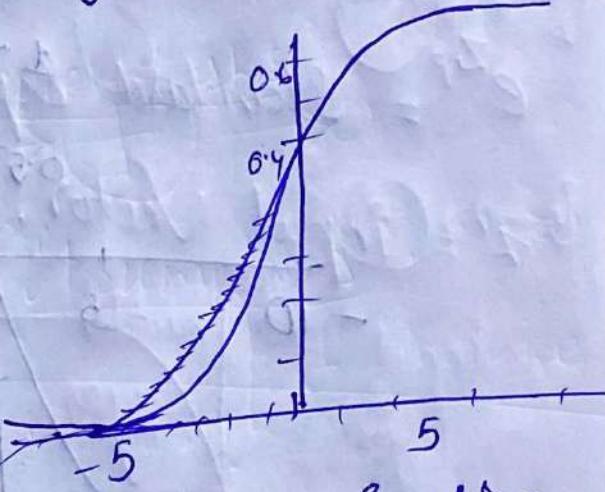
~~The input of~~ The function will be non linear

1) Sigmoid :-

A Sigmoid function is graphed in "S" shape.

It uses the function:

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



* When the logit is very small, the output of a logistic function is very close to "0".

* When the logit is very large, the output of the logistic function is close to "1".

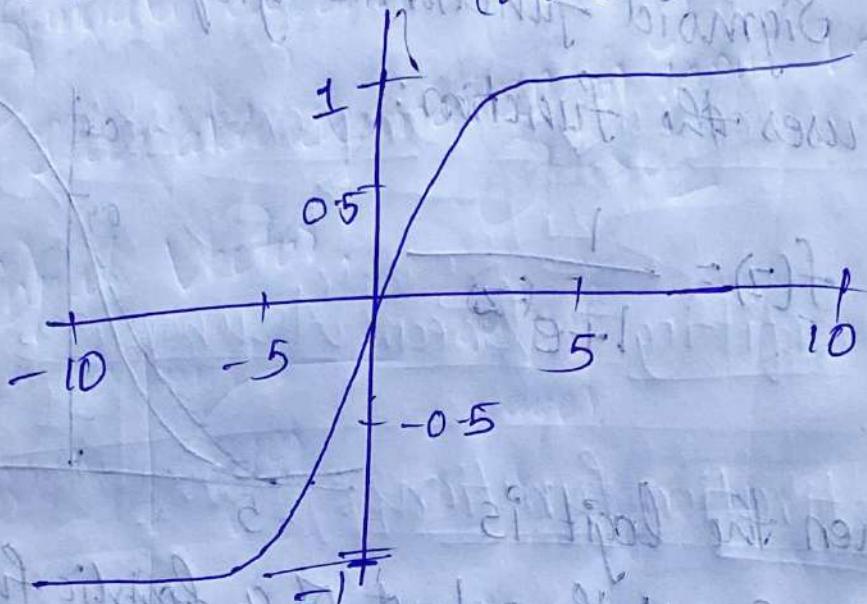
* In b/w them, the function assumes an S-shape.

② Tanh function:

The Tanh function has a similar kind of S-shaped non-linearity, but instead ranging from 0 to 1, the output of tanh function ranges from -1 to 1.

If uses:

$$f(z) = \tanh(z) = \frac{e^z}{(1 + e^{2z})}$$

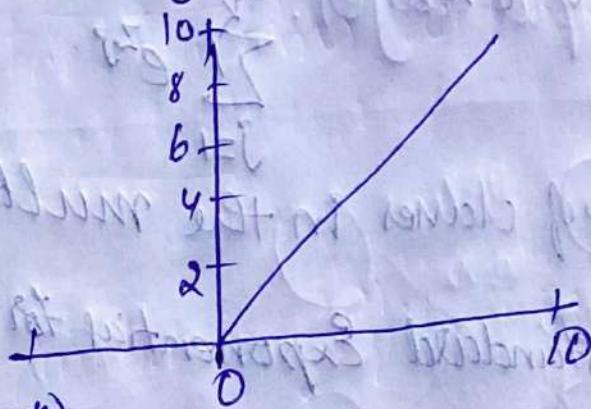


- The function is differentiable.
- The function is monotonic and its derivative is not.
- It is mainly used in between layer of.

3) RELU activation function

It is Rectified Linear Unit

It uses the function $f(z) = \max(0, z)$, resulting in a characteristic hockey-stick-shaped response.



It ranges from 0 to ∞ .

The function is monotonic.

4) Softmax Function:

- The Softmax Function can be used for multiclass classification problems.
- It is described as combination of multiple Sigmoids.
- Sigmoid returns values b/w 0 and 1, which can be treated as probabilities of a data point belonging to a particular class.

This function returns the probability for a data point belonging to each individual class.

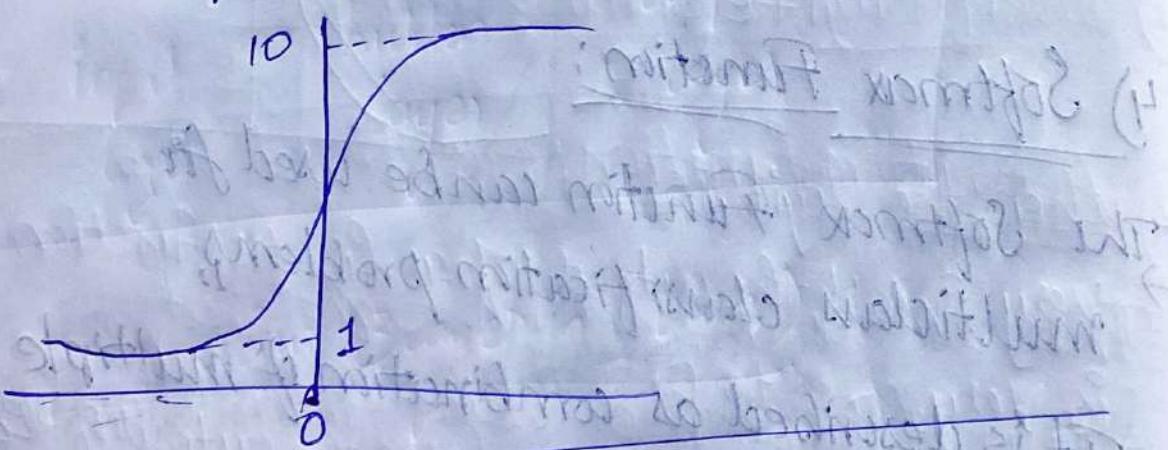
$$\sigma(\vec{z}) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

k = no of classes in the multiclass classifier

e^{z_j} = standard Exponential for output vector

e^{z_p} = standard Exponential for input vector.

\vec{z} = input Vector.



Hyperparameters: (not important)

- 1) The k in KNN
- 2) Learning rate for training a neural Network
- 3) Train-test split ratio
- 4) Batch size
- 5) Number of Epochs
- 6) Branches in Decision Tree
- 7) Number of clusters in clustering algorithms

Feed-forward Neural Network:

The process of receiving an input to produce some kind of output to make some kind of prediction is known as feed forward. The feed-forward neural network is the core of many other types of neural network such as convolution neural network.

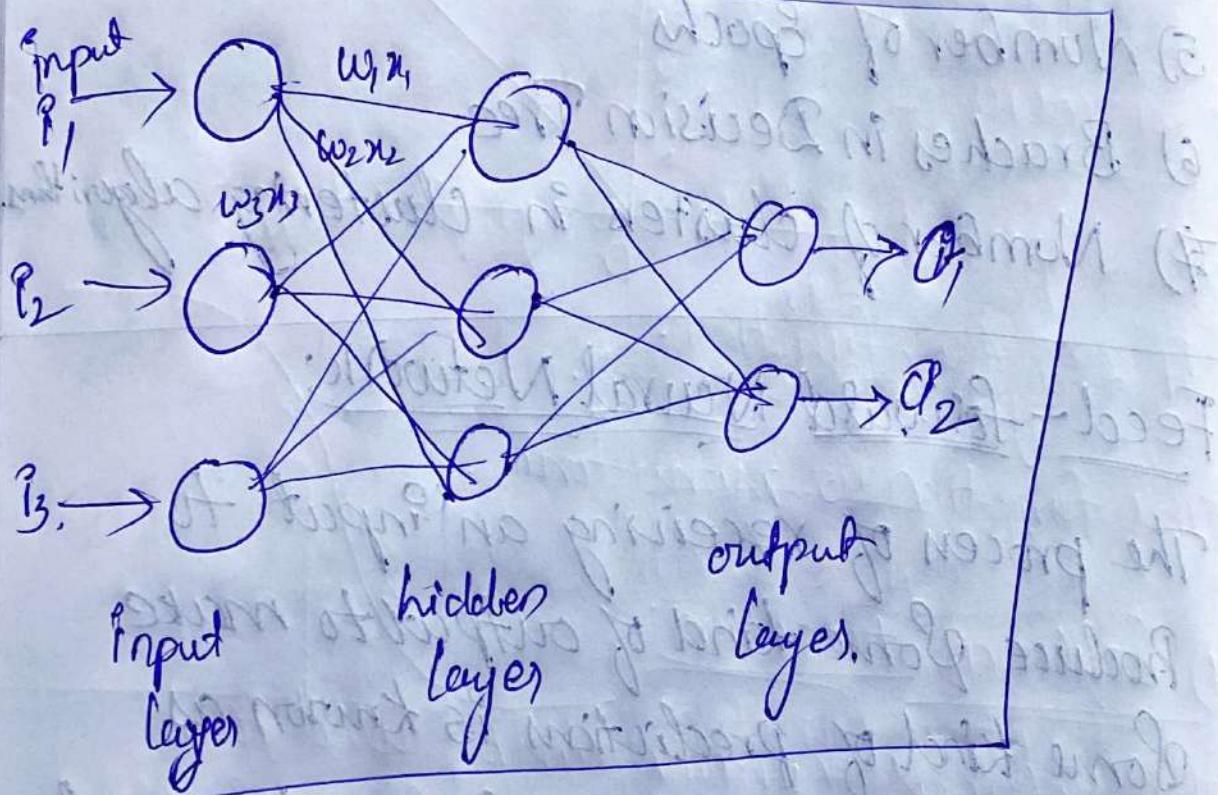
In the feed-forward neural network, there are not only any feedback loops or connections in the network. Here is an example.

there are 3 layers:

1) Input layer

2) Hidden layer

3) Output layer.

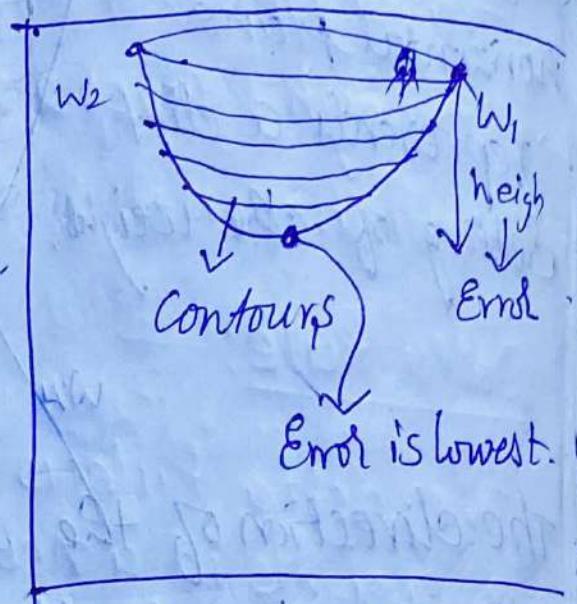


(Explain about 3 layer in own.)

UNIT 1 Chapter 2

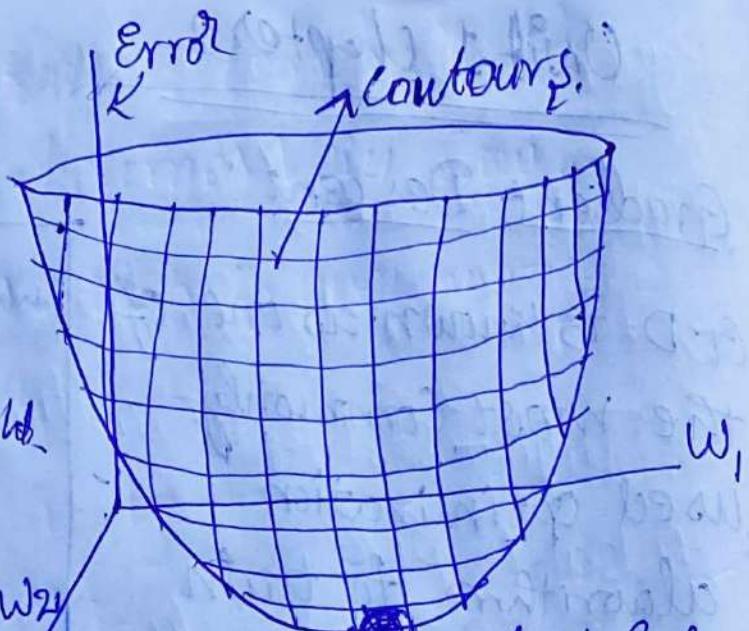
Gradient Descent:

G.D. is known as one of the most commonly used optimization algorithms to train ML models by means of minimizing errors between actual & expected results.



- It helps to find the local minimum of a function.
- The main objective of using a gradient descent algorithm is to minimize the cost function using iteration.
- The 3D space where the horizontal dimensions corresponds to the weights w_1 and w_2 ; and the vertical dimension corresponds to the value of the error function E .
- As shown in figure the circles are elliptical "Contours". Where the minimum error is at the center of the ellipse.

- * Each point on the horizontal plane represents a cliff.
- Setting of the weight.



The direction of the steepest descent is always \perp to these contours.

- We start by initializing the weights (w_1 & w_2) randomly.
- Evaluate the gradient at this starting position to find the direction of steepest descent.
- Take a step in the direction of this gradient.
- After the step, you will be at the new position with a potentially lower error.
- Reevaluate the gradient at this new position and again move in the direction of steepest descent.
- This iterative evaluating the gradient and taking steps in its direction, we

- Progressively move closer to the point of minimum Errors.
 - This Iteration process is the main in gradient descent.
-
- The Descent Rule & Learning Rule :-
- In neural networks, apart from weights, there are other parameters known as hyperparameters that help guide the training process. one such hyperparameter is "Learning Rate".
 - The Learning Rate controls how much we adjust the weights of our Networks wrt to the gradient.
 - When using the gradient descent the learning rate determines how big a step we take at each iteration.
 - The steepness of the error surface can indicate how close we are to minimum error.
 - when surface is flat → we are closer to minimum. Should take smaller steps.
 - Surface is steep → take larger steps.

- If the learning Rate is too small, the training process becomes very slow.
- If too high, there's a risk of overshooting the minimum.

Delta Rule:

The delta Rule is the method of updating the weights in a neural network.

It involves calculating the gradient of the error function wrt. each weight.

By applying the gradient and the learning rate, the delta rule adjusts each weight to reduce the error.

We know that

$$\Delta w_k = -\eta \frac{\partial E}{\partial w_k} \quad \text{--- (1)}$$

We also know that

$$E = \frac{1}{2} \sum (t_k - o_d)^2 \quad \begin{matrix} \text{act. predict.} \\ \uparrow \end{matrix} \rightarrow \text{--- (2)}$$

(2) in (1)

$$\begin{aligned}\Delta w_k &= -\epsilon \frac{\partial}{\partial w_k} \left[\frac{1}{2} \sum (t_i - o_d)^2 \right] \\ &= -\epsilon \frac{1}{2} \cancel{\frac{\partial}{\partial}} (t_i - o_d) \frac{\partial}{\partial w} (t_i - o_d) \\ &= -\epsilon (t_i - o_d) \frac{\partial}{\partial w} \left(\frac{t_i - (w_i \cdot x)}{1} \right) \quad \boxed{\because o_d = w_i \cdot x} \\ &= -\epsilon (t_i - o_d) (-x) \\ \boxed{\Delta w_k = \epsilon x (t_i - o_d)}\end{aligned}$$

Applying this method of changing the weights at every iteration, we are finally able to utilize gradient descent.

Gradient Descent with Sigmoidal Neurons

We

Let's recall the mechanism by which logistic neurons compute their output value from their inputs:

$$z = \sum_k w_k x_k$$

$$y = \frac{1}{1 + e^{-z}}$$

The neuron computes the weighted sum of its inputs, the bias "z", It then feeds it bias into the input function to compute y , its final output.

To do so we start by taking the derivatives of the bias w.r.t input & the weight:

$$\nabla_{w_k} E = \frac{\partial E}{\partial w_k}$$

~~$\nabla_b E = \frac{\partial E}{\partial b}$~~

$$\therefore E = \frac{1}{2} (f(x) - y)^2$$

$$\therefore f(x) = \frac{1}{1 + e^{-(wx+b)}}$$

$$\frac{\partial E}{\partial w_k} = \frac{1}{2} \left[\frac{1}{1 + e^{-(wx+b)}} - y \right]^2$$

$$= \frac{1}{2} \left[2 \times (f(x) - y) * \frac{\partial}{\partial w} (f(x) - y) \right]$$

$$\frac{\partial E}{\partial w_k} = (f(x) - y) * \frac{\partial}{\partial w} (f(x) - y) \rightarrow \text{Q}$$

$$\therefore \frac{\partial}{\partial w} (f(x) - ey)$$

$$= \frac{\partial}{\partial w} \left[\frac{1}{1 + e^{-(wx+b)}} - ey \right]$$

$$= \frac{-1}{(1 + e^{-(wx+b)})^2} \frac{\partial}{\partial w} (e^{-(wx+b)})$$

$$= \frac{-1}{(1 + e^{-(wx+b)})^2} * e^{-(wx+b)} \frac{\partial}{\partial w} (-\underbrace{(wx+b)}_{1 \downarrow 0})$$

$$= \frac{-1}{(1 + e^{-(wx+b)})^2} * e^{-(wx+b)} (-x)$$

$$= (+x) \frac{-1}{(1 + e^{-(wx+b)})^2} * e^{-(wx+b)} \frac{z}{z}$$

$$= \frac{-1}{1 + e^{-z}} * \frac{e^{-z}}{1 + e^{-z}} (-x) \quad (z = wx + b)$$

$$\frac{\partial f}{\partial w_k} = f(x) * (1 - f(x)) (-x) \rightarrow ②$$

keep ② in ① $f(x) = y$

$$\therefore [f(x) - y] [f(x) * (1 - f(x)) (-x)]$$

$$= -x \frac{f(x)(f(x)-y)}{1-f(x)} (1-f(x)) \frac{\partial y}{\partial w_k}, \Delta w_k = -e \frac{\partial f}{\partial w_k}$$

$$\Delta w_k = -x f(x) (f(x)-y) (1-f(x))$$

$$w_{t+1} = w_t - \eta \nabla_{w_t}$$

$$= w_t - \eta \sum \text{error} (1-y)(t-y)$$

$$w_{t+1} = w_t - \eta (f(w-y) f'(y) (1-f(y)) x)$$

As you notice, the new modification rule is just like the delta rule, except with extra multiplicative terms included to account for the logistic component of the Sigmoidal neuron.

Backpropagation algorithm: for reducing Err.

algorithm:

Step 1: calculate the Error at the Output
* Start with the output layer of the neural network

* for each neuron i in the output layer, calculate the Error (E) as the sum of the squares of the difference between the Predicted output (y_i) and the true target (t_i)

$$\boxed{E = \frac{1}{2} \sum_i (t_j - y_j)^2}$$

* compute the derivative of the Error wrt each output neuron's activity:

$$\boxed{\frac{\partial E}{\partial y_j} = \frac{1}{2} \times 2(t_j - y_j) \frac{\partial}{\partial y_j} (t_j - y_j)}$$

$$\frac{\partial E}{\partial y_j} = \frac{1}{2} \times 2(t_j - y_j) \frac{\partial}{\partial y_j} (t_j - y_j)$$

$$\frac{\partial}{\partial y_j} (t_j - y_j) = -1$$

Step 2: Back propagate the Error.

- * Presume that the Error derivatives for layer j are known.
- * To calculate the Error derivatives for the preceding layer i ; accumulate info on how each neuron in layer i influences the neuron in layer j .
- * This can be done as follows using the Partial derivative of the bias wrt incoming output data from the layer beneath is merely the weight of the connection w_{ij} .

$$\frac{\partial E}{\partial y_i} = \frac{\partial E}{\partial z_j} \frac{\partial z_j}{\partial y_i}$$

$$= w_{ji} \frac{\partial E}{\partial z_j} \rightarrow ① \left(\frac{\partial z_j}{\partial y_i} = w_{ji} \right)$$

$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial z_j}$$

$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_i} y_i (1-y_i) \rightarrow ②$$

③ ④ ⑤ ⑥

$$\frac{\partial E}{\partial y_i} = \sum_j w_{ji} y_i (1-y_i) \frac{\partial E}{\partial z_j}$$

Step 3: update the weights

* After calculating the partial derivatives
for all units, update the weights to
minimize the error, using the learning

rate "η"

$$\Delta w_{ji} = -\eta \frac{\partial E}{\partial w_{ji}} \rightarrow ③$$

We know

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial z_j}{\partial w_{ji}} \frac{\partial E}{\partial z_j}$$

From ②

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial y_j} y_j(1-y_j) \times \frac{\partial z_j}{\partial w_{ji}}$$

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial y_j} y_j(1-y_j) \times y_i \quad \left[\because \frac{\partial z_j}{\partial w_{jj}} = y_j \right]$$

~~Now keep the~~

Now for all the training examples, sum up the partial derivatives to adjust the weights:

$$\Delta w_{ji} = -\eta \frac{\partial E}{\partial y_j} y_j(1-y_j) \times y_i.$$

$$\boxed{\Delta w_{ji} = -\eta \sum_{\text{dataset}} y_i y_j (1-y_j) \frac{\partial E}{\partial y_j}}$$

Step 4: Iterate the process:

Repeat steps 1 to 3 for each layer in the network, moving from the output layer towards the input layer.

Iterate the process across all layers
for multiple epochs or until the network
converges to a satisfactory level of
Error.

Stochastic and Minibatch Gradient

after Evaluating each individual training
Example.

→ It is a compromise b/w batch gradient
descent & stochastic gradient descent.

It computes the gradient against a subset
of the data known as mini batch.

$$\frac{1}{B} \sum_{i=1}^B \nabla J(\theta) = \nabla J(\theta)$$

out of individual set of a epoch based
about input digits at most common showed
fewer digits with

Different techniques to prevent overfitting.

1) Regularization

Overfitting :-

the model becomes very good at predicting or fitting the data it has been but fails to generalize to data that hasn't been trained on.

1) Regularization :-

* Regularization is a technique to prevent overfitting by adding a penalty for larger weights in the objective function.

* The modified objective function becomes:

Error + $\lambda f(\theta)$, where $f(\theta)$ ↑ with larger weights
if λ is the regularization strength, a hyperparameter.

* choosing λ is crucial : $\lambda=0$ means no Regularization while a large λ might prioritize small weights over fitting the training data.

L2 Regularization (Weight Decay) :

The most common regularization in ML, L2 Reg. Adds $\frac{1}{2} \lambda W^2$ to the Error function for every weight "W".

It penalizes large weights, Preferring more evenly distributed weight values and Prevents overfitting by encouraging the whole network to use all inputs rather than a few inputs a lot.

During gradient descent, L2 causes weights to decay linearly towards zero.

3) L1 Regularization :

- * Adds the term $|W|$ for each weight in the network.
- * Leads to Sparsity in the weight vector, making the model less sensitive to noise in the inputs.
- * useful for feature Selection as it indicates which features ~~descent~~ are most important for predictions.

4) Max Norm Constraints :-

- * Enforce an absolute upper bound on the magnitude of the weight vectors.
- * If the weight exceeds the bound during the gradient descent update, it is projected back onto a predefined bound.
- * Prevents weights from growing too large, ensuring stability in the learning process.

5) Dropout :-

- * A technique where each neuron is kept active with a certain probability P during training, & set to zero otherwise.
- * makes the network less reliant on any single neuron, effectively combining many different neural network architectures.
- * Dropout can prevent overfitting by ensuring the network remains accurate. Even when certain neurons are dropped.

6) Inverted Dropout

- * A variation of dropout where the neuron's output is scaled at training time, rather than test time.
- * Ensure that at test time, the neuron's output do not need to be scaled, maintaining the expected output level.

→ The Autoencoder can effectively summarize data points when all relevant info is present in the individual data points themselves.

Chapter 1

1) Vanilla Deep Neural networks

Deep learning aims to automate feature selection, crucial for Computer Vision tasks.

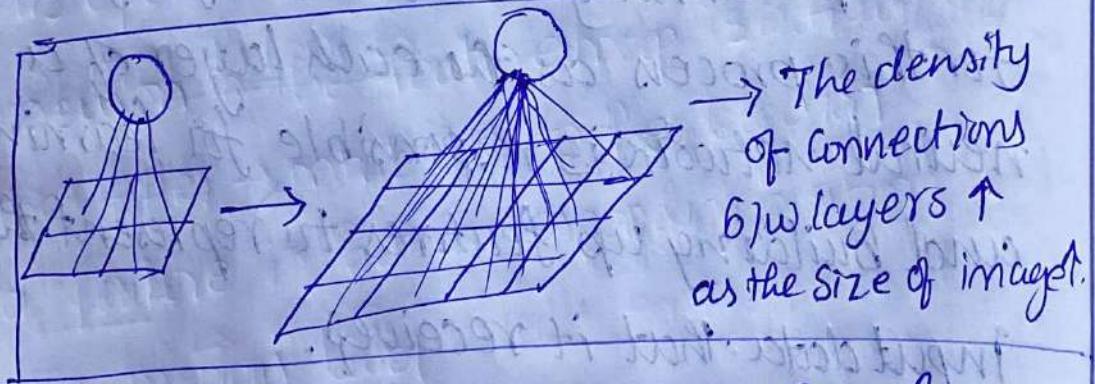
The Deep learning network are prefects for this process as the each layer of a neural network is responsible for learning and building up features to represent the input data that it received.

The Vanilla Deep neural network are not scalable for image classification as image size increases due to an exponential increase in number of weights.

In the MNIST dataset, the problem appears in manageable with 28×28 Pixel images, but Scaling up to larger images rapidly increases the complexity.

for large image, such as 200×200 full-color image, The no of weights required for just the input layer becomes high, leading to high risk of overfitting and inefficiency.

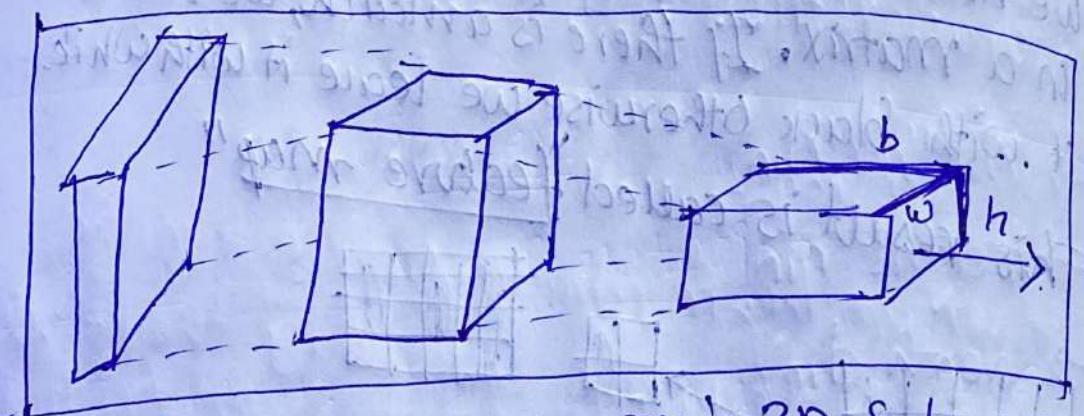
The CNN address this Scalability issue by arranging neurons in 3D and connecting each neuron to only a small, localized region of the previous layer.



This architecture mimics human visual processing and reduce the no of parameters significantly.

CNN manage the complexity by processing info in a 3D; So layers have width, length & depth, so new volumes of processed info layer by layer.

This design of CNN layer avoids the problem of full connectivity and reduce the potential for overfitting while handling larger images more efficiently.



CNN layer arrange neurons in 3D, so layers have width, height; depth.

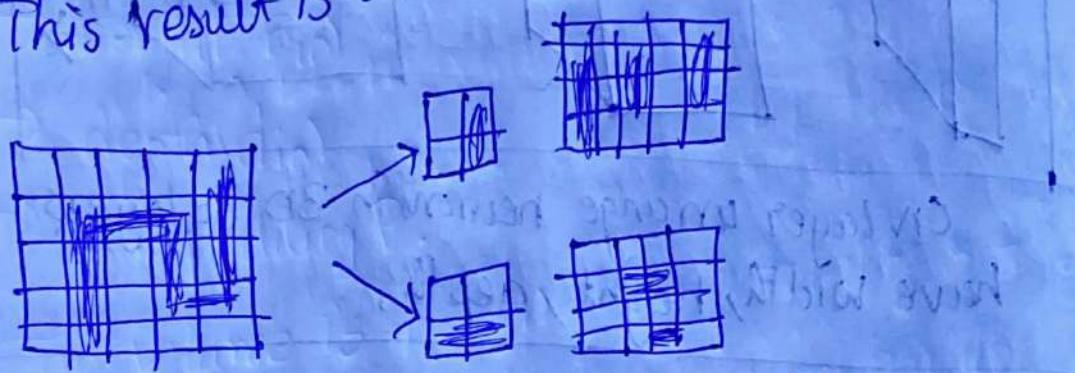
Filters and Feature Map:

The visual cortex layers structures, where each layer builds upon the features detected by the previous one, inspired the design of CNN layers in neural networks.

The filters is essentially a feature detector that slide across the image to identify features like edges, lines & angles.

~~See~~

To detect Vertical lines, we would use the feature detectors on the top, slide it across the image and at each and every step check if we have a match. We keep track of answer in a matrix. If there is a match, we shade it with white otherwise we leave it with black. This result is called "feature map".

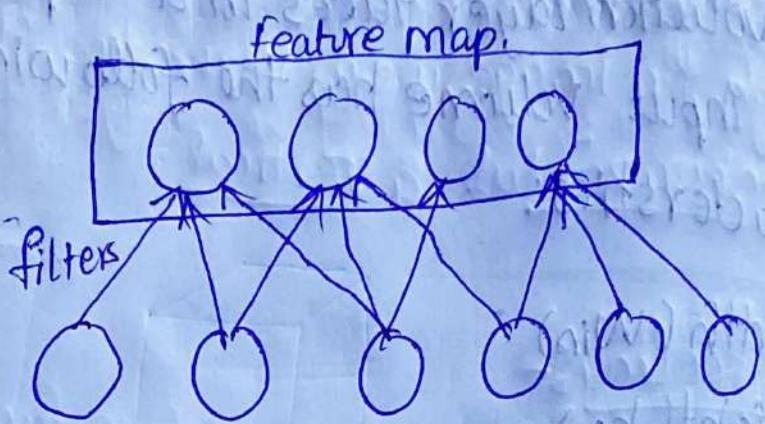


This operation is called Convolution. A neuron in the feature map is activated if the filter contributes to the detection of a feature in the corresponding position in the previous layer.

Expressing the Convolution Operation as neurons in a network

Layers of neurons in a feed forward neural net represents either the original

image or a feature map that get replicated across the entirety of the input.



Let's denote the k^{th} feature map in layer m as m^k . $W \rightarrow$ corresponding filters weight. $b^k \rightarrow$ bias. we mathematically Express feature map as:

$$m_{ij}^k = f((W * x)_{ij} + b^k)$$

The fundamental operation of Slide filters across the input data to produce feature maps, capturing more complex features.. They capture patterns like edges or textures.

Filters in CNN have depth, allowing them to process multi-channel input, such as colour images.

The CNN converts one volume of values into another volume of values.

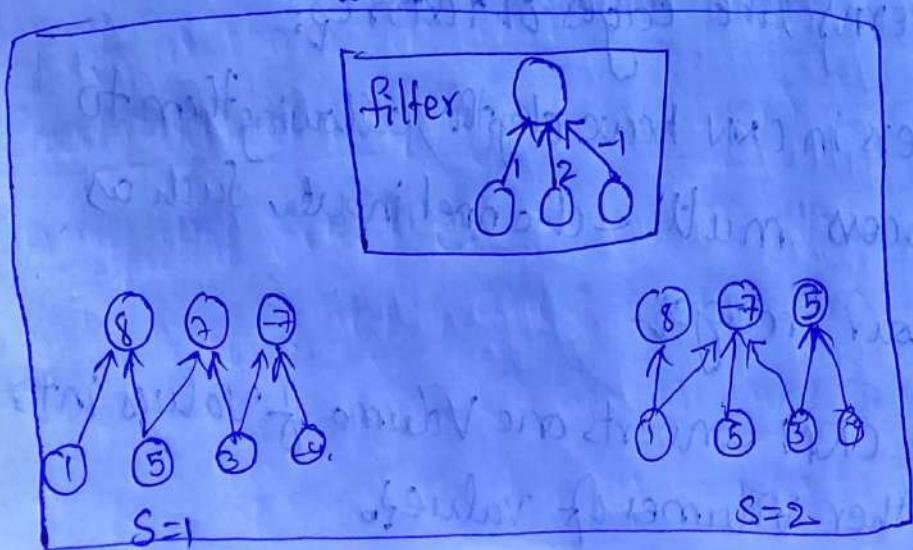
Full Description of Convolution layer

A Convolution layer takes in a Input Volume.
This input volume has the following
Characteristics:

- Width (Win)
- height (hin)
- depth (din)
- Zero padding P.

This volume is processed by a total of K filters,
This filters has many no. of hyperparameters:

- Spatial Extent e. → which is equal to the filters height & width.
- Stride S → distance b/w consecutive applications of the filters on the input volume.
- bias b, → which is added to each component of the convolution.



The Result in the Output Volume has :

→ function f ,

$$\rightarrow \text{Width } W_{\text{out}} = \left\lceil \frac{W_{\text{in}} - e + 2p}{s} \right\rceil + 1$$

$$\rightarrow \text{height } h_{\text{out}} = \left\lceil \frac{h_{\text{in}} - e + 2p}{s} \right\rceil + 1$$

$$\rightarrow \text{depth } d_{\text{out}} = k$$

$$\begin{aligned} w &= 5 \\ h &= 5 \\ d &= 3 \\ \text{ZeroPad} &= 1 \\ w &= 3 \\ h &= 3 \\ d &= 2 \end{aligned}$$

Padding:

If the input is $n \times n$ and filter size is $f \times f$
the output size will be $(n-f+1) \times (n-f+1)$

There are two disadvantages:

- Every time we apply a convolution Operation, the Image Size Shrinks
- Pixels present in the corner of the image are used only a few number of times during convolutional as compared to center pixels.

To overcome the issue we can Pad the image with additional borders. We add one pixel around the edges. This means input size will be 9×9 instead of 7×7 .

input : $n \times n$

Padding : p

filter size : $f \times f$

output : $(n+2p-f+1) \times (n+2p-f+1)$

Striding

No of steps taken for convolution.

Suppose we choose a stride of 2. So while convoluting through the image, we will take two steps: both in horizontal & vertical.

Max Pooling:-

It returns the maximum value from the position of the image covered by the kernel.

~~we can describe~~

The pooling layer is responsible for Reducing the spatial size of the convolved feature.

This is to decrease the computational power required to process the data. through dimensionality reduction. and hence speed up the computation.

two parameters:

1) Spatial Extent 'C'.

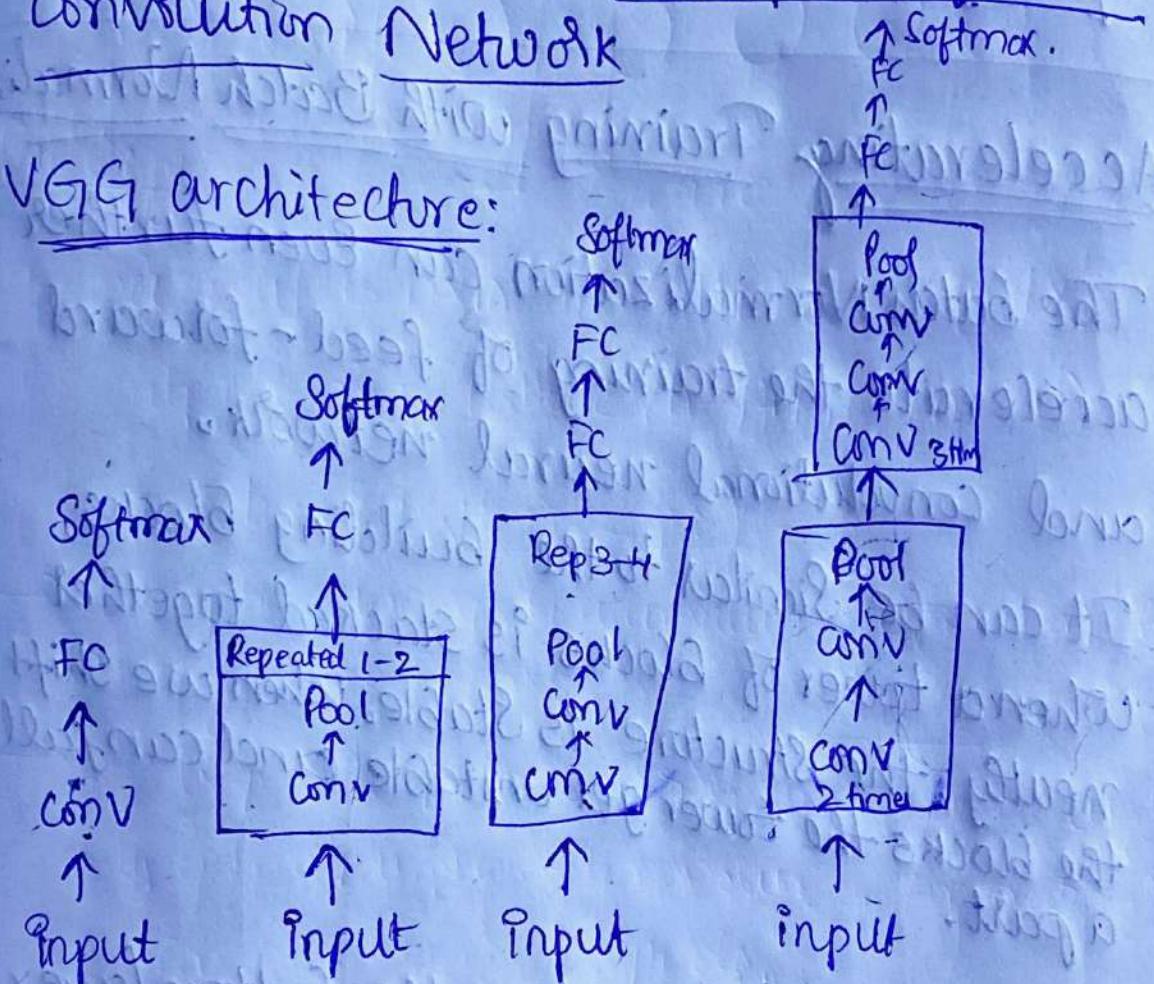
2) Stride 'S'.

$$\text{Width } W_{\text{out}} = \left\lceil \frac{W_{\text{in}} - e}{s} \right\rceil + 1$$

$$\text{height } H_{\text{out}} = \left\lceil \frac{H_{\text{in}} - e}{s} \right\rceil + 1$$

Full Architectural Description of Convolution Network

VGG architecture:



FC → Flatten

The VGG architecture is a deep convolution neural network known for its simplicity and its deep stack of convolution layers with small receptive fields.

Small filters.

Input layer: The input layer would be size $28 \times 28 \times 3$, assuming the images are RGB.

Link

Accelerating Training with Batch Normalization

The batch Normalization can even further accelerate the training of feed-forward and convolutional neural network.

It can be similar to the building blocks, when a tower of blocks is stacked together neatly, the structure is stable. When we shift the blocks the tower get unstable and can fall apart.

A similar phenomenon where the bottom layer (first layer) and top layer (last layer) happens, when we give data to the bottom layer the data keep changing while passing from bottom layer to upper layers. So the last layer has to keep adapting to the changes in order to give

UNIT 10

⑫ Collect output so the more the layers we have the more problem we get.

Normalization of image inputs helps out the training process by making it more robust to variations. Batch Normalization takes this a step further by normalizing inputs to every layer in our neural network.

We modify our network to include the below operations for batch Normalization:

- 1) Grab the vectors of bias which is incoming to a layer before they pass through the non-linearity.
- 2) Normalize each component of the vector of bias across all examples of the minibatch ~~by subtracting the mean~~.
- 3) Given normalized input \hat{x} , use an affine transform to restore representation power with 2 vectors of parameters: $\gamma \hat{x} + \beta$.

Batch Normalization also allows us to significantly \uparrow the learning rate.

UNIT-II

Questions:

- 1) Principal Component Analysis / Diff b/w PCA & Autoencoder
- 2) Autoencoder Architecture.
- 3) Denoising of Autoencoder
- 4) Sparsity in Autoencoder
- 5) Lower-Dimensional Representation
- 6) Diff b/w AE & PCA.
- 7)
- 8)
- 9)
- 10)
- 11)
- 12)
- 13)
- 14)
- 15)

Chapter 9:-

1) Learning Lower Dimensional Representation

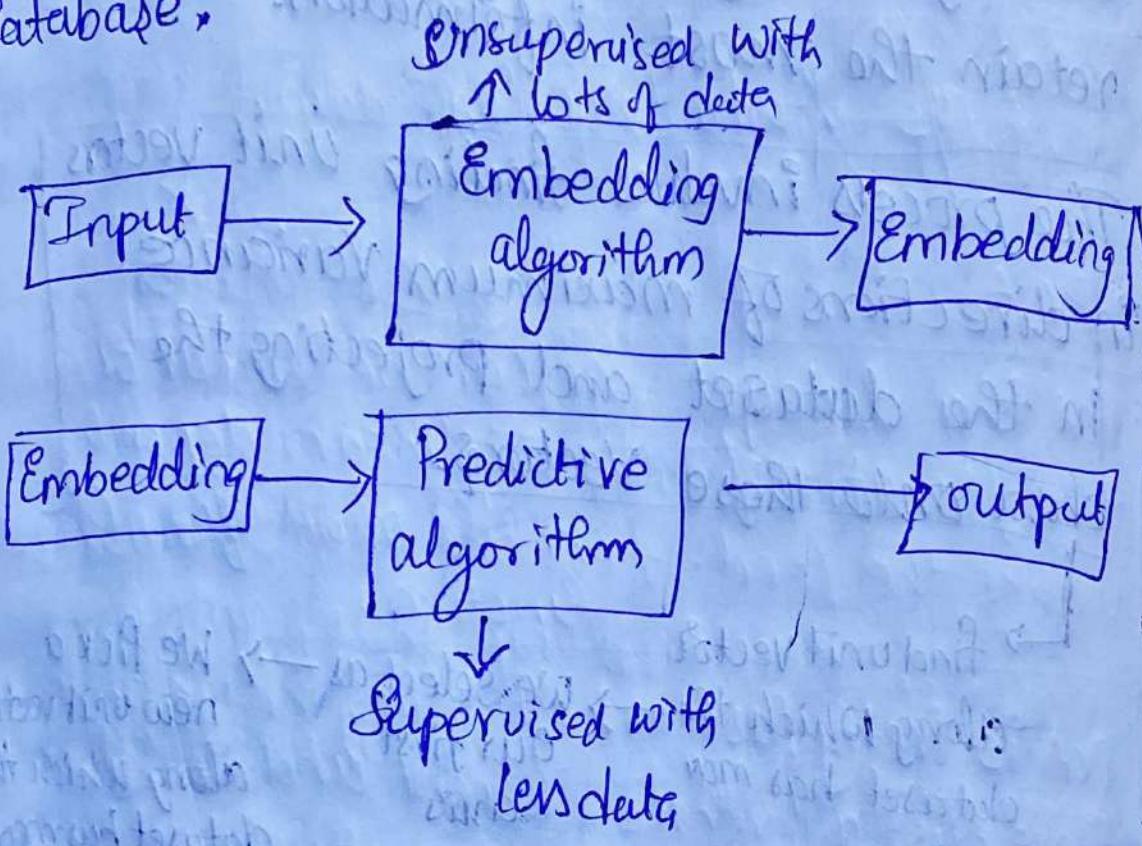
Large input vectors leads to larger models with many parameters, which can be Expressive but are also data-hungry and Prone to Overfitting.

The convolutional architectures address the curse of dimensionality by Reducing the number of Parameters, maintaining Expressiveness without need for Extensive training data.

There is often a Scarcity of labeled data, which is necessary for training Convolutional networks, leading to the needs for methods that can work effectively with less labeled data.

One solution is to develop models that learn from unlabeled data by creating lower-dimensional embeddings in an unsupervised manner.

These embeddings capture essential features and can be used with smaller models, reducing the requirements for large labeled database.



Principal Component Analysis:

Principal Component Analysis is a statistical procedure that aims to transform a set of possibly correlated variables into a set of linearly uncorrelated variables known as principal components.

The goal of PCA is to retain as much information as possible.
Seeks to reduce the dimensionality of data by finding the new axes that retain the most information.

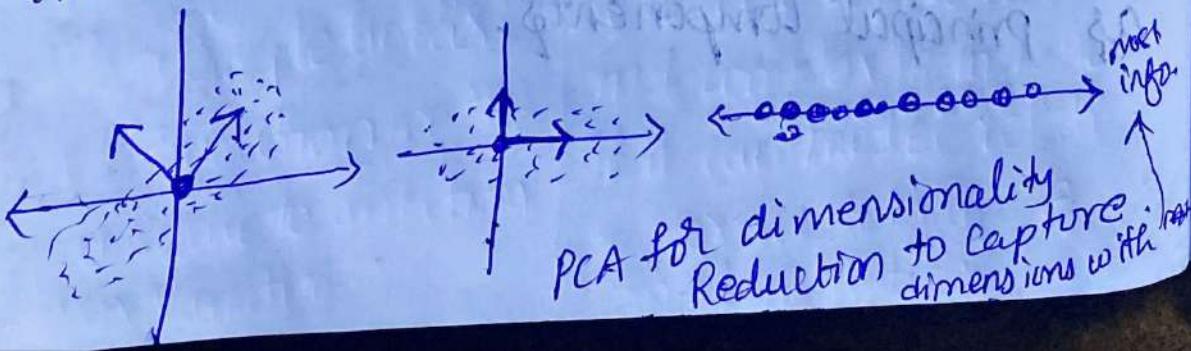
The process involves finding unit vectors in directions of maximum variance in the dataset and projecting the data onto these vectors.

→ find unit vector along which the dataset has max variance → we select as our first axis → we pick a new unit vector along which the dataset has max variance.

Project our data onto this new set of axes.

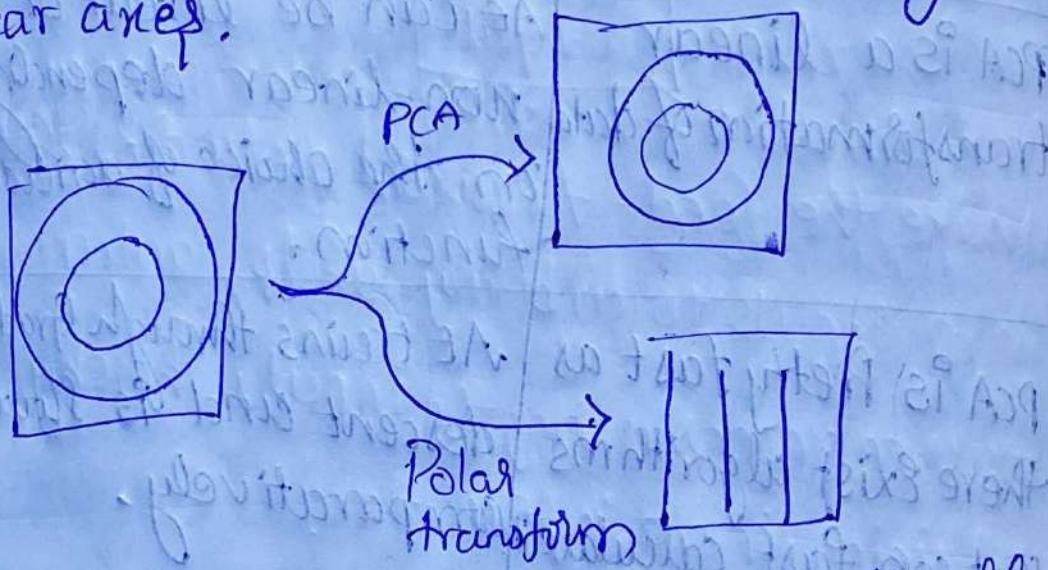
Continue this process until we have found a total of d new vectors.

This result is shown in below fig.



The dataset is represented as a matrix, and an embedding matrix is created that reduces dimensions while preserving as much info as possible.

PCA can fail with complex datasets, particularly where important patterns are nonlinear, as it primarily capture variance along linear axes.



A situation in which PCA fails to optimally transform the data for dimensionality Reducing

The nonlinear Relationships in data, such as the above concentric circles, may require diff approaches, like polar transformation, to Separate features Effectively.

The limitations of PCA in dealing with non linear relationships highlight the needs

for nonlinear dimensionality reduction techniques, which are addressed by neural models in deep learning.

PCA vs Auto Encoders

PCA

PCA is a linear transformation of data

PCA is pretty fast as there exist algorithms that can fast calculate.

PCA Projects data into dimensions that are orthogonal to each other resulting in very low or close to zero correlation in the projected data.

PCA is a simple linear transformation on the input space to directions of max variations

AE

AE can be linear or non-linear depending on the choice of activation function.

AE trains through Gradient descent and is slower comparatively.

AE transformed data doesn't guarantee that because the way it's trained is merely to minimize the reconstruction loss.

AE is more sophisticated and complex technique that can model relatively complex relationships & non-linearities.

one rule to thumb could be
the size of Data. PCA for
small datasets

AE for larger
dataset.

PCA hyperparameter is
"K"

AE it is the architecture
of the neural network.

AE with single layer.
& linear activation
has similar performance
as PCA.

AE with multiple layers
and non-activation function
termed as Deep Autoencoder
Prone to Overfitting.

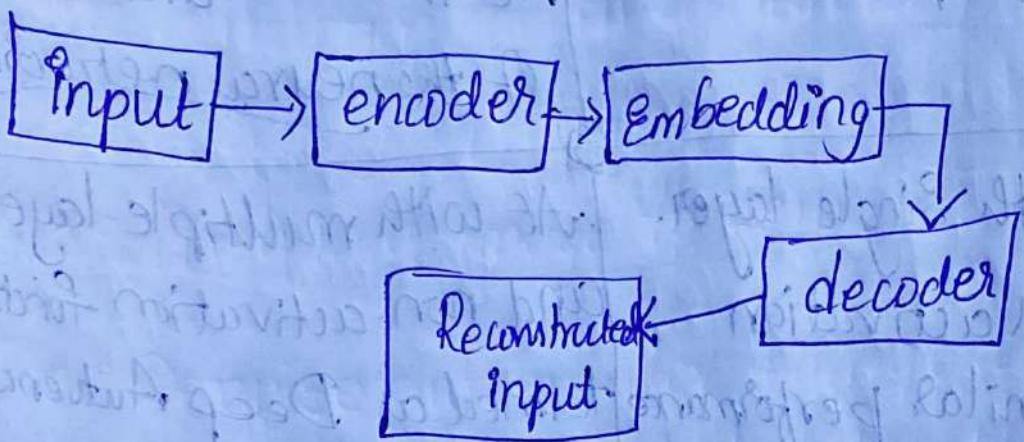
Autoencoder Architecture:

The feed-forward networks learn increasingly refined representation of input data, with the final convolutional layer output serving as a compressed representation.

Autoencoders are introduced as an architecture designed to address the loss of information by capturing and compressing input data into a lower-dimensional space.

An autoencoder consists of two main components:

an encoder that compresses the input into a code, and a decoder that reconstructs the original input from this code.

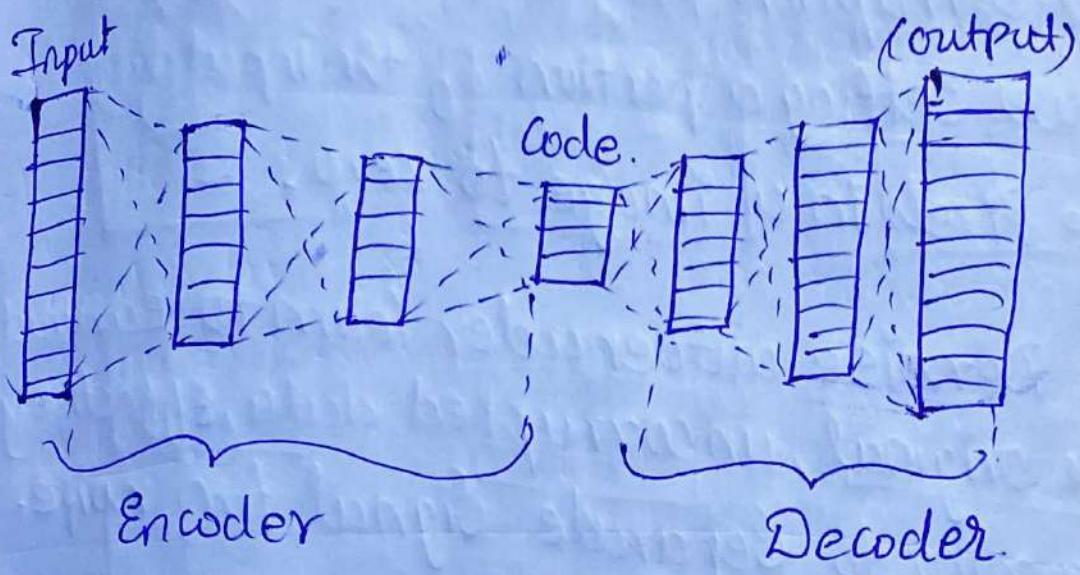


Unlike the feed-forward networks that map inputs to arbitrary labels, autoencoders focus on reconstructing the input, thereby retaining more information.

The encoder part of the autoencoder learns to identify and encode the most critical features of the input data.

The decoder part learns to reverse the encoding process to reconstruct the input data from the encoded representation.

Autoencoders can be used to generate lower dimensional data representation without the supervision required in traditional methods.



Denoising of autoencoder

Denoising:

Denoising is used to enhance autoencoders by generating noise-resistant embedding.

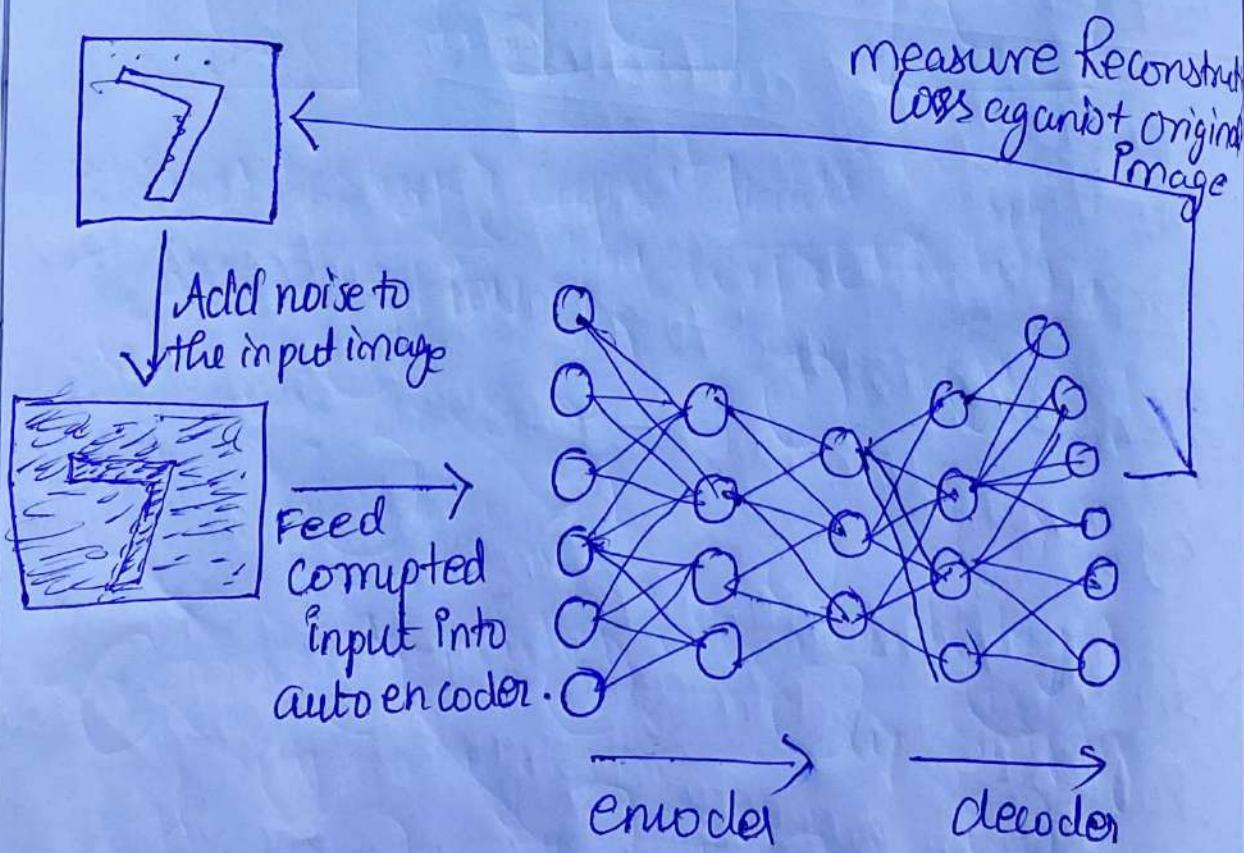
In the human perception we can recognize patterns and images despite significant noise, a property desirable for ML models.

We can even distinguish α and γ if all the pixels got corrupted easily.

The Denoising encoders are trained using corrupted inputs, forcing them to learn robust features.

The corruption process in denoising involves setting a portion of the input data, like pixels in an image to zero.

The Denoise autoencoder learns to reconstruct the original, uncorrupted data, effectively learning to separate signal from noise.



The geometric interpretation involves the concept of a manifold, the shape within the data that the autoencoder learns to reconstruct.

An autoencoder must identify whether a data point belongs to one manifold or another during the reconstruction process.

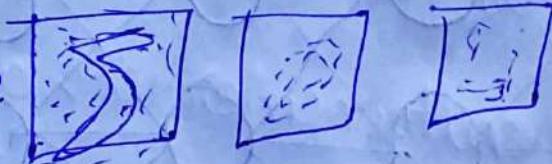
The denoising process conceptually expands the data space to include points around the manifold, with the autoencoder then collapsing these points back onto the manifold.

The denoising objective helps the autoencoder learn to distinguish b/w generalizable features and noise.

Original Image



Corrupted Image

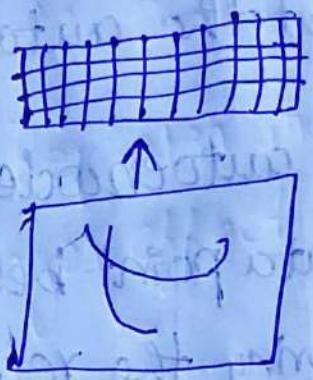
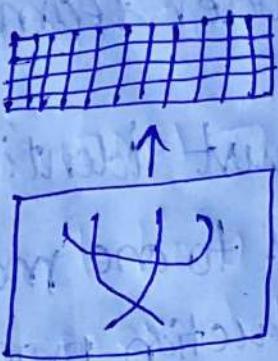
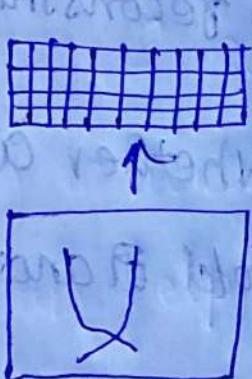


Reconstructed Image

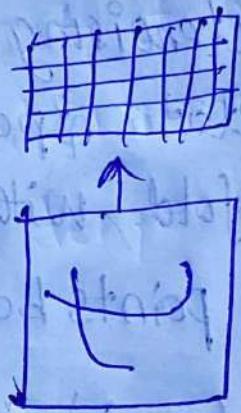
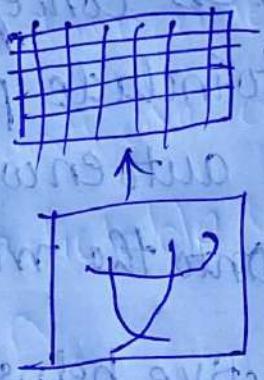
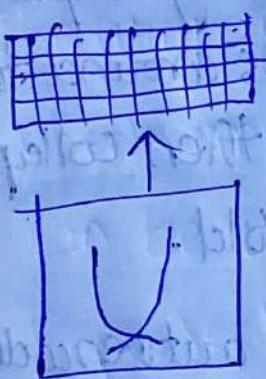


Sparsity in Autoencoders

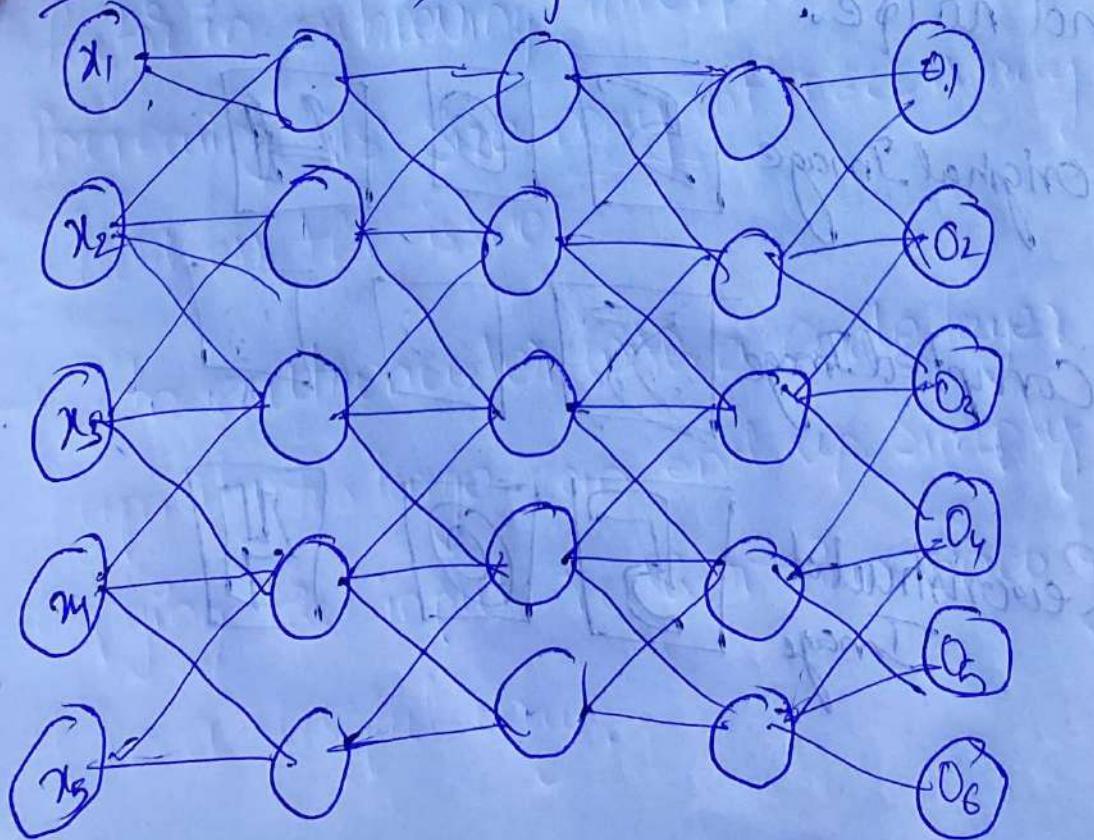
A



B



Input layer $\xrightarrow{\text{encoder}}$ hidden layer $\xrightarrow{\text{decoder}}$ output layer



Deep learning models struggle with interpretability due to their complexity and nonlinearity. Deep learning models are generally very difficult to interpret because of the nonlinearities and massive number of parameters that make up a model.

The interpretability can be solved by autoencoder. In general, an autoencoder representations are dense, which means that many features are combined in ways that are difficult to separate or interpret.

In fig A. The activation of a dense representation combine and overlay information from multiple features in ways that are difficult to interpret.

Introducing Sparsity into autoencoder representation can help make the model output more interpretable.

The goal is to achieve a one-to-one correspondence b/w the high level features in the input

data and the components in the autoencoder code.

→ The Fig B. colour - coding the contribution of individual features to the representation can help understand how changes effect the overall representation.

→ The way an autoencoder's output changes when components are added or removed can be unpredictable with dense representation.

→ The code's layer capacity can be limiting factor in achieving interpretability, increasing its size may not necessarily address the issue.

→ In an autoencoder with a code layer capacity too large can end up simply copying the input rather than encoding it in a useful way.

→ The Autoencoder can effectively summarize data points when all relevant info is present in the individual data points themselves.
