

Machine learning

Govash Shenode

20051627



Machine Learning

Module - 1 : Introduction to Machine Learning

- What is Machine learning
 - What is the need, motivation for ML
 - Well defined ML Problem
 - What are the 3 main features
 - What are the Key concepts of Machine Learning
- > Machine learning → The capability of a machine to imitate intelligent human behaviour
- > Need and motivation →

The combination of computing and communication has given rise to a society that feeds on information.

→ Info is categorized on the basis of nature of experience



- Examples
- Sample
- measurements
- patterns

- Experience
- expertise
- mainly expressed in linguistic form
(Follows IF-THEN RULE)
- ↳ used in SHL.

•> Definition of Machine Learning & 3 main features

→ A computer program is said to learn from experience with respect to some class of task and performance measure, if the performance and task as measured by the performance measure, improved with the experience.

→ 3 feature

- 1> The learning task. → 'E' → Experience.
- 2> The measure of performance. → 'P' → Performance.
- 3> The task experience. → 'T' → Task.

12/1/2023

- > Feature / attributes → Inputs are the feature of ML.
- The key concept is "Learning from the experience"
- Important aspects of "Learning from the experience" are
 - 1) Remembering and Adapting
 - Recognising that last time in a similar situation a certain action was attempted which resulted in the output so it should be tried again or the same action failed in a similar situation so something different should be tried
 - 2) Generalizing
 - This aspect is recognizing similarity between different situations so a decision can be taken for a new situation.
- > Application of ML
 - 1) Social media (Sentiment analysis, filtering spam)
 - 2) Transport (Safety Monitoring, Air traffic control etc)
 - 3) Finance (Fraud detection, Portfolio management)
 - 4) Health care centre (Drug Discovery, Disease Diagnosis)
 - 5) eCommerce (Customer support, product recommendation etc)
 - 6) Virtual assistants (Intelligent agent, NLP)

Name → Vrushali Shende

Roll no → 20051627

Section → CS-11

Class Activity - 1

- Write in about 200-300 words about one application of machine learning in any field.

→ SPACE EXPLORATION:

→ One application of machine learning in the field of space exploration is the use of autonomous systems for spacecraft navigation and control. These systems use machine learning algorithms to analyze data from onboard sensors and make decisions about how to navigate and control the spacecraft. For example, NASA's Jet Propulsion Laboratory has developed an autonomous navigation system for the Mars rovers that uses machine learning to analyse data from cameras and other sensors to identify landmarks and obstacles in the rover's path. The system then uses this information to plan a safe and efficient path to the rover's next destination.

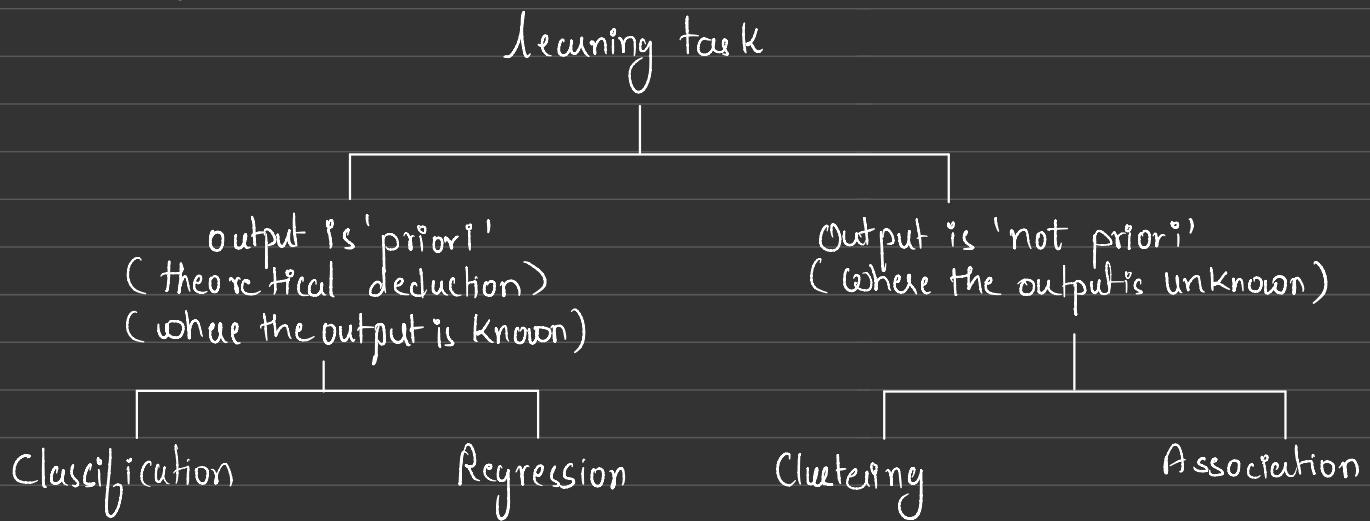
Machine learning can also be used to improve the efficiency of communication between spacecraft and ground control. For instance the European Space Agency's Hera mission, a proposed asteroid deflection test mission plans to use machine learning to optimize the data compression and transmission from the spacecraft to the ground, allowing for more data to be transmitted in a shorter amount of time.

Another example is the use of machine learning for the analysis of data from scientific instruments on space crafts. For instance NASA's Kepler project used machine learning to analyse data from the Kepler spectrometer on Cassini spacecraft, which orbited Saturn for over a decade. The machine learning algorithms helped to identify and classify different types of particles in the Saturnian system, enabling scientist to better understand the composition and structure of the planet and its moons.

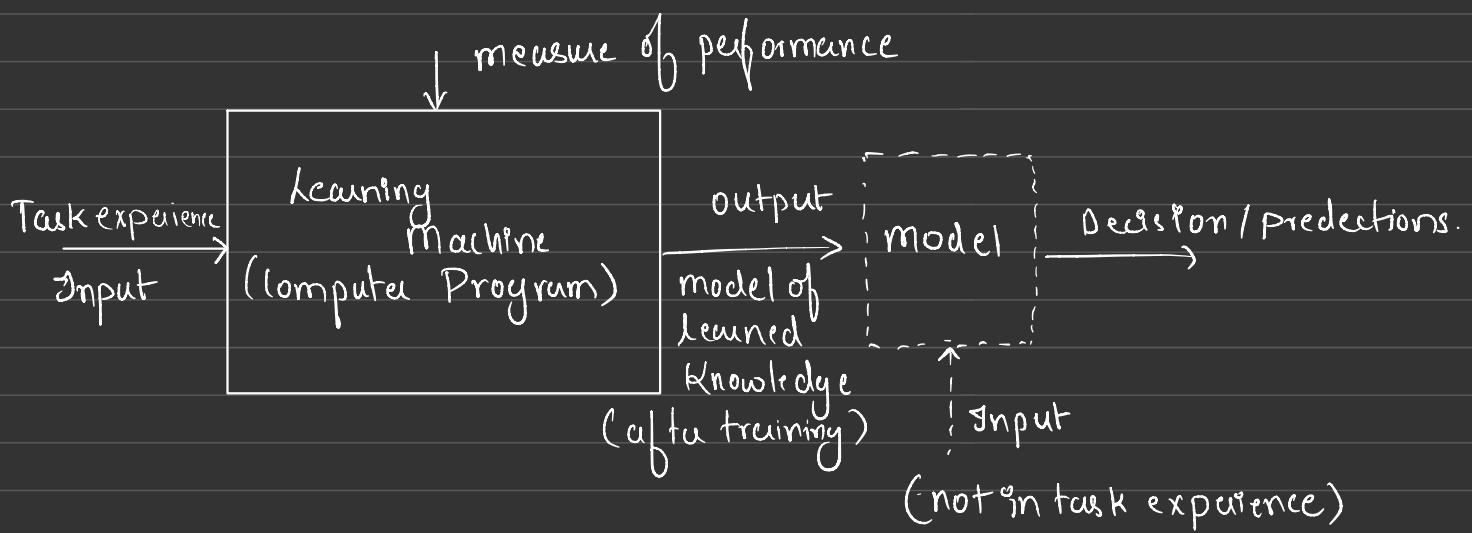
Overall, machine learning plays a crucial role in space exploration by enabling spacecraft to navigate and control themselves more efficiently, analyze large amounts of scientific data, and improve communication between space craft and ground control, with increasing amount of data being generated by spacecraft, machine learning will continue to be an important tool for making new discoveries in the study of our solar system and beyond.

→ learning task

⇒ learning task is divided into 2 types



•) Block diagram Representation of Machine Learning



Input → Training Data (experimental data / Human knowledge / data)
Testing Data

Step used in Machine Learning

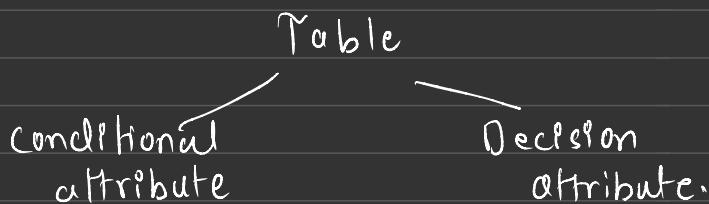
1) Collecting Data (raw data from excel, access, text files etc)

2) Preparing the data : Preprocessing of data

- 3) Training the model : There are two types of inputs 1st is training then it goes to learning machine then it goes to model where testing input is
- 4) Evaluating the model : To test the accuracy , the second part of the data (Test data) is used.
- 5) Improving the performance : This step might involve choosing a different model altogether or introducing more variables.

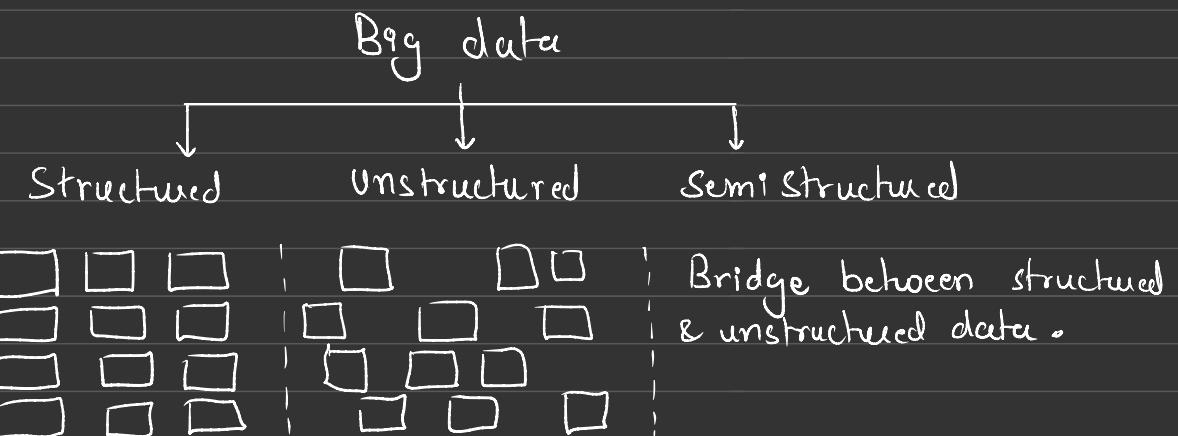
•) Data Representation

- Raw data requires some preprocessing w.r.t class of tasks
- This leads to an information system may be stored in the data warehouse used for decision making. (knowledge in the raw data.)
- Information is divided into table in two ways.



•) Diversity of Data.

- Data in petabytes (10^{15} byte) is called Big data.



- The data which is lacks any predefined to the point & highly organized is referred to as Structured data.

Activity (Tabulate the difference between structured and unstructured data using some parameter .)

Parameter	Structured Data	Unstructured Data
1) Scalability	It can be easily scaled up or down as the volume of data grows	It can be more difficult to scale as the volume & variety of data increases
2) Data Quality	High , as data is well organised & follows pre-defined rules	low, as data may be inconsistent or contain errors
3) Data Accessibility	Easily accessible using standard database management systems	Can be more difficult to access and may require specialized software .
4) Storage	Can be stored in variety of formats , such as spread sheets or databases which are optimized for structured data.	Can be more challenging to store as unstructured data can come in many different formats & may require specialized storage soln
5) Searchability	Can be easily searched and filtered using specific criteria	Can be more difficult to search & filter , often requiring NLP or other advanced Techniques
6) Format	It has predefined format for example consider a XL sheet which is predefined	It has variety of formats ie it comes in a variety of shapes & sizes.

14/11/23

Krishnayuk for ML topics

- > The process of teaching machines can be broken down into 3 parts.

Data as Input

Text files, spreadsheet
SQL Databases

Abstracting the data

Representation of the
data in a structured
format

Generalization

The practical application
happens here where the
learning from the
previous step is used
to develop an insight.

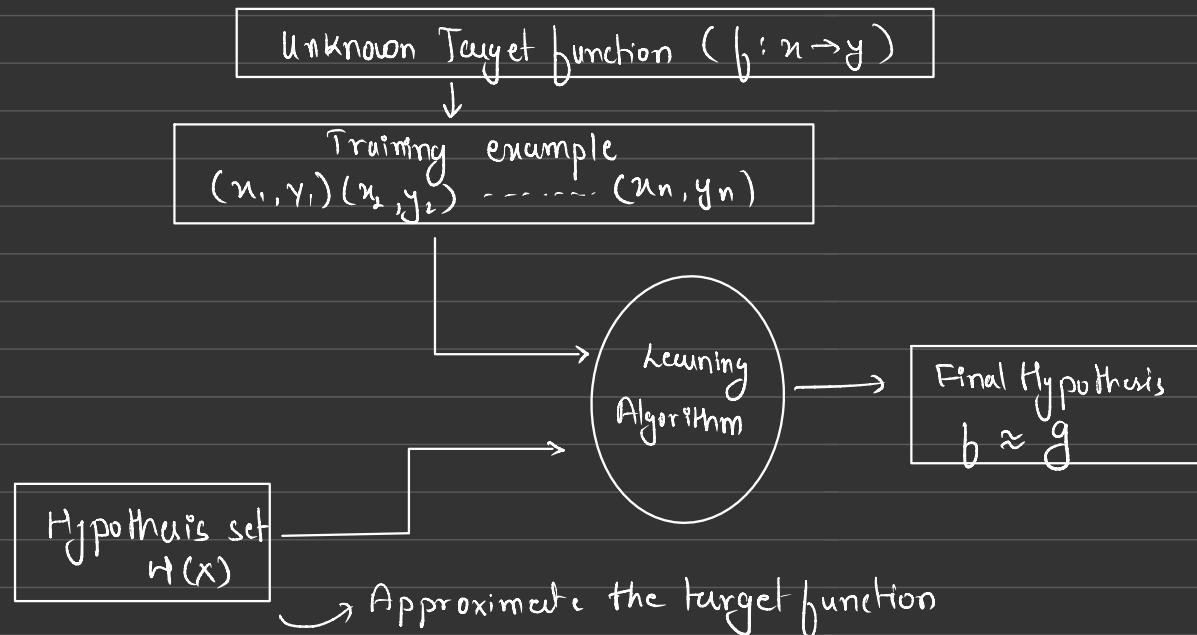
20/11/23 (Lect 4)

- > Simple learning process.
- > Design of a learning process.
→ checkers game

Well defined ML problem (T, E, P)

*> Simple learning process:

Target functions (Mapping function) \rightarrow maps relation between input & output.



•) Design of a learning process

1) Choosing a training experience

- Direct training experience
- Indirect training experience

2) Choosing the target function (mapping: $f(x) \rightarrow y$)

3) Representation of Target function:

$$\text{Output: } w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots$$

4) Error Calculation / Approximation Algorithm.

$$E = \text{Actual value} - \text{predicted value}$$

•) Types of Machine learning algorithms.

1) Supervised Learning

- Classification
- Regression

2) Unsupervised Learning

- Clustering → Targeted marketing, Recommender System, customer segmentation
- Dimensionality Reduction → Meaningful

3) Reinforcement Learning

- Game AI
- Skill application
- learning task
- Robot navigation

•> Supervised Learning.

→ The machine is designed by making use of the prior known information in the form of direct training examples consisting of observe values of system state. and the response of each state is output vector.

$$\begin{aligned} x &\rightarrow \text{(input vector)} \\ y &\rightarrow \text{(output vector)} \end{aligned}$$

•> Regression → Prediction Analysis. (Output is continuous)

→ Futuretic Data analysis (Forecasting)

→ Deals with huge amount of data

Application :

- : population growth analysis.
- : estimating life expectancy.
- : market forecasting
- : advertisement popularity prediction.

•> Classification: Application: Diagnostic.

(output is discrete) : Customer Retention

: Spam mail Detection.

: Image classification.

•> Unsupervised : It is not associated with any kind of guidance.

: No concrete data.

: Random data is allocated.

:

1) Clustering: Not at all structured data is available.
→ Creation of sets by segregating data by some certain attributes is clustering.

2) Dimensionality Reduction: Data reduction. (an application of data)

Application → Structure Discovery.

•) Reinforcement learning.

→ Instant decision making comes under reinforcement learning

•) UNIT 2 (Fundamentals of Machine learning)

→ Three Basic method of Machine learning.

- 1) Least Square method
- 2) KNN method ()
- 3) Distance based method

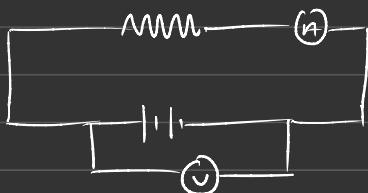


1) Least Squared Method

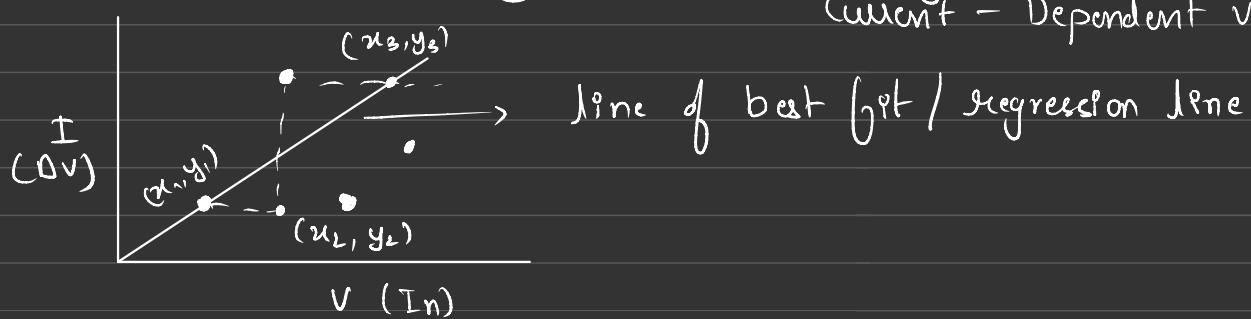
→ Used to find best fit line

Example:

Ohm's law

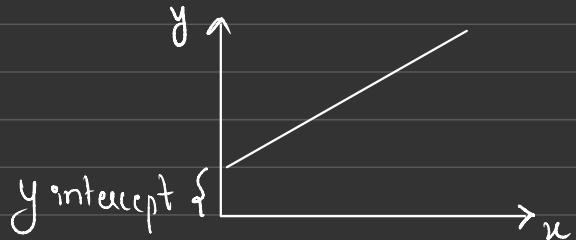


Voltage - Independent variable
Current - Dependent variable.



• Best fit line $\rightarrow y = mx + c$

dependent variable Slope Independent variable
 target variable ↓ Predictor variable



\Rightarrow Definition: It is based on the idea that the square of the errors obtained must be minimized to the most possible extent hence the name is given Least square method.

formula:

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n [y_i - f(x_i)]^2 \text{ should be least.}$$

Simple Linear Regression Model.

The least square line $y = a + bx$ for the given data

We can determine the values of a and b in the simple linear regression equation using the give data by minimising the error function. which is

$$E(a, b) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$E(a, b) = \sum_{i=1}^n (y_i - a + b x_i)^2$$

Differentiating the error function with respect to a and b and equating them to zero constitutes a set of two equations

$$\Rightarrow \frac{\partial E(a, b)}{\partial a} = -2 \sum_{i=1}^n (y_i - a - bx_i) = 0$$

$$= \sum_{i=1}^n y_i - \sum_{i=1}^n a - \sum_{i=1}^n bx_i = 0$$

$$= \sum y - na - b \sum x = 0$$

$$= na + b \sum x = \sum y \rightarrow \textcircled{1}$$

$$\Rightarrow \frac{\partial E(a, b)}{\partial b} = -2 \sum_{i=1}^n x_i (y_i - a - bx_i) = 0$$

$$= \sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i - \sum bx_i^2 = 0$$

$$= \sum xy - \sum ax - \sum bx^2 = 0$$

$$\Rightarrow \sum xy = \sum ax + \sum bx^2 \rightarrow \textcircled{2}$$

$\textcircled{1}$ & $\textcircled{2}$ are called normal equation.

Q1) Find the least square line $y = a + bx$ for the given data

x	-2	-1	0	1	2
y	1	2	3	3	4

x	y	x^2	xy
-2	1	4	-2
-1	2	1	-2
0	3	0	0
1	3	1	3
2	4	4	8

$$\sum x = 0 \quad \sum y = 13 \quad \sum x^2 = 10 \quad \sum xy = 7$$

$$\begin{aligned} 5a = 13 \\ 0 \cdot 0 + 10b = 7 \end{aligned} \quad \begin{aligned} a = 13/5 = 2.6 \\ b = 10/7 = 0.7 \end{aligned}$$

$$y = 2.6 + (0.7)x$$

Q)	Year	1958	1959	1960	1961	1962
	Sales	65	95	80	115	105

deviation ($n = \text{odd}$)

Year	Sales	x	x^2	xy	
1958	65	-2	4	-130	
1959	95	-1	1	-95	225
1960	80	0	0	0	
1961	115	1	1	115	
1962	105	2	4	210	325

$$\text{Total} \quad \sum y = 460 \quad \sum x = 0 \quad \sum x^2 = 10 \quad \sum xy = 100$$

$$Na + b\sum n = \bar{y}$$

$$5a = 460$$

$$a = 92$$

$$\sum xy = \sum ax + \sum bn^2$$

$$100 = b \cdot 10$$

$$a = 92, b = 10$$

$$y = 92 + 10n$$

• Equation for Best fit line ($y = a + bn$)

Normalize equation

$$na + b\sum n = \sum y$$

$$a\sum x + b\sum x^2 = \sum xy$$

Example 3

Find the trend values

1953.5

→ Year : 1950	1951	1952	1953	1954	1955	1956	1957
→ Value : 346	411	392	512	626	640	611	796

Year	New value ($\tilde{y} = y - 1953.5$)	y	n^2	ny
1950	-3.5	346	12.25	-12.11
1951	-2.5	411	6.25	-102.75
1952	-1.5	392	2.25	-588.0
1953	-0.5	512	0.25	-256.0
1954	0.5	626	0.25	313.0
1955	1.5	640	2.25	960.0
1956	2.5	611	6.25	1527.5
1957	3.5	796	12.25	2786.0
$\sum n = 0$		$\sum y = 4334$	$\sum n^2 = 42$	$\sum ny = 2504.0$

$$y = 541.75 + (59.6) n$$

Trend values

333.15

392.75

452.35

$$\begin{aligned} \text{Min } & \sum b^n = \sum y \\ & \sum b^n \times \sum a^n \\ & \sum a^n \times \end{aligned}$$

•> Nearest Neighbour Method. / Non Parametric Method.
 Instance based / Memory based
 Learning method / Distance based.

The steps in KNN classification

Step 1: Look at the data

Step 2: Calculate the distance.

Step 3: Find neighbour

Step 4: Vote on labels

Example

•> Predict the weight of the person with ID 11 based on his height and age.

Height	Age	Weight	
5	45	77	
5.11	26	47	
5.6	30	55	
5.9	34	59	
4.8	40	72	
5.8	36	60	
5.3	14	40	
5.8	28	60	9d 1 5 6
5.5	23	45	
5.6	32	58	$\frac{77 + 72 + 60}{3} = 65.6 \text{ kg}$
5.5	38	69.6 ?	

$$k=3 \quad \text{ID } 11 \text{ weight} = 69.6 \text{ kg}$$

$$k=5 \quad \text{ID } 11 \text{ weight} = 65.2 \text{ kg}$$

KNN works best when data is not too large.

$$k=5$$

$$I_{D1} = \frac{5.5}{(x_2)} \quad \frac{38}{(y_2)} \quad \text{Unknown}$$

$$I_{D1} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} = \sqrt{(5.5 - 5)^2 + (38 - 45)^2} = 7.017$$

$$I_{D2} = \sqrt{(5.5 - 5.11)^2 + (38 - 26)^2} = 12.04$$

$$I_{D3} = 8.0006$$

$$I_{D4} = 4.019$$

$$I_{D5} = 2.011$$

$$I_{D6} = 2.022$$

$$I_{D7} = 19.001$$

$$I_{D8} = 10.0044$$

$$I_{D9} = 15$$

$$I_{D10} = 6.0008$$

•> Distance Based Learning Method.

→ Distance can be defined in multiple ways

Non-negativity : $d_{ii} \geq 0$

Self-proximity : $d_{ii} = 0$

Symmetry : $d_{il} = d_{li}$

Triangle Inequality : $d_{il} \leq d_{ik} + d_{kl}$

•> Properties of Distance

1) If the two objects are not same then distance between them should be greater than zero. (Positivity.)

2) Symmetry

3) Self proximity.

4) Distance between two object is less than the summation of these distance where k is common to both i & l is called triangle inequality.

5) If a distance satisfies all its properties then it is called metric

Types of Distances

↳ Euclidean Distance

$$d_{ii} = \sqrt{(x_1^{(i)} - x_1^{(i)})^2 + (x_2^{(i)} - x_2^{(i)})^2 + \dots + (x_n^{(i)} - x_n^{(i)})^2}$$

•> Statistical Distance

$$d_{ii} = \sqrt{(u^i - u^j)^T \Sigma^{-1} (u^i - u^j)}$$

u^i & u^j are n -dimensional vectors of measurement value of pattern i & j . Σ is covariance matrix of these vectors.

•> Manhattan Distance

$$\rightarrow d_{ii} = \sum_{i=1}^n |u_i - y_i|$$

•> Minkowski

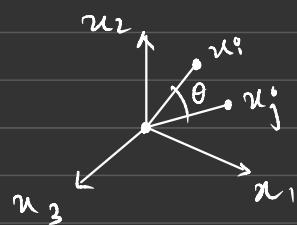
$$\rightarrow L_p(u^i, u^j) = \left(\sum_{j=1}^n |u_j^{(i)} - u_j^{(j)}|^p \right)^{1/p}$$

•> Hamming Distance

Using XOR operation

Calculate distance between two binary vectors.

•> Cosine Similarity :-



The Cosine Similarity between two vectors u^i & u^j is defined as cosine of the angle θ between them.

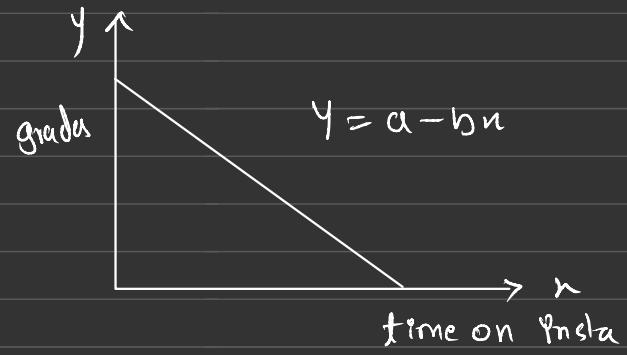
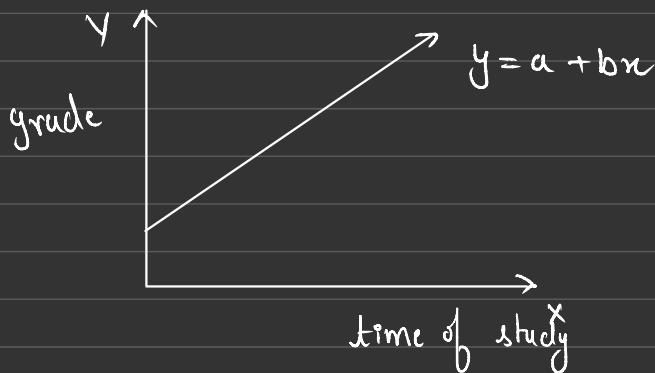
$$\cos \theta = \frac{u_1^i u_1^j + u_2^i u_2^j + \dots + u_N^i u_N^j}{\sqrt{\sum_{j=1}^N (u_j^i)^2} \times \sqrt{\sum_{j=1}^N (u_j^j)^2}}$$

→ For a given new data point $B = (1.4, 1.6)$ find the cosine similarity wrt the data given.

Cosine similarity

Sol	SL NO	A ₁	A ₂	$\frac{1.5 \times 1.4 + 1.7 \times 1.6}{\sqrt{1.5^2 + 1.7^2} \times \sqrt{1.4^2 + 1.6^2}}$	= $\frac{4.82}{4.8200414936} = 0.99999$
1	1.5	1.7			
2	2	1.9			
3	1.6	1.8			
4	1.2	1.5			

Unit-3 (Linear Regression)



We will study this.

- Regression - predicts values of responsible variable from attribute variables.
- Variable - continuous numeric values

Linear Regression

Simple Linear Regression

$$y = \underbrace{w_0}_{\text{Intercept}} + \underbrace{w_1 x}_\text{weight function / regression coefficient}$$

One-one

Multiple Linear Regression.

$$y = w_0 + w_1 x_1 + w_2 x_2$$

Many-one

The linear regression analysis is used to predict the value of response variable based on the value of attribute variable. where the output or the response variable is continuous in nature. here the relation between dependent & independent variable is linear.