

21/7/24

Machine Learning

• Artificial Intelligence (AI) : The ability of machine to think, reason and act like a human. (mimic a human).

- The term was coined by 'John McCarthy' (american) in 1956 at the Dartmouth Conference.
- A British mathematician and logician, 'Alan Turing' introduced the 'Turing Test', originally called 'imitation game', was proposed as a way to assess a machine's ability to exhibit intelligent behaviour indistinguishable from that of a human (when put in a dark room).

• Machine Learning (ML) : A subset of AI, that allows computers to learn from data and improve their performance over time without explicit programming. (glass / white box).

- T : Task then only you can say (is)
- P : Performance that the machine is
- E : Experience - learning.

31/7/24

• Neural Networks (NN) : They are computational models inspired by the human brain that learn from data through interconnected layers of neurons. Effective for image and speech recognition. (Artificial Neurons).
(shallow learning)

➢ Architecture of a neuron :-

- Synapsis - data collection
- Soma - able to process something
- Dendrites - data giving (pass the data to another neuron).
- Axon - a thread which carries the process data to dendrites.

* 60% of ML can be solved by : $y = mx + c$

ML - { Probability
Statistics }

PAGE NO.: 2.

Deep learning (DL) : A ML technique that uses deep neural networks with multiple layers to learn complex patterns from large datasets.

* Overfitting : Occurs when a ML model becomes excessively complex, memorizing the training data rather than learning underlying patterns. This leads to poor performance on new, unseen data. (Basically fitting only type of data, excessively).

* I/P ^{Hidden layer} O/P
layer ||||| layer Activation function : which fires the neurons.

* Neural Network / Deep learning \rightarrow Black box
as a type of testing.

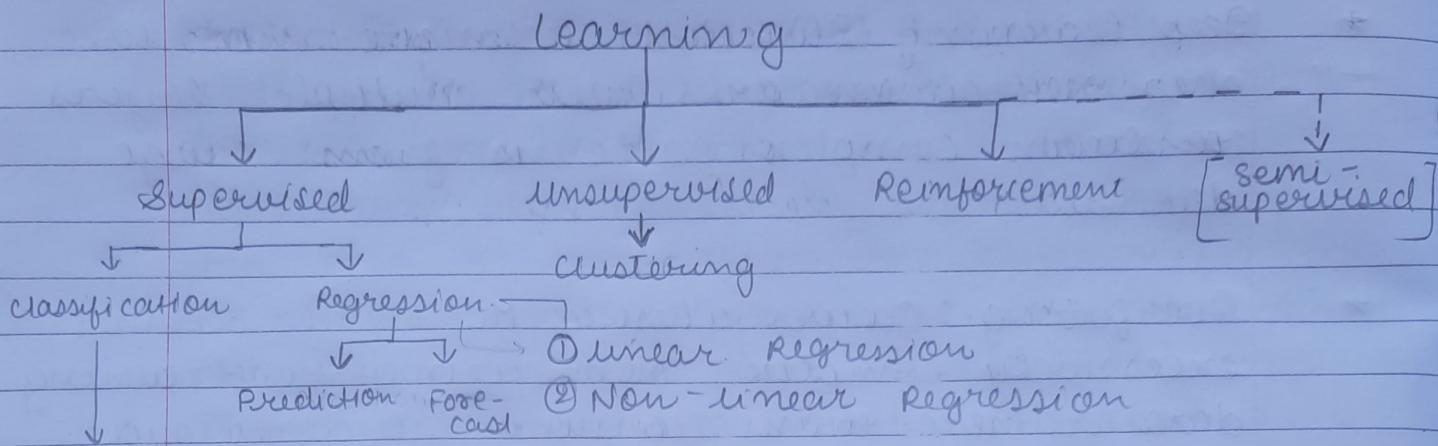
* Internal behaviour of the system is always different because no. of weights will be a learning parameter which will always change; that would be the input. And the output will be learning parameter.

Types of System learning in ML :

* ChatGPT - generative pre-trained transformer (2017)

* NLP - Natural language Processing.

9/7/24



① Binary (2 class)

② Three class

③ Multi class

- **Supervised learning :** A model learns from labeled data, where input data is paired with desired output. (like teacher's guidance)
- **unsupervised learning :** A model learns from unlabeled data, discovering patterns and relationships within the data itself. (like learning independently)
- **Reinforcement learning :** A model learns by interacting with an environment, receiving rewards or penalties for its actions. (like learning through trial and error).

(optional)

- **Semi-Supervised learning :** A model learns from a combination of labeled and unlabeled data, leveraging both to improve its performance. (like learning with partial guidance).



Discriminant Functions : A discriminant function takes an input vector ' x ' and assigned it to one of ' k ' classes; where the class is denoted by C_k .

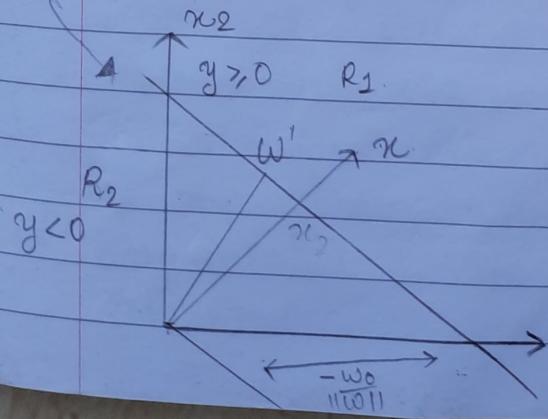
(i) Two class classification

The simplest representation of a linear discriminant funcn is obtained by taking a linear function of input vector ' x '.
 ↗ recurring parameters.
 ↘ Bias

$$\therefore Y(x) = \underline{w^T} \underline{x} + \underline{w_0} \quad \text{--- (1)}$$

weight vector.

- The negative of the bias is called threshold.
- An input vector ' x ' is assigned to some class C_1 , if $Y(x) \geq 0$ and assigned to C_2 otherwise.
- The Corresponding decision boundary is therefore defined by the relation of $Y(x) = 0$.
- $Y(x) = 0$, corresponds to ' $D-1$ ' dimensional hyperplane within D -dimensional input space

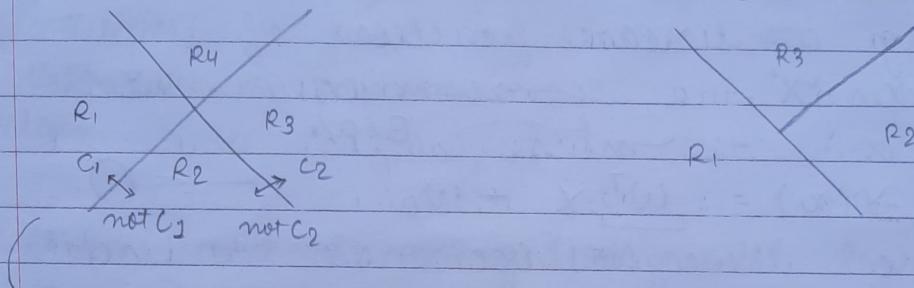


$\frac{y(x)}{\|w\|}$ → The orthogonal distance from point x to decision boundary

Displacement from the weight

ii) Multi-class Classification

It is the extension of $K > 2$ classes. The use of $K-1$ classifier, where each of the classifier solves a two class problem of separating points in a particular class C_k from the points not available in that class is known as 'One - versus - the - rest' classifier.



→ Here, due to the classification arrangement we would get the fourth region and that is an unknown and un-identified region, so it is considered as to be the worst case arrangement. There are only two possible ways of this arrangement i.e. super accurate classification or problem with arrangement.

In Example involving three classes where this approach leads to a region of input space i.e. ambiguous unclassified

$$K(K-1)/2$$

Dimensionality Reduction:

- i) PCA (Principle Component Analysis)
- ii) ICA (Independent Component Analysis)
- iii) FA (Factor analysis)

The primary reasons to simplify the datasets are:-

- (i) Making a dataset to use
- (ii) Reducing the computation cost
- (iii) Reducing the noise inside the dataset for better generalization
- (iv) Making the result easier to understand.

(v) To fit in,

⇒ Principle Component Analysis

A dataset is transformed from its original co-ordinate system to the new co-ordinate system.

The initial new Axis are chosen in the direction of higher variance of the data and later axis are the orthogonal projection of the initial axis.

The same procedure that's repeated for as many features as we have in the original data.

The PCA, reduces the complexity of the data and identifies the most important features which are known as Principle Component Feature.

~> Independent Component Analysis.

The ICA, assumes that the data is generated out of 'n' number of sources.

Here, the data is assumed to be the mixture of observation of the available sources. These sources are assumed to be statistically independent which represents that the data is uncorrelated with each other.

If there are fewer no. of sources than the amount of observed data, then we can achieve the Dimensionality Reduction.

29/7/24

ML activities / General Process of ML :-

In order to develop ML system, multiple processing pipeline require along with some typical preparation activity, once the input data becomes available.

Few of the required activities consist of :

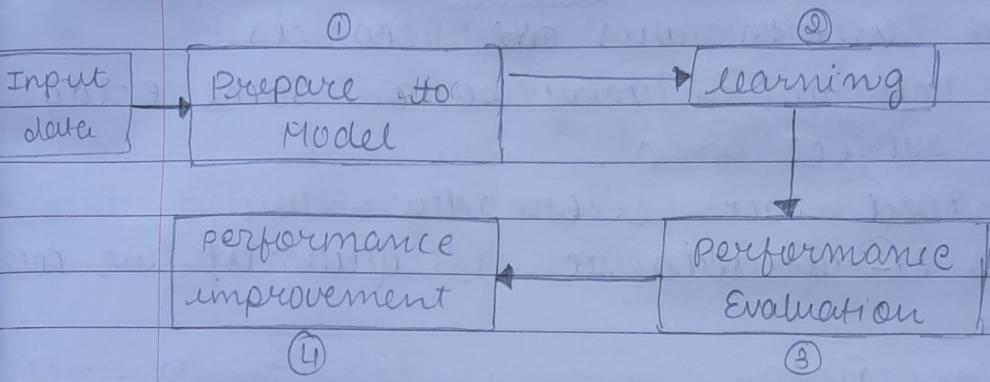
- i) Understand the type of data in the given data set.
- ii) Explore the data to understand the nature, the quality and quantity.
- iii) To understand inter-feature relationship
- iv) Identify the potential issues in the

available dataset, if any.

to perform the division of the dataset

used for evaluation }
 of the model {
 performance. }

-> Training Set	80%
-> Validation Set	(10%)
-> Testing Set	20% (10%)



31/7/24

Step-1: Understand the type of data in the given dataset, which includes pre-processing part of dimensionality reduction and feature selection.

Step-2: learning. In the ML pipeline, the learning step includes the partition of the data, selection of model and cross validation.

Step-3: Performance Evaluation. Here, we conduct the evaluation of ML model with various performance matrices. Here, we should visualize the performance and the behaviour of the model.

Step-4: Performance Improvement. This include

training the model and ensembling, bagging and boosting.

Types of Data in Machine Learning :

In ML, the dataset is a collection of related informations and records.

Information may be on some entity or some subject area.

Data can be broadly divided into two types : ① Qualitative ② Quantitative data

1. Qualitative Data.

It provides information about the quality of an object or information which cannot be measured. (eg: Good or Bad)

It is also called Categorical data, which can be further divided into two types : ① Nominal ② Ordinal

-4

The Nominal Data does not contain any numerical value but a named value. This values cannot be further quantified. (eg: Nationality, blood groups, gender)

-4

The Ordinal Data, in addition to processing of nominal data can also be naturally ordered. (eg: Customer satisfaction, Grades of Student, hardness of material).

2. Quantitative Data

It relates to information about the quantity of an object and hence it can be measured. (we can take interval & ratio as well)

* Explaining the numerical data :-

	10,000
--	--------

0 - 255

* If bunch of quantities are given and we need to segregate the data, then we can see by taking mean, median and variance & look up on the histogram.

• Understanding the Central Tendency :-

To understand the nature of numeric variables, we can apply the measures of central tendency like mean and median.

• Understanding the data spread :-

After central tendency of numerical attributes, we have a clear idea of which attributes have a large deviation between mean and median. Here, we can apply two approach : (1) Dispersion of data (2) Position of the data values.

SVM (Support Vector Machine) is highly accurate given the appropriate size of data and SVM can be used for numerical prediction and classification.

In SVM, the separative hyperplane can be written as,

$$\underset{\text{weight vector}}{w} \cdot \underset{\text{input tuples}}{x} + b = 0 ;$$

where w is weight vector, $w = \{w_1, w_2, \dots, w_n\}$ is a no. of attributes and b is the scalar.

→ x is input vectors, $x = \{x_1, x_2, \dots, x_n\}$, where x_i is the values of the given attributes.

→ Any points ^{that} lies between the separative hyperplane is denoted by,

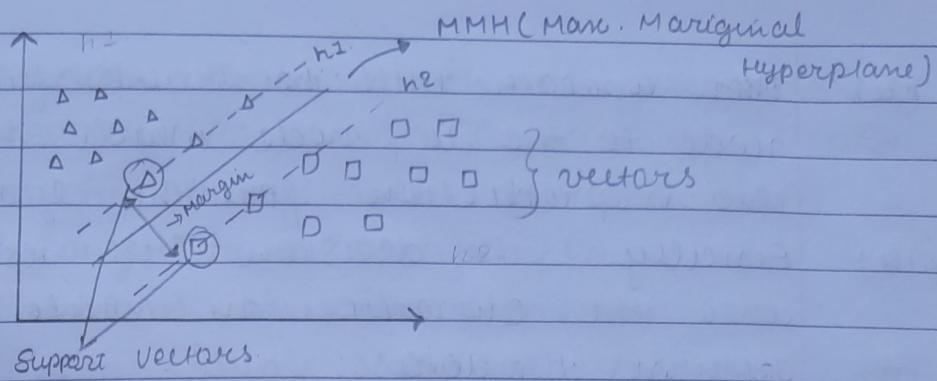
$$w \cdot x + b > 0.$$

→ Similarly, that lies below the hyperplane is denoted by,

$$w \cdot x + b < 0.$$

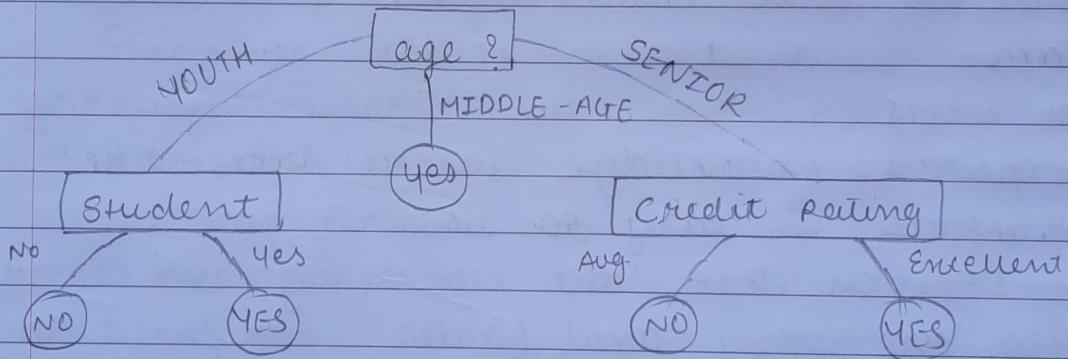
→ Any training tuples that fall on the hyperplanes ^{$h_1 \& h_2$} are called support vectors.

The Support Vectors are most difficult tuples to classify and gives most crucial information regarding classification.



* Decision Tree Classifier :- (Approach)

A Decision tree is a structure where each internal node denotes a test on an attribute. Each branch represents the output of the test and each leaf node holds the label of the class.



In decision tree, the classification algorithm will generate a binary tree structure, which can lead the testing data to the appropriate class label.

-4

How the decision tree classifier are used for classification?

(i) Given a tuple 'x' for which the associated class label is unknown, the attribute vector of the tuple are evaluated against the decision tree



- (ii) Then a path will be traced from the root node to the leaf node which holds the class prediction for the given tuple.
- (iii) Finally, the decision tree is converted into the classification rule which can return the class label:
- 4 Why the decision tree classifier are so popular? (Significance)
- The construction of the decision tree classifier does not require any domain related knowledge.
 - Decision tree can handle multi-dimensional data.
 - In Decision tree, the representation of acquired knowledge in a tree like structure is easy to understand.
 - In Decision tree, the learning and classifications step are easy and faster, compared to other similar techniques.
 - Decision tree classifier, offers good accuracy.

21/8/24
#

Bayes' Classifier : $P(H|X)$

Posterior \rightarrow
 $P(H|X)$ \leftarrow Prior

- They are statistical classifier.
- They can predict class membership probabilities such as the probability that a given tuple belongs to a particular class.

- (3) Let ' x ' be a data tuple, which is considered as an evidence in Bayesian Terminology and ' H ' be some hypothesis such as the data tuple ' x ' belongs to a specific class ; say ' C_1 '.
- (4) For the problem of classification, we want to determine ; $P(H|X)$, where the probability that hypothesis ' H ' holds given the data tuple ' X '.

$$\therefore P(H|X) = \frac{P(X|H) \times P(H)}{P(X)} ; \text{ Bayesian Equation.}$$

Likelihood Calculation

SR.	Weather	Play
1.	Sunny	NO
2	Rainy	YES
3.	Sunny	YES
4.	Overcast	YES
5.	Rainy	NO
:	:	:
14.		

Pop(YES)	Pop(NO)	Weather	YES	NO	
4/9	0/5	Overcast	4	0	$= 4/14 = 0.29$
3/9	2/5	Sunny	3	2	$= 5/14 = 0.36$
2/9	3/5	Rainy	2	3	$= 5/14 = 0.36$
		Total =	9	5	
					$9/14 = 0.64 \quad 5/14 = 0.36$

$$(ii) P(\text{YES}/\text{overcast}) = \frac{P(\text{Overcast}/\text{YES}) \times P(\text{YES})}{P(\text{Overcast})}$$

$$= \frac{4/9 \times 9/14}{4/14} = \frac{0.44 \times 0.64}{0.29} = 0.97$$

$$(iii) P(\text{NO} \mid \text{Overcast}) = \frac{P(\text{Overcast} \mid \text{NO}) \times P(\text{NO})}{P(\text{Overcast})} = 0.03$$

* Steps of Bayes' Classification :- (Single Feature)

- Step-1: Calculate the prior probability for given class ^{label}.
- Step-2: Find likelihood probability with each attribute with each class.
- Step-3: calculate posterior probability , putting the values of likelihood probability.
- Step-4: Check and select the class / hypothesis having highest probability

* Steps of Bayes' Classification :- (Multiple Feature)

- Step-1: Calculate prior probability for given class label
- Step-2: Calculate conditional probability with each attribute with each class.
- Step-3: Multiply same class conditional probability
- Step-4: Multiply prior probability with the conditional probability of (Step-3).

Step-5: Check and select the class label or hypothesis for highest probability.

- $P(\text{YES} | w = \text{overcast}, \text{Temp} = \text{Mild}) = P(w = \text{overcast}, \text{Temp} = \text{Mild}) \cdot P(\text{YES})$ ↳ ①
- $P(w = \text{overcast}, \text{Temp} = \text{Mild} | \text{YES}) = P(\text{overcast} | \text{YES}) \cdot P(\text{Mild} | \text{YES})$ ↳ ②
- The goal in classification is to take an input vector ' x ' and to assign it to one of the hypothesis based on the probability calculated with an evidence.
- In Bayes' Theorem, we are looking at the probability that a given tuple ' x ' having certain no. of feature belongs to one of the ' C ' classes having the prior knowledge of attribute description of ^{given} tuple ' x '.

28/8/24

Clustering in

1. Partition Based Methods

This category of methods uses mean values of the available data sample to represent the center of ^{the} cluster.

2. Hierarchical Methods.

This category of methods propose hierarchical or tree like structure through

decomposition or merging of data samples.

3. Density Based Methods

This category of methods are applicable for identifying arbitrarily shaped clusters.

* K-means Clustering :-

In supervised learning, we try to match input to some existing pattern while in unsupervised learning we will try to discover the patterns in raw and unlabeled data.

Clustering is the process of grouping similar type of data by isolating these similar data.

We want the data points in a cluster with a common property that separates the data points of other clusters.

Clustering is used to identify the potential group in a dataset while classification is used to match an input to an existing group.

We can define the task of clustering as :

Input: Training Pairs

$$\mathcal{D}_{\text{train}} = [x_1, x_2, \dots, x_m]$$

Output: Assignment of each point to cluster

$$z = [z_1, z_2, \dots, z_m]; \text{ where } z_i \in \{1, \dots, k\}$$

The task of clustering is take a set of data points as input and return a partitioning of the points into k-clusters.

- Step-1: Select k-points in the data space and make them as initial centroids.
- Step-2: Assign each point in the data space to the nearest centroid to form k-cluster.
- Step-3: Measure the distance of each point in the cluster from the centroid.
- Step-4: Calculate the error to measure the quality of the cluster.
(with SSE)
- Step-5: Identify new centroids of each cluster on the basis of distance between the points.
- Step-6: Repeat the steps ② to ⑤ until centroids needs no changes.

Probability of Machine Learning :-

In ML, we train the system by using a limited set of data called the training data. And based on the confidence level of the occurring data, we expect the ML algorithm to depict the behaviour of the

larger set of actual data.

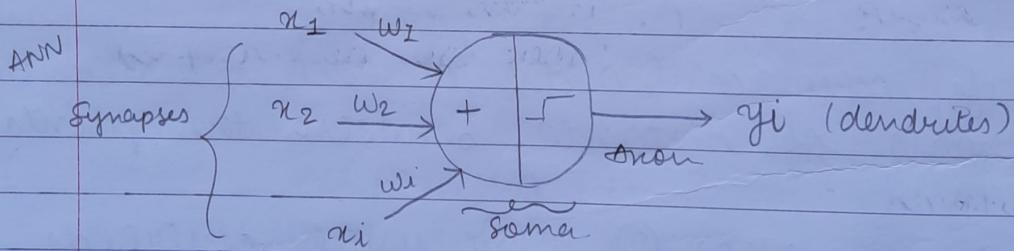
If we have an observation on a subset of events called a sample, then there will be some uncertainty in attributing the sample results to the whole set or population.

Question arises, how a limited knowledge of a sample set can be used to depict the behaviour of the real set with the same confidence?

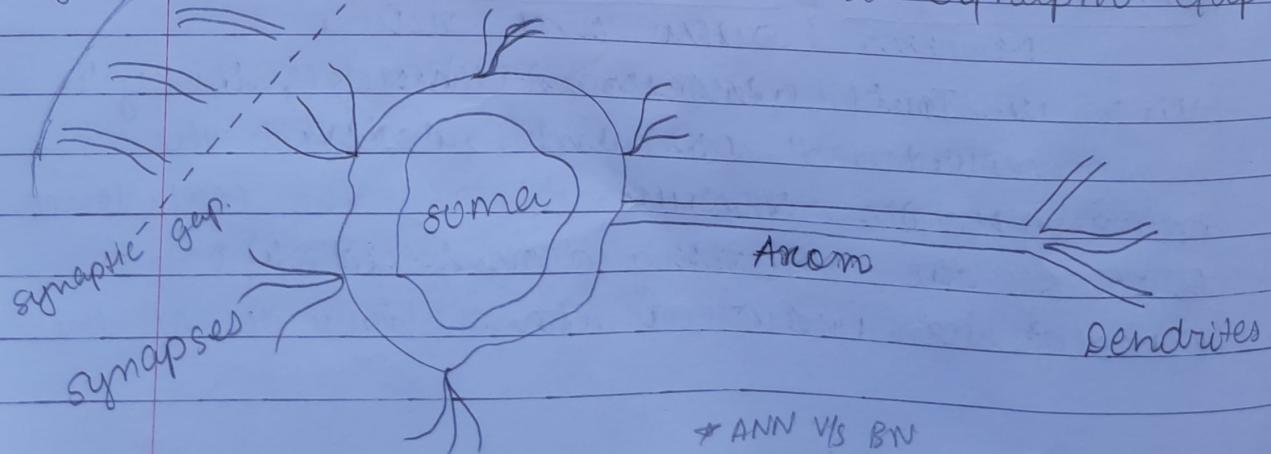
$$\boxed{\text{Knowledge of Sample}} + \boxed{\text{Known uncertainty}} = \boxed{\text{useful Knowledge.}}$$

18/9/24

Artificial Neural Networks :-

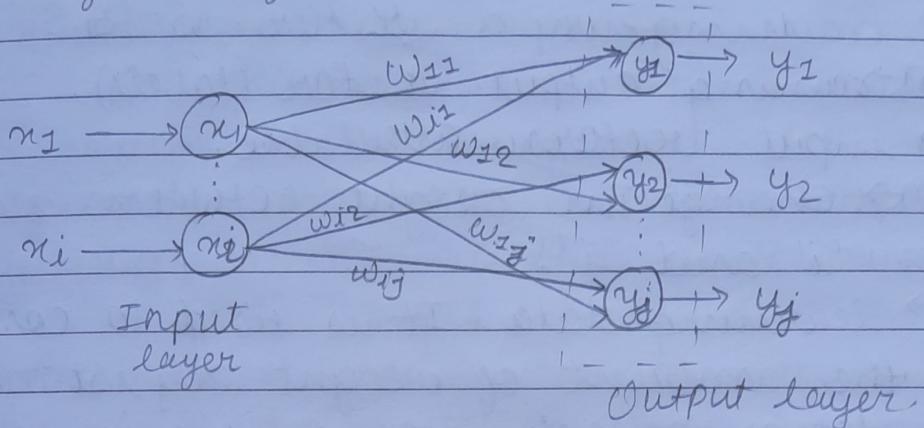


- * There is always a gap between two interconnected neurons and it is called 'Synaptic Gap'.

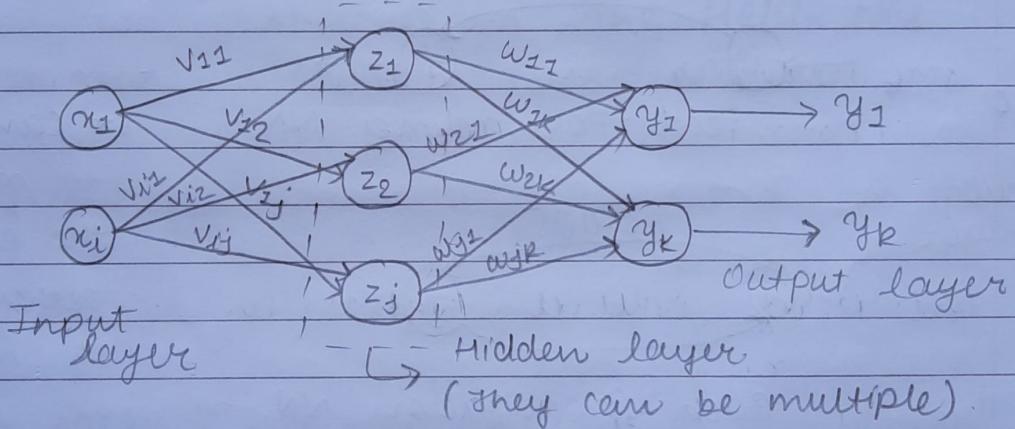


* ANN v/s BN

i) Single layer ANN



ii) Multi-layer ANN

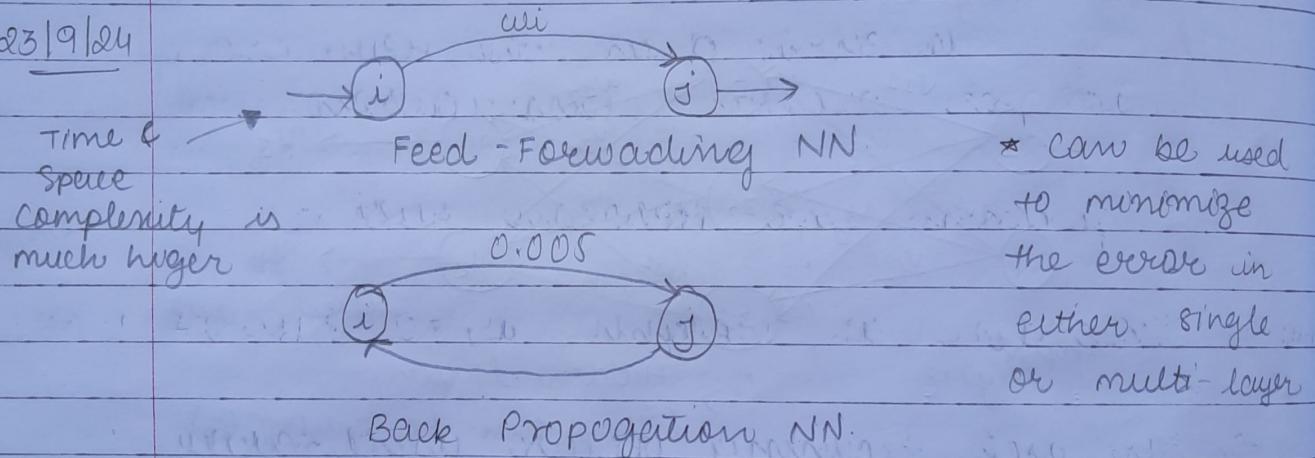


* Notations :-

- $x_i \rightarrow$ inputs signal
- $y_j \rightarrow$ output signal (produced by $f(y_in_j)$)
- $w_{ij} \rightarrow$ weight between 'i' and 'j'. usually between hidden and output
- $v_{ij} \rightarrow$ weights between 'i' and 'j', usually between input and hidden
- $b_j \rightarrow$ Bias on Neuron 'j'
- $W \rightarrow$ weight matrix formulated, $W = \{w_{ij}\}$
- $y_in_j \rightarrow$ Net input to unit ' y_j '.

- $\theta_j \rightarrow$ threshold on the unit 'j'
- $S \rightarrow$ training input vectors, $S = \{S_1, S_2, \dots, S_{12}\}$
- $t \rightarrow$ training output vector (labels)
- $x \rightarrow$ input vectors available.
- $\Delta w_{ij} \rightarrow$ change in weight between units 'i' and 'j'.
- $\eta \rightarrow$ learning data. It is used to control (eta) the amount of weight adjustment during each step of training.

23/9/24



Hebb Net \approx (\rightarrow Hebbian learning rule
 \rightarrow extended Hebbian rule)

The earliest and simplest learning rule for a NN is generally known as Hebb Rule.

Hebb proposed that the learning occurs by modification of synaptic strength or weights.

In a manner, such that if two interconnected neurons are in activated state at a same time then the synaptic

strength between those neurons should be increased.

However, the strongest form of learning occurs if we determine the synaptic strength between neurons during their inactive state.

This result, extended Hebbian rule which is a single layer feed-forward NN architecture, popularly known as Hebb Net.

In early days, the Hebb Net was used for pattern classification.

* Algorithm for Hebbian Rule :-

Step -1: Initialize all weights, $w_i = 0$ for $i = 1 \text{ to } n^3$.

Step -2: For each training vector and target output pair 's:t' perform step ③ to ⑤

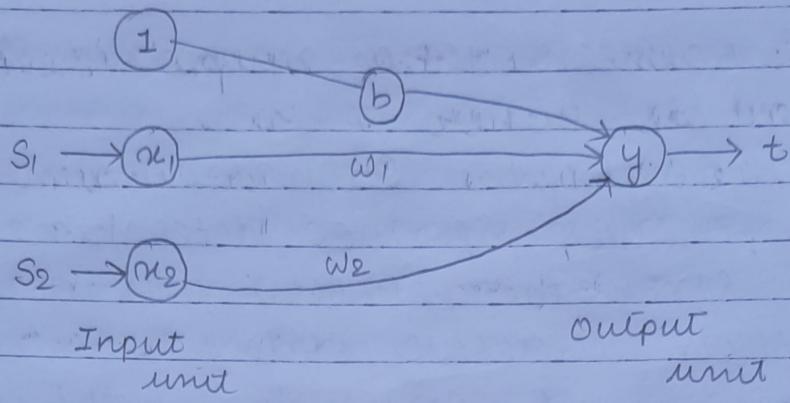
Step -3: Set activations for input unit,
 $x_i = s_i$ (for $i = 1 \text{ to } n$).

Step -4: Set activations for output unit,
 $y_j = t_j$ (for $j = 1 \text{ to } m$)

Step -5: Adjust the weights,
 $w_i(\text{new}) = w_i(\text{old}) + x_i \cdot y_j$.
 (for $i = 1 \text{ to } n$)

Always activated;
 $b_i(\text{new}) = b_i(\text{old}) + y_j$.

Step -6: For the stopping condition.



Perception :-

After the Hebb Net, the perceptron are having far reaching impact on NN community in the early days.

The perceptron learning rule is a more powerful rule compared to the Hebbian Rule under the specific assumption its iterative learning procedure can be proved ^{with} convergence.

are the no. of different types of perception proposed, since 1962, having the linear capabilities. The Perceptron can be decomposed into 3 major blocks:

- i) The Sensory Unit
- ii) The Associatory Unit
- iii) The Response Unit

The activation fnⁿ for each associator unit is the binary step fnⁿ with the arbitrary but fixed threshold '0'.

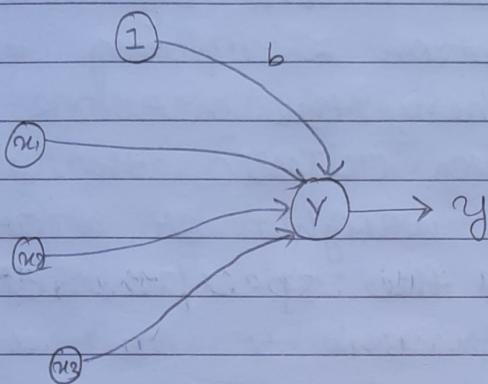
So, the signal sent from the

associator units to the output / response unit, will be either '1' or '0'.

The Output of the Perceptron is

$$y = f(y_{-in}) ;$$

$$f(y_{-in}) = \begin{cases} 1 & ; y_{-in} > 0 \\ 0 & ; -\theta \leq y_{-in} \leq 0 \\ -1 & ; y_{-in} < -\theta \end{cases}$$



The weights from associator unit to the response unit are getting adjusted using 'perceptron learning rule'. For each training input, the net would calculate the response of the output unit. Then the net would determine whether an error occurred for the given pattern or not, comparing the output of the net with the targeted value of the given input pattern.

The Perceptron Learning Rule includes an additional known as 'learning Rate' and the rule can be state as,

$$\therefore w_i (\text{New}) = w_i (\text{old}) + \eta t n_i$$

If the error occurred during learning a specific pattern then the

weights of the network are getting updated using the given Perceptron learning rule.

The training of the network would continue until no error remains or the possible error minimization.

The Perceptron learning rule converges only when the network meets a specific learning criteria that results the best possible weights or learning parameters, given the specific set of input sequences.

* Perceptron Algorithm :-

Step -1: Initialize the weight and bias. For simplicity set weights and bias to '0'. Set the learning rate between, ($0 < \eta < 1$) [we can set the learning rate to 1].

Step -2: While stopping condition is false, repeat step ③ to ⑥.

Step -3: Set activations of input unit, $a_i = s_i$

Step -4: Compute the response of output unit,
 $y_{-in} = b + \sum x_i \cdot w_i$

Step -5: Update the weights and bias, if the error occurred for the given pattern

$$\text{If } y \neq t; \quad w_i(\text{New}) = w_i(\text{Old}) + \eta \cdot \text{tri}$$

$$b(\text{new}) = b(\text{old}) + \eta \cdot t$$

else ; $w_j(\text{new}) = w_j(\text{old})$
 $b(\text{new}) = b(\text{old})$

Step - 6: For the Stopping Condition.

30/9/20

Back-Propagation ANN :-

↳ Nomenclatures.

$$\textcircled{1} \quad n = (n_1, \dots, n_i, \dots, n_n)$$

$$\textcircled{2} \quad t = (t_1, \dots, t_i, \dots, t_n)$$

$$\textcircled{3} \quad \delta_k = w_{jk} \cdot y_k$$

↳ proportion of error correction

weight adjustment for w_{jk} i.e. due to an error at the output unit, y_k . Also the information about the error at the unit y_k is propagated back to the hidden units that fit into the y_k .

$$\textcircled{4} \quad \delta_j = v_{ij} \cdot z_j$$

↳ Proportion

for v_{ij} i.e. due to the back propagation of error information from the output layer to the hidden unit z_j .

$$\textcircled{5} \quad \eta = \text{learning rate}$$

$$\textcircled{6} \quad n_i = s_i$$

Input unit 'i' for an input unit, the input signal and the output signal are the same.

$$\textcircled{7} \quad v_{oj} = \text{Bias on hidden } j'$$

Bias on the hidden unit 'j'.

* Activation function. $f(x) = \frac{1}{1 + e^{-cx}}$

Sigmoidal

⑧ $z_j = z_{\text{in } j}$

The hidden unit 'j', the net input to 'z_j' is denoted by 'z_{in}_j' and will be calculated by,

$\therefore z_{\text{in } j} = b_{0j} + \sum_i w_{ij} \cdot u_{ij}$
and the output signal or activation of 'z_j' is denoted by 'z_j', (small 'z').

$$\therefore z_j = f(z_{\text{in } j})$$

⑨ w_{0k} = Bias on output unit 'k'.

⑩ y_k = Output of output unit 'k'; k is no. of possible outcomes and the net input to 'y_k' is denoted by, 'y_{in}_k'.

$$\therefore y_{\text{in } k} = w_{0k} + \sum_j z_j \cdot w_{jk}$$

and

$$\therefore y_k = f(y_{\text{in } k})$$

* Back Propagation algorithm:-

Step-1: Initialize the weights and bias.

Step-2: While stopping condition is false, perform step ③ to ⑩

Step-3: For each available training pair, feed forward, perform step ④ to ⑨.

Step-4: Each input unit 'x_i' receives input signal 'x_i' and broadcast this signal to all the units in layer above i.e. hidden units

Step-5: Each hidden unit ' z_j ' sums the weighted input signal given by,

$$\therefore z_{-in\ j} = v_{oj} + \sum_{i=1}^m w_{ij} \cdot v_{ij}$$

& It applies activation fnⁿ to compute its output signal, given by

$$\therefore z_j = f(z_{-in\ j})$$

& sends this output signal to the unit in the layer above i.e. Output units.

Step-6: Each output unit ' y_k ' sums its weighted input signal, given by,

$$\therefore y_{-in\ k} = w_{ok} + \sum_{j=1}^m z_j \cdot w_{jk}$$

& applies its activation fnⁿ to compute its output signal, given by,
 $\therefore y_k = f(y_{-in\ k})$

Step-7: Each output unit ' y_k ' receives a target pattern for corresponding to the input training pattern. We can compute the Error Information Term by,

$$\therefore \delta_k = (t_k - y_k) \cdot f'(y_{-in\ k})$$

To calculate the weight correction term, we have,

$$\therefore \Delta w_{jk} = n \cdot \delta_k \cdot z_j$$

To calculate its bias correction term, we have,

$$\therefore \Delta w_{ok} = n \cdot \delta_k$$

Step-8: Each hidden unit ' z_j ' sums up Δ inputs from layer above hidden layers,

$$\therefore \delta_{\text{in } j} = \sum_{k=1}^m \delta_k \cdot w_{jk}$$

Multiply the quantity by the derivative of its activation function to calculate its error information term.

$$\therefore \delta_j = \delta_{\text{in } j} \cdot f'(z_{\text{in } j})$$

& To calculate its weight correction term, we are having,

$$\therefore \Delta w_{ij} = n \cdot \delta_j \cdot n_i$$

& To calculate its bias correction term, we have,

$$\therefore \Delta v_{oj} = n \cdot \delta_j$$

Step-9: Each output unit ' y_p ' updates its weight and bias using,

$$\text{O/P unit: } \therefore w_{jk}(\text{new}) = w_{jk}(\text{old}) + \Delta w_{jk}$$

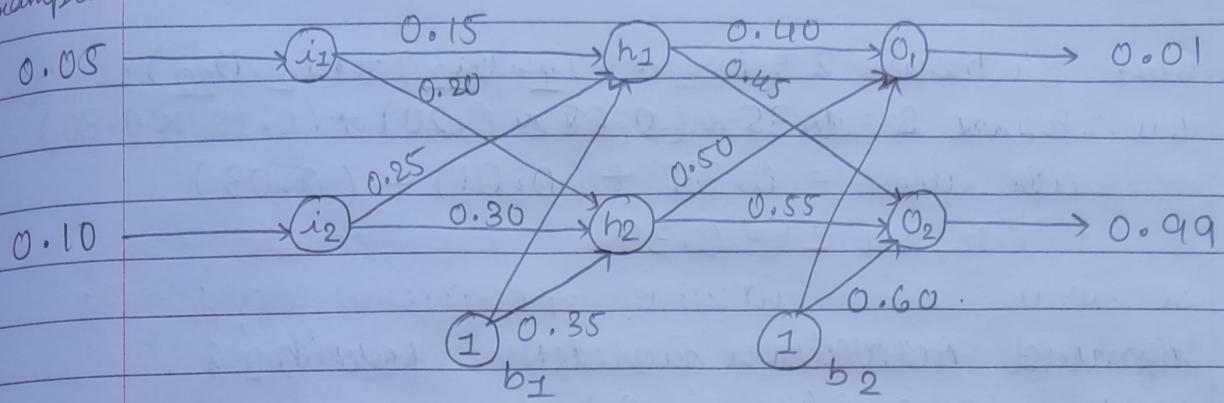
Each hidden unit ' z_j ' updates its bias.

$$\therefore v_{ij}(\text{new}) = v_{ij}(\text{old}) + \Delta v_{ij}$$

Step-10: It is for the stopping condition; it can be more than one.

$$\text{Sigmoidal} \Rightarrow f(x) = \frac{1}{1 + e^{-bx}}$$

Example:-



Given,

→ (i) Inputs : $i_1 = 0.05$ $i_2 = 0.10$.

weights (input - hidden) : $w_{11} = 0.15$ $w_{21} = 0.25$
 $w_{12} = 0.20$ $w_{22} = 0.30$

Biases (hidden) : $b_{h1} = 0.35$ $b_{h2} = 0.35$.

weights (hidden - output) : $w_{11} = 0.40$ $w_{21} = 0.50$
 ~~$w_{12} = 0.45$~~ $w_{22} = 0.55$

Biases (output) : $b_{o1} = 0.60$ $b_{o2} = 0.60$.

Targets : $o_1 = 0.01$ $o_2 = 0.99$.

Forward Pass :-

(For input - hidden layer) $[z_j - \mu_{ij} = v_{oj} + \sum_{i=1}^n x_i \cdot w_{ij}]$

$$h_{\text{in } 1} = b_{h1} + (i_1 \cdot w_{11}) + (i_2 \cdot w_{21})$$

$$\therefore h_{\text{in } 1} = 0.35 + (0.05 \times 0.15) + (0.10 \times 0.25)$$

$$\therefore h_{\text{in } 1} = 0.35 + (0.0075) + (0.025)$$

$$\therefore h_{\text{in } 1} = 0.3825.$$

→ Sigmoid Activation Function for ' h_1 ' : $[z_j = f(z - \mu_j)]$

$$h_1 = f(h_{\text{in } 1})$$

$$\therefore h_1 = \frac{1}{1 + e^{-(0.3825)}} \approx 0.5944759307$$

> $h_{\text{in} 2} = b_{h2} + (i_1 \cdot w_{12}) + (i_2 \cdot w_{22})$
 $\therefore h_{\text{in} 2} = 0.35 + (0.05 \times 0.20) + (0.10 \times 0.80)$
 $\therefore h_{\text{in} 2} = 0.35 + (0.01) + (0.08)$
 $\therefore h_{\text{in} 2} = 0.39$

→ Sigmoid Activation Function for ' h_2 ' :

$$h_2 = f(h_{\text{in} 2})$$

$$\therefore h_2 = \frac{1}{1 + e^{-(0.39)}} \approx 0.5962826992$$

(For hidden-output layer)

> $O_{\text{in} 1} = b_{o1} + (h_1 \cdot w_{11}) + (h_2 \cdot w_{21})$
 $\therefore O_{\text{in} 1} = 0.60 + (0.5944759307 \times 0.40) + (0.5962826992 \times 0.50)$
 $\therefore O_{\text{in} 1} = 1.18893172188$.

> $O_{\text{in} 2} = b_{o2} + (h_1 \cdot w_{12}) + (h_2 \cdot w_{22})$
 $\therefore O_{\text{in} 2} = 0.60 + (0.5944759307 \times 0.45) + (0.5962826992 \times 0.55)$
 $\therefore O_{\text{in} 2} = 1.195469653375$

→ Sigmoid Activation Function for ' O_1 ' :

$$O_1 = f(O_{\text{in} 1})$$

$$\therefore O_1 = \frac{1}{1 + e^{-(1.18893172188)}} \approx 0.7569319154$$

→ Sigmoid Activation Function for ' O_2 ' :

$$O_2 = f(O_{\text{in} 2})$$

$$\therefore O_2 = \frac{1}{1 + e^{-(1.195469653375)}} \approx 0.7677178797$$

16/10/24

Convolutional Neural Network :~ (Feed Forward Architecture)

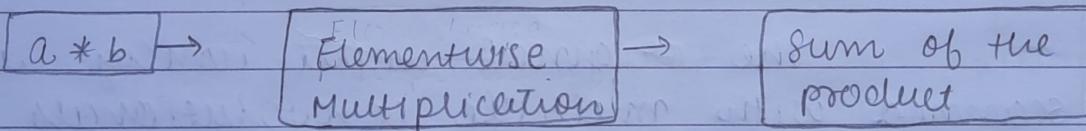
CNN is a type of feed forward NN architecture. Usually applicable for image classification.

It is also known as ConvNet.

The convolutional operation forms the basis of CNN.

In convolutional operation, the arrays are multiplied by element wise and the product is sum together to create a new feature array.

$$\begin{matrix} \textcircled{1} & \textcircled{2} \\ a = [5, 3, 2, 5, 9, 7] & b = [1, 2, 3] \end{matrix}$$



$$[17, 22, 39, 44]$$

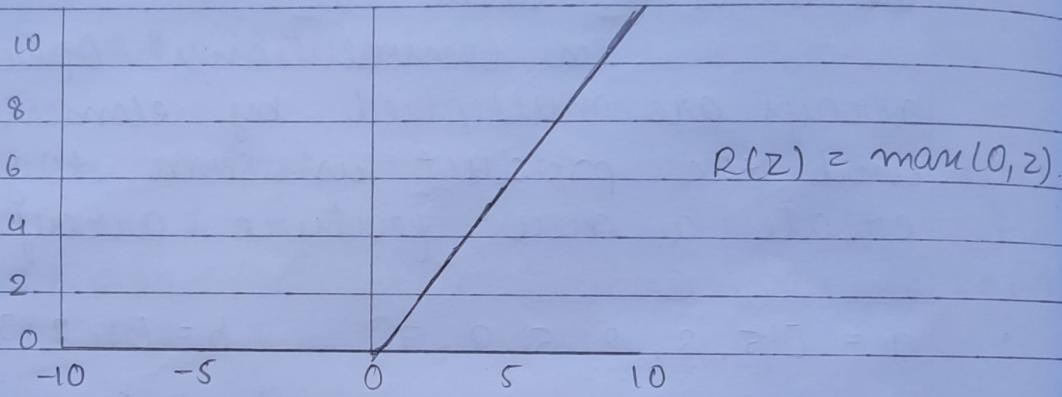
There are 7 possible layers in CNN generally, which are:-

- ① Convolutional layer
- ② ReLU
- ③ Pooling Layer
- ④ Fully connected layer
- ⑤ Activation layer
- ⑥ Flattening layer
- ⑦ Output layer

The ReLU stands for rectif

Once we have ^{the} feature maps extracted from the input then this vector is forwarded to ReLU layer.

The ReLU is responsible to perform elementwise operation and sets all the negative pixel to 0. Hence, the ReLU introduces non-linearity to the



The original input image is scanned with multiple convolutionals and ReLU layers for locating larger set of feature.

→
the
pooling
layer

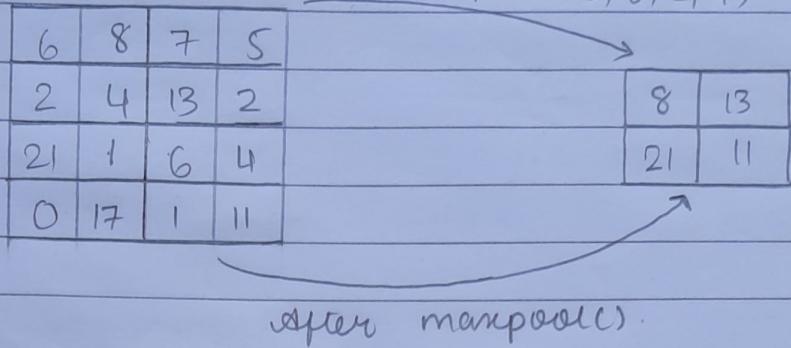
The Pooling is down sampling, that reduces the dimensionality of the feature map. The Pooling layer uses variety of layers to filter the variety of features.

After the pooling operation, we apply the flattening to convert the resultant 2-D array into single linear vector.

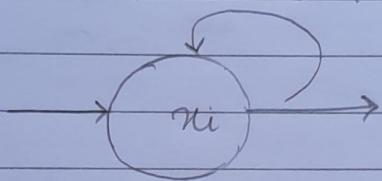
The Flattened vector is further fed as an input to fully connected

network to classify the image.

maxpool(6, 8, 2, 4)



Recurrent Neural Network :



LSTM
LSTM cell ↓
RNN Recurrent cell ↓

LSTM cell

↓
RNN
↓

Seq-to-seq learning
↓

Seq-to-seq with Attention
↓
Global
Local

Transformer