

UNIT - 1 Introduction to machine learning

Evolution of machine learning, Paradigms for ml, learning by Rote, learning by Induction, Reinforcement learning, Types of Data, matching, stages in machine learning, Data Acquisition, feature Engineering, Data Representation, model selection, model learning, model Evaluation, model Prediction, search and learning, data sets.

UNIT - 2 Nearest Neighbour - Based models

Introduction to Proximity measures, Distance measures, Non-metric similarity functions, Proximity b/w binary patterns, Different classification Algorithms based on the distance measures, K-Nearest Neighbor classifier, Radius Distance Nearest Neighbor Algorithm, KNN Regression, Performance of classifiers, Performance of Regression Algorithms

UNIT - 3 Models Based on Decision Trees

Decision trees for classification, Impurity measures, Properties, Regression Based on Decision trees, Bias-variance Trade-off, Random forests for classification and Regression.

The Bayes classifier: Introduction to the Bayes classifier, Bayes' Rule and Inference, The Bayes classifier and its optimality, multi-class classification / class conditional independence and Naive Bayes classifier (NBC)

UNIT - 4 Linear discriminants for machine learning

Introduction to linear discriminants, linear

Discriminants for classification, Perceptron classification, Perceptron learning Algorithm, Support vector machines, linearly Non-separable case, Non-linear SVM, Kernel Trick, Logistic Regression, Linear Regression, multi-layer Perceptrons (MLPs), Back Propagation for training an MLP

UNIT - 5 clustering

Introduction to clustering; Partitioning of data, matrix factorization / clustering of patterns, divisive clustering, Agglomerative clustering, Partitional clustering, k-means clustering, soft partitioning, soft clustering, fuzzy c-means clustering Rough clustering, Rough k-means clustering Algorithm, Expectation maximization - based clustering, Spectral clustering.

on successful completion of the course, students will be able to		Pos
CO1	Identify machine learning techniques suitable for a given problem	PO6, PO7, PO8, PO12
CO2	Solve real-world problems using various machine learning techniques	PO6, PO7, PO8, PO9, PO12
CO3	APPLY Dimensionality reduction techniques for data preprocessing.	PO6, PO7, PO8, PO9, PO12
CO4	EXPLAIN what is learning and why it is essential in the design of intelligent machines	PO6, PO7, PO8, PO9, PO12
CO5	Evaluate Advanced learning models for language, vision, speech, decision making etc.	PO6, PO7, PO8, PO9, PO12

Introduction to machine Learning

Evolution of machine learning :-

Machine learning is a process of a learning model that can be used in prediction based on data

- * Prediction involves assigning a data item to one of the data (classification) or associating the data item with a number (regression)
- * ML gain its importance because of the input processing space and storage space of computer availability of large data state for experimentation.
- * Deep learning is an offshot of ML Perceptron is the popular ML tool. It is basic building log of various architecture such as multilayer perceptron, convolution neural network (CNN) Recurrent neural network (RNN).
- * The early base of AI it was assumed that mathematical logic is used for the AI system.
- * Some of the contribution, general problem solver (GPS) automatic trainer, pruner, rule based system, programming language like Prolog & LISP were outcomes of this view during 20th century.

AI researcher were of the view logic is AI and AI is logic the reasoning system were developed based on this view.

- * In early 21st century the view is that AI is DL and DL is AI. The advent of graphical processing unit tensorflow, pytorch along with CNN, RNN, have changed every aspects of science & engineering activities across the globe. High level view of AI

Data structures & Algorithm of basics to AI systems:-

The logic and describe system. Play an imp role in the analysis synthesis of AI system. In ML the data input may be viewed as a matrix, called a data matrix. In there are n data items each represented as L dimensional vector then the corresponding data matrix is of size $n \times L$. Linear algebra is useful in analysing the weights associated with the edges in a neural network. The clustering may be viewed as data manipulates data factorization. The Probability & statistics helped in estimation the distribution analysing the data. Optimization is essential in training neural networks. Information theory concepts like entropy, mutual information, divergence are important to understand. Topics such as decision tree, classifier, feature selection, deep neural networks.

Paradigms for ML:- There are basically 5 types.

There are different ways Paradigms for ML such as

- 1) Learning by rote
- 2) Learning by deduction
- 3) Learning by abduction
- 4) Learning by induction
- 5) Reinforcement Learning

Learning by rote:-

This involves memorization in an effective manner. Memorizing alphabets and numbers, multiplication tables. Examples of rote learning. In the case of data checkend the computer values are stored so that the values no need to recompute later. The checking is implemented by such engines and it is also

- example of rote learning the when computation is more expensive then recall this rote learning cancel the a significant large amount of time

chess master :- spend a lot of time memorizing the great gains of the past that is called rote learning the teacher them how to think in the chess.

Learning by deduction :-

Deductive learning deals with exploitation of deductions may explain this type of learning is based on reasoning that is truth preserving.

* given A & $A \rightarrow B$ we can deduce B we can use B along with $B \rightarrow C$ to deduce C . Note that A & $A \rightarrow B$ are true, then B is true. that is the truth preserving nature of learning by deduction.

consider the following statement :-

1. It is raining
 2. It rains, the roads are wet
 3. If a road is wet, it is slippery
- from 1 & 2 we can infer using deduction that
4. the roads are wet.
- this deduction can then be with 3 to deduce the
5. the roads are slippery

* The statements 1, 2 & 3 are true then statement 4 & 5 are automatically true

* deductive learning is apply in game playing.

Learning by abduction :-

Here we can infer A from B & $(A \rightarrow B)$ this is not

truth preserving as both $B \wedge (A \rightarrow B)$ can be true and A can be false.

→ for example

1) An Aeroplane is a flying object
Aeroplane \rightarrow flying object

2) A is a flying object
A is an Aeroplane

A can be a bird or a kite

from i & 2 we infer using abduction that A is an Aeroplane
this kind of reasoning may lead to incorrect conclusion.

"A can be a bird or a kite"

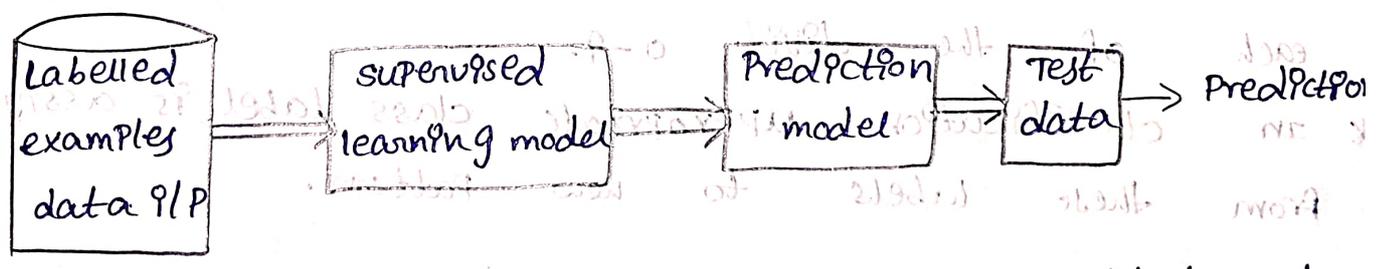
Learning by induction:-

* Here learning by achieved with the help of examples are observations it can be may be categorized as.

i) Learning from examples (supervised learning)

ii) Learning from observations (unsupervised learning)

Learning from examples:-



* Here a collection of label examples are provided and machine learning system uses these examples to make a prediction on a new data pattern.

* In supervised classification are learning for examples there are two machine learning models that is

classification and Regression.

classification:-

consider hand written digits shown in figure

0 0 0

1 1 1

2 2 2

3 3 3

4 4 4

5 5 5

6 6 6

7 7 7

* Here each row has three examples of each of the digit the problem is to learn on machine learning model using this data to classify new data pattern this is called supervised learning as the model learning with help of exemplar data.

* These exemplar data is provided by an expert for example a medical doctor provide examples of normal patient and patients infected by covid based on test results.

* In the case of hand written digits we have ten class labels. one class label correspond to each of the digit 0-9.

* In classification appropriate class label is assigned from these labels to new pattern.

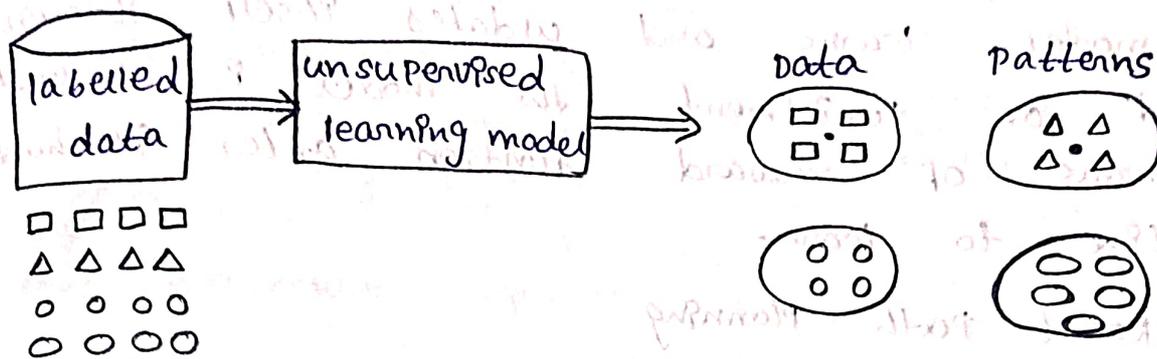
Regression:-

* There are several prediction applications where the labels come from an infinite set.

EX:- share value of a stock could be positive real number.

* The stack may have different values at a particular time and each of these values is a real number

Learning from observations :-



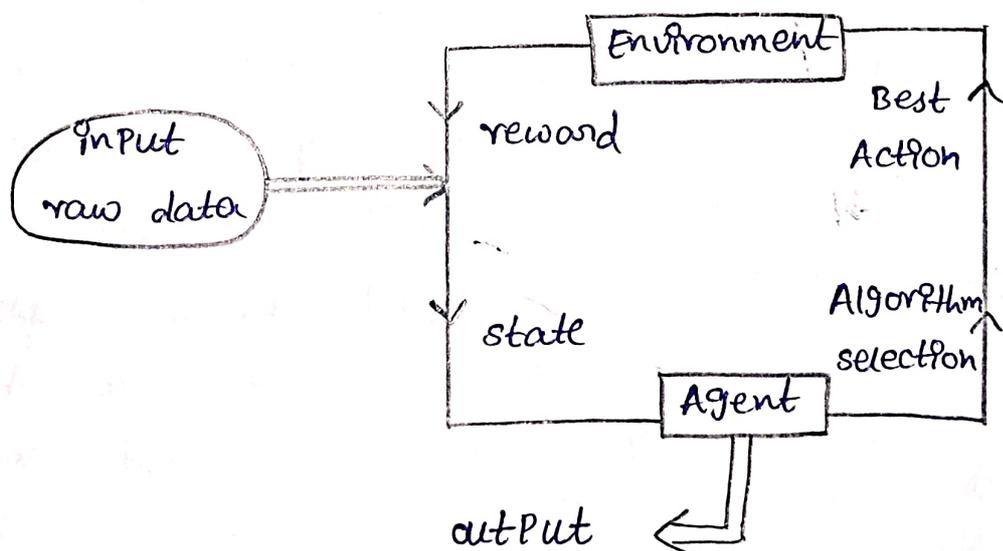
* observations are the instance that need not be labelled in this case i.e. cluster (or) group the observations into a smaller number of groups.

* such grouping is performed with the help of clustering algorithms that assigns similar patterns to the same group (or) cluster.

* each cluster is represented by centroid or mean for example x_1, x_2, \dots, x_p are p elements are cluster then the centroid is defined by

$$\frac{1}{p} \sum_{i=1}^p x_i$$

Reinforcement Learning :-



* In reinforcement learning the agent learns optimal policy to optimize some reward function the learned policy helps the agent in taking an action based on current state of the problem

* the model learns and updates itself through reward or punishment the model is evaluated by means of reward function after it had some time to learn.

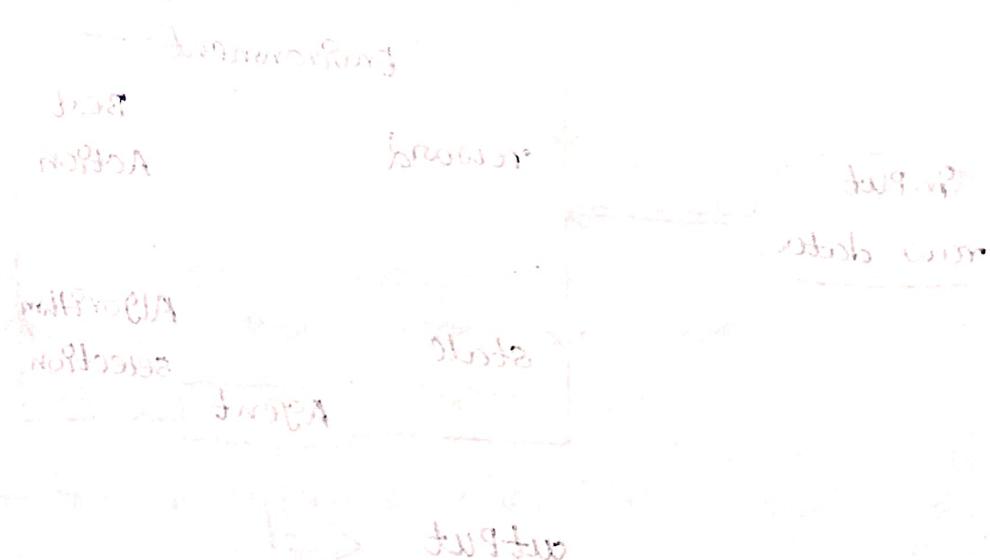
Ex:- Robot Path Planning

write the difference between supervised and unsupervised reinforcement learning.

Supervised learning is a type of machine learning where the model is trained on a set of input-output pairs. The goal is to learn a function that maps the input to the output. This is done by minimizing the loss function. Supervised learning is used for tasks like image classification, text classification, and regression.

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

Reinforcement Learning



Supervised	unsupervised	Reinforcement
<ul style="list-style-type: none"> * This type of learning is used when you know how to classify a given data, or in other words classes or labels are available. * labelled training data is needed. model is built based on training data. * The model performance can be evaluated based on how many misclassifications have been done based on a comparison b/w predicted and actual values. * There are two types of supervised learning problems classification and regression. * simplest one to understand 	<ul style="list-style-type: none"> * This type of learning is used when there is no idea about the class or label of a particular data. The model has to find pattern in the data. * Any unknown and unlabelled data set is given to the model as input and records are grouped. * Difficult to measure whether the model did something useful or interesting homogeneity of records grouped together is the only measure. * There are two types of unsupervised learning problems clustering and association * more difficult to understand and implement than supervised learning. 	<ul style="list-style-type: none"> * This type of learning is used when there is no idea about the class or label of a particular data. The model has to do the classification. It will get rewarded if the classification is correct, else get punished. * The model learns and updates itself through reward punishment. * Model is evaluated by means of the reward function after it had some time to learn * No such types * Most complex to understand and apply.

standard algorithms

include

- * Naive Bayes
- * K-nearest neighbour (KNN)
- * Decision tree
- * Linear regression
- * Logistic regression
- * Support vector machine SVM etc

practical applications

include

- * Handwriting recognition
- * stock market prediction
- * Disease prediction
- * fraud detection etc

standard algorithms

are

- * K-means
- * Principal component Analysis (PCA)
- * self-organising map (SOM)
- * A priori algorithm
- * DBSCAN etc.

practical applications

include

- * market basket analysis
- * Recommender systems
- * customer segmentation, etc

standard algorithms

are

- * on-learning
- * sarsa

practical applications

include

- * self-driving car
- * intelligent robots
- * Alpha Go zero (the latest version of DeepMind's AI system playing Go)

Types of data sets:-

The data can be divided into 2 types.

1. Qualitative data (Categorical data)
2. Quantitative data (Numerical data)

Qualitative Data:-

* Qualitative data provides information of the quality of object which cannot be measured.

Ex:- Name, Roll No etc

* It can be further divided into two types

1. Nominal data
2. Ordinal data

Nominal data:-

* It has no numeric values, but a named value.

* Nominal values cannot be measured.

Ex:- 1. Blood group: 'A', 'B', 'AB', 'O'

2. Gender: 'male', 'female'.

3. Nationality: 'Indian', 'American', 'British'...

* Mathematical operations +, - etc cannot be performed on nominal data that is why mean variance can not be applied. where as a basic count is possible.

* So mode most frequently occurring value can be identified for nominal data.

Ordinal data:-

* These are named objects and can be arranged a sequence of increasing or decreasing value, so that we can say whether the value is better than.

another value.

Ex:- 1. customer, satisfaction: 'very good', 'Good', 'Average', 'poor'.

2. Grades of student: 's', 'A', 'B', 'C'---

3. Hardness of metal: 'very hard', 'hard', 'soft'

* mode can be identified.

* since ordering is possible median can be identified.

quantitative data:-

* quantitative data relates to information about the quality of object hence it can be measured.

Ex:- marks of subjects.

* quantitative data can be divided into

1. Interval data

2. Ratio data

Interval data:-

* The interval data is numeric data for which the order is known and the exact differences between values is also known.

Ex:- celsius, temperature.

* The difference between each value remains same in celsius temperature.

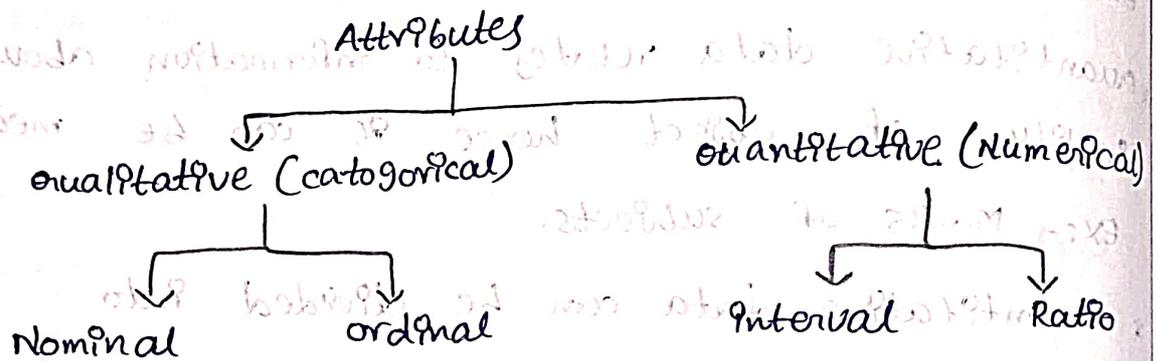
* The interval data do not have a true zero value there is very low zero temperature. Hence, only

addition and subtraction can be applied for interval data.

* The statistical function such as mean, median, mode, standard deviation, can be calculated.

Ratio data:-

- * The Ratio data represents numeric data for which exact value can be measured.
 - * absolute zero is available for ratio data,
 - * all mathematical calculation applied and all statistical calculations are applied.
- Ex:- height, weight, age, salary etc



- * Attributes can be classified as discrete or continuous based on the number of values that can be assigned.

Discrete Attributes:-

- * It can take a finite or countable number of values.

Ex:- Pincode, Rank etc.

- * discrete Attribute which can assume two values only is called binary attribute.

Ex:- yes / no, true / false etc.

- * continuous attribute it can assume any possible values which is a real number.

Ex:- length, height etc.

- Matching :-
- * Matching in machine learning involves techniques to find similar items (or) records across datasets (or) with in a single data set.
 - * It is used in both supervised and unsupervised learning.
 - * Matching is found by using proximity measured such as distance / dissimilarity measure (or) similarity measure.
 - * To data items u & v represented as 1 dimensional vectors and these are matched when the distance between them is smaller (or) when the similarity between them is larger.
 - * The distance measure is Euclidean distance and a popular similarity measure is the cosine of the angle between vectors.

$$d(u, v) = \sqrt{\sum_{p=1}^l (u(p) - v(p))^2}$$

cosine similarity

$$\cos(u, v) = \frac{u^t v}{\|u\| \|v\|}$$

$\|u\|$

where, $u^t v$ is the dot product between the vectors u & v .

$\|u\|$ is Euclidean distance between u and origin it is

also called Euclidean Norm.

Applications of matching in machine learning :-

Finding the nearest neighbor of a pattern

let x be 1 dimensional vector and

$x = \{x_1, x_2, \dots, x_n\}$ is a collection of

n data vectors. the nearest neighbour of x from X denoted by $NN_x(x)$ is x_j . if

$$d(x, x_j) \leq d(x, x_i), \forall x_i \in X$$

* Assigning a set with the nearest representative:

centroid :-

* let $x_1, x_2, x_3, \dots, x_k$ are k sets with x_1, x_2, \dots, x_k are there centroids (representative) the vector x is assign to c_j if

$$d(x, x_i) \leq d(x, x_j) \text{ for } j = \{1, 2, \dots, k\}$$

where c_j is j th group (or) cluster x_j is centroid of c_j this will be useful in clustering or learning from observation:

stages in machine learning :-

The figure shows number of steps involved in machine learning system.

Application domain \rightarrow Data Acquisition



feature Engineering = Preprocessing + Representation

model selection \rightarrow choose a model \leftarrow Domain Knowledge

model learning \rightarrow train the model \leftarrow Training data

model evaluation \rightarrow validation the model \leftarrow Evaluation data (validation data)



model prediction \rightarrow learn the model \leftarrow Test data

model explanation \rightarrow explain the model \leftarrow Expert feedback.

* The available data is split into training and validation and test data.

* Training data is used in model learning are Training.

* Validation data is used to tune the machine learning model.

* Test data is used to examine how the learnt model is performing.

Data Acquisition: -

* It is the process of gathering and collecting relevant data from various sources and storing the data.

* It is important for obtaining high quality data for model training and optimizing the performance.

* This data make machine learning organisms to learn and perform (method).

Ex: - * To distinguish between adults and children the measurement of height and weight are sufficient.

* To distinguish between covid and non covid patient the body temperature and chest congestion are important.

Feature engineering (methods)

This step involves the combination of data PreProcessing and data representation.

1) Data PreProcessing

The raw data available needs to be updated before it can be used by an ML model. The common problems that occur with raw data or missing values, different ranges for different variables, and

the presence of outliers.

a) missing values

occurs as a consequence of the inability to measure a feature value or due to unavailability are erroneously data entry. some ML algorithms can work whenever there are less number of missing data values in this case there is no need for preprocessing. where as in most of the cases ML models can not work with missing values therefore there are different methods for predicating the missing values.

use the nearest neighbour

x is vector and its q th component $x(q)$ is missing

* let $x = x_1, x_2, \dots, x_n$ are the set of n training vectors.

* let $x_j \in X$ is the nearest neighbour of x based on the remaining $(l-1)$ components i.e, the predicted values x_j is $x_j(q)$ $x_j = x_j(q)$

use a larger neighbour

use the k -nearest neighbour of x to predict the missing value of $x(q)$ let the k NN of x using the remaining $l-1$ components from X or x_1, x_2, \dots, x_k Now the predicted value $x(q)$ is the average of q th component of this k NN.

$$x(q) = \frac{\sum_{j=1}^k x_j(q)}{k}$$

Ex:- the data vector is

$(1, 1, 1) (1, 1, 2) (1, 1, 3) (1, -, 2) (1, 1, -) (6, 6, 1)$

$k=3$

missing data vector is $(1, -, 2)$

k NNs are $(1, 1, 1) (1, 1, 2) (1, 1, 3)$

$$\frac{1+1+1}{3} = "1"$$

the data vector = (1, 1, 2)

cluster the data & locate the nearest cluster

this approach is based on clustering the training data and locating the cluster to which $x \in$ based on the remaining $d-1$ components. let x with its i th value its missing belongs to cluster c if u^c is the centroid of c then the predicted value of $x(i)$ is $u^c(i)$, $x(i) = u^c(i)$.

Ex:-

cluster 1: { (1, 1, 1) (1, 2, 2) (1, 3, 2) } centroid 1 = (1, 2, 2)

cluster 2: { (3, 4, 3) (3, 5, 3) (3, 3, 3) } centroid 2 = (3, 4, 3)

cluster 3: { (6, 5, 6) (6, 8, 6) (6, 7, 6) } centroid 3 = (6, 7, 6)

the missing data vector is

$$(1, -, 2)$$

is near to centroid 1 = (1, 2, 2)

the predicted value is "2"

the data vector is (1, 2, 2)

* the mean square error is used to find deviation of predicted value from original value. let the true values are y_1, y_2, \dots, y_n and predicted values are $y_1^{\wedge}, y_2^{\wedge}, \dots, y_n^{\wedge}$ and mean square value

$$\sum_{i=1}^n (y_i - y_i^{\wedge})^2$$

$$MSE = \frac{\sum_{i=1}^n (y_i - y_i^{\wedge})^2}{n}$$

b) different ranges for different variables :-

The scales of value of different features could be very different this will effect the process to depend more upon the features

that has larger values. that will reduce the contribution of the features with smaller values.

EX:- classification of adult (or) child Here height is measured in meters and weight in grams. the adult is represented by $(1.6, 75000)$ the child is represented by $(0.6, 5000)$ the range of height is $[0.5, 2.5]$ range of weight $[2000, 200000]$

the Euclidean distance between adult & child is

$$\text{Euclidean distance} = \sqrt{(1.6 - 0.6)^2 + (75000 - 5000)^2} \approx 7000$$

$$\text{The cosine similarity} = \frac{\begin{bmatrix} 1.6 \\ 75000 \end{bmatrix} \cdot \begin{bmatrix} 0.6 \\ 5000 \end{bmatrix}}{\sqrt{(1.6)^2 + (75000)^2} \sqrt{(0.6)^2 + (5000)^2}}$$

The Euclidean distance and the cosine of angle between the two vectors depend more upon the weight feature while the contribution of height is negligible. This can be handled by scaling different components differently and this process is called scaling. There are two important normalization methods.

scaling using the range:-

for any categorical features the values of two features either match or miss match and the contribution to the distance is either zero (match) or 1 (miss match) in the case of numerical features the contribution should be in the range $[0, 1]$ this can be obtained by scaling the differences by the range of values of the feature. If p th component contribution to the distance between two objects x^i & x^j is

$$|x^i(p) - x^j(p)|$$

Range P

where Range P is Range of i th feature the value of this term is in the range $[0, 1]$ the value of 1 is obtain when $|x^i(p) - x^j(p)| = \text{Range P}$ and zero is obtain $x^i(p) = x^j(p)$.

standardization :-

The data is normalized so that it will have mean and unit variance this is the property of standard normal distribution

Ex:- Consider the 5th component of the data vector is 60, 50, 20, 100 & 40

$$\text{mean} = \frac{60 + 50 + 20 + 100 + 40}{5} = 60$$

zero mean data = 0, 20, -40, 40 & -20

$$\text{variance } \sigma^2 = \frac{0^2 + 20^2 + (-40)^2 + 40^2 + (-20)^2}{5} = 800$$

$$\text{standard deviation } \sigma = \sqrt{\sigma^2} = \sqrt{800} = 28.284$$

$$\frac{0}{28.28} \quad \frac{20}{28.28} \quad \frac{-40}{28.28} \quad \frac{40}{28.28} \quad \frac{-20}{28.28}$$

zero mean $\Rightarrow 0, 0.707, -1.414, 1.414, -0.707$

& unit variance data

This data is standard normal distribution data.

(c) outliers in the data

The outliers in data item that is either noise or are erroneous noise measuring instrument are

error data recording are responsible for presence of outliers. A data item is called an outlier if its values are away from average of data items or it is not similar to other items in terms of its characteristics. outliers occur because of different reasons.

1) data entry errors

2) measurement error

due to faulty devices or experiment set up problem.

3) Experimental error:-

due to faults in experimental design

4) intentional outliers

5) data performing error

6) natural variances

It is possible to use clustering to locate search outliers.

Data representation:-

each data item is represented by a vector in the partial applications the dimensionality of data or the size of the vector can be very large the high dimensionality data is common in bioinformatics information retrieval.

satellite imaging

the difficulties in using high dimensionality data vectors are

* computation time increases with the dimensionality.

* storage space are also increases

* performance of the model high dimensionality data performance of the model

* there is a problem called peaking phenomenon that is in dimensionality increases the accuracy of the ML model increases until some values and beyond that value the accuracy starts decreasing. the model will also called overfitting that is the model will perform well for training data and fails to perform well of validation data.

the dimensionality reduction is used in ML model so that the model does not over fit for available data the dimensionality reduction approaches.

* feature selection :-

let $f = \{f_1, f_2, \dots, f_L\}$ are L features in the feature selection approach will select only P no. of feature $P < L$ such that these P no. of features will maximize the performance of ML model.

Feature extraction :-

let $f = \{f_1, f_2, \dots, f_L\}$ are L features a set

$H = \{h_1, h_2, \dots, h_l\}$ of $l < L$ features is extracted

so this is possible with following methods.

1. linear method :-

In this $h_j = \sum_{i=1}^L \alpha_{ij} f_i$, $j=1, 2, \dots, l$ that is each element of "H" is a linear combination of original feature. some of the methods in this category are

a. Principal component

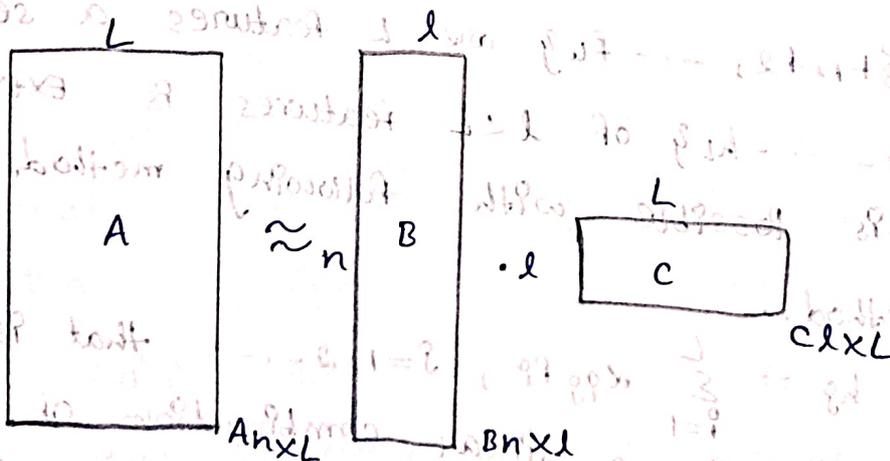
* consider data of n vectors in L -dimensional space that is $A_{n \times L}$ the covariance matrix of "A" is of size $L \times L$ that is $E_{L \times L}$ and the eigen vectors of "E" form the principal components.

* the eigen vector corresponding to the largest eigen value is the first principal component similarly the eigen vector corresponding to second largest eigen value is the second principal component. both the original features and principal component are used to represent any data vector.

b. Non negative matrix factorization (NMF)

NMF :- is a method to break down larger dataset into smaller meaningful parts and ensuring that all the values are non-negative this is useful in extracting important features from the data and these features are easy to analysis and process it.

NMF is factorization of matrix $A_{n \times L}$ into a product of $B_{n \times l}$ & $C_{l \times L}$



Non-linear feature extraction:-

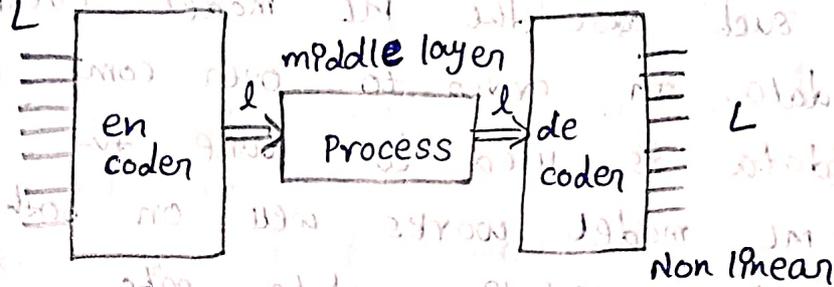
$H = \{h_1, h_2, \dots, h_k\}$ is represented using $h_i = t(f_1, f_2, \dots, f_k)$

with t is a non-linear function.

Ex: - $f = \{f_1, f_2\}$ linear

$$h_1 = a_1 f_1 + a_2 f_2 + a_3 f_1 f_2 \text{ Non-linear}$$

This is a non linear because of the term $f_1 f_2$ in h_1
Autoencoder is a Non-linear feature extraction tool or method



Model selection

selection of the model depends upon the nature of the data and knowledge of application domain. If some features or numerical and others are categorical that classifiers based on Perceptron (or) support vector machine are not suitable because the dot product between non-numerical data is not possible.

Therefore Bayesian models and decision tree models are suitable for categorical data (Non-numerical)

Model learning

This step depends on the size and type of training data. A subset of labelled data is used as training data for learning the model. Another subset is used for model validation or model evaluation. Some of the ML models are highly transparent while others are opaque or black box models.

Ex: - A decision tree models or transparents
Neural networks are opaque as the output of

Intermediate layers may not offer transparency.

Model Evaluation (Model Validation)

This step requires validation data of the ML model works well on the training data then we say that the model is well training some times it may not work well on validation data in such case the ML model overfits the training data. In order to overcome over fitting validation data is used to tune the ML model so that ML model works well on both the training and validation data sets.

Model Prediction

This step deals with testing the model this step uses test data set for the prediction is real world the model is used for prediction as new input keep coming in.

EX:- An ML model built for medical diagnosis should predict like doctor an means that diagnosis when a new patient uses it.

Model explanation :-

This step is important to explain an expert or a manager. why a decision was taken by the ML model. explanation is become very important in deep learning because deals system used neural networks that are opaque or black box. such opaque behaviour has created the need for explainable AI.

search and learning:-
search is the fundamental operation in both ML & AI. search was used in problem solving theorem proving, planning, and knowledge based systems. search is used in several computer science applications.

Ex:-
search is used in databases for answering any queries. In ML search is used in learning a classification model, a proximate measure for clustering and classification and learning model for regression.

* Inference is search in logic and probability.

* In algebra matrix factorization is search in optimization regularizer is used to simplify the search in finding a solution. In information theory the search for low entropy.

* optimization, inference, matrix factorization are important in ML and all are based on search.

* learning itself is search.

Explanation offered by the model:-

conventional AI systems are logic based, are rule based so the reasoning process has transparency and explainability both forward and backward. is used in conventional AI systems. some knowledge based systems are used in diagnosis and in teaching because of this flexibility.

Ex:-

The knowledge base used by the "MYCIN" expert system used in teaching medical students using another expert system "GUIDON". the problems

with conventional AI systems are.

1. There was no general framework for building AI systems. The experience in building one AI system did not simplify the building of another AI system.

2. Acquiring knowledge is a great challenge because different experts differed in their conclusions meeting to inconsistency.

current AI systems:-

current AI systems based on deep learning and "DL" is large data dependent they learn representations automatically they used multilayer neural networks and back propagation algorithm in training the models. The difficulties with "DL" systems are.

1. They are data dependent where performance improves as the size of the data set increases.

2. Learning in "DL" systems involves a simple change of weights in the neural networks this is done with the help of back propagation.

3. "DL" systems are black box system and does not have explainable capability therefore the current "AI" researches working on explainable "AI".

Data sets used:-

data sets for char classification:-

1. MNIST Hand written digit dataset

There are 10 classes corresponding to digits 0, 1, ..., 9 and each digit is an image of size 28×28 pixels and each pixel having

values in the range 0 to 255 (0 - black, 255 white).

* there are around "6000" digits hand training patterns and around 1000 text sample in each class the class label is also provided for each of the digit.

2. Fashion MNIST Data set :-

* It is a data set of Zalando's article images and consisting of a training set of 60,000 examples and test set of 10,000 examples. each example is 28x28 gray scale image and associated with a label from 10 classes.

* It is used as replacement for original MNIST data set for ML models.

* This data set is taken from in built data set of kaggle.

3. Olivetti face data set

* It consist of 10 different images for each of 40 distinct objects.

* for some objects the images was taken at different times, varying the lighting.

* facial expressions (open eye, close eye, smiling, not smiling) and facial details that is glasses and not glasses.

* all images was taken against dark homogeneous background.

* each image is of the size 64x64.

* It is available scikit-learn platform.

4. Wisconsin Breast cancer data set

* we consist of 569 examples and each is a

30 dimensional vector there are two classes Benign and malignant the number of samples from benign is 357 and the number of malignant patterns is 212.
* It is available on scikit-learn.

Data sets for Regression

1. Boston Housing datasets :-
It has 506 examples each example is 30 dimensional vector. It is available on scikit-learn platform.

2. Airline Passengers data sets

This data set provides monthly totals of US Airline Passengers from 1949 to 1960. This data set is taken from in built data set of "kaggle".

3. Australian weather data set

It provides various weather record details for cities in Australia. The features include location, minimum and maximum temperature. This data set is taken from an in built data set of "kaggle".

Nearest Neighbour - Based Models

Introduction to Proximity measures:-

Proximity measures are used to measure the degree of similarity or dissimilarity between two or more data vectors. These data vectors can represent documents, images, or audio or video files. Proximity measures are used by machine learning algorithms to compare and classify or group or make predictions using these vectors. The important proximity measures are distance measures.

1. Distance measure:-

Distance measures are mathematical functions used to find how similar or dissimilar two objects are based on their features. The distance measure is used for clustering, classification, and information retrieval.

1) Euclidean distance

It represents the shortest distance between two vectors. It is used in clustering analysis.

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

It is used for continuous numerical data.

2) Manhattan distance

It determines the absolute difference, the pair of coordinates in a plane. Coordinates are

(x_1, y_1) & coordinates (x_2, y_2) the Manhattan distance between P and Q is

$$|x_1 - x_2| + |y_1 - y_2|$$

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

* It is used for high dimensional data. It is also called city block distance.

3) Minkowski distance

* It is a generalized distance measure that includes both Euclidean and Manhattan distance as special cases and is controlled by a parameter p .

$$d(x, y) = \left[\sum_{i=1}^n |x_i - y_i|^p \right]^{1/p}$$

* It is used for flexible distance calculations. For $p=1$ it becomes Manhattan and for $p=2$ it becomes Euclidean.

4) Hamming distance

* Number of positions where two vectors differ and it is used for error detection and sequence comparison.

$$d(x, y) = \sum_{i=1}^n [x_i \neq y_i]$$

where, $[x_i \neq y_i] = 1$ if symbols differ, $= 0$ else

* It is used for binary string DNA sequences error correction.

EX: - 1. Karolin

2. Kathryn

$$d(x, y) = 3$$

Similarity measures

It is used to calculate how similar the objects are. The higher the value of similarity, the more similar the object.

1. cosine similarity
- * measure cosine distance of angle between two vectors
 - It focuses on orientation rather than magnitude.

$$\text{cosine}(x, y) = \frac{x \cdot y}{\|x\| \|y\|} \rightarrow \text{similarity formula}$$

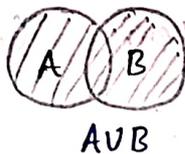
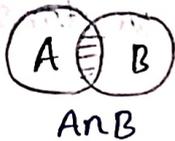
$$\text{distance formula} = 1 - \frac{x \cdot y}{\|x\| \|y\|}$$

- * It is used for text mining "NLP" It also used in recommendation systems.

Natural language Process

Jaccard index

Intersection union



$$* J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- * The Jaccard distance is set based distance that compares this similarity by the ratio of intersection to union.

$$d(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

- * It is used for binary or categorical data specially for sets.

Pearson correlation

It measures linear relation between variables

distances measure

- * Distances measure is used to find the dissimilarity between patterns represented as vectors.

- * Patterns which are more similar could be closer the distance function (d) is a metric or non-metric

The most popular distance metric is min Kowsky metric and a metric is type of a measure that satisfies three attributes:

1. Positive reflexivity

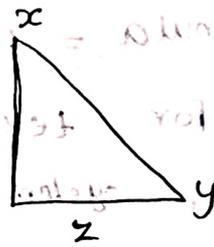
$$d(x, x) = 0$$

2. Symmetry

$$d(x, y) = d(y, x)$$

3. Triangular inequality

$$d(x, y) \leq d(x, z) + d(z, y)$$



where x, y, z are three patterns.

* The different types of dissimilarity or similarity measures between pattern (data, P/P) vectors.

1. Minkowski distance

The Minkowski metric can be defined as $d_r(P, Q) =$

$$\left(\sum_{k=1}^L (|P_k - Q_k|^r) \right)^{1/r}$$

where,

* r is a parameter that determines the type of metric

* P and Q are L -dimensional vectors based on the value of r .

* L_∞ norm: - Here $r = \infty$, $d(P, Q) = \max_k (|P(k) - Q(k)|)$

Here $r = \infty$ and $d(P, Q) = (|P(k) - Q(k)|)_{k \in \{1, \dots, L\}}$

* L_2 norm: - Here $r = 2$ and $d(P, Q) = \left(\sum_{k=1}^L (|P(k) - Q(k)|^2) \right)^{1/2}$

is Euclidean distance

* L_1 norm: - Here $r = 1$ and $d(P, Q) = \left(\sum_{k=1}^L (|P(k) - Q(k)|) \right)$

is the city-block distance.

fractional norm :- It is possible that R is a fraction the resulting distance is called fractional norm. It is not a metric as it does not satisfy the triangular inequality. Normalization of feature values is used to keep all the features in the same range. The Mahalanobis distance is a popular distance measure that is used in classification and it is computed using the covariance matrix. The squared Mahalanobis distance is given by

$$d^2(x, y) = (x - y)^t \Sigma^{-1} (x - y) \text{ where } \Sigma \text{ is covariance matrix.}$$

EX:- If $x = \begin{pmatrix} x_1 & y_1 & z_1 \\ 5 & 2 & 4 \end{pmatrix}$
 $y = \begin{pmatrix} x_2 & y_2 & z_2 \\ 3 & 4 & 2 \end{pmatrix}$

Euclidean distance = $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$
 $= \sqrt{(3 - 5)^2 + (4 - 2)^2 + (2 - 4)^2}$
 $= \sqrt{(-2)^2 + (2)^2 + (-2)^2} = \sqrt{4 + 4 + 4} = \sqrt{12}$
 $\Rightarrow 3.46$

Manhattan distance

$$d(x, y) = \sum_{i=1}^n |x_i - y_i| = |2| + |2| + |2| = 6$$

2. Non-metric similarity function

- * This category includes the similarity functions that do not obey the triangular inequality of similarity.
- * These are commonly used for image or string data and they do not include outliers or noisy data.

EX:- A non-metric similarity function k-median

A non metric similarity function is k -median distance between 2 vectors.

$$x = \{x(1), x(2), \dots, x(n)\}$$

$$y = \{y(1), y(2), \dots, y(n)\}$$

the formula called k -median distance is

$$d(x, y) = k\text{-median} \{ |x(1) - y(1)|, |x(2) - y(2)|, \dots, |x(n) - y(n)| \}$$

where

k -median operator returns the k th value of the ordered difference vector.

* the another similarity measure is the cosine of the angle between vectors $(x \& y)$. It is symmetric because $\cos(\theta) = \cos(\theta)$. the cosine similarity measure is

$$s(x, y) = \cos(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

the distance function corresponding to cosine similarity

is $d(x, y) = 1 - s(x, y)$, $d(x, y)$ is also symmetric

& it does not satisfies triangular inequality.

Ex:- let x, y, z are 3 vectors, the angle between

x and z is 45° the angle between z & y is

45° . show that it does not satisfies triangular

inequality.

$x, y \& z$

$x \& z$ is 45° , $z \& y$ is 45°

$$d(x, z) = 1 - \cos(45^\circ)$$

$$= 1 - \frac{1}{\sqrt{2}} = 0.2928$$

$$d(z, y) = 1 - \cos(45^\circ) = 1 - \frac{1}{\sqrt{2}} = 0.2928$$

$$d(x, y) = 1 - \cos(90^\circ) = 1$$

$$d(x, z) + d(z, y) = 0.2928 + 0.2928 = 0.5856$$

$$d(x,y) \leq d(x,z) + d(z,y)$$

$$1 \neq 0.586$$

3. Proximity between binary patterns:-

let P & Q are l 'l' bits, binary vector

(1) Hamming distance (HD):-

If $P(i) = Q(i)$ then P & Q are matched at i th bit

else if $P(i) \neq Q(i)$ then P & Q mismatched on i th bit,

the hamming distance is no. of mismatched bits of the l bit locations.

Ex:- $P = [1\ 1\ 1\ 0\ 1\ 1\ 1\ 1\ 1\ 1]$ hamming distance
 $Q = [1\ 0\ 0\ 1\ 0\ 0\ 1\ 1\ 0]$ is "7"

because they mismatched in 7 locations

2. simple matching co-efficient (SMC)

let us define the following:

where p is 0 and q is 1, m_{10} is No. of bpts

where p is 1 and q is 0, m_{00} is No. of bpts where

$p=0$ and q is 0, m_{11} is No. of bpts where p is 1

and q is 1

$$SMC(p, q) = \frac{m_{11} + m_{00}}{m_{00} + m_{01} + m_{10} + m_{11}}$$

Jacard co-efficient

Jacard co-efficient is defined as

$$JC(p, q) = \frac{m_{11}}{m_{10} + m_{01} + m_{11}}$$

$$P = 1001001101000$$

$$Q = 011100101110000$$

$$m_{01} = 3$$

$$m_{10} = 2$$

$$m_{00} = 2$$

$$m_{11} = 3$$

$$SMC(p, q) = \frac{m_{11} + m_{00}}{m_{00} + m_{01} + m_{10} + m_{11}}$$

$$= \frac{3 + 2}{2 + 3 + 2 + 3}$$

$$= \frac{5}{10} = 0.5$$

$$JC(p, q) = \frac{m_{11}}{m_{10} + m_{01} + m_{11}}$$

$$= \frac{3}{2 + 3 + 3}$$

$$= \frac{3}{8}$$

$$= 0.375$$

Different classification algorithm based on distance measure

1. Nearest Neighbor classifier (NNC)

* Here a test pattern T is classified based on its nearest neighbor in the training data. let

* let $X = \{(x_1, l_1), (x_2, l_2), \dots, (x_n, l_n)\}$ is labeled training data set of n patterns. the first component is pattern vector. and second component is its class label.

* l_i is class labels if there are P classes with their labels coming from the said $C = \{c_1, c_2, \dots, c_P\}$
 $l_i \in C \quad i=1, 2, \dots, n.$

* The nearest neighbor of test pattern T is given by

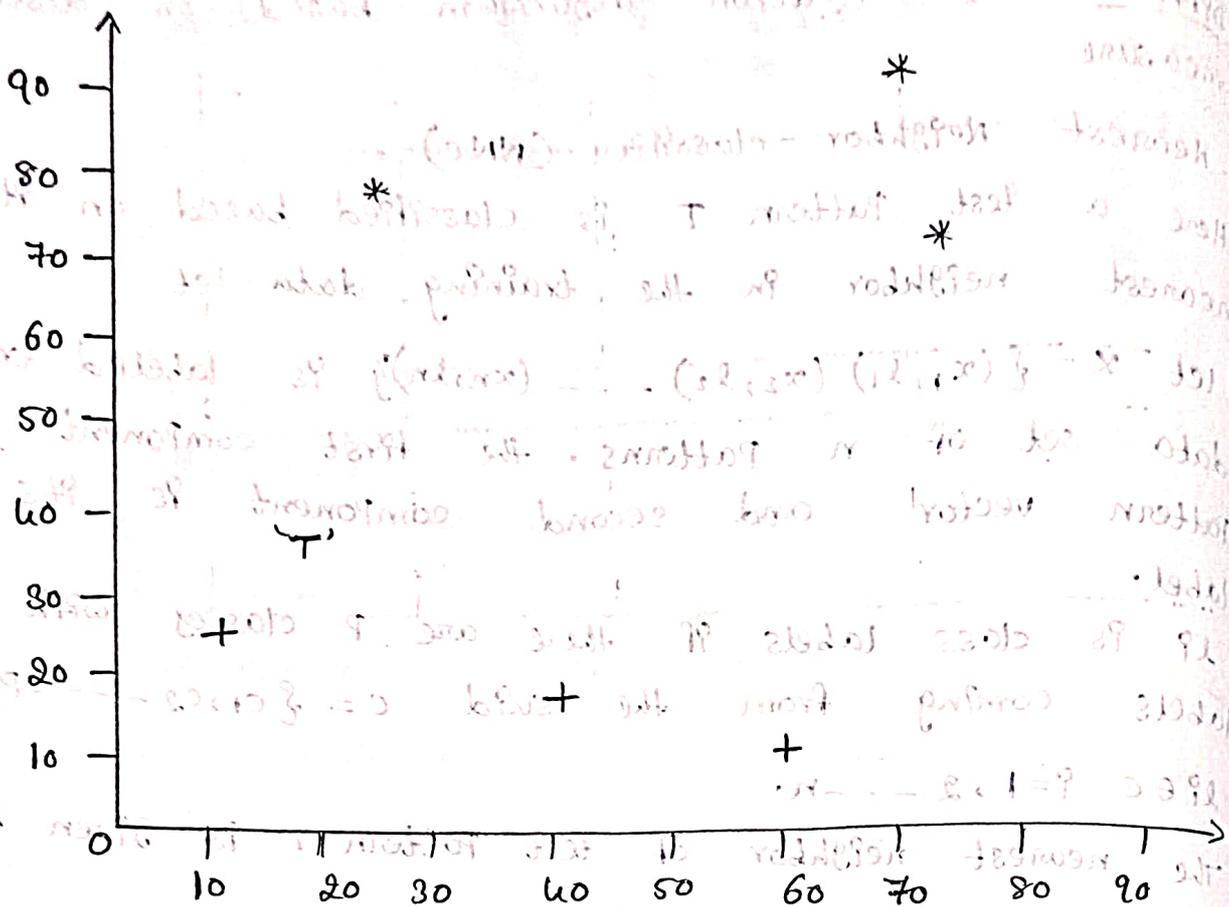
$$NN(T) = \arg \min d(x_i, T)$$

where x_i is i th training pattern and $d(x_i, T)$ is the distance between x_i and T , $NN(T)$ is at the minimum distance from T and it is maximum similar to T . the NN rule assigns the T with the class label of $NN(T)$.

Ex: - find the class of test pattern using Nearest Neighbor classifier.

	Brightness	saturation	class
x_1	40	20	Red
x_2	50	50	Blue
x_3	60	90	Blue
x_4	10	25	Red
x_5	70	70	Blue
x_6	60	10	Red
x_7	25	80	Blue
	20	35	?

Red '+'
Blue '*'



$$dTx_1 = \sqrt{(20-40)^2 + (35-20)^2} = \sqrt{20^2 + 15^2} = 25$$

$$dTx_2 = \sqrt{(20-50)^2 + (35-50)^2} = \sqrt{30^2 + 15^2} = 33.54$$

$$dTx_3 = \sqrt{(20-60)^2 + (35-90)^2} = \sqrt{40^2 + 55^2} = 68.01$$

$$dTx_4 = \sqrt{(20-10)^2 + (35-25)^2} = \sqrt{10^2 + 10^2} = 14.14$$

$$dTx_5 = \sqrt{(20-70)^2 + (35-70)^2} = \sqrt{50^2 + 35^2} = 61.03$$

$$dTx_6 = \sqrt{(20-60)^2 + (35-10)^2} = \sqrt{40^2 + 25^2} = 47.16$$

$$dTx_7 = \sqrt{(20-25)^2 + (35-80)^2} = \sqrt{5^2 + 45^2} = 45.27$$

$NN^i(T)$ is x_4

The class of x_4 is assign to T

class of T is "Red"

Brightness	saturation	class
20	35	Red

let the training

	brightness	saturation	class	distance
	40	20	Red	25
	50	50	blue	33.54
	60	90	blue	68.01
x_u	10	25	red	14.14
	70	70	blue	61.03
	60	10	Red	47.16
	25	80	blue	45.27

smallest

* let the training set consist of the following two dimensional patterns with associated labels.

$$x_1 = (0.7, 0.7), l_1 = 1;$$

$$x_2 = (0.8, 0.8), l_2 = 1;$$

$$x_3 = (1.1, 0.7), l_3 = 1$$

$$x_4 = (1.1, 1.1), l_4 = 1$$

$$x_5 = (3.7, 2.7), l_5 = 2$$

$$x_6 = (4.1, 2.7), l_6 = 2$$

$$x_7 = (3.7, 3.1), l_7 = 2$$

$$x_8 = (3.4, 0.6), l_8 = 3$$

$$x_9 = (3.1, 0.8), l_9 = 3$$

$$x_{10} = (3.1, 0.6), l_{10} = 3$$

$$1 = "\Delta"$$

$$2 = "+"$$

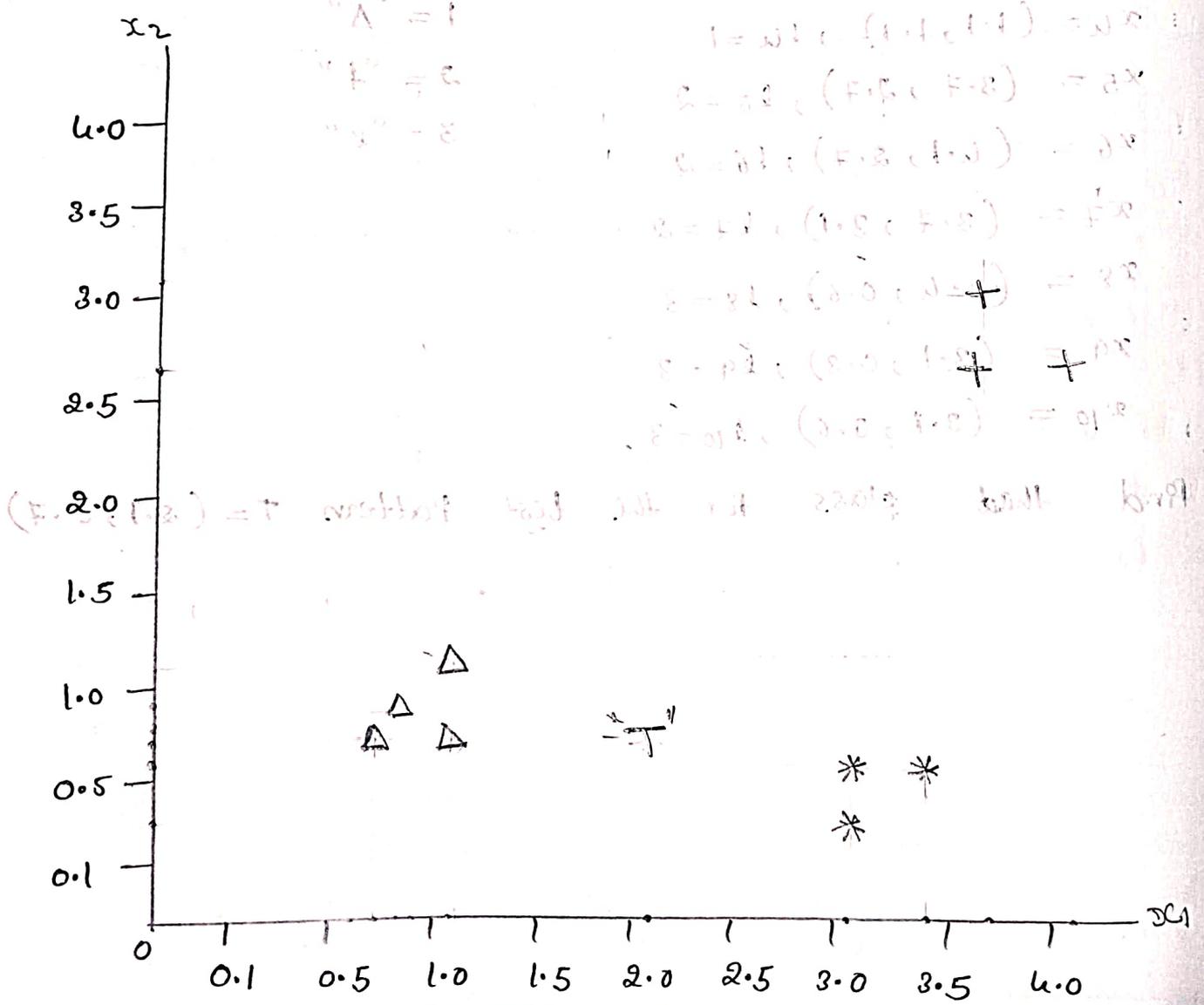
$$3 = "*"$$

Find the class for the test pattern $T = (2.1, 0.7)$

	Brightness	saturation	class
x_1	0.7	0.7	1 Δ
x_2	0.8	0.8	1
x_3	1.1	0.7	1
x_4	1.1	1.1	1
x_5	2.7	2.7	2 \times
x_6	4.1	2.7	2
x_7	2.7	1.1	2
x_8	2.4	0.6	3 $*$
x_9	3.1	0.3	3
x_{10}	3.1	0.6	3

l_1
 l_2
 l_3
 l_4
 l_5
 l_6
 l_7
 l_8
 l_9
 l_{10}

Brightness	saturation	class
2.1	0.7	?



$$dTx_1 = \sqrt{(2.1 - 0.7)^2 + (0.7 - 0.7)^2} = \sqrt{(1.4)^2 + (0)^2} = 1.4$$

$$dTx_2 = \sqrt{(2.1 - 0.8)^2 + (0.7 - 0.8)^2} = \sqrt{(1.3)^2 + (0.1)^2} = 1.303$$

$$dTx_3 = \sqrt{(2.1 - 1.1)^2 + (0.7 - 0.7)^2} = \sqrt{(1)^2 + (0)^2} = 1$$

$$dTx_4 = \sqrt{(2.1 - 1.1)^2 + (0.7 - 1.1)^2} = \sqrt{(1)^2 + (0.4)^2} = 1.077$$

$$dTx_5 = \sqrt{(2.1 - 3.7)^2 + (0.7 - 2.7)^2} = \sqrt{(1.6)^2 + (2)^2} = 2.561$$

$$dTx_6 = \sqrt{(2.1 - 4.1)^2 + (0.7 - 2.7)^2} = \sqrt{(2)^2 + (2)^2} = 2.828$$

$$dTx_7 = \sqrt{(2.1 - 3.7)^2 + (0.7 - 3.1)^2} = \sqrt{(1.6)^2 + (2.4)^2} = 2.884$$

$$dTx_8 = \sqrt{(2.1 - 3.4)^2 + (0.7 - 0.6)^2} = \sqrt{(1.3)^2 + (0.1)^2} = 1.303$$

$$dTx_9 = \sqrt{(2.1 - 3.1)^2 + (0.7 - 0.3)^2} = \sqrt{(1)^2 + (0.4)^2} = 1.077$$

$$dTx_{10} = \sqrt{(2.1 - 3.1)^2 + (0.7 - 0.6)^2} = \sqrt{(1)^2 + (0.1)^2} = 1.004$$

Brightness	Saturation	class	distance
0.7	0.7	1	1.4
0.8	0.8	1	1.303
1.1	0.7	1	1
1.1	1.1	1	1.077
3.7	2.7	2	2.561
4.1	2.7	2	2.828
3.7	3.1	2	2.884
3.4	0.6	3	1.303
3.1	0.3	3	1.077
3.1	0.6	3	1.004

→ smallest

NN (T) x_3

the class of x_3 is assign to T

class of T is "1"

Brightness	saturation	class
0.7	0.7	1

2. k-Nearest Neighbor Classifier

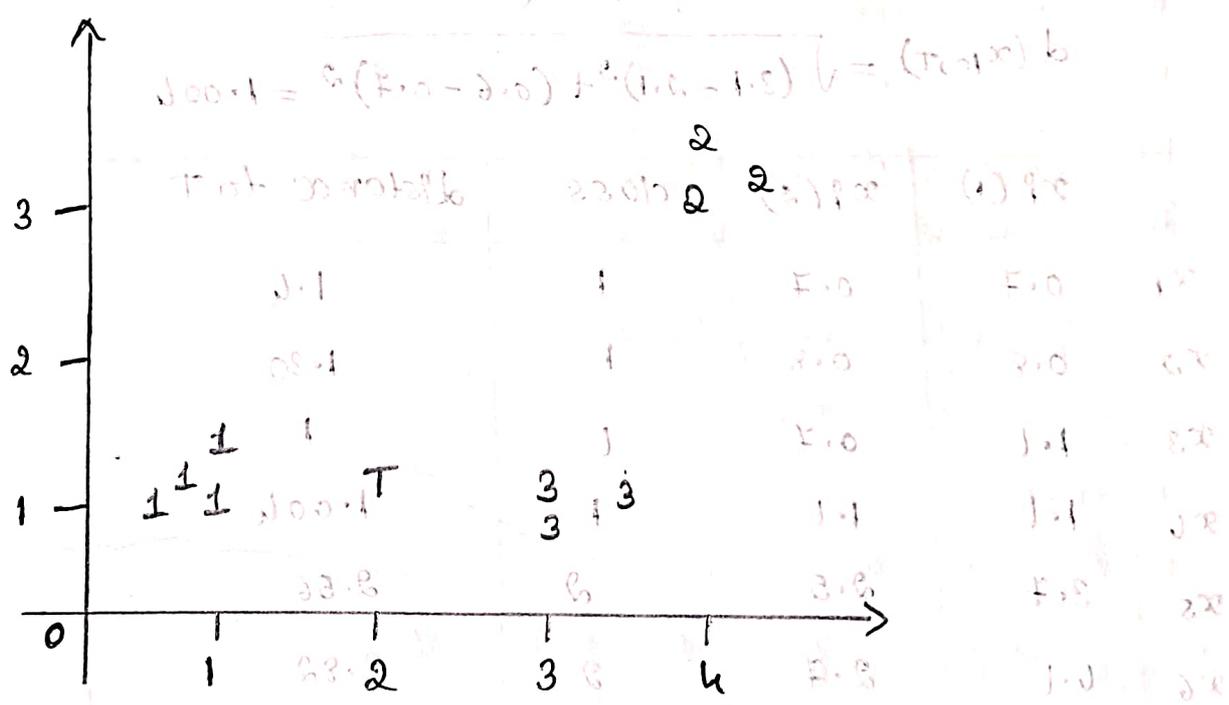
- * In k-Nearest Neighbor (KNN) algorithm we find the k-nearest neighbor for the test pattern "T" from the training data "X" then assign majority class label among k-neighbors to the "T".
- * To determine the k-nearest neighbor of "T" it is necessary to calculate the distance b/w "T" and each of the "N" training patterns.
- * The distance metric used depends on specific problem the distance metric can be Euclidean distance, Manhattan distance or cosine distance.
- * The class label of "T" is determined based on the majority class label among its k-nearest neighbor.
- * By using this method the errors in classification can be reduced when ever the training patterns are noisy.
- * The case of large data set "k" can be increased to reduce the error.
- * The value of "k" can be selected such that the lowest error occur in classification the value of "k" usually selected of odd number ($k=5$).

EX:-

- * Let training set consist of following data points with the label

- $x_1 = (0.7, 0.7), l_1 = 1$
- $x_2 = (0.5, 0.5), l_2 = 1$
- $x_3 = (1.1, 0.7), l_3 = 1$
- $x_4 = (1.1, 1.1), l_4 = 1$
- $x_5 = (3.7, 2.5), l_5 = 2$
- $x_6 = (4.1, 2.7), l_6 = 2$
- $x_7 = (3.7, 3.1), l_7 = 2$
- $x_8 = (3.4, 0.6), l_8 = 3$
- $x_9 = (3.1, 0.3), l_9 = 3$
- $x_{10} = (3.1, 0.6), l_{10} = 3$

Find the label of test pattern $T = (2.1, 0.7), l = 2$ using KNN for $k=5$



To find the class of test pattern we have to calculate the distance from test pattern to each of the training data using euclidean distance.

$$d(x_i, T) = \sqrt{(x_i(1) - T(1))^2 + (x_i(2) - T(2))^2}$$

$$d(x_1, T) = \sqrt{(0.7 - 2.1)^2 + (0.7 - 0.7)^2} = 1.4$$

$$d(x_2, T) = \sqrt{(0.8 - 2.1)^2 + (0.8 + 0.7)^2} = 1.30$$

$$d(x_3, T) = \sqrt{(1.1 - 2.1)^2 + (0.7 - 0.7)^2} = 1$$

$$d(x_4, T) = \sqrt{(1.1 - 2.1)^2 + (1.1 - 0.7)^2} = 1.07$$

$$d(x_5, T) = \sqrt{(3.7 - 2.1)^2 + (2.5 - 0.7)^2} = 2.56$$

$$d(x_6, T) = \sqrt{(4.1 - 2.1)^2 + (2.7 - 0.7)^2} = 2.82$$

$$d(x_7, T) = \sqrt{(3.7 - 2.1)^2 + (3.1 - 0.7)^2} = 2.88$$

$$d(x_8, T) = \sqrt{(3.4 - 2.1)^2 + (0.6 - 0.7)^2} = 1.30$$

$$d(x_9, T) = \sqrt{(3.1 - 2.1)^2 + (0.3 - 0.7)^2} = 1.07$$

$$d(x_{10}, T) = \sqrt{(3.1 - 2.1)^2 + (0.6 - 0.7)^2} = 1.004$$

	$x_i(1)$	$x_i(2)$	class	distance to T
x_1	0.7	0.7	1	1.4
x_2	0.8	0.8	1	1.30
x_3	1.1	0.7	1	1
x_4	1.1	1.1	1	1.004
x_5	3.7	2.5	2	2.56
x_6	4.1	2.7	2	2.82
x_7	3.7	3.1	2	2.88
x_8	3.4	0.6	3	1.30
x_9	3.1	0.3	3	1.077
x_{10}	3.1	0.6	3	1.004

	$x_p(1)$	$x_p(2)$	class	distance to T
x_3	1.1	0.7	①	
x_{10}	3.1	0.6	3	1.006
x_4	1.1	1.1	①	1.06
x_9	3.1	0.3	3	1.077
x_2	0.8	0.8	①	1.30

since the value of $k=5$ we will consider first 5 rows the majority class with in the first 5 nearest neighbors to the test pattern p_s therefore the label of test pattern p_s is 1.

Test Pattern

$$T = (2.1, 0.7), d=1$$

find the class of test pattern using nearest neighbor classifier.

Brightness	saturation	class	distance
40	20	Red	25
50	50	blue	33.54
60	90	blue	68.01
10	25	Red	16.14
70	70	blue	61.03
60	10	Red	47.16
25	80	Blue	45.27

Brightness	saturation	class	distance
10	25	Red	16.14
40	20	Red	25
50	50	blue	33.54
25	80	blue	45.27
60	10	Red	47.16

$T = (20, 35)$ class = Red

consider a scenario where a realstate agency wants to predict the price of a house based on features such as area (square feet) and age of the house (yes) explain how kNN recursion can be used. In this case given a small data set of 5 houses demonstrate the steps to predict the price of a new house using $k=3$ include the distance calculation, selection of neighbors and competition of predicted price discuss how different values of k might effect the prediction outcome.

House ID	Area (sq.ft)	Bedrooms	Age (year)	Price (in 1000s)
H1	2100	3	10	1500
H2	1600	2	5	330
H3	2400	4	20	550
H4	1410	3	15	299
H5	1800	2	8	410

New house to predict

* Area = 1900 sq.ft

* Bedroom = 3

* Age = 10 years

3. Radius distance Nearest Neighbour:-

* this algorithm consider all the neighbours within a specified distance "r" of the test pattern this algorithm can be described as.

1. Given a point T identify the data points that fall within the radius centered at "T" denoted by

$$Br(T) = \{x \in X \mid \|T - x\| \leq r\}$$

2. If $B_r(T)$ is empty the output is the majority class of entire dataset
3. If $B_r(T)$ is not empty - the output is the majority class of class of the data points with in $B_r(T)$
4. this algorithm is useful for identifying outliers the choice of value of radius (r) is important because it can effect performance of the algorithm.
5. this algorithm is more efficient than K-NN algorithm.

* for the training data set given below

$$x_1 = (0.7, 0.7), l_1 = 1;$$

$$x_2 = (0.8, 0.8), l_2 = 1;$$

$$x_3 = (1.1, 0.7), l_3 = 1;$$

$$x_4 = (1.1, 1.1), l_4 = 1;$$

$$x_5 = (3.7, 2.7), l_5 = 2;$$

$$x_6 = (4.1, 2.7), l_6 = 2;$$

$$x_7 = (3.7, 3.1), l_7 = 2;$$

$$x_8 = (3.6, 0.6), l_8 = 3;$$

$$x_9 = (3.1, 0.3), l_9 = 3;$$

$$x_{10} = (3.1, 0.6), l_{10} = 3;$$

find the class of test pattern $T = (1.7, 2)$ using radius distance nearest neighbor algorithm with the radius of 1.65

To find the data points with in the radius of 1.65 we need to find the Euclidean distance between each of the data point is the training set and test pattern. the euclidean distance

$$d(x_i, T) = \sqrt{(x_i(1) - T(1))^2 + (x_i(2) - T(2))^2}$$

Now find out distance

$$d(x_1, T) = \sqrt{(0.7 - 1.7)^2 + (0.7 - 2)^2} = 1.64$$

$$d(x_2, T) = \sqrt{(0.8 - 1.7)^2 + (0.8 - 2)^2} = 1.58$$

$$d(x_3, T) = \sqrt{(1.1 - 1.7)^2 + (0.7 - 2)^2} = 1.63$$

$$d(x_4, T) = \sqrt{(1.1 - 1.7)^2 + (1.1 - 2)^2} = 1.08$$

$$d(x_5, T) = \sqrt{(3.7 - 1.7)^2 + (2.7 - 2)^2} = 2.11$$

$$d(x_6, T) = \sqrt{(4.1 - 1.7)^2 + (2.7 - 2)^2} = 2.29$$

$$d(x_7, T) = \sqrt{(3.7 - 1.7)^2 + (3.1 - 2)^2} = 2.28$$

$$d(x_8, T) = \sqrt{(3.4 - 1.7)^2 + (0.6 - 2)^2} = 2.20$$

$$d(x_9, T) = \sqrt{(3.1 - 1.7)^2 + (0.3 - 2)^2} = 2.20$$

$$d(x_{10}, T) = \sqrt{(3.1 - 1.7)^2 + (0.6 - 2)^2} = 1.97$$

	$x^p(1)$	$x^p(2)$	label	distance to T
x_1	0.7	0.7	1	1.64
x_2	0.8	0.8	1	1.58
x_3	1.1	0.7	①	1.63
x_4	1.1	1.1	①	1.08
x_5	3.7	2.7	2	2.11
x_6	4.1	2.7	2	2.29
x_7	3.7	3.1	2	2.28
x_8	3.4	0.6	3	2.20
x_9	3.1	0.3	3	2.20
x_{10}	3.1	0.6	3	1.97

From the table the patterns located within a radius of 1.65 or x_3 and x_4 the class of these patterns are 1. therefore the text pattern assigned with

$$T = (1.7, 2), \text{ class} = 1$$

KNN Regression

* let x is set of n label data

$$x = \{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_n, y_n)\}$$

where,

x = is a data vector

y = scalar values

Here $x_p, p=1, 2, \dots, n$ is a data vector

y_p, p a scalar

The regression model needs to use x to find the value of y for a text vector T . In the case of regression based on KNN the following steps are followed:

* find k -nearest neighbor of " T " from n data vectors

let KNN are x_1, x_2, \dots, x_k

* consider y values associated with this k nearest neighbors. let them y_1, y_2, \dots, y_k

* take the average value of this y values and assign this average value as predicted value of y to " T "

the predicted value of y is $y^{\wedge} = \frac{1}{k} [y_1 + y_2 + \dots + y_k]$

Ex:- consider the data shown in table find the predicted target value for $T = (0.3, 0.4)$ using KNN regression.

	$x_p(1)$	$x_p(2)$	target y_p
x_1	0.2	0.4	8
x_2	0.4	0.2	8
x_3	0.6	0.4	12
x_4	0.8	0.6	16
x_5	1.0	0.7	19
x_6	0.8	0.4	14
x_7	0.6	0.2	10
x_8	0.5	0.5	12
x_9	0.2	0.6	10

To find k nearest neighbors we apply Euclidean distance b/w data vectors and test pattern

$$d(x_1, T) = \sqrt{(x_1(1) - T(1))^2 + (x_1(2) - T(2))^2}$$

$$d(x_1, T) = \sqrt{(0.2 - 0.3)^2 + (0.4 - 0.4)^2} = 0.1$$

$$d(x_2, T) = \sqrt{(0.4 - 0.3)^2 + (0.2 - 0.4)^2} = 0.223$$

$$d(x_3, T) = \sqrt{(0.6 - 0.3)^2 + (0.4 - 0.4)^2} = 0.3$$

$$d(x_4, T) = \sqrt{(0.8 - 0.3)^2 + (0.6 - 0.4)^2} = 0.52$$

$$d(x_5, T) = \sqrt{(1.0 - 0.3)^2 + (0.7 - 0.4)^2} = 0.76$$

$$d(x_6, T) = \sqrt{(0.8 - 0.3)^2 + (0.4 - 0.4)^2} = 0.5$$

$$d(x_7, T) = \sqrt{(0.6 - 0.3)^2 + (0.2 - 0.4)^2} = 0.36$$

$$d(x_8, T) = \sqrt{(0.5 - 0.3)^2 + (0.5 - 0.4)^2} = 0.223$$

$$d(x_9, T) = \sqrt{(0.2 - 0.3)^2 + (0.6 - 0.4)^2} = 0.223$$

$x_i(1)$	$x_i(2)$	Target y_i	distance
0.2	0.4	8	0.1
0.4	0.2	8	0.223
0.6	0.4	12	0.3
0.8	0.6	16	0.52
1.0	0.7	19	0.76
0.8	0.4	14	0.5
0.6	0.2	10	0.36
0.5	0.5	12	0.223
0.2	0.6	10	0.223

x_1	0.2	0.4	8	0.1
x_8	0.5	0.5	12	0.223
x_2	0.4	0.2	8	0.223
x_9	0.2	0.6	10	0.223
x_3	0.6	0.4	12	0.3

If $k=1$ then the Nearest neighbor is x_1
 If $k=5$ then the 5 Nearest neighbor the text
 pattern x_1, x_8, x_2, x_9, x_3 and the labels of these
 neighbors are 8, 12, 8, 10, 12 the Predicted target value

$$y^{\wedge} = \frac{1}{5} (8 + 12 + 8 + 10 + 12)$$

$$= \frac{1}{5} (50) = 10$$

the target value

$$T = (0.8, 0.6), y^{\wedge} = 10$$

Performance measure

Performance of classifier:-

i) Classification Accuracy

let n the total no. of patterns that are classified
 by a classification algorithm. let nc is the no. of
 correctly classified patterns. Classification Accuracy

$$= \frac{nc}{n}$$

ii) confusion matrix:-

it is used to store to result simplified matrix formate
 this data is further simplified to analyze the
 result.

Ex :-

let they are three classes c_1, c_2 and c_3 there are
 200, 100 and 50 patterns from classes c_1, c_2 and c_3
 respectively. let a classifier classify this patterns
 into three classes as shown in the table.

True / Predicted	c_1	c_2	c_3
c_1	180	15	5
c_2	5	85	10
c_3	3	2	45

* The first row of table shows that 200 patterns from C_1 are classified as

* 180 are assigned to C_1 , 15 are assigned to C_2 and 5 are assigned to C_3 . Similarly second row for 100 patterns from C_2 and third row shows the assignment of 50 patterns from C_3

* The classification Accuracy is obtained by considering the correctly classified patterns that is 180 patterns from C_1 , 85 patterns from C_2 , and 45 patterns from C_3 therefore no. of correctly classified patterns.

$$= 180 + 85 + 45$$

$$= 310$$

total no. of patterns = $C_1 = 200$
 $C_2 = 100$
 $C_3 = 50$

350

the classification Accuracy = $\frac{n_c}{n} = \frac{310}{350} = 0.8857$

the percentage classification Accuracy is 88.57%

* It is possible to simplify the confusion table with respect to C_1 the compact table is of size 2×2

True \ Predicted	C_1	\bar{C}_1 (Not C_1)
C_1	TP 180	FN 20
\bar{C}_1	FP 8	TN 162

* The first row of the table shows that 180 patterns of C_1 are correctly classified therefore these are called True Positives (or) $TP \cdot TP = 180$

* The second column and first row has 20 patterns these patterns are from C_1 that are not classified

as c_1 they are assigned to \bar{c}_1 (Not c_1) they are called false negative (or) FN. $FN = 20$.

* the first entry in (second row) shows that 8 patterns from other class are assigned to c_1 so they are called false positive (or) FP. $FP = 8$.

* The second column in the second row shows that 162 patterns from \bar{c}_1 (not c_1) assigned to \bar{c}_1 they are called true negative (or) TN. $TN = 162$ the confusion matrix is given by

True / Predicted	c_1	\bar{c}_1
c_1	TP	FN
\bar{c}_1	FP	TN

this matrix is used to find following measures

* Precision :-

It is the ratio of TP to total no. of predicted positives $(TP + FP)$

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{180}{180 + 8} = \frac{180}{188} = 0.9574$$

* Recall :-

It is the ratio of TP to total no. of positives in the training data $(TP + FN)$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{180}{180 + 20} = \frac{180}{200} = 0.9$$

* F1-score

It is the harmonic mean of precision and recall

$$f_1\text{-score} = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}$$

$$= \frac{2}{\left(\frac{1}{0.9574}\right) + \left(\frac{1}{0.9}\right)}$$

$$= 0.9278$$

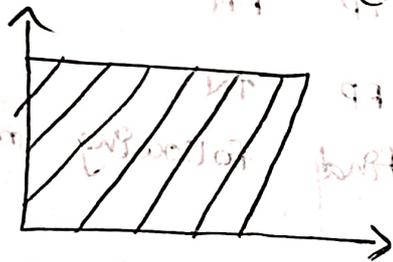
* True Positive Rate or TPR

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

* False Positive Rate or FPR

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

* Area under the curve (AUC)



AUC is obtained by plotting the graph b/w FPR on X-axis and TPR on Y-axis calculate the area under this is also called the receiver of (receiver operator characteristics) or (AUC)

Performance of Regression algorithm:-

1. Mean Squared Error (MSE)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

where n is no. of patterns

y_i is the target value for i th pattern

\hat{y}_i is the predicted value of i th pattern using regression model.

2. Mean absolute Error: - (MAE)

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

It is average of difference b/w target and predicted value.

Ex:- If $TP=2$, $FN=0$, $FP=1$, $TN=1$ find Precision, Recall, f_1 -score, True Positive Rate, False Positive Rate and classification accuracy.

True / Predicted	C_1	\bar{C}_1
C_1	$TP=2$	$FN=0$
\bar{C}_1	$FP=1$	$TN=1$

Precision:-

$$\frac{TP}{TP+FP} = \frac{2}{2+1} = \frac{2}{3} \Rightarrow 0.6666$$

Recall:-

$$\frac{TP}{TP+FN} = \frac{2}{2+0} = \frac{2}{2} = 1$$

f_1 -score:-

$$\frac{2}{\frac{1}{0.6666} + \frac{1}{1}} = \frac{2}{1.5 + 1} = \frac{2}{2.5} = 0.8$$

True Positive Rate:-

$$TPR = \frac{TP}{TP+FN}$$

$$= \frac{2}{2+0} = \frac{2}{2} \Rightarrow 1$$

False Positive Rate:-

$$FPR = \frac{FP}{FP+TN} = \frac{1}{1+1} = \frac{1}{2} \Rightarrow 0.5$$

Classification Accuracy

$$\frac{nc}{n}$$

$$nc = TP+TN = 2+1 = 3$$

$$n = TP + TN + FP + FN$$

$$= 2 + 1 + 1 + 0$$

$$= 4$$

$$\frac{nc}{n} = \frac{3}{4} = 0.75$$

The Percentage of Accuracy is 75%

2) $TP = 2, FN = 0, FP = 0, TN = 2$

calculate classification

Performance measures

True/Predicted	CI	\bar{C}_I
CI	TP = 2	FN = 0
\bar{C}_I	FP = 0	TN = 2

Precision:-

$$\frac{TP}{TP + FP} = \frac{2}{2 + 0} = \frac{2}{2} = 1$$

Recall:-

$$\frac{TP}{TP + FN} = \frac{2}{2 + 0} = \frac{2}{2} = 1$$

F1 score:-

$$\frac{2}{\frac{1}{Precision} + \frac{1}{Recall}} = \frac{2}{\frac{1}{1} + \frac{1}{1}} = \frac{2}{2} = 1$$

True Positive rate:-

$$TPR = \frac{TP}{TP + FN} = \frac{2}{2 + 0} = \frac{2}{2} = 1$$

False Positive rate:-

$$FPR = \frac{FP}{FP + TN} = \frac{0}{0 + 2} = \frac{0}{2} = 0$$

classification Accuracy

$$\frac{nc}{n}$$

$$= \frac{3}{4} = 0.75$$

$$n_c = TP + TN = 2 + 2 = 4$$

$$n = TP + TN + FP + FN$$

$$= 2 + 2 + 0 + 0$$

$$= 4$$

$$\frac{n_c}{n} = \frac{4}{4} = 1$$