# INTRODUCTION TO MACHINE LEARNING.

## MACHINE LEARNING :-

Machine learning is said as a subset of AI that is mainly concerned with the development of algol- which allow a computer to learn from the past data experience on their own. The term Machine learning was first introduced by Arthur samuel in 1959.

## DEFINITION :-

ML enables a machine to automatically learn from data, improve performance from experiences and predict things without being explicitly Programmed.
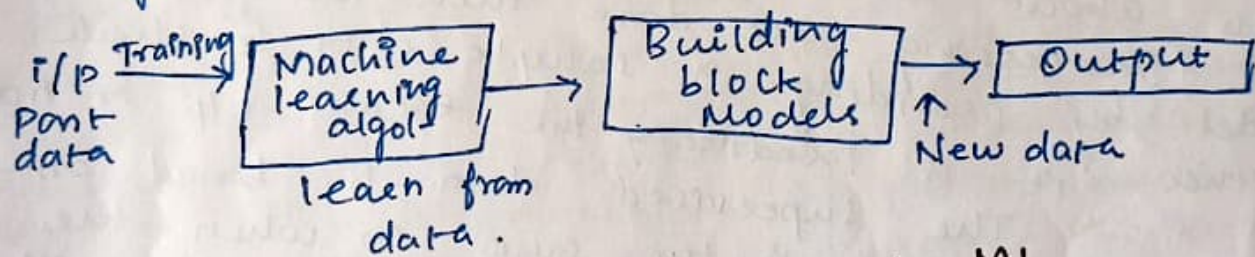


Fig: Working of ML.

A Machine s/m learns from historical data, builds the prediction models and whenever it receives new data and predicts the o/p for it.

## Features of Machine learning :-

Machine learning uses data to detect various patterns in a given set.

It can learn from past data and improve automatically.

It is a data driven technology.

Machine learning is much similar to data mining as it also deals with huge amount of the data.

Classification of Machine learning :-
It is classified into three types. They are

1) Supervised learning.
2) Unsupervised learning.
3) Reinforcement learning.

## SUPERVISED LEARNING :-
It is a type of ML Method in which we provide sample labelled data to the ML s/m in order to train it and on that basis it predict the o/p. The s/m creates a model using labelled data to understand the datasets and learn about each data once the training and processing are done then we test the model by providing a sample data to check whether it is predicting the exact o/p or not. The supervised data is based on supervision and it is the same as when there student learns things in the supervision of the teacher. The example of supervised learning is spam filtering. Supervised learning can be grouped further into two categories of algo-

1) Classification
2) Regression.

## UNSUPERVISED LEARNING :-
It is a learning Method in which a Machine learns without any supervision.

The training is provided to the machine with the set of data that has not been labeled classified or categorized and the algo- needs to act on the data without any supervision. The goal of unsupervised learning is here to restructure the i/p data into few features or a group of objects with similar patterns.

In supervised learning we don't have a predetermined result. The machine tries to find useful insights from the huge amount of data. It can be further classified into two categories of algo.

1. Clustering.

2. Association.

Reinforcement learning :-

It is a feedback based learning method, in which learning agents gets a reward for each right action and gets a penalty for each wrong action. The agent learns automatically with the feedbacks and improve its performance. In reinforcement learning, the agent interacts with the environment and explores it. The goal of an agent is to get the most reward points and hence it improve its performance.

LINEAR ALGEBRA :-
Linear algebra is an essential field of mathematics, which defines the study of vectors, matrices, planes, mapping and lines required for linear transformation. It plays a vital role and key foundation in machine learning and it enables machine learning algo- to run a huge number of datasets. The concepts of

linear algebra are widely used in developing algol- in machine learning. It can also Perform the following tank:

1) Optimization of data.
2) Applicable in loss functions, regularisation, Covariance Matrices, singular value decomposition (SVD) matrix operations and support vector Machine classifications.

Implementation of linear Regression in ML. The linear algebra is also used in neural nlws and data science field.

The general linear equation is represented as

$$a_1 x_1 + a_2 x_2 \ldots \ldots + a_n x_n = b$$

where,

a = Represents the co-efficients
x = Represents the unknowns
b = Represents the constant.

Transposes and Inner Products :-

A collection of variables may be treated as a single entity by writing them as a vector.

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$x^T = [x_1 \quad x_2 \quad x_3].$$

The transpose operator also turn row vectors into column vectors. By defining the inner product of two vectors

$$x^T y = [x_1 \quad x_2 \quad x_3] \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$= x_1 y_1 + x_2 y_2 + x_3 y_3$$

$$= \sum_{i=1}^{3} x_i y_i .$$

The outer Product of two vectors Produces a Matrix

$$x y^T = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} [y_1, y_2, y_3]$$

$$= \begin{bmatrix} x_1 y_1 & x_1 y_2 & x_1 y_3 \\ x_2 y_1 & x_2 y_2 & x_2 y_3 \\ x_3 y_1 & x_3 y_2 & x_3 y_3 \end{bmatrix}$$

When applying the transpose matrix the $i^{th}$ row becomes the $i^{th}$ column. That is, if

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}.$$

Then

$$A^T = \begin{bmatrix} a_{11} & a_{21} & a_{31} \\ a_{12} & a_{22} & a_{32} \\ a_{13} & a_{23} & a_{33} \end{bmatrix}$$

Outer Product :-

The outer product of two co-ordinate vectors is a matrix. If the two vectors have dimensions $n$ and $M$, then their outer product is an $n \times m$ matrix.

$$c = a_{*1} B_{1*} + a_{*2} B_{2*} + a_{*3} B_{3*} + \cdots + a_{*n} B_{n*}.$$

Inverse :-

A matrix $X$ its inverse $x^{-1}$ is defined by the properties

$$x^{-1} x = I$$
$$x x^{-1} = I.$$

Eigen values and Eigen vectors :-

The eigen vectors $x$ and eigen values $x$ of a matrix $A$ satisfy $\rightarrow Ax = \lambda x.$

# Singular Values and Singular Vectors :-

A Singular value and pair of singular vectors of a square or rectangular matrix A are a nonnegative scalar $\sigma$ and two non-zero vectors u and v so that

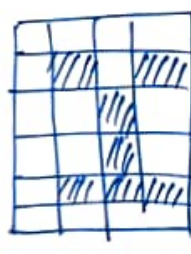$$Av = \sigma u,$$

$$A^H u = \sigma v.$$

# Examples of linear algebra in Machine learning :-

## Data sets and Data files :-

Each ML Project works on the dataset and we fit the Machine learning model using this dataset.

$$dog = \begin{bmatrix} 0 \\ 7 \\ 0 \\ 0 \\ 5 \\ \vdots \\ 0 \end{bmatrix}$$

## Images :-



$$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 0 \\ \vdots \\ 1 \end{bmatrix}.$$

## YES/NO or Ratings :-

Given users and items (eg. Movies), vectors can indicate if a user has interacted with the item (1 = yes, 0 = no) or the users ratings. Say a number b/w 0 and 5.

$$User\ 1 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \begin{matrix} NO \\ YES \\ NO \\ NO \\ \\ YES \\ NO \end{matrix} \qquad or \qquad USER\ 4 = \begin{bmatrix} 0 \\ 5 \\ 0 \\ 3 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} ? \\ Love \\ ? \\ Like \\ \\ ? \\ Dislike. \end{matrix}$$

# One Hot Encoding :-

How to represent non-numerical data such as language. One way is to apply one hot encoding.

Assign to each word a vector with one 1 and 0s elsewhere

$$apple = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad Cat = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad home = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

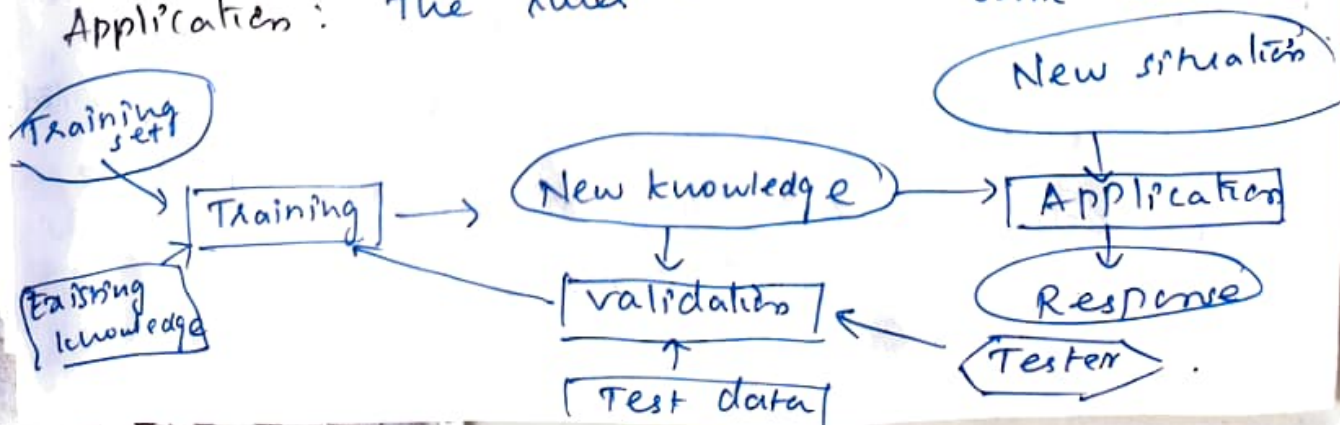$$tiger = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

# How Machines learn ?

Machine learning typically follows three phases :

1. Training
2. Validation
3. Application.

**Training:** A training set of examples of correct behavior is analyzed and some representation of the newly learnt knowledge is stored.

**Validation:** The rules are checked and if necessary, additional training is given.

**Application:** The rules are used in responding to some new situation.

# Why Machine learning is important?

ML algol- can figure out how to perform important tasks by generalizing from examples.

Machine learning provides business insight and intelligence. Decision makers are provided with greater insights into their organizations. This adaptive technology is being used by global enterprises to gain a competitive edge.

## Ingredients of Machine Learning :-

The ingredients of ML are as follows:

1. Tasks: The problems that can be solved with Machine learning. A task is an abstract representation of a problem. The standard methodology in Machine learning is to learn One task at a time. Large problems are broken into small, reasonably independent sub-problems that are learned separately and then recombined.

2. Models: The o/p of Machine learning. Different Models are geometric Models, probabilistic Models, logical Models, grouping and grading.

3. Features: The workhorses of Machine learning. A good feature representation is central to achieving high performance in any Machine learning task.

## Examples of Machine learning Applications :-

Optical character recognition: Categorize images of handwritten characters by the letters represented.

Face detection: Find faces in images (or indicate if a face is present).

**Spam filtering:**
Identify email messages as spam or non-spam topic spotting: Categorize news articles (say) as to whether they are about politics, sports, entertainment etc---,

Spoken language understanding: Within the content of a limited domain, determine the Meaning of something uttered by a speaker to the extent that it can be classified into One of a fixed set of categories.

## Face Recognition and Medical Diagnosis:-

### Face Recognition :-
It is effortlessly and every day we recognize our friends, relative and family Members. We also Recognition by looking at the Photographs. In Photographs, they are in different Pose, hair styles, background light, Makeup and without Makeup.

### Medical diagnosis :-
The i/p are the Relavant information about the patient and the Classes are the illness. The i/p contains the age of patient's, gender, Past Medical history and current symptoms.

## Google Home and Amazon Alexa !-

### Amazon Alexa/ siri :-
Everytime Alexa or Siri Make a mistake when responding to our request, it uses the data it receives based on how it responded to the original query to improve the next time.

# Google Home :-

Google services such as its image search and translation tools are sophisticated machine learning which allow computers to see, listen and speak in much the same way as human do.

# Unmanned Vehicles :-

An unmanned Aerial Vehicle (UAV), sometimes known as a drone, is an aircraft or airborne s/m that is controlled remotely by an onboard computer or a human operator. The ground control station, aircraft components, and various types of sensors make up the UAV s/m.
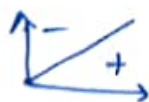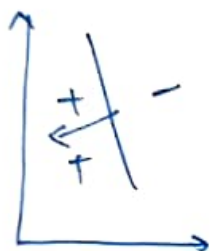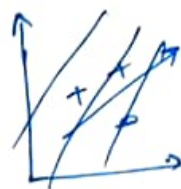
# VAPNIK - CHERVONENKIS (VC) DIMENSION :-

The Vapnik - Chervonenkis dimension, more commonly known as the VC dimension, is a model capacity measurement used in statistics and machine learning.

# Shattering :-

It is the ability of a model to classify a set of points perfectly. The model can create a function that can divide the points into two distinct classes without overlapping. It is different from simple classification because it considers all possible combinations of labels upon those points. ($2^N$ Possible ways).
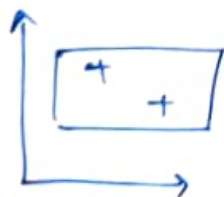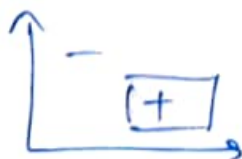
# Find VC Dimension :-

$2^4 \rightarrow$ is not Possible

In VC dimension $2^3$ is Possible

Maximum no. of points in $R^2$ that be shattered by straight line is 3.

VCD (Axis aligned Rectangle)



$\rightarrow$ outside the rectangle negative
Inside the rectangle Positive

until 4 point. Possible

5 $\rightarrow$ not Possible

$$VCD = 4 \quad (\text{axis aligned})$$

# PROBABLY APPROXIMATELY CORRECT LEARNING :-

PAC learning is a framework for the Mathematical analysis of Machine learning.

Goal of PAC :- With high probability ("Probably"), the selected hypothesis will have lower error ("Approximately Correct")

In the PAC model, we specify two small Parameters, $\varepsilon$ and $\delta$, and require that with Probability atleast $(1-\delta)$ a system learn a Concept with error at Most $\varepsilon$.

## $\varepsilon$ and $\delta$ Parameters :-

$\varepsilon$ gives an upper bound on the error in accuracy with which h approximated (accuracy : $1-\varepsilon$).

$\delta$ gives the probability of failure in achieving this accuracy (confidence : $1-\delta$).

X instance
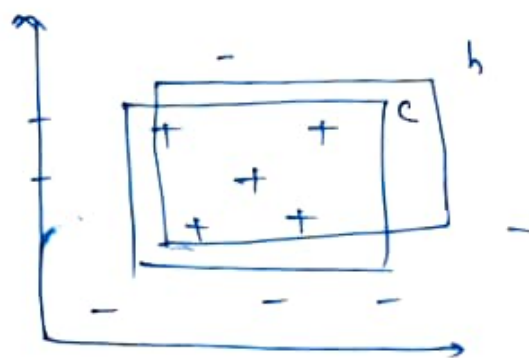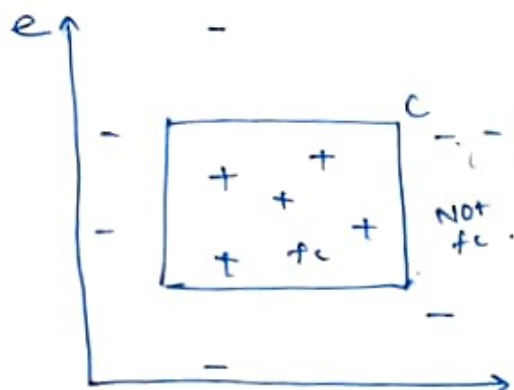C Target
H Hypothesis
D Training data

## Example :-

N number of Car having Price and Engine Power, as training set $(P, e)$, find the Car is family car or not.

An algol- gives answer whether the Car is family care or not.

C - Target function.

Instances within rectangle represents family Cars and outside are not family Cars.

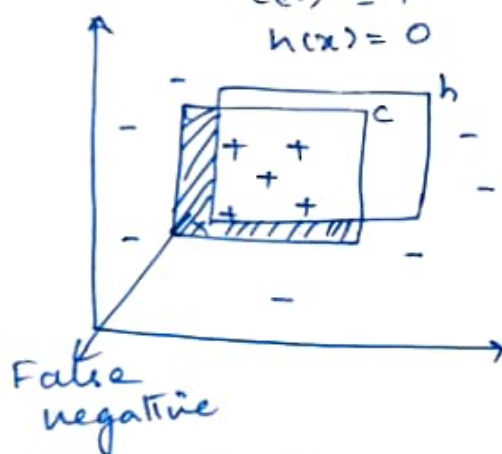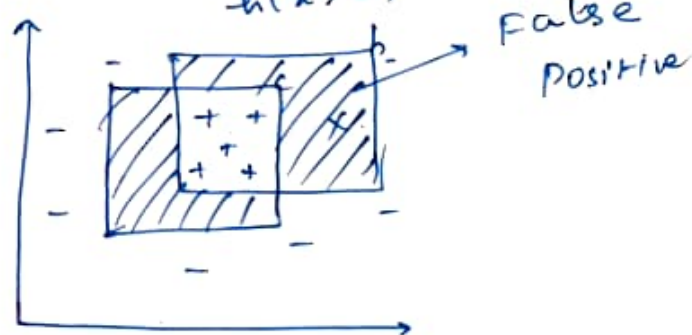Hypothesis h - closely approximate C, and there May be error region.

## False Negative & False +ve :-

Instances lies on shaded region are +ve/-ve according to our actual function 'c'. but those are -ve/+ve based on the hypothesis h. Hence it is called as false -ve or false +ve.

$c(x) = 1$
$h(x) = 0$

$c(x) = 0$
$h(x) = 1$

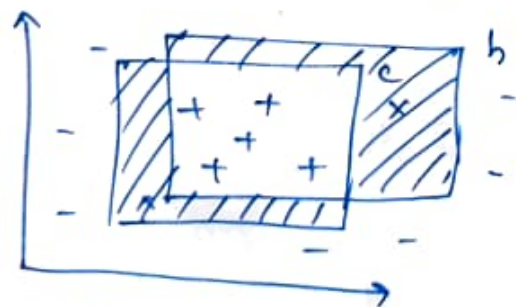False Positive



False negative

## Error Region :-

The Probability of error region to be Small. The error region : $P(c \text{ XOR } h) <= \varepsilon$.

$P(c \oplus h)$.

Error region : $c \text{ XOR } h$



## Approximately Correct :-

the hypothesis h, that approximately correct, and error is less than or equal to $\varepsilon$

where $0 <= \varepsilon <= 1/2$ . (i.e) $P(c \text{ XOR } h) <= \varepsilon$

## Probably Approximately Correct :-
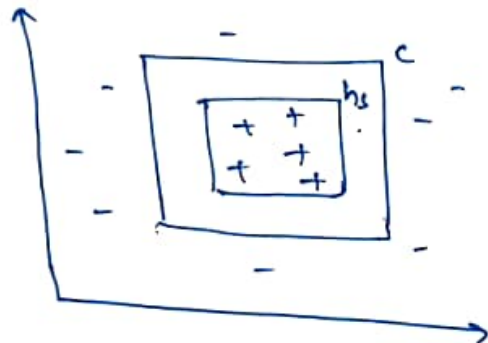
low generalization error with high Probability

$$[P( error (h) <= \varepsilon)] < = 1- \delta$$

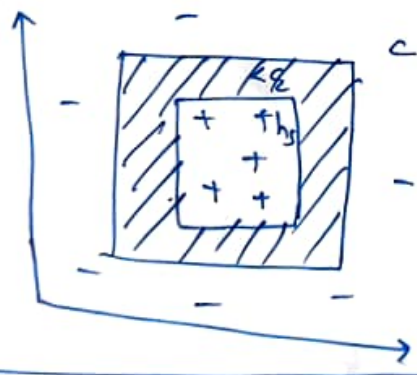$$P[P(c \text{ XOR } h) <= \varepsilon] <= 1-\delta .$$

PAC learnability for axis - aligned rectangle :-

## Specialization :-

$h_s$ is the tighest possible rectangle around a set of +ve training.
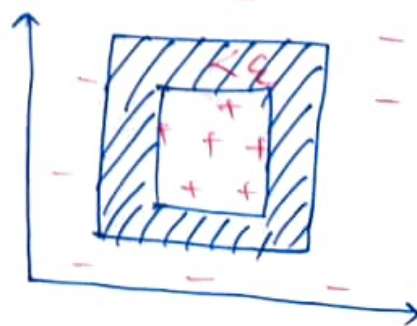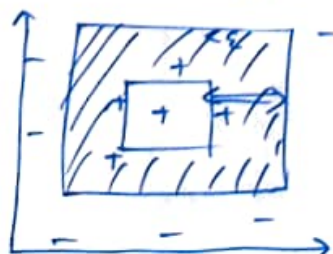
$h_s$ is subset of c, Hence error region $= c - h$.



If an hypothesis lies blw h and c (shaded region) then it is approximately correct.



If the generated hypothe does not touch any of these region.

Error region is greater than $\varepsilon$ and not approximately correct, because the error region got increased

Atleast one +ve example at each side of the rectangle.

# Error Region:-

Error Region = Sum of four rectangular strips $< \varepsilon$.
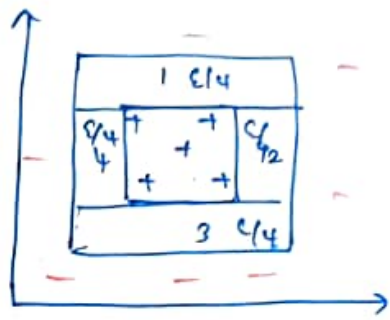
Each strip is at most $\varepsilon/4$.

Probability of +ve example falling in any one of the strip (error region = $\varepsilon/4$)

Probability that a randomly drawn +ve example misses a strip = $1 - \varepsilon/4$ . ⊗ .



$P(m$ instance miss a strip) $\neq \varepsilon/4$
$$= (1 - \varepsilon/4)^m .$$

$P(m$ instance miss any strip) $< 4(1 - \varepsilon/4)^m .$

Finally we get $m > 4 / \varepsilon \log 4 / \delta$.

## Example 1:-

Hypothesis $h1$ generated the errors with respect to price and engine power of given to samples.

Given $\varepsilon = 0.05$     $\delta = 0.20$

$P(h_i) >= 1 - \delta$

$P(h_1) = 8/10 = 0.80$

($3^{rd}$ and $8^{th}$ values are greater than $\varepsilon$)

$\therefore 0.80 >= (1 - 0.20)$ i.e. $0.80 = 0.80$ ⊗

$H1$ is probably approximately correct.

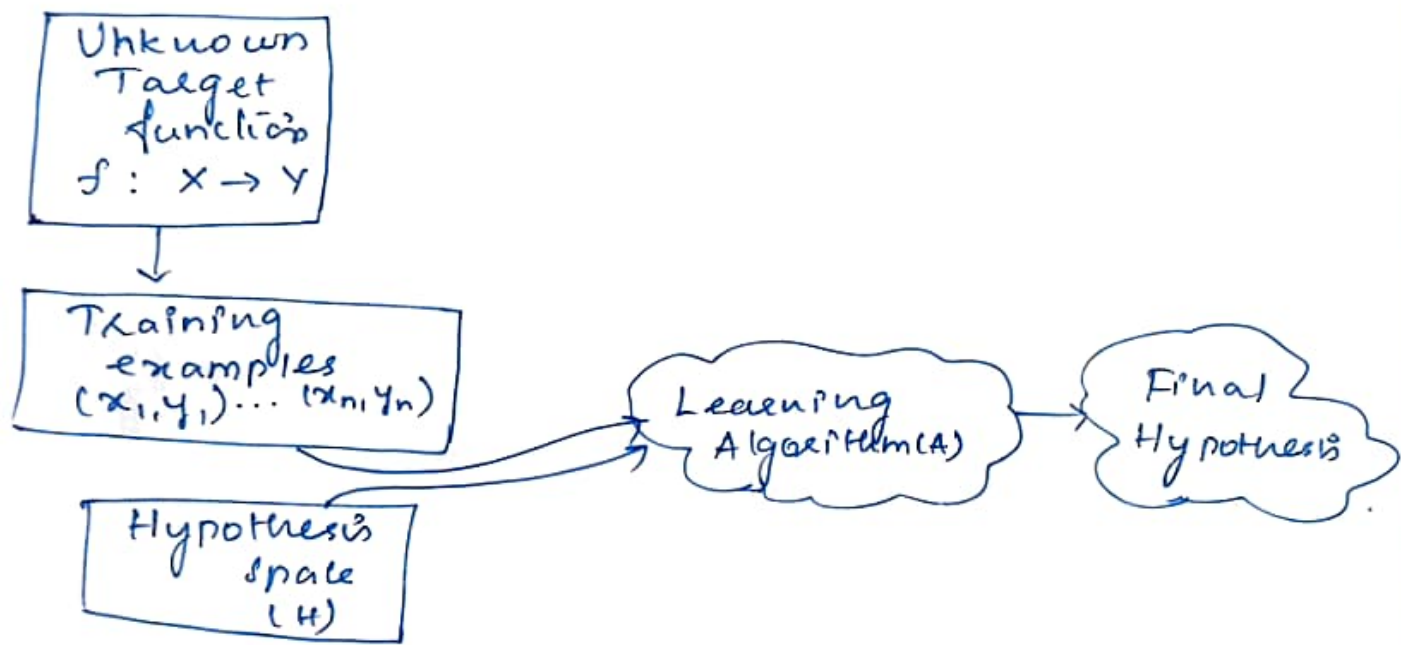| Sl. No | Error(hi) |
|---|---|
| 1 | 0.001 |
| 2 | 0.025 |
| 3 | 0.07 |
| 4 | 0.003 |
| 5 | 0.035 |
| 6 | 0.045 |
| 7 | 0.027 |
| 8 | 0.065 |
| 9 | 0.012 |
| 10 | 0.036 |

## HYPOTHESIS SPACES:-

The hypothesis space is defined on the supposition or proposed explanation based on the insufficient evidence or assumptions. It is just as guess based on some known facts but has not yet been proved. A good hypothesis

- is testable, which content is not contained in false. which results is either true or false.

1 | **Hypothesis in Machine learning (ML) :-**

The hypothesis is one of the commonly used concepts of statistics in ML. It is specifically used in supervised learning where an ML Model learns a function that best Maps the I/p to corresponding o/ps with the help of an available dataset.



**CANDIDATE ELIMINATION ALGORITHM :-**

It is a method of learning concepts from data that is supervised. Given a hypothesis space H and a collection E of instances, the candidate Elimination procedures develop the version space continuously.

| TIME | WEATHER | TEMP | COMPANY. | HUMIDITY | WIND | GOES |
|------|---------|------|----------|----------|------|------|
| MORNING | Sunny | Warm | Yes | Mild | Strong | yes |
| EVENING | Rainy | Cold | No | Mild | Normal | No |
| Moring | Sunny | Moderate | Yes | Normal | Normal | yes |
| Evening | Sunny | Cold | yes | High | Strong | yes |

$G_3 = \langle \text{Morning} ? ? ? ? ? \rangle \langle ? \text{ sunny} ? ? ? ? \rangle \langle ? ? ? ? \text{yes} ? ? \rangle$

$x_3 = \langle \text{evening}, \text{sunny}, \text{cold}, \text{yes}, \text{High}, \text{strong} \rangle$.

+ve → specific

$S_4 = \langle ?, \text{sunny}, ?, \text{yes}, ?, ? \rangle$

$G_4 = \langle ? \text{ sunny} ? ? ? ? \rangle \langle ? ? ? ? \text{yes} ? ? \rangle$.

## Specific hypothesis

$\langle ?, \text{sunny}, ?, \text{yes}, ?, ? \rangle$

## General hypothesis

$\langle ? \text{ sunny} ? ? ? ? \rangle \langle ? ? ? ? \text{yes} ? ? \rangle$

## Sample Error and True Error :-

The Sample error $(err_s(h))$ of $h$ with Respect to target function $(f)$ and data sample $(s)$ is the proportion of examples $h$ Misclassifing. The sample test error is the Mean error over the test Sample.

$$err_s(h) = \frac{1}{n} \sum_{i=1}^{n} L(f(x_r), h(x_r)).$$

The true error of hypothesis $h$ with respect to target function $f$ and Distribution $D$ is the probability that $h$ will Misclassify an instance drawn at random according to $D$.

$$Err(h) = E\left[L\{f(x)\}, h(x)\}\right].$$

## INDUCTIVE BIAS :-

The Candidate Elimination Algol- will Converge toward the true target Concept Provided

— it is given accurate training examples →

— its initial hypothesis space contains the target Concept.

Bias ↓ average Squared diff. b/w predictions and true values

Generalize

Solve on what days the person likes to go on walk using Candidate elimination algo-.

**Sol**

## CEA :-

- extended form of Find S algo-
- considers both +ve and -ve instances
- finds all the hypotheses that Match all the given training examples.

Attributes :- Time, weather, Temp, company, Humidity, wind.

Target : Goes for a walk or not.
$$(+ve = yes ; -ve = No).$$

S ( specific hypothesis) $= < \phi, \phi, \phi, \phi, \phi, \phi >$

G ( General hypothesis) $= < ?, ?, ?, ?, ?, ? >$

1) $x_0 = <$ Morning, Sunny, warm, yes, mild, strong $>$
   +ve → Specific (s)
   $(x_0, S_0)$ $S_1 = <$ Morning, Sunny, warm, yes, Mild, strong $>$
   $G_1 = < ?, ?, ?, ?, ?, ? >$

② $x_1 = <$ Evening. ~~Rainy~~, cold, No, Mild, Normal $>$
   -ve → General (G)
   $S_2 = <$ Morning, Sunny, warm, yes, Mild, strong $>$
   $(x_1, S_1)$
   $G_2 = <$ Morning, ? ? ? ? ? $>$ $<$ Sunny. ? ? ? ? $>$
   $< ? ?$ warm ? ? ? $>$ $< ? ? ?$ yes ? ? $>$
   $< ? ? ? ? ?$ strong $>$.

③ $x_2 = <$ Morning, Sunny, Moderate, yes, Normal, ~~Strong~~ $>$
   +ve → Specific (s)
   $(x_2, S_2)$
   $S_3 = <$ Morning, Sunny, ?, yes, ?, ? $>$ .

# BIAS XXXX VARIANCE TRADE OFF :-

- What if the target concept" is not contained in the hypothesis space? experimental Practice we observe the bias

- Can we avoid this difficulty by using a hypothesis space that includes every possible hypothesis?

## A Biased Hypothesis space :-

Support we wish to ansure that the hypothesis space contains the unknown target concept. The obvious solution is to enrich the hypothesis space to include every possible hypothesis. To illustrate, consider the Enjoyspert example. in which we restricted the hypothesis space to include only conjunctions of attribute values.

< Sunny, High. Normal, Steang, Cool, same >.

| Example | Sky | Aretemp | Humidily | Wind | water | Forecant | Enjoysport |
|---------|-------|---------|----------|-------|-------|----------|------------|
| 1 | Sunny | Waem | Normal | Strng | Cool | change | Yes |
| 2 | Cloudy | Waeue | Normal | Strong | Cool | change | Yes |
| 3 | Rainy | Waem | Normal | Strong | Coul | change | No. |

So : < ⌀, ⌀, ⌀, ⌀, ⌀, ⌀ >

$S1$ : < Sunny, warm, Normal, Strong, Cool, change >
$S2$ : < ?, warm, Normal, String, Cool, change >
$S3$ : < ?, warm, Normal, Strong, Cool, Change >
$G_2$ : < ?, ?, ?, ?, ?, ? >
$G1$ : < ?, ?, ?, ? ? ? ? >
$G_0$ : < ? ? ? ? ? ? ? >

Because of the restriction, the hypothesis space is unable to represent even simple disjunctive target concepts such on "sky = Sunny or sky = Cloudy".
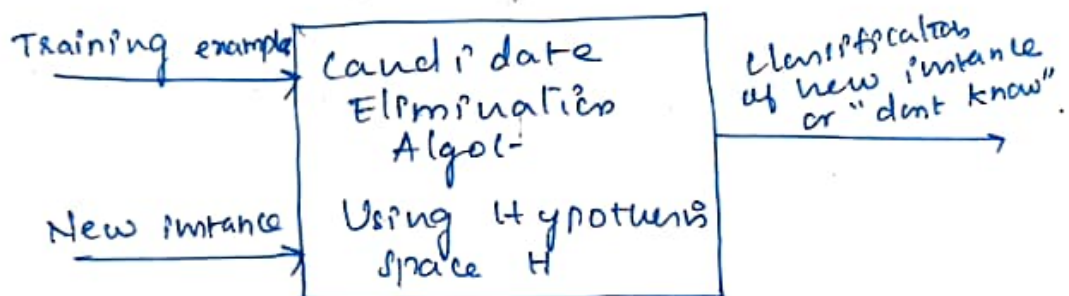
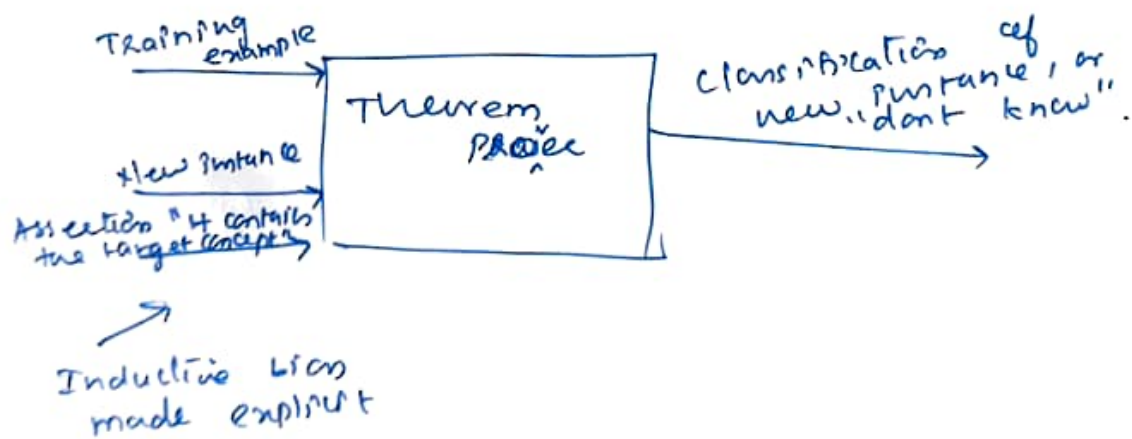< "sky = Sunny or sky = Cloudy"? warm, Normal, Strong, Cool, change >

# Unbiased Learner :-

The obvious solution to the problem of ansuring that the target concept is in the hypothesis space H is to provide a hypothesis space capable of representing every teachable concept; that is, it is capable of representing every possible subset of the instances.

In general, the set of all subsets of a set X is called the powerset of X.

# Inductive sim :-



Training example → [ Candidate Elimination Algol- Using Hypothesis Space H ] → Classification of new instance or "dont know".

New instance →

# Equivalent Deductive sim :-



Training example → [ Theorem Prover ] → Classification of new instance, or "dont know".

New instance →

Assertion "H contains the target concept"

→ Inductive bias made explicit

# Generalization :-

Generalization is a definition to demonstrate how well is a trained model to classify or forecast unseen data. Training a generalized machine learning model means, In general, it works for all subset of unseen data.

# BIAS ~~AND~~ VARIANCE TRADE OFF :-

In the experimental Practice we observe an important Phenomenon called the bias variance dilemma.
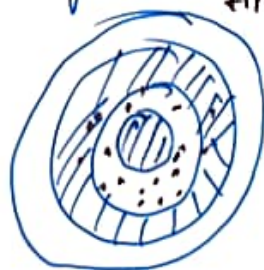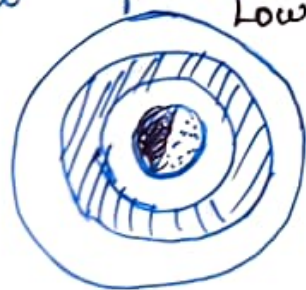
In Supervised learning, the clans value assigned by the learning model build based on the training data may differ from the actual clans value. This error in learning can be of two types, errors due to 'bias' and error due to 'variance'.

The bias - variance dilemma is the Problem of simultaneously Minimizing two source of error that Prevent Supervised learning algol- from generalizing beyond their training set.

1. The bias is error from erroneous assumptions in the learning algol-. High bias la came an algol- to miss the relevant relations b/w features and target Olp.

2. The variance is error from sensitivity small fluctuations in the training set. High varie can cause an algol- to miss the relevant Relation b/w features and target O/Ps.

Low variance        High Variance

Bias - assumption
variation - deviate (how much)

Low bias

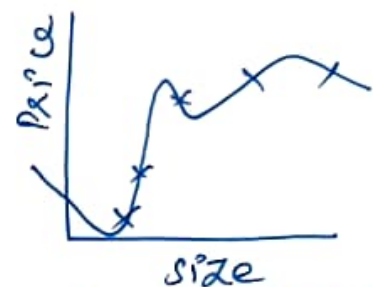High bias

# Underfitting ( High bias and low variance) :-

A statistical model or a Machine learning algol- is said to have underfitting when it cannot capture the underlying trend of the data. It can be avoided by using more data and also reducing the features by feature selection.



$Q_0 + Q_1 x$

High bias (underfit)

$Q_0 + Q_1 x + Q_2 x^2$

High bias (underfit)

$Q_0 + Q_1 x + Q_2 x^2 + Q_2 x^2_4$

$Q_2 x^2$

High variance (overfit).

# Overfitting (High variance and low bias)

A statistical model is said to be overfitted, when we train it with a lot of data.

When a Model gets trained with so Much of data, it starts learning from the noise and inaccurate data entries in our data set.

Then the Model does not categorize the data correctly, because of too Many details & noise.

A solution to avoid Overfitting is using a linear algol- if we have linear data or using the parameters like a Maximal depth if we are using decision trees.