Department of Artificial Intelligence and Machine Learning

II B.Tech II Sem

Machine Learning

UNIT-1

Introduction to Machine Learning

By

Mrs.K.Reddemma,

Assistant Professor,

Department of AI&ML.

Machine Learning

T-TINU

Introduction to Machine learning

- 1 Evolution of ML
- @ paradigms: learning by Rote, Introduction, Reinforcement, etc.
- 3 Types of Data, Stages in ML
- A Data Acquisition, Feature Engineering, Reportsentation
- 5 Data Acquisition, Model selection, Learning, Evaluation, prediction
- 6 Search & leavining, Data sets

UNIT-11

Nearest Neighbor - Based Models

- 1 proximity & Distance Measures
- 2 Non-Metric Similarity Functions
- 3 Binary pattern proxumity
- (A) classification Algorithms (Distance-Based)
- (5) K-NN classifier, Radius Distance NN Algorithm
- 6 K-NN Regoussion, periformance of classifiers & Regoussion Algorithms.

Introduction to Machine Learning

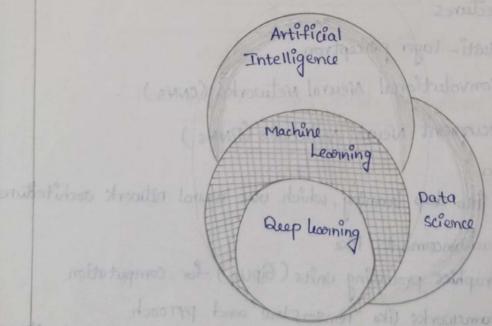
What is Machine Learning:

Machine Lewining is a branch of Artificial Intelligence that develops algorithms by Leaving the hidden patterns of the datasets used it to make peredictions on new similar type data, without being explicitly programmed for each task.

Totaditional machine leavining combines data with statistical tools to possolict an output that can be used to make actionable insights.

Machine learning is used in many different applications, From image and speech recognition to natural language processing, recom mendation systems, fraud detection, portfolio optimization, automated wisk and so on.

Machine learning models on also used to power autonomous vehicles, dorones and vobots, making them more intelligent and adaptable to changing environments.



Venn diagram for AI, ML, Deeplearning and Data Science

Evolution of Machine Learning

Machine Learning

Machino learning is the process of learning models that can predict outcomes based on data.

It involves prodicting classes (or) regression values for given

Early AI and symbolic systems:

AIs beginnings focused on symbolic logic

- * General problem solver (Gips) and programs using prolog were
 - * Relied on rule-based systems and automated reasoning.
 - * AI wor often Seen on "logic-bared"

shift to Data-Driven learning:

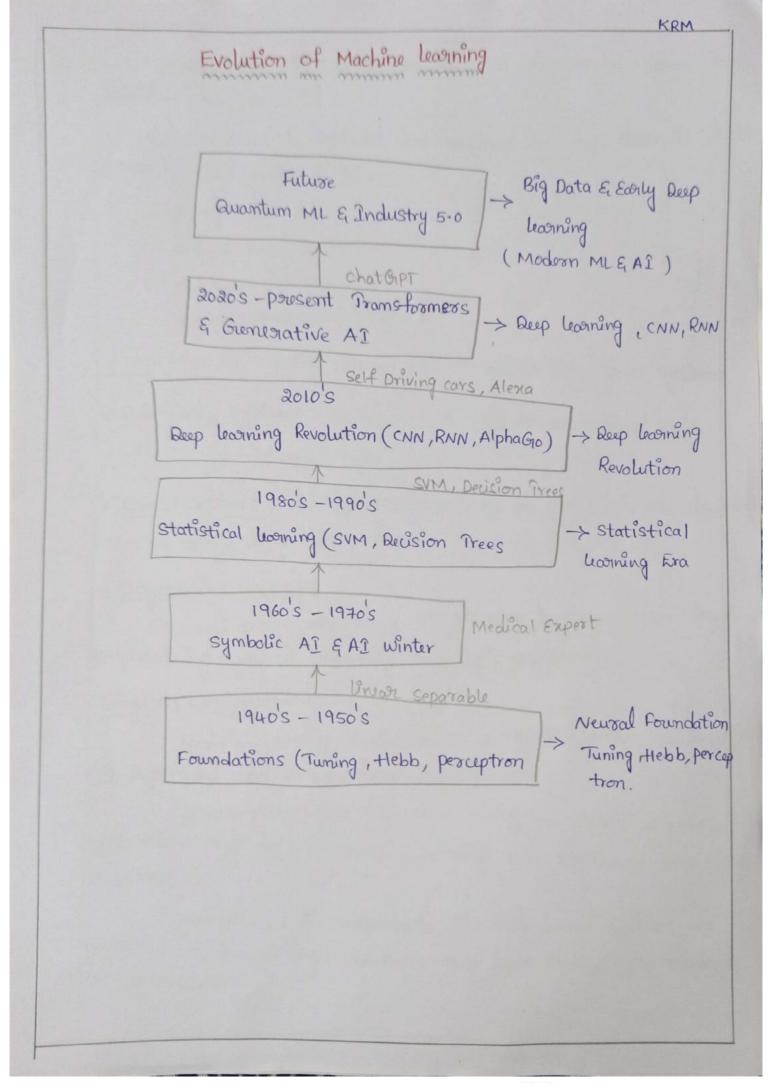
The 1980's saw challenges in symbolic A? due to complex real-world problems.

- * Focus shifted to data-driven approaches with the rise of neural network architectures
 - * Multi- layer perceptron
 - * convolutional Neural Networks (CNNs)
 - * Recurrent Neural Networks (RNNs)

Deep Learning Eva:

- * AI Evolved into deep learning, which uses neural network architectures.
- * Enabled by advancements like
 - 1. Graphics processing units (Gipus) for computation
 - 2. Frameworks like TensorFlow and pyrrooch

Impact of AI Today: Deep bearing, especially convolutional neural networks and generative models, influences vociour scientific and engineering fields globally.



conventional AI

Conventional AI is a type of AI that focuses on narrow, taskspecific intelligence.

Conventional As systems are torained on large datasets to loom patterins and make decisions.

- * Logic for AI
- * Rule-Based Systems
- * Expert Systems
- * programming languages.
- * Logic for A? : Using formal logic to design intelligent systems.
- * Rule-Based systems: systems that vely on predefined vules. Ex: if-then logic.
- * Expert systems: AI systems that reday on poor mimic the decisionmaking ability of a human Expert.

* programming Languages:

Lisp and poolog were early programming languages used for AI, emphasizing symbolic reasoning and logic programming.

Diagram Explanation:

Machine Learning is decicted are a key area within AI, Supported by supervised and unsupervised learning tasks.

pattern recognition and data mining one shown as confical applications that vely on machine learning and structured data from databores.

conventional AI components like rule-bard systems are persented on foundational concepts that have evolved into modern AI approaches. char internation sufers to interface (or) systems that allow

Main components:

1. Machine Learning (ML):

A subfield of AI focused on enabling systems to learn from data. It is further divided into

* learning from Examples - (Supervised learning) involves tasks like classification and Regoussion

classification: predicting categorical outcomes

Ex: Mails spam (or) Not spam

Regardsion: predicting continuous values

Ex: House posses

* learning from observations-(unsupervised Learning)

Leconing from observations involves clustering clustering meant refers to the Grouping similar data points.

Ed: Customer Segmentation

2. pattern Recognition:

Focuses on identifying patterns in data, often linked to tasks in machine borning such ou image (or) speech recognition.

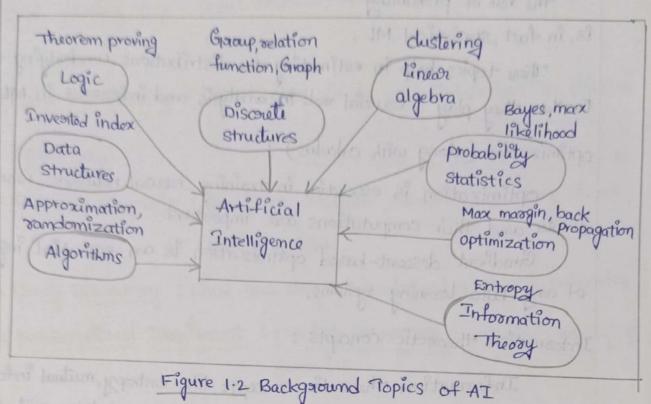
3. Data Mining

The process of efficiently organizing and analyzing large datasets (databous) to uncover useful patterns and insights.

This fooms the backbone of AI by providing structured data -for analysis.

Databases and Human-Machine Interaction

Databores one confical for storing and managing data, while human-Ma chine interaction refers to interfaces (or) systems that allow humans to interact with AI effectively



Data structures and algorithms are basic to both conventional and Cussent AI systems. Logic and Discorete structures played an important

role in the analysis and synthesis of conventional AI system!

In ML, we deal with vectors and vector spaces and these topics one better appreciated through linear algebra. The data input to an ML sy-Stem may be viewed on a matrix, popularly called the data matrix.

If there are n data items, each represented ar an L-dimensional vector, then the corresponding data matrix A is of size n xL. Linear algebra is useful in analysing the weights associated with the edges in a neural network. Matrix multiplication and eigen analysis are important in initializing the weights of the neural network and in weight updates. It can also help in weights normalization.

The whole activity of clustering may be viewed ar data matrix factorization.

* probability and statistics:

The role of probability and statistics need not be explained on ML is, in fact, statistical ML.

There topics help in estimating the distributions underlying the data. Further, they play a coucial role in analysis and inference in ML.

* optimization (along with calculus):

optimization is essential in training neural networks where gradients and their computations are important.

Gradient descent-bored optimization is an essential ingredient of any Data learning systems.

* Information theoretic concepts:

Information theoretic concepts like entropy, mutual information and Kullback-Leibler divergence one essential to understand topics Such or decision tree classifiers, feature selection and deep neural networks.

- * self-driving cars: Detect objects, follow lames, poreducts trafic movements.
- * Don've Monitoring : Detect donowsiness (or) distraction in drivers

5 Healthcore !

- * Diseases Riagnosis: poudicts diabetes, Cancer, and more from patient data:
- * Doug Discovery: Find new drug molecules using neural networks.
- * personalized Treatment: predict responses to medications using genomics.

6. Finance :

- * Fraud Retiction: Betict unusual transactions (or) Hacking
- * Algorithmic Totading: Make Trading decisions using historical data.
- * coudit scoring: Assess risk levels for loan applications.

7. E-Commerce & Recommendation systems:

- * products Recommendations: Suggest îtems bored on user behavior. Eg: Amazion, Net-flix
- * Search Ramking: Improve relevance in product search.

& Gaming & Simulation:

- * Grame AI: Bots that learn to play games better than humans Eq: AlphaGio
- * Simulation of Real-World Scenarios: Weather, traffic, robotics

- 9. Robotics
 - * perception: Robots understand their environment using sensors and vision.
 - * Motion planning: Navigate around obstacles.
 - * Human-Robot Interaction:

use voice, vision, and touch to interact with people.

10. Big Data and Anomaly Detection:

Detect forand, network attacks, (or) equipment failures using large-scale pattern analysis.

Machine learning Basics:

- 1 Supervised And Unsupervised learning
- 1) Superivised Learning:

Supervised learning Algorithms are those that learn a mapping from inputs of to outputs y using a labeled dataset. The dataset contains examples of inputs along with the correct outputs (labels), often provided by a human supervisor.

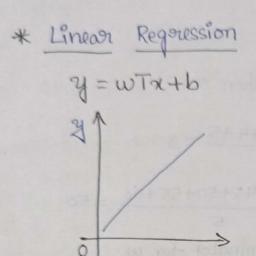
Even when outputs of one generated automatically, the team superivised learning still applies.

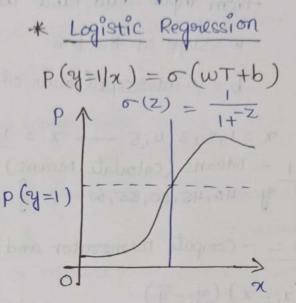
(0V)

Superivised learning is a type of machine learning where a model learns to map inputs x to outputs y using a labelled dataset. In Deep learning this mapping is learned using deep neural networks that automatically discover useful features from data.

Binoxy Cross-Entropy (log loss) = - [ylog(y)+(1-y) log(1-y)

Ex: Spam detection, disease prediction (yes/No) Sentiment classification (positive / Negative)





Example problem too linear Regoussion:

To predict the salary of an employee bared on their years of expenience. Find aline y=wx+b.

Years of Experience (n)	salary (y in \$ 1000)
I	40
2	45
3	50
4	55
5	60

Answer: (0) 1. Use least squares (or) gradient descent to fit the best line. 2. After Braining data from input data table w=5, w = slope of the line b = y - Intercept (values of y when x = 0) $x = 1,2,3,4,5 \rightarrow x = 1+2+3+4+5 = 3$ step 1 - Means (calculate means) 5 y=40,45,50,55,60 -> y=40+45+50+55+60 = 50 step 2 - Compute numerator and denominator for w Σ (x; - \) (y; - \) =(1-3)(40-50)+(20-3)(45-50)+(3-3)(55-50)+(5-3)(60-50) = (-2)(-10) + (-1)(-5) + 0 + (1)(5) + (2)(10)= 20+5+0+5+20 $\sum (9^{\circ}-7)^{2}=(-2)^{2}+(-1)^{2}+0^{2}+1^{2}+2^{2}$ = 4+1+0+1+4 \[\(\(\alpha \)^2 = 10 y = wx+b [: (x=6 It is a new input to posedict salony

So y = wx + b [: (x = 6) It is a new input to powdict salony) y = 5x + 35 $y = 5x + 35 \Rightarrow 30 + 35 \Rightarrow 65$

Stages in machine learning!

The process of building a machine learning (ML) System involves several critical steps, as highlighted figure 1.6. The steps emphasize data in the form training, validation, and test datasets Below is 01 explanation of each step. an

Application domain -> Data acquistion , 6 Extertainment and media @ Finance- Fraud detection Feature engineering = preprocessing + Ex . Healthcare and Medicine

Representation (3) E-commerce

(Natural Longuage processing (NLP) Model Selection -> choose a model <-- Domain knowledge

model learning - train the model - Training data

Model evaluation -> varidate the model (varidation data) cours - 19 patients from healthy

model prediction -> Learn the model -- Test data

dividuals. Explain the model -- Expert feedback model explanation ->

Stages of a machine learning System:

1. Application:

3

3

4

3

0

0

0

3

3

9

3

understand the probelm or context where

ML System will be applied. the

Examples:

Identifying adults vs. children

Diagnosing covid-19 through specific metrics (eg. temperature, chest congestion)

2. Data Acquistion:

Data is collected based on the application domain

Key conditions:

The type of data depends on the problem being solved.

Examples:

Height and weight to distinguish adults from

children.

Body temperature and chest congestion to differentiate covid-19 patients from healthy individuals.

3. Feature Engineering:

This combines data preprocessing and data representation to ensure the ML Model can use the ML data effectively.

6

6

6

6

Steps in Feature Engineering!

Data preprocessing: prepare raw data by

handing issues like:

Missing Values

Noise or irrelevant information.

4. Model Selection:

choose an ML model that first the probelm.

considerations:

99999999

9

3

9

9

9

0

0

0

0

9

3

3

3

9

3

the complexity of the probelm.

The need for domain knowledge to enchare model Selection.

5. moder learning:

Train the chosen model using the training data.

expert elections has velines the masei

the model learns patterns and relationship in the data.

6. Model Evaluation:

Validate the models performance using vali-- dation data.

This Step ensures the model generalizes well and avaids overfiting.

7. Model prediction:

the trained model is tested on unseen test data to predict outcomes and evaluate its real-world performance.

8. Model Explanation:

Interpret and Explain the model's predictions. Expert feedback can help refines the model.

Data acquistion: Importance: Data collection is the first step and heavily depends on the probelm domain.

Examples:

For distinguish adults from Children. metrics like height or weight are Sufficient.

For diagnosing covID - 19:

Measurements like body temperature and chest congestion are more relevant than height or weight.

Feature Engineering:

Feature Engineering is a two-part process.

1. Data preprocessing:

Handless issues in raw data such as:

Missing values: Chaps in the data that must be filled or imputed.

outliers: Extreme values that can distort the

Moise: unwanted variations or errors in the data.

2. Data Representation:

Transform raw data into a format Sutiable for the ML model.

Examples:

9

9

9

9

9

Encoding categorical variables (e.g., converting "male (Female" into numerical values) Normalizing numerical data to Standardize Scales.

Datasets in Mc:

1. Training data:

used for training the model to learn patterns.

2. validation Data: " De la suzzi zzalbandi

Helps fine-tune the model and validate its performance during training, but again to bourge outliers: Extreme values that can distort the

3. Test Data:

Evaluates the model's final performance on on bosnoones : solo unseen data.

22/1/25 Learning from observations: * observations are also instances like examples but they are different beaute observations need not be labelled. * In this case, we cluster (or) group the observations into a Smaller number of goloups. * Such grouping is performed with hulp of a clustering algori. thm that assigns simillar patterns to the same group/cluster. * Each cluster could be represented by it's centroid (or) mean * let M, M2 --- sop be P elements of a cluster * Controld of the cluster is defined by 2 thing 1 2 90; P 1=1 2 will be to the transport (mean to) 3 Formula for controld of a cluster is

Centrold = $\frac{1}{p} \sum_{i=1}^{p} \pi_i$ 3 5 0 0 0 0 0 0 0 0 0 0 0 0 9 9 9999999999 8 Examples of Handwritten digits labelled 0 To 9

Handworften Rigits:

* This figure shows examples of handwarten digits lebeled from

* Each row contains multiple examples of a single digit. For instance, the first row contains several samples of "o" The Second row contains "1" and so on "9"

* The labels (0-9) are used on ground truth for supervised learning tasks.

Centroids in clustering

The text introduces the concept of a certifoid is the average (or mean) reportesentation of a cluster of data points.

For given cluster, if there one p elements x, x2--xp the centroid is calculated on.

* In this figure, handwortten digits are gorouped into clusters board on their class labels

En: Digits 0, 1 and 3 are clustered Separately.

* The Centroids of each cluster are represented or

There centroids are averages computed from all the samples within each class.

Sum of coordinater Sum of all the x-coordinates of the data points. Sum of n-coordinates = = = n; * Rivided by the Total Number of points. Average for each dimension by dividing the summed values by the total number points (p). Controid x-coordinate = = x; 1. compute controid in 20: (2,u) (4,5) (6,7) Sum of 91-coordinates 2+4+6=12, 3+5+7=15 2. Divide by the Number of points: Total no. of points (p) = 3 controid x-coordinate 12=4 controid y-coordinate 15 = 5 Corrtroid is (4,5) Centroid: A controid is a point in the feature space where the coordinater one the average of the coordinater of all points in the cluster. Centroid (C) = 1 \(\sum_{i=1}^{p} \) \(\text{2} \) \(\text{2} \)

Scanned with OKEN Scanner

casine Similarity

Casine similarly is the measure of similarity between two non-zero vectors widely applied in many machine learning and data analysis applications.

The Formula for cosine similarity is $\cos(0) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$

- * A.B is the dot product of vectors A and B
- * 11A11 and 11B11 one the magnitudes of vectors A and B
- similarity between two vectors in the machine learning
- * It's often used in text analysis and recommendation system
- How is cosine similarities used in machine learning?

 In the realm of machine learning, cosine similarity plays a pivotal role in embansing text analysis processes. By measuring the similarity between textual data points, this metric emables efficient document clustering and information retrieval systems.
- what is the cosine algorithm in orachine learning?

 cosine similarity is used as a metric in different machine

 learning algorithms like the for determing the distance between the

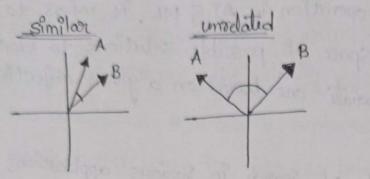
 neighbors, in recommendation systems, it is used to recommend

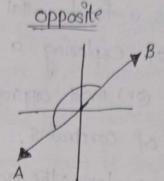
 neighbors, in recommendation systems, it is used to recommend

 movies with the same similarities and for texual data, it is

 movies with the same similarity of texts in the document.

The smaller the angle between the two vectors, the more similar they are to each other.





The cosine similarity is described mathematically on the division between the dot product of vectors and the product of the euclidean nooms (or) magnitude of each vector.

$$||a|| ||b||$$

$$||a|| ||b||$$

$$||a|| = \sqrt{a_1^2 + a_2^2 + a_3^2 + - + a_n^2}$$

$$||b|| = \sqrt{b_1^2 + b_2^2 + b_3^2 + - + b_n^2}$$

Where, a and b one vectors in a multidimensional space.

Since the cos(6) value is in the range [-1, 1]

-1 value is indicates strongly opposite vectors no similarity

* o indicater independent vectors (or orthogonal).

* I value is indicate high similarity between the vectors.

Sewich and learning in AIEML:

Search is a fundamental operation in AIG ML. It refers to the process of emplosing a space of possible solutions to identify the best (or) most appropriate one board on a given objective (or) set of contrains.

We explore the vole of Sewich in Various applications of computer science and it's importance in ML and AI and along with relevant examples.

Ex: 1. problem-solving _ Sudoka (or) finding shortest path inagraph =

2. Theorem-proving-Automated systems like prolog use search to verify logical proofs. Robots use search algorithms like

(BFS) (or) (DFS) to plan pla paths in navigation.

depth-first-search

Breadth-First-Search

Examples of BFS algorithm:

Adjacency lists

A - BID

B-CIF

D-DE

E-BIF

F-A

G-E

Matching: Matching is an impostant activity in ML. It is used in both supervised learning and in learning from observations.

Matching is cardied out by using a proximity measure which was can be a distance / dissimilarity measure or a similarity measure.

Two data items, a u and v, represented on L-dimensional vectors, match better when the distance between them is smaller (or) when the similarity between them is larger.

A popular distance moorure is the Euclidean distance and a popular similarity measure is the cosine of the angle between vectors.

Fuclidean Distance:

Euclidean Distance 18 the straight-line distance between two points in a muttidimensional space. It is calculated using the formula

Euclidean Distance =
$$\begin{bmatrix} n \\ \frac{1}{2} = 1 \end{bmatrix} (x_i^2 - y_i^2)^2$$

using this formula to measure the absolute distance between two vectors (or) points.

Example: let's consider two points in a 2D space

$$A = (112)$$
 and $B = (416)$

Euclidean Distance =
$$\sqrt{(4-1)^2 + (6-2)^2}$$

= $\sqrt{(3)^2 + (4)^2}$ = $\sqrt{9+16}$ = $\sqrt{25}$

If the Euclidean distance is small, the points one closer together.

2. cosine similarity:

Cosine similarity measures the cosine of the angle between two vectors. It indicates how similar two vectors are in terms of direction, regardless of magnitude.

quies similarity proxime is the collins

Example: let's consider two vectors

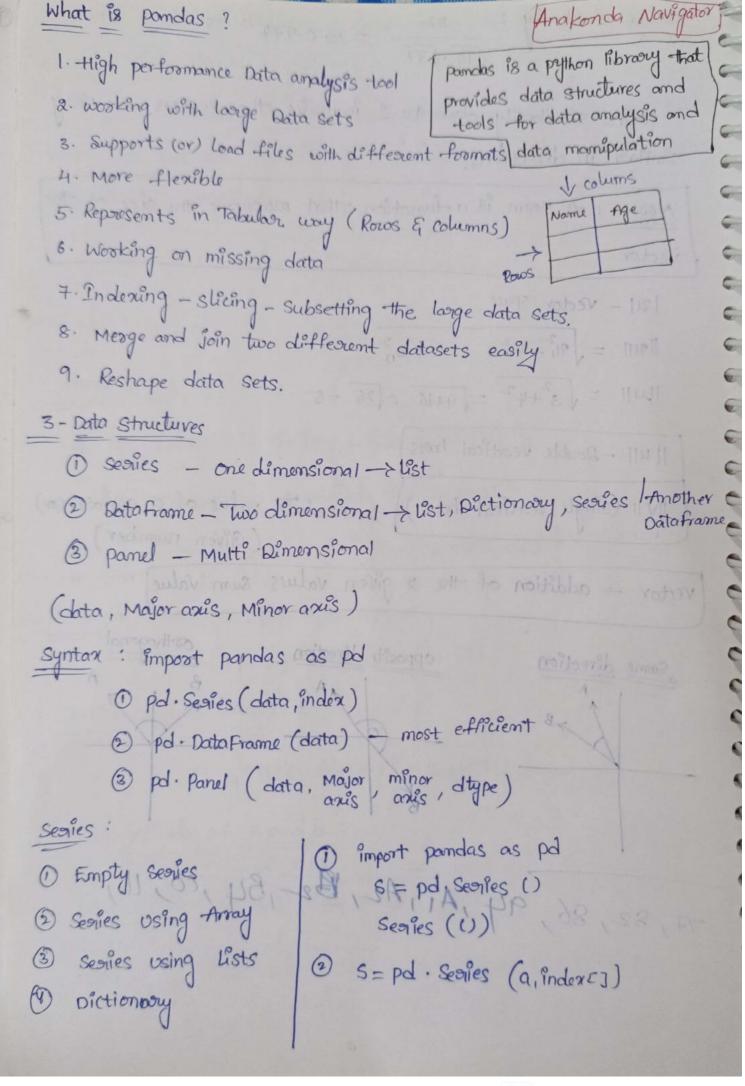
I. Dot product (A.B):

$$(1.4) + (2.5) + (3.6) = 4 + 10 + 18 = 32$$

$$|A| = \sqrt{1^2 + 3^2} = \sqrt{1 + 4 + 9} = \sqrt{14}$$

$$181 = \sqrt{4 \times 2 + 5 \times 2 + 6 \times 2} = \sqrt{16 + 25 + 36} = \sqrt{77}$$

√14·√77 = ≈ 0.974 cosine similarity: cosine similarity close to I indicates the vectors are very similar in directions Noom: A noom is a function that measures the size of the vector [21] - Vector space 11211 = 21+42+ -- + un | | | = \32+42 = \9+16 = \25 = 5 11 411 - Double resitical books 11VII (single vertical boors) > absolute value of a scalar (or) (Given number) vector - addition of the 2 given values sum value opposite direction appropriate contraction Same direction Man ales (data, and) (3) pd. Parul (data, Major, major, dupe) 1 import pandas as pa 94 (d) 201802 B2 184, (8) 111, 100 77,82,86,94(1) 201802 Sessies using Lists (a, Index (a, Index (a) orctionson



impost pandar on pd data = [10, 20, 30, 40] pandas servies with labels = ['A', 'B', 'c', 'p'] Seglies = pd. Seglies (data, index = labels) point ("pandas series with labels:") paint (series) D Rictionary impost pandas con pol data-dict = {[A':10, 'B':20, 'C':30, 'D':40], A Senies = pd. Senies (data_dict) posint ("pandas seosies from Dictionary:") dtype: int 64 point (Senies) olp: Name Age City per-tramence con unsum deta. pandas data frame import pandas as pd Jata = {

'Name': ['Alice', Bob', 'charlie', 'David'], 3 data = { : [24,27,32,29], o mano 2010) city : [New York, los-Angeln', chicago', 'Houston', phonin'] df = pd. Data Frame (data) posint (" pandas DataFrama: ") point (df)

Learning by Rote:

* Learning by Rote involves memorization in an effective manner

* It is a form of lowning.

Ex: That is a foom of learning numbers in elementary schools. where the alphabets and numbers are memorized.

* Memorizing the simple additions and muttiplications, tables are also examples for the Rote learning.

31/1/25 11, 64, Absentes

In Machine learning (ML), learning by Rote refers to a model memorizing training data instead of generalizing from it. This means the model recalls specific examples rather than identifying underlying patterns, leading to poor performance on un seen data.

* learning by memorization * We store the computed values so that we do not have

to recompute them later.

En: tables, square and cube, Additions, muttiplication * It is a simplest type of learning-without any modifications is simply copied into the knowledge bare.

(" consider Diderions :")

19 - pd. Data France (deta)

Data sets in ML is a collection of data used to train and test Data Sets machine loanning models. Inductive - harning by Inductive; Inductive teaching and learning activities promt students to generate knowledge through manisy reasoning, observation (or) experience, rather than receive it through direct instruction Ex: images, objects, data There are 4 booming styles - 1. Visual 2. auditory 3. Reading 4. Wouting There are two types of Induction process 1. Mutual 2. Self Introduction. Inductive Reasoning Pattern specific Recognition observation To identify the difference blu 2 values and products the next value number.

Learning by Induction:

The process of looming general patterns (or) rules from a set of specific examples.

Someone is formary introduced into a place of work (or) an organization.

There are 4 type of data - Reading Writing visuals

Audios

classification: Hand wouthen digites - (0,9)

Real-world Applications

- Omedical Diagnosis
- 3) spam Detection
- 3) face Recognition
- (object Detection

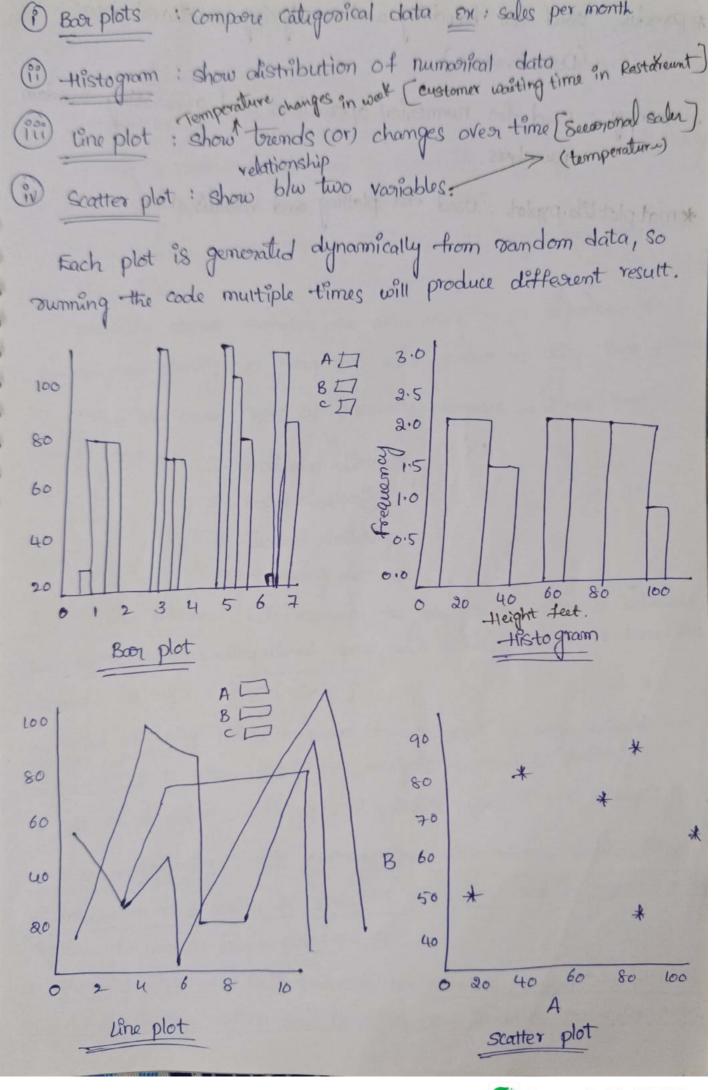
Regoression: Stock markets

Real -woold -Applications

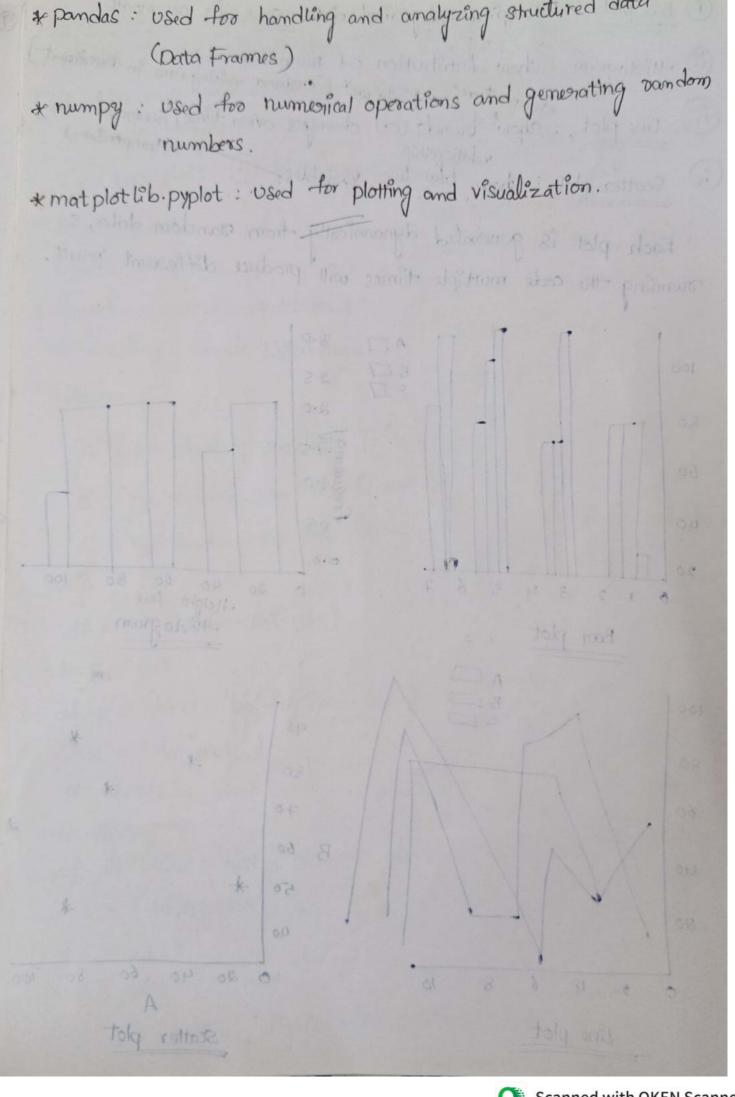
- 1 House price prudiction
- @ Weather Fore casting
- 3 Saler prediction
- Energy consumption Estimation

Benchmark dataset for testing Data Sets Used Pata sets for classification hearing model! I MNIST Handwortten Rigits Data set - used for Handwortten digit Recognition (0-9) 28x28 = 784 (pixels per 2. Fashion MNIST Data set image) A modern alternative to MNIST used for Each fixel value o torss black white clothing item classification. provided by zalando, an online fashion Training sets - 60,000 images Test lets = 10,000 images retailer. (Rang 0 70255) classes: 1.T-shot /top Modified National Institute of standards 2. Trouser and Technology MIST [MNIST] 3. pullarer Grayscale - 3 (Red, green, Blue) 4. Dress 5. Coat 83,86 LEII 6. Sandal of momory usage less 7. shist * processing speed - faster 8. Sneaker 9. Bag 10 Ankle boot € face Recognition - 40 people ×10 images - 64 × 64 gray scale - 569 samples - 30 features -cell size. Boston Housing - 506 - 13 features ted boils) tob 9 (5) Airline passongers - time servies forecasting monthly record. 6 Australian weather - weather preduction told and the H. plot (Kind = 11ing , -figsize = (815), title =

Designing and coealing easy-to-communicates cut pandas library: visualization using graphs and plots. @ Write a program which use pandas inbuilt visualization to plot (1) Bor plots (ii) Histograms (iii) line plots (iv) scatter plots. -following graphs: - Scatter plot df. plot (Kind = 'scatter', x='A'; impost pandas as pol Y='B', figSize = (8,5), title = impost numpy as np import matplotlib. pyplot as pit "Scatter plot") # creating a Sample Data Frame ptt. show () data = § A B, c random integers) (1,100,10) A: np. random. randint (1,100,10), : np. random. randint plot: Sequence of the (1,100,10), : np. random · randint data framing df = pd. Data frame (data) =# Bor plot of plot (kind = bin, figsize = (8,5), > plot size title = 'Bor plot") plt. show () # Histogram df. plot (kind = hist', bins = 10, alpha = 0.7, figsize = (8,5), title = "Histogram") plot. show () # Line plot of plot (kind = 'line', figsize = (8,5), title = "line plot") plt. Show ()



Scanned with OKEN Scanner



Department of Artificial Intelligence and Machine Learning

II B.Tech II Sem

Machine Learning

UNIT-2

Nearest Neighbor-Based Models:

By

Mrs.K.Reddemma,

Assistant Professor,

Department of AI&ML.

The Jaccard similarity fromula

The Jaccard similarity (or Jaccard Index) measures the similarity blue

two sets. It is defined on J (A,B) = [AnB]

* ANB = Number of elements common in both sets.

* AUB = Number of elements in both sets.

Example:
$$A = 1,2,3,4$$
 $B = 3,4,5,6$

* ANB = 3,4 (common elements) -> size = 2

* ADB = 1,2,3,4,5,6 (all unique elements) -> size =6

Formula
$$J(A_1B) = \frac{2}{6} = 0.33$$

4 Hamming Distance

The Hamming distance measures the number of positions at which two strings (or binary seawonce) of equal length differ. It is mainly used in error detection, error correction, and text similarity.

Formula:
$$H(A_1B) = \sum_{i=1}^{n} [A_i \neq B_i]$$

* A And B one two strings (or) bit sequences of the same length.

* [Ai + Bi] is 1 if the characters at position i are

different, otherwise o

* The Sum gives the number of differing positions

Example. I Binary String let's say we have two 8-bit binary numbers: 44444444 B: 1001001 comparing each bit : 10 11101 Hamming distance = 2 (Since two bits differ) Example : 2 For two woods of equal length 18 3 A: "karolin" B: "kathrin" 3 comparing characters 3 karolin kathrin -Hamming distance = 3 (difference at positions 2,4 and 6) 2.2 Distance Measures A distance measures is used to find the dissimilarity between patterns represented as vectors. patterns which one more similar should be closer. The distance function (d) could be a metric (or) a non-metric. The most popularly used family of distance metrics is called entere: P.9. one vectors of the Minkowski metric. metric A metric is a type of measure that possesses three key attributes Measurements used to evaluate the performance of a model. Error identification positive seflexivity, symmentary and triangular inequality

* positive reflexivity: d(a, x) =0

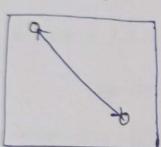
* Symmentry : d (ny) = d (y,x)

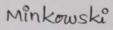
* Taiangular inequality: d (x,y) < d(x,z)+d(z,y)

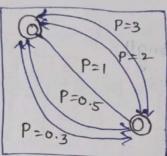
Where n, y, z one any three patterns

The following sections describe different types of dissimilarity /similarity (proximity) measures between patterns vectors.

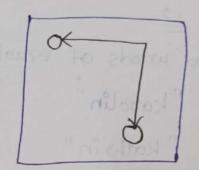


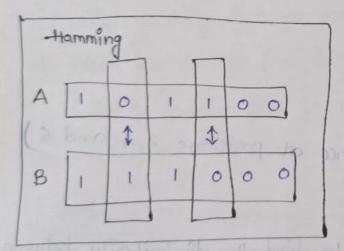






Manhattan





Hamming distance

Minkowski Distance

The Minkowski distance is a generalized metric used to calculate the distance blu two vectors. It is given by

 $d(p,q) = \left(\sum_{k=1}^{L} | p_k - 2k|^r\right)^{r/4}$ Absolute difference blue

Where: P.2 one vectors of length L.

or is a parameter that determinen the type of metric

Weighted distance Measures are used in Machine locating to calculate the distance between data points while assigning different impos-

tance (weights) to each features (or) dimension.

This is oseful in scenarios where certain features one more

there are types: 1. Weighted Manhattan Distance

2. Weighted Euclidean Distance

3. Minkouski Distance (Generalized Distance)

1 Weighted Manhattan Distance:

Manhorttan distance is a distance metric used to calculate the distance between two points in a good-like path.

It is also called LI-distance (or) taxicab distance becourse it represents the total sum of absolute difference between co-ordinates.

For two points X and Y in an L-dimensional space.

where

* Absolute difference | nk-yk | ensures no negative distance.

Example calculation:

Suppose we have two points in a 20 space

9 = 315

Using Manhattan distance formula: $d(n_1y) = (x_1 - y_1) + (n_2y)$ $d(n_1y) = (3-1) + (5-9) = 2+u = 6$ distance b|w= there points is 6

Use cases 1 K-Newst Neighbors (KNN): used to measure similarity blu points 2. K- Means clustering : When using II hoom for cluster assignments. 3. Recommendation systems: To find the closest matches. 4. Anomaly Detection: Helps in identifying outliers board on feature differences. Example - program from scipy spatial distance impost cityblock y = [1,9] district (00) smoth of bottom solo distance = cityblock (n,y) point ("Manhattan Distance: ", distance) Manhattan Distance : 6 the values of the K-IL teather for points as difference our of onsures no regative distance Suggest We have two points in a 20 space thetome formula , of (my) = (21-41) + 90 9 (det) = (3-1) + (0-0) + 5 + 11 = 6

JBFs algorithm Bacadth - Frost Sewich (BFS) Algorithm is a fundamental graph traver Sal algorith (used to explore nodes and edges of a graph systematical It is porticularly usaful for finding the shortest path in an unweighte graph. * Finding shootest path in an unweighted graph * Network broadcasting [finding all reachable nodes) * Solving puzzles * Web corawling * Detecting cycles deque - Double- ended queue (deaue) 18 faster than a list for appending and popping elements Remove - from Right FIFO - Smitialize Loop until Queue is Empty Implementation of BFS Algorithm 1 collections - extend the functionality & build Sin type, list, tupple 3 @ deane - appending And popping alements from both sides. 3 Graph - Direct graphs, undirect graphs, Weighted graph, in weighted graphs, cyclic graph, Acyclic graph. @ defautalet ML-knowledge graphs Scanned with OKEN Scanner 1 Collections - extends graph fundionally using with build in fundionally using with build in fundionally strains

3 deque - appending And popping element,

@ Graphs - Direct, undirect, weighted, unweighted graphs

- 3 de-Pault
- O Distance Measures O Weighted Euclidean
 - @ Weighted Manhattan
 - 3 Minkowski Distance
- @ weighted Distance measures minkowski Distance

8 = parametres value.

P.2 = vector elements

L- diamenssion

$$n = (5,2,4)$$
 $y = (3,4,2)$

 $w_1 = 0.3$, $w_2 = 0.5$, $w_3 = 0.2$

= 0.3 $\times (5-3)^2 + 0.5 \times (2-4)^2 + (4-2)^2 \times 0.2$

 $= 0.3 \times (2)^{2} + 0.5 \times (2)^{2} + 0.2 \times (2)^{2}$

= 0.3 x4 + 0.5 x4 + 0.2 x4

= 1.2+2+0.8

= 4

Non-Metric Similarity Measures in Machine Learning:

Non-Metric similarity Measures do not satisfy the properties of a Metric distance function.

data and they are resistant to orelivers (or) extreamly Noice data one example of a Non-Metric similarity function is the k-median distance between two vectors

Given
$$\alpha = (x_1, x_2 - -x_1)$$
 $y = (y_1, y_2 - -y_1)$
 $d(x_1y) = k$ -median $\{(x_1 - y_1) - - (x_1 - y_1)\}$

Similarity measures $s(x_1y) = \frac{x^ty}{1100 \text{ try II}}$

d(ny) = 1 - s(ny)(symmetric)

Ex: any and 2 be three vectors in two dimensional space

77777799999999999999

$$d(x_1y) = 1 - cos(T/2)$$

$$= 1 - cos(-1/2)$$

$$d(a_{1}z) + d(z_{1}y) = 1 - \cos(\pi/4) + 1 - \cos(\pi/4)$$

$$= 2 - \frac{2}{\sqrt{2}} \quad \cos(\pi) \text{ value} = -0.6536$$

$$= 0.586$$

Buta Defugliation :

The levenshtein Distance (Edit distance) is a metric used in Natural language processing (NLP) and string matching tasks to Measure the difference between two Sequences (Usually woods (or) sontences) It counts the minimum number of operations required to con-Vert one string into another string.

- 1. Insertion (adding a character)
- 2. Relition (Removing a character)
 - 3. Substitution (Replacing one character with another)

usages of levenshtein Distance

1. Text Similarity & spell checking:

used in seasich engines and outocorrect to find the closest matching words.

2. plagiarism Detection:

compares document similarity by checking word variations.

3. DNA Sequence - Analysis

used in bioinformatics to compare genetic sequences.

4. Speech Recognition & chatbots : 10000

Helps recognize variations in words (or) phoases.

5. Data Deduplication

Identifies similar entries in databares.

```
Example for levenshtein Distance;
 Command for Levenshtein - pip install python-Levenshtein
  import levenshtein
  wood1 = "machine"
  wood = " maching " & Rotine two woods
 distance = levenshtein. distance (word 1, word 2)
 # calculate Levenshtein Distance
 posint (+ "levenshtein Distance between \quandif and \quand 23 is:
                      Edistance? )
output :
   levenstien Distance between 'machine' and machine is :1
Explaination:
    This means only one edit operation (substituting "e" with g")
is needed to convert "machine" into " maching.
Levenshtein Distance: levenshtein distance between two woods is the
minimum number of single character edits required to change one
word into the other.
     * Edit Distance (levenshtein Distance)
    * Edit distance is a way to measure the similarity of two
      * Insertion, deletion, replace.
EX: FXOUS FLOMAX ? F->V - Substitution
            VDLMAX J L-> 0 - Substitution
                          0-> L - substitution
      GILY ? E-> Substraction
                    E> addition.
```

9 9

-5

9

9

9

-

-

0

9

-5

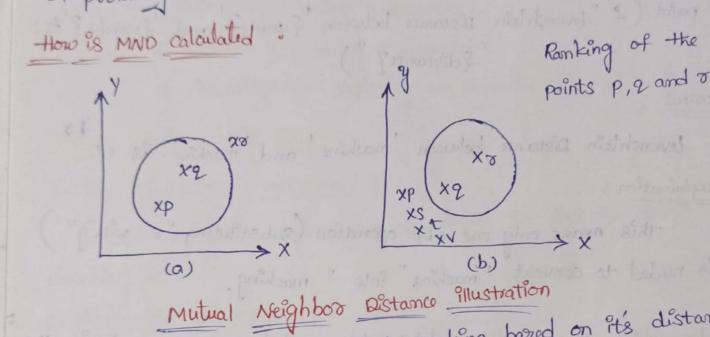
5

5

-

Mutual Neighbook hood Distance (MND): Mutual Neighborhood Distance (MND) is a function that measures the Similarity between two points (or) patterns. X and y bored on their neighborhood relationships. Instead of relying on traditional distance metrics (like Euclidean distance)

It considers how the two points rank each other in terms of proximity to other points in the dataset.



* Each point hour a neighborshood roomking borred on it's distance

from other points.

* NIN (AIY) reporesents the rank of y when ordered by it's

* N,N (y,x) represents the rank of x when ordered by it's

pronumity to y.

* The Mutual Neighborhood Distance (MND) is then computed or MND (My) = NN (My) + NN (y,x)

Relative positional vanking:	sulla suit abient
1 2 Mutual Distance:	Laborator of 9 15
P 9 8 MND (P,2) =2	Library 83 p. s.
9 P 8 MND (210) =3	I have S as to be
8 9 P MND (PID) = 4	on debalar or apple
1 2 3 4 5 Mutual Distan	co .
(00)=	5
PStV28 MND (PA) =	3 at both .
MND (210) =	7
2 P & S t V MND (PIS) =	to pot and other
8 a p s t v (BM n (Q))	at milhoods
properties of MND:	stelman at +
1. Symmetry _ MND (ny) = MND (y1x)	to Compare
2. Reflexivity-If a point is compared wit	th itself $(x=x)$ then
2. Reflexivity-If a point is confi	NIND (MIX) =0,
3. Not a Metric - While MND is Usefull.	1844 annot
It does't satisfy the triangle inece	and, manuf
be used on a strict mathematical distan	ce metric.
Formula for MIND:	MU (3) M] T;
	Example Let's CA!
MND (AIB) =	1
Where * A and B are the two items you	are composing
where * A and B are the two rems you	Silvi to Por Stance
* N(A) And N(B) one the Sets of	neighbores Too Tooms
A and B, respectively!	*
1 w wal so the number of	metrial neighbours
* N(A) n N(B) is the number of	Trace of the second
between A and B.	o proto
If there one no mutual neighbours, the dis	tance is intina,
"roucaling that	(e) unoulated.

consider three individuals A,B and c (or) p, and or * p is connected to 9, 8 and 8 * 9 is connected to pro and t of 8 is connected to PI2 and U Steps to calculate Mutual Neighboron Distance (MND) 1. Identify the neighbours of P and 2 [N(P), N(2)] 2. Find the mutual neighbours. This done by finding the intersection of the two sets of neighbows. N(P) n N(2) * The number of mutual neighbours is represented ar N(P) nN(2) | the size of the intersection. 3. calculate the mutual Neighborier Distance. MND (P12) = [N(P) n N(2)] if |N(P) NN(a) | =0 there is no mutual neighboros. 0 Example: let's say we have 4 Hems AIBIC and D * A is connected to Bic and D * B 18 connected to AIC and D * c is connected to A and B * D P8 connected to A and B. ATT 23 (8) M A (A) A out no million relighbours, the distance is infinite

N(A) = {B,C,D} N(B) = {AcciD} => MND (AB) = [N(A) n N(B)] N(A) 1 N(B) = C,D Mutual distance between - A And B ison [N(A) N N(B)] = 2 N(c) = {AIB? N(D) = {A1B} MND for each pair : MND (AIB) : $N(A) n N(B) = \{c_1 D\} \Rightarrow |N(A) n N(B)| = 2$ MND = 1/2 = 0.5 MND (AIC) : $N(A) \cap N(C) = \{B\} \Rightarrow |N(A) \cap N(C)| = 1$ MND = 1/ =1 MND (AID) : N(A) n N(D) = {B} => |N(A) n N(D) = 1 MND = 1/1 =1 N(B) NN(D) = {A} => |N(B) NN(D) |=1 MND (BIC) : MND = 1/1 = 1 N(C) n N(D) = {A18} => |N(C) n N(D) | =2 MND (CID) : MND = 1/2 = 0.5

Determine the Relative position of Each Hem! 1. Calculate the Sum of MNDs too each item. 2. Compare Sum of MNDs Sum of MNDs for A, B, C, D A = MND (AIB) + MND (AIC) + MND (AID) = 0.5+1+1 = 2.5 · B = MNO (BIA)+MND (BIC)+MNO (BIO) = 0.5+1+1 = (DAMINIA) - 10.5 - LOWN (A) M = 2.5 C = MND (C,A)+MND (C,B) + MND (C,D) MAD (AIC) : = 1+1+0.5 = OMACAMA - 147 - OMACAMA D = MND (D,A) + MND (D,B) + MND (D,C) 1 = 2.5 : (OIA) GAM = 1+1+0.5 = (a) MACAM = 181= (O) MACAM = 2.5 Mutual Neighboron Distance Ps 2.5 0 Rank the îtems sequential order dépendance on the Distance 0 1st place : A 3rd place : C MND = 1 = E.S 4th place : D

proximity Between Binory patterns:

when this section describes different ways to measure similarity

(or) distance between binary patterns (binary strings).

We can view L-dimensional binory partierns on binory strings

of length 1. let p and 2 be two 1-bit binary strings.

some of the popular proximity measures on such binary patterns

ore: 1. Hamming Distance

a. Simple Matching coefficient (SMC)

3. Jaccord Coefficient (Ic)

1. Hamming Distance (+10):

Hamming Distance measures how different two bindry strings one by counting the number of bit positions where they differ.

* P(i) = 2(i) -> P and 2 match on their ith bit

* p(i) # 2(i) -> p and 2 mis match on their ith bit

(Hamming distance is the number of mismatching bits of the 1-bit locations)

Example: consider the two 10-bit patterns P and 2 given by.

The Hamming distance is 2 at they mismatch in bit positions I and 10. 2. Simple Matching Coefficient (SMC):

SMC Measures the proportion of Matching bits (both o and 1) out of the total bits.

1. Hamming Distance (-410) :

ne define bits where

* Moo = number of times both p and 2 have o

* M11 = " " pand 2 have 1

* Mol = " " pis o and 2 is1

* MID = " " P is 1 and 2 is 0

Formula:

SMC (P,Q) = MII+MOO

MOO+MOI+MIO+MII

* MII = 1 (Bit position 7)

* Moo = 7 (Bit positions: 2,3,4,5,6,8,9)

*MOI = I (Bit position 10) (enotional tidal

*MIO = 1 (Bit position 1)

 $SMC(PA) = \frac{M_{II} + M_{00}}{10} = \frac{1+7}{10} = \frac{8}{10} = 0.8$

3. Jaccord Coefficient (Jc):

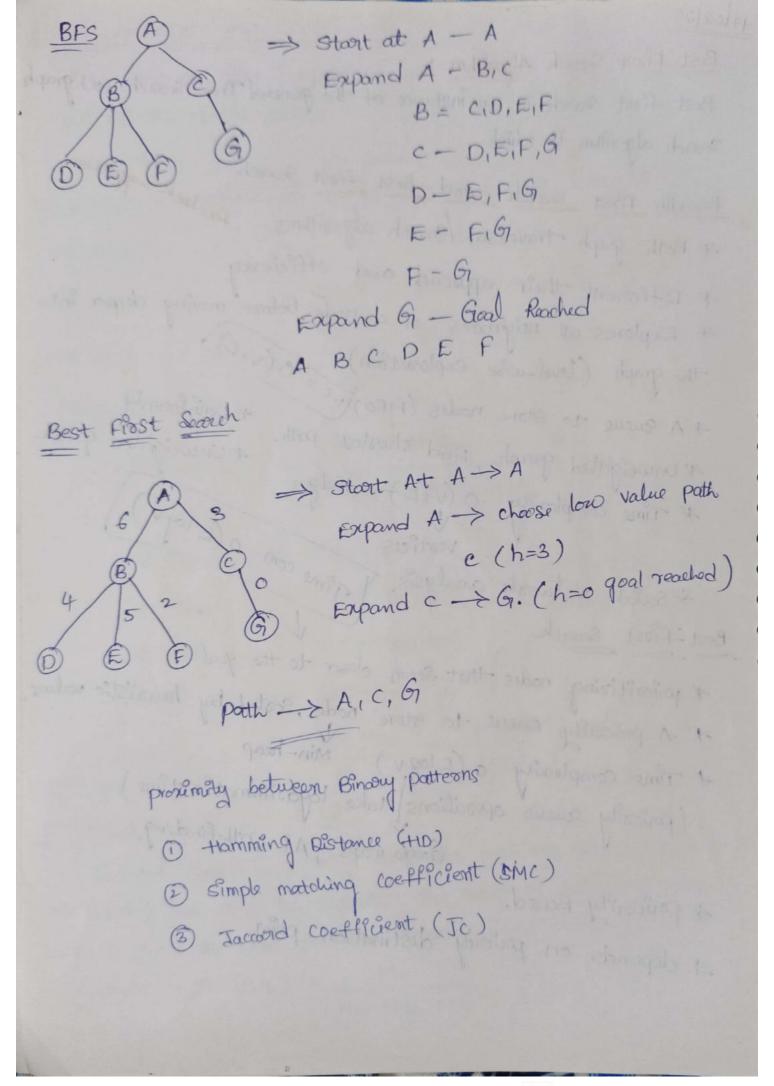
Jaccord Coefficient only considers positions where at least one storing how a 1. It measures the natio of positions where both one 1 to the total positions where at least one is 1.

Foomula: Ic (P19) = MII

Moi + Mio + MII

17/02/25 Best First Second Algorithm: Best first search is an instance of the general Tree search (or) graph Search algorithm is which * Both graph traversal / Search algorithms for book suggestions * Different their approach and efficiency * Explores all neighbors of a node before moving deeper into * A queue to store nodes (FIFO) (CT) * uniformly * unweighted graph, find shortest path. * unweighted graph * Time complosity O (V+15) > Edges. *Social Network analysis. The con- (ElogV)

t-Flost Search -3 Best-First Socorch * pricontizing nodes that seem closer to the goal 3 * A periosity and to store nodes, sosted by heuristic values -3 -0 * Time complexity o (£logv) Min-heap -0 (priority avene operations take logarithmatic time) 9 Google maps, Al path-finding. 3 -3 * poucetty Booked. -0 * depender on policity destination point.



For the given Example * Mn =1 $Tc(p_1q_1) = \frac{1}{1+1+1} = \frac{1}{3} \approx 0.33$ * MOI =1 * MID = 1 consider again : P=1000001000 9 =0000001001 MII = 1, MOO = 7, MOI = MIO = 1

$$SMC(p_{1}, 2) = \frac{1+7}{7+1+1+1} = \frac{8}{10} = 0.8$$

Differente classification Algorithms Based on the Distance Measures

Many dassification algorithms vely on distance measure to determine the similority (or) dissimilarity between data points. This algovittems classify new data points bared on their proximity to existing labeled data.

Below one some key classification algorithms that utilize distance Measures:

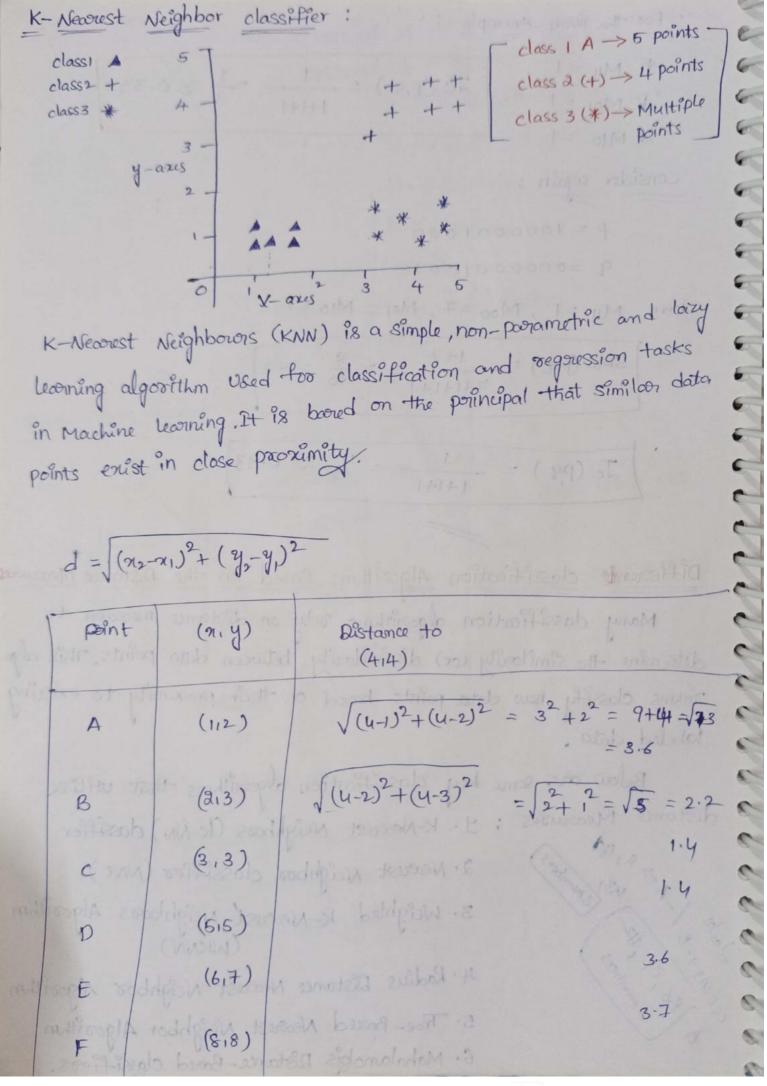
I. K-Newrest Neighbors (K-NN) classifier 2. Newest Neighbor classifier (NNC)

3. Weighted K-Newsest Neighbors Algorithm (WKNN)

4. Radius Distance Nearest Neighbors Algorithm

5. Tree-Bould Nearest Neighbor Algorithm

6. Mahalamobis Distance-Booled Classifiers



Diffesions dassification algorithms should on the standing Message.

22/02/25

KNN - Regoussion:

K-Newest Neighbor Regoression is a Supervised borning algorithm used to paedict continuous values instead of catagorical labels. It works by finding the K-Newvest data points to a given input and averaging their corresponding output values to make a preduction KNN-Regaussion works:

1. Store all training data - No explicit model & built 2. Find the k-Nearest Neighbors of a given input data point borred

on a distance metric

Ex: Fuclidean distance.

3. Extract the output values (y-values) of there nearest neighbours, 4. Comput the predicted output by taking the average of there Y-values.

Mathematical Formula:

For a given input &, the poudicted value of is Ay = 1 K Say

1. Find the K-nearest neighbors of x from n data vectors. let them be x1, x2, --- xk

2. Consider the 1x-y values associated with there x's. let them be y, y2, ---, yk.

3. Take the average of there y's and declare this average values to be the powdicted value of y associated with x.

so the powdicted values of y, call it of is

	-		I - Thistop of breeze
Ex:	Number (i)	pattern (ni)	Target (yi)
	1	(0.2,0.4)	8
	2	(0.4,0.2)	0/18 atal priming to motert
	3	(0.6,0.4)	12 08 -
	4.	(0.8,0.6)	16 06-1 *
	1 5 1	(1.0,0.7)	19 0.4 - 62,000
	6.	(0.8,0.4)	14 0.2 (o 4,0.2) ×
	7.	(0.6,0.2)	10 01 02 03 04 0.5 06 07 08 09
	8.	(0.5,0.5)	12 le tors more o sor
	9.	(0.2,0.6)	10
			12 1-7 14 12

New pattern be x = (0.3, 0.4) k = 3

* The 3 Newvest Neighbours of on from the patterns in the table are (0.2,0.4), (0.4,0.2), and (0.5,0.5)

* The corresponding target values observed in the table 8.8 and 12

* The average of there values is $\frac{8+8+12}{3} = 9.33$ * so predicted target value for x = (0.3, 0.4) is 9.33.

Composison with a linear Model The data can also be modeled using a linear equations. yo = aχo (1) + b 90 (2) +c 1. Boston Housing Dataset Example * A real-woold dataset containing 13 features Er: coûme rate 200 training samples and 50 test samples The goal is to posedict house pouces using k-NN regoussion. Mean Stuared Error (MSE) is used to measures prediction acturacy 50 45 40 35 30 25 20 15 10 30 Test pattern number MSE values on 24 the Boston Hou-23 sing data KNN Regoussion. 21 20 of Neighbours.

per-tormance Measures:

There are several measures used to evaluate the performance of ML models. We will consider some of the popular measures in this sections.

periformance of classifiers

* dassification Accuracy

classification accuracy is a basic performance metric that calculations how many posedictions were correct out of the total posedictions made.

Given

a = Total number of patterns classified.

he = Number of correctly classified patterns.

Formula

classification Accuracy = $\frac{hc}{a}$ x100

If a classifier correctly classifies 280 out of 350 patterns.

350 x100 = 80% accuracy.

* Confusion Matrix: Confusion Matrix provides a detailed boreakdown cof a classifier's performance by showing how many samples were correctly (or) incorrectly classified for each class.

* let there be three classes ci, c2 and c3

* Let there be 200, 100 and 50 patterns from classes c, , re and co

respectively.

Tome/predicted	CI	C2	C3
C,	180	15	5
C2	5	85	10
C3	3	2	45

0

0

0

	KNN- Keg	oussion -> cla	extination accuracy	EX - House Pouce Pro	
	art the		broug on the sy	July sidem omanding	
•	. 0.0	Making Al	array In TIP	The second	
	(on-tusio	D. Manney	TOATA	1+FP+FN 35+50+10+5	3
,			- Partie	The state of the s	
9			ingipules	35 John 35 = 0.78	
		= 85 =	85% = 0.85		
>		the same of	Librar hyp it	35 0.78	
>	DOSHIVE	e values pried	uetion =	35+10 US 1	
>				35 35 = 0.78 35+10 us	
			to the skullenil	STINEM HOMMED TOURS	-
			dlone = TA	= 3545	2
6	NAMINA	values poredu	dions = TP+FA	3575	
	The state of the s	and the second	16M - / 7 a	was a second a room.	
9		an Rey X	- (218, 31) + (21)	ty2) (xnign)	
>	THE P	The state of the s	in south forther bone	belilberg metalist mores	
>		to at the same		his of pros si st and to his	
,				01 0 15002	
	As	tial values =	X = 1 1 - 1 =	1 - 100	
	410	have your	N		
,			9: - 2	क्षेत्र कर्म के पुष्ट : कर्मान	
/	125-		O STANK IX	moo out so he	
(2)	2/02/25-4-	1 perforn	nance motrics for	the Regoussion Algorithm	
	Eli				
	15 Ith	(a) Hea	n Absolute Error (M	MAE) - 1 5 (90-90)	
	2 mambers	(b) Meo	in squared Error (MSE) _ I h	
	0 /		The state of the s	n 5 (0) 01. 12	
	8			2 (91 - 91)	nor.
	V		The second secon	$MSE) - \frac{1}{h} \sum_{i=1}^{n} (y_i - \hat{y_i})^2$	mor)
	House	Actual value	predicted values	Absolute squared	mor)
	V		The second secon	Absolute squared	more)
	V	Actual value	The second secon	Absolute squared values values	eros;
	House	Actual Value 4;	predicted values	Absolute squared	and the
	House	Actual value	predicted values	Absolute squared values values	aon all
	House	Actual Value 4;	predicted values 210 210	Absolute squared values values 10	aon and
	House 2	Actual Value 41; 200 250	predicted values 41	Absolute squared values values	aon a
	House 1	Actual Value 41; 200	predicted values 41 210 215	Absolute squared values 9:-9: -10 5 -10 5	aon and and and and and and and and and an
	House 2	Actual Value 41; 200 250 220	predicted values 210 210	Absolute squared values 9:-9: -10 10 5 -10	aon and and and and and and and and and an
	House 1	Actual Value 41; 200 250	predicted values 41 210 215	Absolute squared values 9:-9: -10 5 -10 5	aon and and and and and and and and and an

performance Metrics for Reguession Algorithms:

performance Metric: performance Metric one used to evaluate how well a model is performing with respect to the task it is designed for. There metrics give you insight into the models accuracy, realiability, and effectiveness in making paeductions.

the choice of metric depends on the type of machine learning task and the specific goals of the model.

the most common metrics include ; I. Mean Absolute Error (MAE)
2. Mean squared Error (MSE)

I Mean Absolute Error (MAE): MAE measures the average magnitude of Errors between powdicted and actual values, without considering their directions. It is easy to interpret and is less sensitive to outliers.

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

whose: yo is the actual value

yo is the possible values

n is the number of data points

Example problem | Real world entity usage Example:

House posice prediction: posedict the posice of a house boosed on it's features

Ex: size - saucore feets

Location

Number of Rooms

After training your model, you want to evaluate its performance asing MSE.

6

6

6

0

6

0

0

MSE	calculations:	the day of	7)	es[02]25	15 70 11.15	
Actua	1 posice	poredicted price	2	-(A2)-A-	M. 15,	
250	,000	260,000	Tatal and	CSE(AI)-A-	(A) (A) (A)	
300,	000	310,000	()	2, 10, 11, 10, 128, 128, 128, 128, 128, 128, 128, 128	491 00,000,00	
350,	000	340,000		0-1	as deliver (1)	
400,	000	590,000			A MAR B	
450,	000	60,000	The second of the second		1.1 5h	
The !	The MSE would be calculated as 12,20,29,32,36,46.54					
MSE	MSE = $\frac{1}{h} = \frac{1}{1} \left(Y_i - \hat{y}_i \right)^2$ 59,60					
MSE			Manus A-		0	
Mean Sa	Mean squared Error (MSE): MSE measures the average squared difference blw actual and					
MSE M	laxures The un					
portdicte	posedicted values. $MSE = \frac{1}{n} \sum (y_i - \hat{y}_i^2)^2$					
Houses	Actual price	pordicted	Error (y, y,)	Absoluti	saucred Error	
1	200	210	-10	10	100	
2	250	aus	185 b	5	85	
3	180	190	-10	10	100	
4	120	215	5	5	25	
5	275	280	-5	6	25	
MAE	= + 1 / 3: - 2	1 - To find Average of	MSE :	100+25+100+	-4;)2 outliers.	
MAE =	10+5+10+5+6	3 SALES OF THE PAS	MSE =	275 = 55	COLUMN TO SERVICE STATE OF THE PERSON SERVICE STATE SERVICE STATE SERVICE STATE OF THE PERSON SERVICE STATE SERVICE STATE SERVICE STATE SERVIC	
THE BUILDING						

Department of Artificial Intelligence and Machine Learning

II B.Tech II Sem

Machine Learning

UNIT-3

Models Based on Decision Trees:

By

Mrs.K.Reddemma,

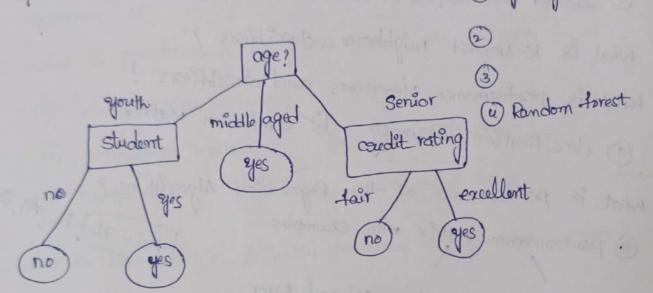
Assistant Professor,

Department of AI&ML.

UNIT-3 Model Based on Decision Trees

Introduction to Decision Prees

A Recision Tree is a Supervised learning model used for classificat ion and Regordssion. It splits data into smaller subsets bord on feature values, counting a tree-like structure, where each internal node reporesents a decision.



income student coudit-rating class: buys_computer goverdy - purity in splitting - each child node dependes on the povern

- It can be expensive no of features depends on the no of patterns

To identify the shortest path of the path is better update the path, given edges. f_score, g_score.

chapter-3 Models Based on Decision Trees

Decision Trees: Decision Tree is a popular structure in ML It is built by splitting the data set associated with a node into Subsets that typically have less entropy (or) better provity. splitting is corried out using the values of a selected teatures this feature-borred decision-making at each of the internal nodes offers transparency and is ideally suited for explanation and for easier understanding by application domain experts.

properties of Recision Tree building algorithms:

- 1. It is gowedy
- 2. It can be expensive
- 3. It can overlit
- 4. Random Forest
- 5 Ada Boost (Adaptive Boosting) / poudictive models using Decision
- 6. Gradient Boost

1. It is goverdy: the best factorie

A Decision Tree is greedy because it selects the best features (in terms of purity) at each step, without considering the overall best tree structure

when the dimensionally is longe.

* So, tests feature used at each of the child nodes will depend on the feature selected at the parent node,

* Such a sequential goverdy process may fail to be optimal on the whole.

Example:

Imagine you are classifying students bared on their grades and study habits.

* Split the study time because it provides the highest pusity at that moment.

* But a better long-term (stratigy) could have been splitting bared on previous exam Scores.

2 It can be expensive:

Selection of a feature depends both on the number of features and the number of patterns. so Tests conducted on the feature values at each node need to be simple.

Even simple tests can become computationally demanding when the dimensionality is longe.

Example:

Imagine a dataset with 10,0000 records and 100 features * Each node how to check 100 different splits.

* If the tree grows deep (Ex depth = 10)

there one 1000 of computations.

solution: To reduce computation expenses

Calculation of the Entropy Values: Doosion in Cossing Entropy; Entropy measures the impurity (or) uncertainty in a data Set. It is used in decision trees to determine the best feature for splitting data. Entropy Formula: H(S) = - E P; log P; Entropy +1(5) for a dataset S with classes C, 1C2 --- Cn is given by HIS) = - Epilog, Pi Where : * Pi = proportion of samples belonging to class i * log_ is the bore-2 logoroithm. Example Calculation: Case 2: perfectly impure Dataset (50-50) split case 1 : pure Dataset Ex: 5 positive, 5 negotive (All in one class) probabilities = P+ = 0.5, P= = 0.5 En: 10 positive, o negative Entropy: probabilities: P+=1, P-=0 Entropy: $H = -(1 \cdot \log_2 1 + 0 \cdot \log_2 0)$ $H = -(0.5 \log_2 0.5 + 0.5 \log_2 0.5)$ H = 0 high purity case 3 : portfally mixed Dataset (7, positive, 3 Negative) 20 9/00 Enthopy (H) = - (0.7 log, 0.7+0.3 log, 0.3) H = 0.88

Entropy Before splitting let's say we have a dataset with 10 samples. * 6 belong to class A (PA = 6ho = 0.6) * 4 belong to class B (PB = 4 = 0.4) -H parent = - (class A log class A + class B log 2 class B) = - (0.6 log_ 0.6 + 0.4 log_ 0.4) Entropy After splitting: Now suppose we split this dataset into two subsets: * left node - 4 Samples (3A, 1B) 00000000 Entropy (H) = 0.81 * Right node : 6 Samples (3A,3B) Entropy = 1 Total Emtropy after splitting +1 children = $\frac{4}{10}(0.81) + \frac{6}{10}(1)$ + = 0.924 Information Gain calculation IG = Howent - Hchildren = 0.97 - 0.924 [IG = 0.046]

values.			Pavadnost	
class	classi	class 2	Entropy	
	(Pi)	(P2)	V	
1	100 (P1 = 1.0) C	(P2 = 0.0)	0.0	
2.	75 (P1 = 0.75) 25	5 (B = 0.25)	0.811	
3.	50 (P1 = 0.5) 50	$(P_2 = 0.5)$	1.0	
4.	25 (P = 0.25) 75	(P2=0.75)	0.811	
			7.3	
<u> </u>	10 (P1=0.0) 100	(Ps=1.0)	0.0	
This is be	symetric meaning	nd = : 10	1 to end	.2.
this is be	cause $(P_1 P_2) = -P_1(\log p_2) = emtroper$	9 P;) - P2(la	og P2) - P2 logs	02 - Pi log
entropa entropa	cause $(P_1 P_2) = -P_1(\log p_1)$ $-p_2(P_1 P_2) = emtrop_1$ $-class problem, em$	py (P2, P1) tropy is ma	og P2) - P2 logs	$P_1 = P_2 = 0.5$
entropa entropa	cause $(P_1 P_2) = -P_1(\log p_1)$ $-p_2(P_1 P_2) = emtrop_1$ $-class problem, em$	py (P2, P1) tropy is ma	og P2) - P2 logs	$P_1 = P_2 = 0.5$
entropa entropa	cause $(P_1 P_2) = -P_1(\log p_1)$ $-py(P_1 P_2) = emtrop$ $-class problem, emblem, emble of em$	pp;) - P2(la py (P2, P1) tropy is ma tropy and	og P2) - P2 logs namum when I equating them	$P_1 = P_2 = 0.5$
entropa entropa	cause $(P_1 P_2) = -P_1(\log p_1)$ $-p_2(P_1 P_2) = emtrop_1$ $-class problem, em$	pp;) - P2(la py (P2, P1) tropy is ma tropy and	og P2) - P2 logs	$P_1 = P_2 = 0.5$
entropa entropa	cause $(P_1 P_2) = -P_1(\log P_1) = \frac{1}{2} = $	pp;) - P2(la py (P2, P1) tropy is ma tropy and P2) = -	og P2) - P2 logs eximum when I equating them $1 - \log(P_1)$	$P_1 = P_2 = 0.5$
entropa entropa	cause $(P_1 P_2) = -P_1(\log p_1)$ $-py(P_1 P_2) = emtrop$ $-class problem, emblem, emble of em$	pp;) - P2(la py (P2, P1) tropy is ma tropy and P2) = -	og P2) - P2 logs namum when I equating them	$P_1 = P_2 = 0.5$

A dataset	used to illustra	ati splitting:	my houself his	1 of distances	8 9
Pattern	Feature 1	Feature 2	class		-
1	undarinet	1005	1		6
2	1	3	1 19		6
3	2	(01 ,0)	(1)		6
4	2	3	(2 a a a) à F		6
5	4	Tan India	1 000	-	0
6	4	3	2		10
7	5	(SE 0 - 4) S	2		6
8	5	(138)	2000		0
		Diame	e seeds out		0
Embore	Calculation:			,	6
	The state of the s	Postus by	moun sidente acoust	Padoulies	60
3	ubset 1			0 040	9
	4/4/09	(1) + 0 log (0))=0	od 85 अंतर	0
19 pol 19 - 6		10g pg) - Bell			100
50		79.01 0			10
	2-41	08(4)-3	10g (3)		96
	The state of the s				6
	4	9 (4) + 4	$og\left(\frac{4}{3}\right)$		6
	2 0 .	244219	to philovious		6
1	(intental) prome6 .		60
L.vo 2	,		96	_	8 2
Feature 2	* * ! # .		-dimensional p		00
3	(3) bol - 1=	dato	set given in	Table 2	0
	* * *	#			0
1 1	2 3 4	5 Feature:	1		0
	2 3 4	Feamle .			6
The state of the s		and the second s	The second secon	THE RESERVE TO STATE OF THE PARTY OF THE PAR	1

3. Decision Prees can overfit overfitting happens when a decision tree borns too many details from training data, making it less generatizable to new da * A deep tree can memorize training data rather than finding general parterns. * Small Variations and noise in the training data one toward at important, reducing performance on test data. pauning: pouring is a technique that removes unneccessary bran chess after training. * pore-poruming - Early stopping * post-pouning - Girow the full tree than remove branches -that do not improve accuracy on a validation set * DT - Accuracy on data Set - low Decision 1 Gireedy Nature - Makes fast decisions but may not be optima! - large data sets and High-dimensional fator Expensive încouare processing time. 3 overfit - Deep trees memorize data but fail on new enu 28.350

Here each of Subsets has 4 elements out of the total of 8 elements pousent. so the weight are equal to 4 - 1/2 so the Weighted value of entropy impurity is 1 x0+1 x0.244 =0.122 so total weighted emtropy associated with the split is 0.1221 Feature 14 Feature 125 class I Feature 2 23 class 2 class 2 class 1 Decision Tree for the data set (0.244+0.301) = 0.2726 (2) Feature 1 (3 Subset \$ 5,6,7,89 Subset finzy

Subset 1 15,6,7,89 Subset : {1,2,3,4} Entropy: 0.244, weight =0.5 Entropy : 0, Weight = 0.5 3.3 Impurity Measures For Decision Tree Construction 1) Gini Index (impurity at a node for a C-class problems) The Gini index is a metric used in ml to measure the impurity (or) disorder of a data set. It is commonly used in Decision Trees to determine the best feature to split the data Formula; Gini = 1- E Pi where * c = Total no. of classes * Pi = probability of a positicular class in a node. => Gini =0 -> The node is pure (all samples belong to the same class) → Gini =1 -> The node is maximally impure (Samples are evenly distributed among classes) =1-0.68 =0.32(If a node contains =1-(P1+P2) 80% of class P, and $=1-((0.8)^2+(0.2)^2)$ 20% of class P2 the 1- (0.64+0.04) Gini Indox Value is . 0.32

Scanned with OKEN Scanner

(2) Misclassification impusity 1 - max(Pg) * P; is the probability of a class at that node. Ex: Suppose a node contains the following class distribution class A = 807- (0.8) class B = 15 %-(0.15) class c = 5%- (0.5) Then the misclassification impurity is 1-max (0.7) = 1-man(P;) - $= 1 - \max(0.8)$ (Misclassification impurity is a metric used in decision trees c to measure how often a sandomly choosen element would? be incorrectly classified if we used the majority class in a given nodo.) Example: To find the outlook, Temperature, Humidity properties of decision Tree properties 2. Temperature 3. Humidity. To determine the data set proposties board on the given data 18/03/25-10 70 11 CSE (AE) - A

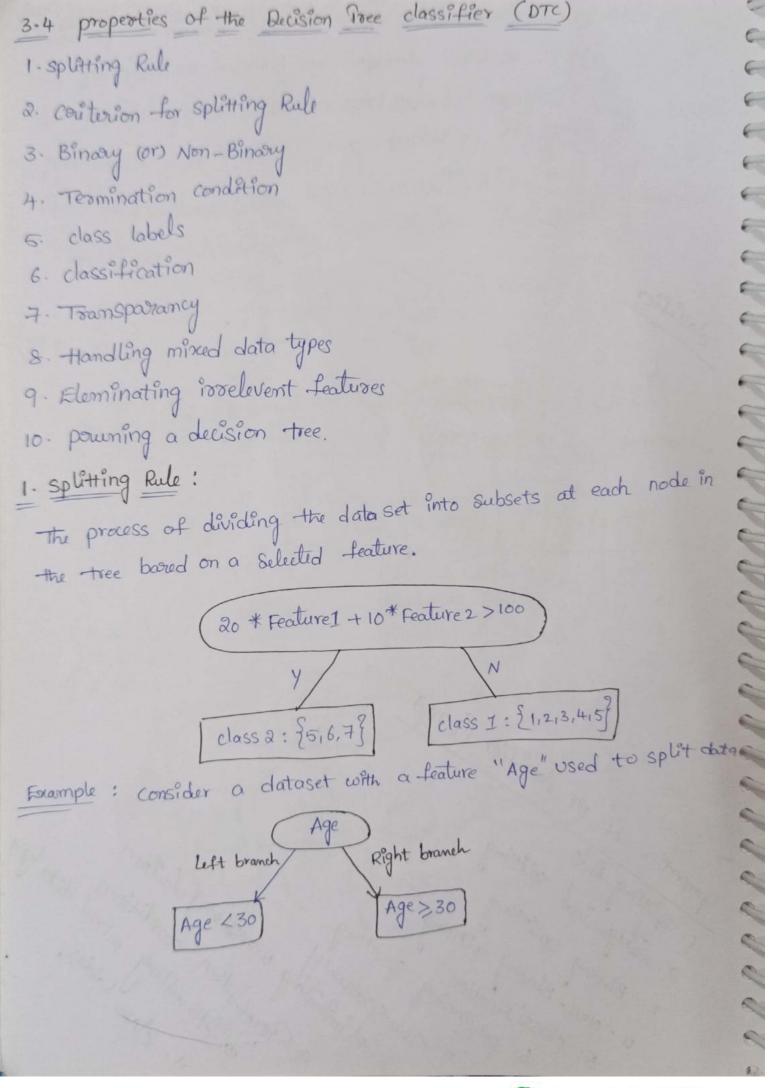
	pecision Tre	e properties Independent	features	
Pattern	outlook	Temperature	+lumidaty (%)	class
1	Sunny	65	91	NP
2	Sumny	85	92	NP
3	overcast	85 /	86	P
4	over cast	70 ×	65	P
5	Rainy	70	65/	NP
6	Rainy	79	62/	P
* subset	5 -> (1,2) (5,6) is pu	$(3,4)$ $(5,6)$ re $P_1=1$, $P_2=$		ex Rang (0,1)
	3,4,5,6) (1,2) and (3,4	$= P_1 = P_2 =$	17	N {1.4.5}
1213 W15,6	ini index =	$= 1 - (P_1^2 + P_2^2)$ $= 1 - (0.5)^2 -$	1641313	Temp >60

Bucksian Pres

using Temperature (F) > 70 subsets Subset I' - (213,6) subset 2 - (1,4,5) pure - Gini index =0 for (2,3,6) -> P, = 2/3, P2 = 1/3 Güni index = $1 - (p_1^2 - p_2^2)$ $=1-(\frac{2}{3})^2-(\frac{1}{3})^2$ G = 4 3 using Humidity < 91% split at Humidity < 91 * Two Subsets: Subset {1} is pure - Gini index =0 Subset { \$13,4,5,6} P1 = 4 , P2 = 1/5 $= 1 - \left(\frac{4}{5}\right)^2 - \left(\frac{1}{5}\right)^2$ Gini index Weighted impusity = 1/6 x0 + 5 x 8 25 * Temperature hor the lowest impurity measure (weighted impurity = 2/9) Best feature for splitting at root is Temperature.

By using Temperature > 70 at the root We obtain ; * pure Subset {2,3,6} -> classified on play * Second Subset {1,4,5}, which needs further splitting Two options for further splitting 1. Temperature > 62 2. outlook-based split Decision Pree for the Data Set @ using Temperature-based split * First split : * Temp >70 -> {2,3,6} (play) and {1,4,5} (needs further split) * Second split Temp >62 -> {4} (play), {1,5} (No play) (b) using outlook-bord split * First split Temp > 70 -> {2,3,6} (play) and {1,4,5} needs further split * Second split: outlook * Sunny -> {1} (No play) * over cast -> {4} (play) * Rainy -> {5} (play)

In Machine learning an impurity measure is a metric used to quantify the disorder (or) randomness in a data set. porticularly in classification problems. It helps in decision tree algorithms (like CART - classification and Regoussion Trees) to determine the best feature for splitting nodes. Temp > 70 Temp > 62 {115} No play Deusion Tree for the data set Temp > 70) {1,2,3,4,5,6? outlook {2,3,6}: play



2. Courtesion for splitting Rule: A mathematical measure used to determine the best way to split data at a node: common contemia include: I. Gini impurity (used in CART) 2. Entropy (information (Used in 203) Ex: For splitting based on Gini impurity It a node contains 80% class A and class B 20% Gini improvity = 1 - (P,2-P,2) = 1- (6.8)2-(0.2)2) (GI = 0.32) 3. Binouy (or) Non-Binouy: * Binory splitting: Each node splits into exactly two brances. * Non-Binory splitting: A node may split into multiple branches boored on categorical values. Example: A mo Binory: Salory < 50K -> Yes (left) /No (Right) Non-Binary: A node may split 1 (A2)-A 12:15 PM LEOS,

Regoussion Board on Decision Trees:

In decision tree-board regoussion, the goal is to powdict a continuous target variable by recursively splitting the dataset board on the feature values. The decision tree choosens the best split by minimizing the Weighted equared Error (SE).

step I : Understanding the given Data :

The dataset consists of two features (frand for) and a torget value. The table represents 9 data points with different values of fr, for and their respective target values.

	- coons	en puno	'n Me pi
Number	-fi	-f2	Tomat (91)
1	0.2	0.4	Target (2)
2	0.4	0.2	3 1 8 8 2 1 2
3.	0.6	0-4	12
4.	0.8	0-6	16
5	1.0	0.7	8 31 19 11
6	0-8	0.4	14
7.	0.6	0.2	10
8.	0.5	0.5	(6) + (6)
9.	0.2	0.6	+ +1 10 + 1

Step 2 : Occision Tree Splitting Conterion :

the decision Tree determines splits by evaluating different values of he and choosing the split the minimizes the squared error.

The sourced error for a split is calculated as:

SEsplit =
$$\frac{n_1}{n}$$
 $\sum_{\chi_i^* \in S_1} (y_i - \mu_1)^2 + \frac{n_2}{n} \sum_{\chi_i^* \in S_2} (y_i - \mu_2)^2$

* S, and So one the two subsets corealed after the split * no and no one the number of elements in each subset * Us and us one the means of the target values in each Subset

The mean of each subset is calculated or:

$$u_1 = \frac{1}{n_1} \sum_{\substack{\chi_1^* \in S_1}} y_1^*, \quad u_2 = \frac{1}{n_2} \sum_{\substack{\chi_1^* \in S_2}} y_1^*,$$

First we have to find the U, and Uz values and split the data sets on a s, and s2.

$$f_1 < 0.8 \longrightarrow \{1,2,3,7,8,9\} \{4,5,6\}$$

$$u_1 \rightarrow 8+12+10+12+10 = \frac{60}{6} = 10$$

$$u_2 \rightarrow 16+19+14 = \frac{49}{3} = 16.33$$

SE (SI) \Rightarrow (8-10)2+(8-10)2+(12-10)2+(10-10)2+(10-10)2 $\Rightarrow (2)^{2} + (2)^{2} + (2)^{2} + (0)^{2} + (2)^{2} + (0)^{2}$ 4+4+4+0+4+0

SE (S1) = 16
SE (S2)
$$\Rightarrow$$
 (16-16.33)²+ (19-16.33)²+ (14-16.33)²
 \Rightarrow (0.33)²+ (2.67)²+ (2.33)²
 \Rightarrow 0.1089 + 7.1289 + 5.428

Weighted (SE) =
$$\frac{6}{9} \times 16 + \frac{3}{9} \times 12.66$$

= $\frac{2}{3} \times 16 + \frac{1}{3} \times 12.66$
= $10.66 + 4.22$
= 14.88
W(SE) = 14.88

step-3: possible splits Bared on fi

To determine the best split the decision tree considers different threshold values for fi

$$f_1 < 0.4$$
 $f_2 < 0.5$
 $f_3 < 0.6$
 $f_4 < 0.8$
 $f_5 < 1.0$

For each split, the data set is divid into two subsets, and the squared error is computed: The split with the lowest squared error is choosen.

step. 4: Evaluating Frach split:

For each split

- I. Divide the dataset into S, (values satisfying the condition) and Sz (Remaining values)
- 2. compute the means U, and U2
- 3. Compute the squared error Contribution from both subsets.
- H. Find the total saucred errors.

The split with the Smallest squared error is the optimal choice for the first day decision node in the tree.

squared	Error (SE) va	lucs for possi	blo aptits
Test	se of si	SE of S2	Weighted SE
f120.4	2	82	64.22
fa <0.5	2.67	52.93	36.11
\$1 <0.6	n	48.8	32.04
f1 <0.8	16	12.66	18/go 14.9 att approvides or
£1 <1.0	55.5	0	49.33
f1 < 1.0	108	0	108
att book	$S_1 = \{1, 2, 3, \dots, 2, \dots, $	7,8,9}	optimal split are
	1, = 10 and	U2 = 16.33	
		(f, <0.8)	N Files does sold
£1,2,3,7	1,8,9} -1,20.5	<	(f2 < 0.7) {415,6}
21/20 21/2		3,718 3:11.33	
ofodo lomb	Deusio		the regrussion data set
			of military one decision to

Step I: Decision Tree Construction . 1. Start at the Root Node: 1) and many topol sunt with * The entire dataset is considered at the root * The model selects the best feature (f, (or) fz) on the best split booled on minimizing the squared error. 2 splitting the data: * Suppose we choose for LO.8 as the first split. * If f, <0.8, the data goes to the left child. * If fi >0.8, the data goes to the right child. 3. Further splitting * For the left child (f, <0.8) another split occurs board on f1 < 0.5 *If fico.5, another split board on fiz is made. 4. Leaf node : * When no further meaningful splits are found, the leaf nodes ore tormed. * Fach leaf nade is represented by the mean of the target Values in that region. step 2 : poudicting a new Data point : let's say we have a new test point (0-3,0-4) Since f1 = 0.3 < 0.8 -> Move to the left child. * stoot at the root * Next, since fi = 0.3 < 0.5 -> Move to the left child. * Finally, the instance lands in a leaf node with an average target value of 8.66 * so, the decision tree predicts 8.66 for (0.3,0.4)

Step. 3: Computing the squared Error:

* The tome target value for (0.3,04) 189

The squared error is calculated on

$$(9-8.66)^2 = (0.34)^2 = 0.111$$

composiison with KNN Regoression:

* In KNN Regoussion, a different approach is used. Ex: Averaging k-nearest neighbors.

* For the same test point (0.3,0.4), KNN Regoression predicted 9-33.

* The sourced error for KNN was $(9-9.33)^2 = (-0.33)^2$

= 0.111

The sancred error for kNN wor O.111

Greneral steps to Build a Quision Tree for Regoussion:

- 1. Start with the Entire rest made dataset at the root node.
- 2. Find the best feature and split point by minimizing the weighted squared error.
- 3. Divide the dataset into two subsets board on the choosen split.
- 4. Repeat the process recursively on each subset until a stopping condition 18 met.

Ex: max depth, minimum samples per node.

- 5. Assign the mean target values to each leaf node.
- 6. Use the tree to make predictions by following decision rules from the root to a leaf:

0

0

3-7 Bias - Variance Trade-OFF:

Bias - Variance Trade-OFF

The bior-variance trade-off is a fundamental concept in machine learning. It explains the balance between a model's simplicity (bias) and complexity (variance).

- * Bias: Foros due to overly simplistic assumptions in the Learning model.
- * Variance: Error due to excessive sensitivity to the training data.

This example illustrates bias and variance through paynomial fitting.

Step I: Understanding the Rata in Table.

The table consists of four examples with features z and target values y:

Example	Z	Tonget (y)	2/51	ypi	JP2
1	0	1	2/3	this of	9 1
2	t	3	11/3	3	9 3
3	2	7	20/3	7	97
4	1/2	(1.75)714	2.16	1.84	1.87

y: Adval Torget value

yal: psedictions using a linear model

ypilype: pordictions using polynomial models of different degrees

Step 2 : Fitting a linear Model A simple linear regoression model is of the form Given the first three data points, we set up the boot-squares equations. (smaller) phiceland bomo (soid) phiseland i=1 3C+m = 3C+m = Nilgmie pleavo at sub social and a $\lim_{n \to \infty} \frac{3}{\sum_{i=1}^{3} x_i^2 y_i^2} = C \sum_{i=1}^{3} x_i^2 + m \sum_{i=1}^{3} x_i^2$ Substituting values : 100 book and startfulli agrees side solving for m and c: m=3, $c=\frac{2}{3}$ Thus, the linear model is: $y=3x+\frac{2}{3}$ Step: 3: Making predictions:

Bias - Variance Trade-OFF Example: There are four Examples with their respective target values consider only for the first three columns in the table, let us selvet the first 8 rows to train the model and the 4th row for testing let us stort with a simple linear model. It is a straight-line fit of the from y=mx+c where m and c are the unknowns. The least-square fit for the first three points gives us the following two equations board on the two unknowns. We of bee nothings took out of polynomial fits to illustrate bias and variance Target (4) yp, Example 2 2/2 111 0 3 3 100 3.67 1.84 7/4 ovenfit ovenfit Step 1: Simple linear Model The first approach is to fit a linear model of the form: y = mx+c -> where-m is the slope, cis the intercept using the least saucones method, two equations are desived. The Summation of target values for training is taken from the first three rows of table 3.7, where the target values (4) are given or : Example

2

0

2

Scanned with OKEN Scanner

Target (y) The total Sum of the targe

> The given linear model is y= mx+c > The Least squares method gives the first summation equation Σ40 = C Σ1+m Σχο ⇒ since there are three training data points: \21=3 And the Sum of the x-values is \(\gamma \gamma^2 = 0 + 1 + 2 = 3 \) ⇒ So, Substituting there into the equation; 11 = 3c+3m This is the first equation used to solve for m and c. y = mx+c 11 = 3x + 3m - 1Step 2: Second Equation: Sum of x14 least squares Equations 000000 Sum of y values: $\sum x_i^2 y_i^2 = c \sum x_i^2 + m \sum x_i^2$ ∑ y: 3C+m ∑ η; x values x target values ∑y; = c ∑1+m ∑ χ; (0X1) + (1X3) + (2X7) = C(0+1+2)+ m (02+12+22) 0+8+14 = c(3)+m(1+4) (Parget Values Sum) 17 = 3C+5m This Second Equation 1+3+7 = 3c + m (0+1+2) (sum of n values) 17 = 3c+5m This gives us the first equation: 11=3c+m3

step 3: solving for c and m Now we have the system of equations 11 = 3c+8m (2) 17 = 3C+5m Substract the first equation from the second equation. (3c+5m) - (3c+3m) = 17-11 am = 6 m=3Substituting m=3 into the first equation 11 = 3c + 3(3)11 = 30 +9 30=11-9 So, the final linear model is: $y = 3x + \frac{2}{3}$ For this data of three points, we can also fit (exertit) a degree-3 polynomial. In general we can have infinite polynomials of degoce 3 that can fit the three training patterns, we consider two such polynomials that capture all the three training examples, that is the 3p=1+ax+bx2+cx3 first 3 rows, perfectly. The polynomials one y, for 1 3 x + 1 x2 + 1 x3 UP = 1+3 2+ 4 27+ 4 x3 3P2 = 1+5 x+ = x3

Depute-3 Polynomial:
$$(9p_1 = \frac{3}{2}x + \frac{1}{4}x^2 + \frac{1}{4}x^3)$$
 $9 = a_0 + a_1x + a_2x^2 + a_3x^3$

We use the three transaming points

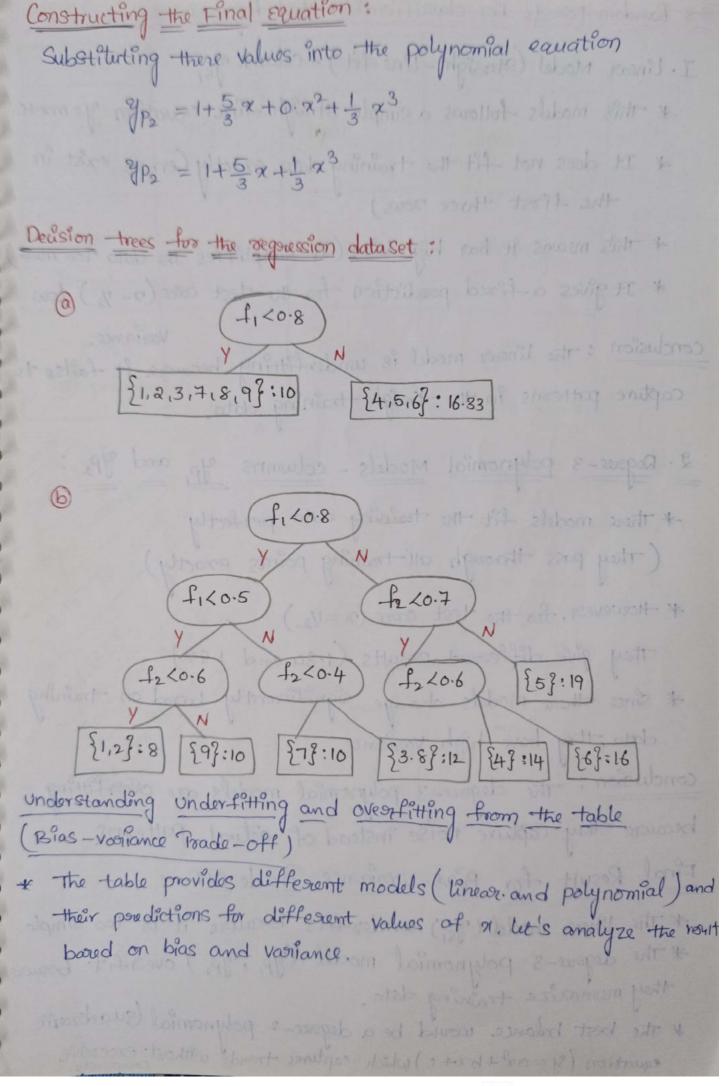
 $0 (x_1, y_1) = (0, 1)$
 $(x_2, y_2) = (1, 3)$

Substituting there values into the polynomial equation

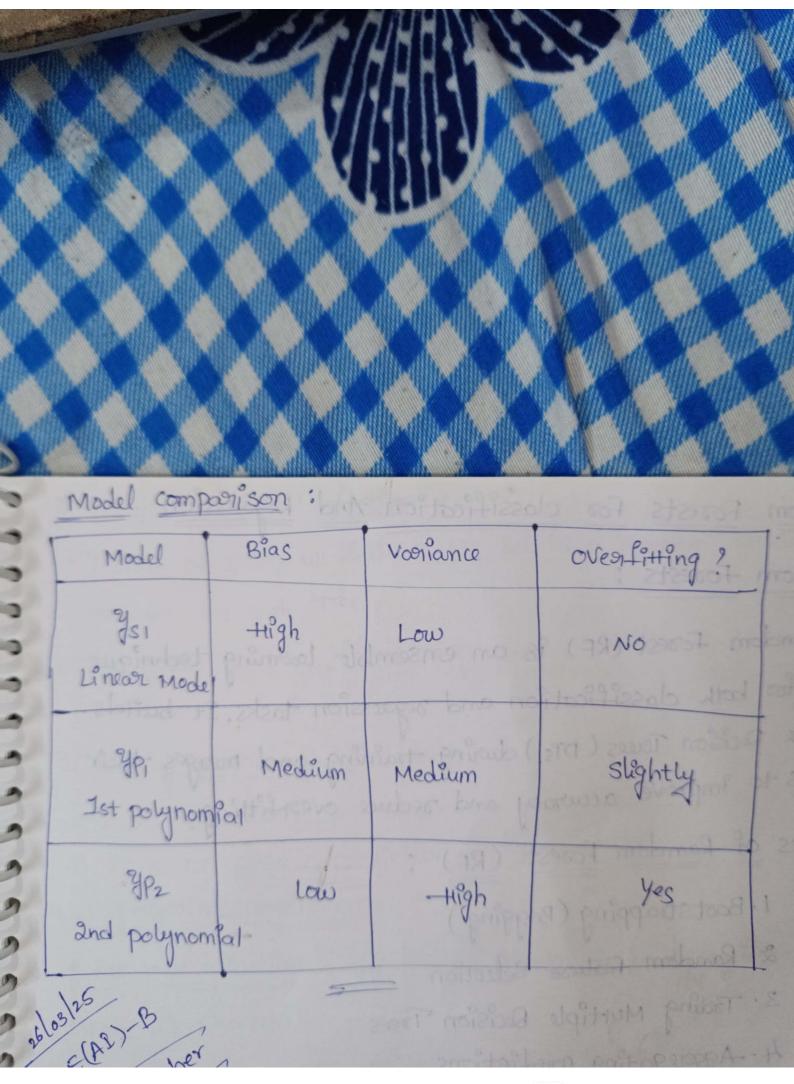
 $(x_1, y_1) = (0, 1)$
 $y = a_0 + a_1x + a_2x^2 + a_3x^3$
 $x = 0, y = 1$
 $1 = a_0 + a_1(0) + a_2(0)^2 + a_3(0)^3$
 $1 = a_0 + a_0 + a_0 + a_0 + a_0$
 $1 = a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0 + a_0 + a_0$
 $1 = a_0 + a_0$

(73, y3) = (2,7) x=2, y=7 y = a0+a1x+a2x2+a3x3 7 = 90 +91(2) +92(2) +93(2)3 7 = 90+ 2(91)+4(92)+8(93) 2(a1)+4a2+8a3=1-7 291+492+893=6 Now solving the system: 1. 91+92+93 =2 a. 201+492+893=6 Divide the second equation by 2: $\frac{2}{3}a_1 + \frac{4}{2}a_2 + \frac{8}{2}a_3 = \frac{6}{2}$ 101+202+403=3 $a_1 + 2a_2 + 4a_3 = 3$ Now Substract the first equation from this (a1+2a2+4a3) - (a1+a2+a3) = 3-2 02+803=1 choosing reasonable values that satisfy the System, we get. $a = \frac{3}{2}$, $a_2 = \frac{1}{4}$, $a_3 = \frac{1}{4}$ Thus, the final equation for yp, is yp, = 1+3x+ 4x+4x This polynomial perfectly fits the three training points.

Degove-2-	polynomi	ials	(3p2=1+5x+1x3)
Example	α	Targetly)	
1	0	,	
2	1	3	2 sot x to t x to too = f.
3	2		
polynomial	Equatio	n y= a0+a	117+9222+0323
(71	7,) =(0), (x2, y2	=(1,3),(913)=(2,7)
(n, y,) =		(x21/2) =	
n=0, 1	<i>‡=1</i>	2=1, 4=	
00=1		a1+92+0	13=2 891+492+893=6
			2- 80,+ 992+893=K
Solving =	for aia	2,93	parde the second concile by ?
		e System	至-10年十四年
(D) a,	1+92+0	13 = 2	
(a) a	9,+402	+892 =6	10, +30, 1103 = 3
Deviding	equatic	m 2 by 8:	9,+292+493=3
Now Sub	tract eq	uation, from	this.
			92+93)=3-2
	And in case of the last of the	-393 =1	
choo	sing the	Value - that	Satisfy the System.
2+L22+L	E-1-1 -	$a_1 = \frac{5}{3}, a_2 =$	=0, 93 = 1/3
		the set with p	The third continue



3.8 Random Forests For classification And Regoussion: I. Linear Model (Straight-Line fit) - column ys, * this madels follows a simple straight-line equation y=mn+c * It does not lit the training data exactly (errors exist in 6 the first three rous) 6 * This means it has high bias (it simplifies the data too much) * It gives a fixed possibilition for the test care (n-1/2) low conclusion: The linear model is underfitting because it fails to 6 capture patterns in the traigh training data. 6 2. Dogowe-3 polynomial Models - columns yp, and yp.: * there models fit the training data perfectly (They pass through all-training points exactly) * However, for the test core (n=1/2) They give different results (1.84 and 1.87) * since there models change significantly bared on training 6 data; they have high variance. conclusion: the degoue-3 polynomial models are overfitting because they capture noise instead of actual patterns. Final Result for Bias-Variance Trade Off: * The Unear model (ys,) under fits because it is too simple * The degree-3 polynomial model (yp, 1 yp2) over fit because they memorize training data. * The best balance would be a degoice-2 polynomial (quadratic equation (y = ax + bx + c) which captures trends without excessive completely.



Random Forests For classification And Regoussion Random Forests:

Random Forest (RF) is an emsemble lowning technique used for both classification and regression tasks. It builds multiple Decision Trees (DTs) during training and merges their autputs to improve accuracy and reduce overfitting.

Features of Random Forest (RF)

- 1. Bootstrapping (Bagging)
- 2. Random feature Selection
- 3. Traing Muttiple Docision Trees
- 4. Aggregating predictions.

1. Bootstrapping (Bagging)

Bootstrapping means sampling multiple sets of training data. In the care of RF, it works on follows.

Given n training data patterns in some L-diamensional space, m samples, each of size n data points, are selected with replacement, it is possible for the same data points to occur more than once in a sample.

* Multiple subsets of training data are trained selected with replacement.

\$ Each Subset is used to train an individual Decision Tree 28/03/2015

2. Random Feature selection * Instead of using all features, RF selects a random subset of features for each tree. * This prevents trees from being overly similar and reduces overfiting 3. Training Multiple Docision Prees: * Each tree learns patterns from it's Subset. * Each DT gives a possibilition independently. 4. Aggregating predictions: * For classification: The class with the majority votes from all trees is choosen. * For Regression: The final prediction is the average of all tree outputs. Advantages of Random Forest: 1 Reduces overfitting 2. Handles High - Dimensional Data 3. Handles Missing Data 4. Robust to Noisy Data. 1. Reduces over fitting The averaging effect of multiple trees prevents over fitting & Handles High-Dimensional Data: can work with datasets that have a large number of features 3. Handles Missing Data: can maintain accuracy even with some missing values. 4. Robust. to Noisy Data: Since it aggregates results, noisy data has less impact.

Differences be	tween classification a	nd Ragoussian in RF:	00
Feature	classification (RF)	Regaression (RF)	0
output Type	catigorical labels	continuous values	6
parediction Method	Majority Voting	Averaging productions	8888
Loss Function Unit-iu-2nd Point	Entropy, Gini Index	Mean squared Errors (MSE) Mean Absolute Error (MAE)	1000
The Bayes classifi	of Boyes classifier	Advantages of Rondon	A A A B B
3. the Baye's old	assifier and it's op	and Naive Bayes classifier	00000
		(NBC)	11
			10
			1110
aniby poli			100
		the principal of the same and t	1

Introduction of Boyes classifier: - A Bayes classifier is a simple but powerful classification algorithm that applies Bayes Theorem to predict the class label of a given instance It is borred on the assumption that the probability of a given Existramer, class given the input data can be completed using posion knowledge about the distribution of data. used applications: Y. spam detection 2. Medical diagnosis 3. Sentiment analysis 4. Document classification Roger Rids one Sufferiors Example: let us assume that 10.000 people in a community had undergone the covid-19 text and 50 of them tested positive, while the remaining 9950 tested negative. In this two-calm (binary class) problem, a simple forequency estimate will give us the following values for the probabilities for the two classes. Understanding prior probabilities: In the example 10,000 people are tested, and the results are * 50 tested positive -> probability. P (positive) = 50 P (Negative) = 9950

classification Board on porior probabilities:

If a new person from the community needs classification without any additional Photomation.

- * since p (negative) = 0.995 is much greater than p (positive) = 0.005, the person will be classified on covid-19, negative.
- * this means that every new person would be labeled or negative because it's the most probable outcome.
- * However, if the person is actually positive, this classification is incorrect with a probability of 0.005.

This method works well only when the perior probability of a class is significantly different, but it fails when dealing with rare events, like covid-19 positive cases in this example.

Relation to k-Newvest Neighbors (KNN) classifier:

In KNN classification, a new Sample is classified bared on the majority class among it's k-nearest neighbors.

If k=n (all training data is used), them KNN simply estimates probabilities ar : $p(positive) = \frac{k_1}{n}$, $p(negative) = \frac{k_2}{n}$

this approach gives the same result on classification bound on posior probabilities.

Additional information Such an: 1. Symptoms

00

probability Rule: p(A) [[0,1] -Addition Rule for Disjoint Events: If two events A and B cannot happen at the same time. (i.e -AnB = 0). then p(AUB) = p(A)+p(B). disjoint. 2. Baye's Rule and Inferience Formula for Baye's - theorem is p(cla) = p(Alc) p(e) P(A)

Where AUB is the union of the sets A and B. This property holds for a countable union of events if they one pairwise

Formula - for Baye's - theorem is
$$p(clA) = p(Alc) p(c)$$

$$p(A)$$

$$p(c_1) = p(c_2) = 0.5$$

$$p(wB|c_1) = \frac{20}{30} = \frac{2}{3}$$

$$p(wB|c_2) = p(RB|c_2) = \frac{15}{30} = \frac{1}{2}$$

$$p(c_1|wB) = \frac{p(wB|c_1)}{p(wB)}$$

$$p(wB) = p(wB|c_1) \ p(c_1) + p(wB|c_2) \ p(c_2)$$

$$p(wB) = (\frac{2}{3} \times 0.5) + (\frac{1}{2} \times 0.5)$$

$$= \frac{1}{3} + \frac{1}{4} = \frac{4}{12} + \frac{3}{12} = \frac{7}{12}$$
applying Bayle's Rule = $p(c_1|wB) = \frac{2}{3} \times 0.5$

$$= \frac{4}{12} + \frac{7}{12} = \frac{7}{12}$$

$$= \frac{7}{12}$$

Example 2: covid-19 Test -Accuracy We are given of covid-19 P(A) = 0.001 not having * probability of covid 19: P(AC) = 1-0.001 P(AC) = 0.999 * Tour positive Rate (sensitivity) P(B/AC) = 0.98 * False positive Rate = p(B/AC) = 1-0.98 = 0.02 We want to find p (AlB), the probability that a patient who tested positive actually har covid-19. By using Boye's theorem P(AlB) = P(BlA) P(A) P(B) To compute P(B), we use the law of total probability. P(B) = P(B/A) P(A) + P(B/AS) P(AS) P(B) = (0.9 x0.001) + (0.02 x0.999) = 0.0009 + 0.01998 P(B) = 0.02088 Now applying Bayes Rule

P(A|B) = P(B|A) P(A)P(B) $=(0.9\times0.001)$ 0.02088 = 0.0009 0.02088 = 0.0431 P(A/B) = 0.0431 01/04/2025 CSE (AI)-A 11:15 AM

The Bayes classifier and they aptimality The Baye's Rule classifier and it's optimality 6 The Bayes Rule classifier is probabilistic model that assigns a List pattern 2 to a class bared on the posterior probabilities of the classes given Z. It is optimal because it minimizes the probability of classification error. Bayes theorem in classification; From Baye's theorem $p(c_1|x) = p(x|c_1) p(c_1)$ p(2/x) = p(a/c2) p(c2) * P (a/x) and P (2/x) one posterior probabilities (probability of a class given the data)

probability of a class given the data)

prior

probability of a class probabilities

(probability of a class occurring before observing the data) * p (elc1) and p (elc2) are probabilities * p(alc,) and p(alcz)

Department of Artificial Intelligence and Machine Learning

II B.Tech II Sem

Machine Learning

UNIT-4

Linear Discriminants for Machine Learning

Ву

Mrs.K.Reddemma,

Assistant Professor,

Department of AI&ML.

UNIT-4 - Linear Discomminants for Machine Learning D'Enear Risconiminants for Machine Leconing Introduction & Linear Risconiminants for classification (3) perceptron classifier Desceptron learning Algorithm @ Support Vector Machines 6 Linearly Non-separable case 7 Non-linear SVM_ 1 Kernel Trick (9) Logistic Regoussion / (10) Linean Regoussion (1) Mutti-layer perceptrons (MLPs) Back propagation for Braining an MLP. O Linear Disconiminants Introduction linear Disconiminants for classification are functions used to separate datapoints into different classes using a straight decision boundary. Ex: (like a line in 20 (or) hyperplane in highter dimensi ons). It works best when the data is linearly separable. The goal is to find a function g (x) = Wx+b where or - Input feature vector W - Weight W Vector b - bias (threshold)

3. classification Rule

If g(n) >0 -> class I

If g(n) <0 -> class 2

g(a) =0 -> point lies on the decision boundary

puspose: To separate data points bared on their features used in binary classification problems.

Applications of linear Discouminants

1. peaception

a. Logistic Regoussion

3. support vector Machines (svm)

4. Linear Disconiminant Analysis (LDA)

visualization: 1. In 2D: Its aline

2. In 3D: It's plane

3. In higher dimensions: It's a hyperplane.

Example problem: Find a linear discouminant-function 9 (7) = wx+b that can separate the two classes.

Two classes of points in 20

class 1 (label = +1):
71=(2,3), 72=(3,3)
100 0 (10/01 = -1)

$$closs_2$$
 (label = -1)

Data points	tabels.	label2
NI	2	3
1/2	3	3
α3	0	1
24	- 4	1)

step-by-step Solution: We want a line that separater the two classes. Assume a disconnant function. g (a) = w1711 + w2 x2 +b let's tony w = (111), b=-4 so, g(n) = n1+x2-4 For class 1: 9 (213) = 2+3-4=1>0 -> class+1 g (3,3) = 3+3-4 = 2>0 -> class+1 For class 2: 9 (0,1) = 0+1-4=-320>class-1 9 (111) = 1+1-4 = -2 <0 -> class -1 A valid linear disconiminant function is: 9 (91) = 91+92-4 9(9)=71+712-4 Decision boundary 3 class I Class 2 class 2 2 ELAN B 25 (41 20 25 (41 20 25 (41 20 25 (41 2) - 4 2 (25 2) 2 (2 711

(3) pesicoptron classifier

The perceptron is a linear classifier that takes to find a staright line to separate two classes of data,

Mathematical formula: g(x) = Wx+b

steps for peaceptron

- 1. Initialize weights W and bias b
- 2. For each data point poudict class using current weights, if paediction is warng, update the weights.

Where:

y-18 actual class

n - is the learning vate (Small positive values)

Example problem: We have the following 20 data points feature x, x2

2(1	22	class (y)
2	1	-+1
1	-1	1
2	3	+1
-1	-2	-1

25/4/2025 B	
25/4/2025 B CSE (A?)-B	
25/2/27-0	
(AZ) AM	
10.00	
1055	
· Is	
a. Me A	
70 / 110	Ų
1 1 CHY	ű
160	
951 //	
121	
6	
10 1.651	
/ L /NV	
21	
62193,8217 3 Member A 62193,8217 12:15 - CSE (AC) - A	
0.10	
19	
9/1	

1 Initial Setup: start with W=[0,0], b=0, looning rate 7=1 point (211), y=+1 porediction -> g (x) = w, x, + w2 x2 +b w=w+nyx = [0,0]+1*1*[2,1] W = [211] b = b+1 *1 b=1 point (1,-1), y=-1 9(1) = 2*1+1*(-1)+1 = 2-1+1 9(1)=2 1 g = +1 point (213), y=+1 9(9)=1*2+2*3+0 = 2+6 = +1

1 Training: we update weights only if the points is misclassified. =0 *2+0*1+0 The perceptron learning Algorithm 18 a simple superivised borning algorithm used for binary classification. It helps find a linear decision boundary that separater data into two classes Ex: class +1 and class -1 point (-1,-2), y=-1 9(9)=1*(-1)+2*(-2)+0 --1-4 9 (9) = -5 Final output : Final Weights w = [12] Final bior b=0 Deusion Boundary: $x_1 + 2x_2 = 0$ (or) $x_2 = -0.5x_1$

Backpropagation for Training an MLP (Muttilayer perceptron): Backpropagation is used to train a newsal network by adjusting the weights to minimize the output error. od () o wing of It works in two main steps: 1. Forward Pass: * you start with random weights. * You apply the input at the input layer. * Each layer processes the inputs and passes the output to next layer * At the end, you get the network's output y_obt (obtained output). & Backward pass (Backpropagation) * Calculate the error: difference between y tar and youte * propagate the error backward through the network. * use gradient descent to update the weights to reduce = The basic weight update onle i's: Wupdated = Woment - n x VE (Woument) * n = learning rate (how big a step you take) where: * VE (Wourset) = gradient of the error with respect to weights. Gradient descent tries to move the weights towards a local minimum where the error is small.

) Forward pan E(W) problem: Forror let warrent = (0.5,0) and the pattern x = (112). Let $E = (\omega_1 \chi_1 + \omega_2 \chi_2 - 1)^2$ >W Here JE (wourrent) is obtained on follows: WZ optimal DE = 2 (w12/1 + w22/2-1)2/1

DE = 2 (w12/1 + w22/2-1)2/2

DE = 2 (w12/1 + w22/2-1)2/2 Gradient descent for minimizing error Given cursent weights Wourrent = (w(1), w(2)) = (0.5, 0) $X = (n_1, n_2) = (112)$ E = (W1 71 + W2 · 72 -1)2 square error loss y = w1 1/1 + w2 1/2 with torget output =1 Goal: compute the gradient of the error with respect to each weight. V E(W) = (DE , DE)

1 Apply chain Rule

$$\frac{\partial E}{\partial \omega_1} = 2(\omega_1 \eta_1 + \omega_2 \eta_2 - 1) \cdot \eta_1$$

Now compute gradients

$$\frac{\partial E}{\partial \omega_1} = 2(-0.5)(1)$$

Gradient vector

$$\nabla E(\omega_{\text{current}}) = (-1, -2)$$

Capally in the contract of the

Kennel Trick

Keemal Trick an important concept in machine learning especi ally in Support Vector Machines.

I. Data is not linearly separable in it's original space

2. We might wount to map data to a higher-dimensional space to make it separable:

3. But computing this mapping explicitly can be hard (or) inefficient.

Instead of mapping vectors emplicitly using a function $\phi(z)$, we compute the dot product in the feature space directly using a kernel function.

$$K(n^i, x^j) = \phi(n^i)^T \phi(x^i)$$

This Lunctions veturn the dot product in higher-dimensional. space without doing the transformation.

polynomial Konnel (degouep):

$$K(n_i, x_j) = (n_i^T x_j)^P$$
 $K(n_i, x_j) = (1 + n_i^T x_j)^P$
 $K(n_i, x_j) = [1 + n_i^T x_j)^P$
 $K(n_i, x_j) = [1 + n_i^T x_j]^P$
 $K(n_i, x_j) = \tanh(\alpha n_j^T x_j^T + b)$

polynomial Konnel of Rogoree 2 To show how the kennel torick works, the book shows,

$$\phi(\alpha) = \int \alpha(0)^{2}$$
 $\sqrt{2} \alpha(0)^{2}$
 $\sqrt{2} \alpha(0) \alpha(0)$

Then

$$K(n, x) = \phi(\alpha)^{T} \phi(\alpha)$$

$$= n(1)^{4} + n(2)^{4} + 2n(1)^{2} n(2)^{2}$$

$$= (n(1)^{2} + n(2)^{2})^{2}$$

This matches (nTn)2, so the kearnel function gives the same result ar explicitly computing in higher-dimensional space.

Mutti-layer perceptions (MLPs)

MLPs are a class of feed-forward artificial neural networks that consist of multiple layers of neurons. stoucture:

Input Layer: Accepts input

Department of Artificial Intelligence and Machine Learning

II B.Tech II Sem

Machine Learning

UNIT-5

Clustering

By

Mrs.K.Reddemma,

Assistant Professor,

Department of Al&ML.

UNITS chapter -7 Raw Data -> Algorithm -> oranges clustering -Introduction to clustering: clustering is a unsupervised learning technique used to group 1 similar data points together board on their features. It is commonly used for pattern recognition, customer segmentation, anomaly detection, 0 and data compression. clustering refers to the process of arranging (or) organizing 6 objects according to specific contenia. clustering involves dividing or grouping the data into smaller data sets based on similarities dissimi-100 locaties. Depending on the requirements, this grouping can lead to 13 various outcomes, such ar particularing of data, data re-organization 0 1. Moorket Segmentation Compression of data and data summorization. -2 statistical data Ex: shopping malls, Super markets. analysis Perovigence 3. Soual network -Fauits regitables -0 analysis 9 4. Image Segmentation Partitioning of Data: (Unstering Method) 5. Anamoly detection Wear -10 -0 portitioning of data in MI refers to the perocess of dividing -0 a dataset into smaller subsets for efficient processing training -0 and evaluation. this helps improve model performance, reduce overfither optimize computational efficiency. 0

Ex: Employee-Table

1. There is an Employee data table containing 30,000 seconds.

2. The table includes a column Dept_code, which has 1,0,00 records distinct values.

3. The records are evenly distributed among there 1,000 deposit

So each department har 30,000 Records 1,000 departments

30 records per deportment

= 450,000 accesses (Without clustering)

With clustering;

Dept - code = 15

1000 departments the Second Uses binary search log_ (1000) accesses.

Approximate value = log (1000) & 10 accesses.

* Without clustering = 45,000 accesses

with clustering = log_ (1000) + 30 2:15 pm To 3:15 pm

02/04/2025

CSE(AI) - A

(3) Matrix Factorization: Matrix Factorization is a technique that can be used to represent clustering mathematically. It allows us to express a dataset on the product of two matrices, [X & Bc] Where: * X - is the original data matrix * B- 18 the cluster assignment matrix (Indicater which data points belong to which clusters) * c - is the representative matrix (contains the cluster centroids or leaders) This approach helps in understanding clustering ar a facto rization problem, where we break down a large data matrix into meaningful components. Example: Consider the dataset $X = \begin{bmatrix} 6 & 6 & 6 \\ 6 & 6 & 8 \end{bmatrix}$ duster I a 4 2 Here, there are 4 data points in a 3-dimensional space, it We use the leader Algorithm with a threshold of 3-units, we obtain two clusters. duster I: (6,6,6) and (6,6,8) with (6,6,6) or the leader.

duster 1: (6,6,6) and (6,6,8) with (6,6,6) or the leader. duster 2: (2,4,2) and (2,2,2) with (2,4,2) as the leader.

Each vow in B reposesents a data paint, and the I's indicates
the assigned cluster. The cluster Reposesentative Matoix C

(leader of cluster) is:

Thus, the matoix factorization: X & BXC

Hard clustering (Vs) soft clustering

Hard clustering: (leader Algorithm)

Each data point belongs to only one cluster (values in B

are either 0 or 1)

Soft clustering: (Fuzzy clustering)

A data point can belong to multiple clusters with different probabilities, (values in B are between o and I, and each

you sums to I)

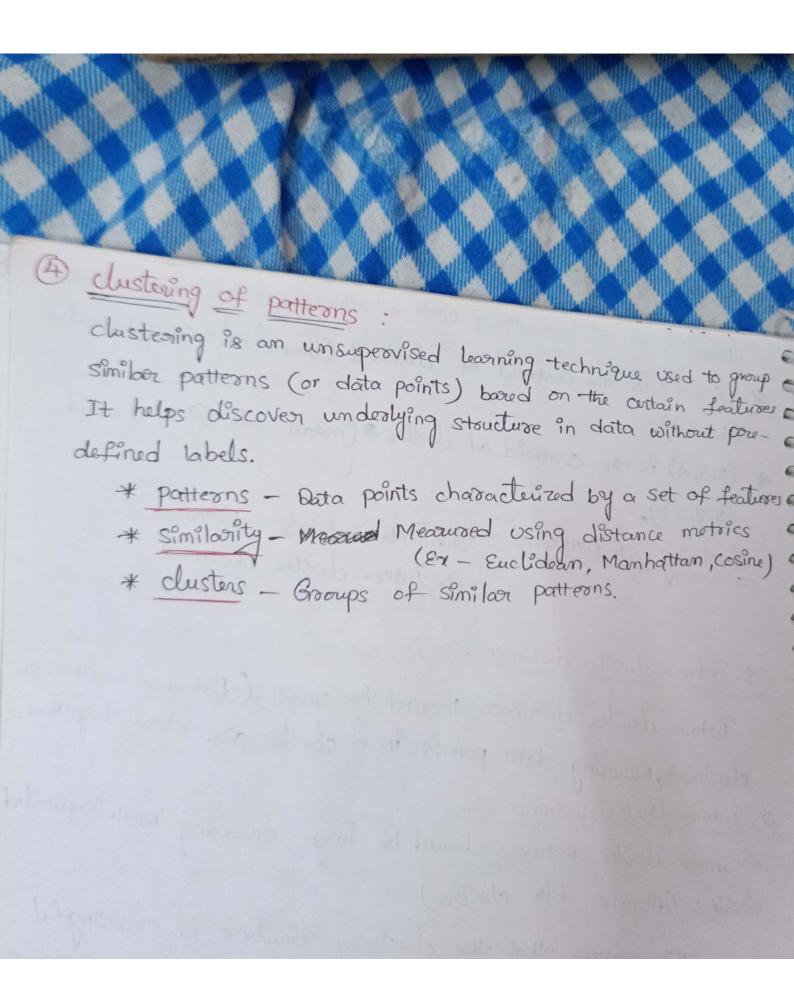
Example of soft clustering: $B = \begin{bmatrix} 0.8 & 0.2 \\ 0.7 & 0.3 \\ 0.1 & 0.9 \end{bmatrix}$

Using centroids Instead of leaders
Instead of using a single leader for a cluster we can compute
the centroid (mean) of points in each cluster. The updated
matoix is

1111

Here * (6,6,7) is the centroid of cluster I (mean of (6,6,6) and (6,618)). * (213,2) is the centroid of cluster 2 (mean of (2,4,2) and (2,2,2)). clustering of patterns : I. Intoa-cluster distance 2. Inter- cluster distance 7. Intra-cluster distance Intra-cluster distance should be small, (distance within a cluster), meaning data points in a cluster are close together. 2. Inter-duster distance Inter-cluster distance should be lærge, ensuring well-separated clusters. (distance blw clusters) This ensure that the clustering structure is meaningful and useful for applications like image processing, customer segmentation and recommendation system. (5/58/62) aloulands CSE (AI) -B

9:15 Am TO 1



Break down : I. Matrix c with Hoord clustering (Leader Algorithm) The representative value (of one directly picked from existing data points. Given that c = (6,6,6) and (2,4,2) this means one cluster is deposemented by (6,6,6) and the other by (2,4,2) 2. Matrix C with Controld-Bared clustering: Instead of choosing a specific data points the representative of a cluster is computed at the centroid (mean of points in the cluster) The Centroid is calculated as follows; Frost cluster: { (6,6,6), (6,6,8) } Controld $\left(\frac{6+6}{2}, \frac{6+6}{2}, \frac{6+8}{2}\right) \Rightarrow \left(\frac{12}{2}, \frac{12}{2}\right) = (6, 6, 7)$ Second cluster: {(2,4,2), (2,2,2)} Centroid $\left(\frac{2+2}{2}, \frac{4+2}{2}, \frac{2+2}{2}\right) = \left(\frac{4}{2}, \frac{6}{2}, \frac{4}{2}\right) = \left(\frac{2}{3}, \frac{3}{2}\right)$

Divisive clustering in Machine Learning ,

Divisive clustering is a top-down hierarchical clustering approach. It starts with all data points in one large cluster and then recursively splits the clusters into Smaller ones, until each point is in Pts own closter or a stopping condition is met.

How Divisive Clustering Morks:

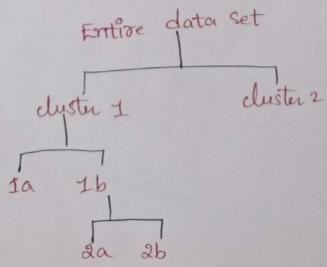
- 1. Start with all data 9n one cluster.
- 2 At each step: & choose the cluster to split Coffen the one with the highest dissimilarity or largest size).

& Split it into two clusters based on a certain criterion (eg. distance, variance)

3. Continue Splitting until:

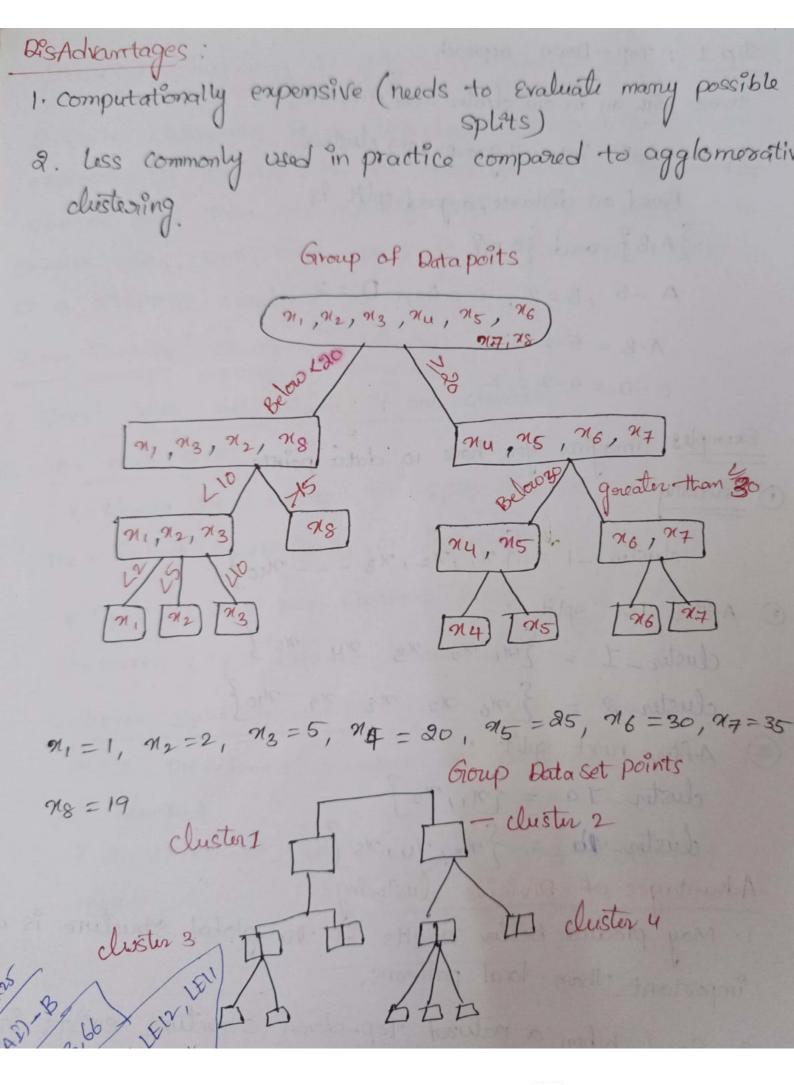
* A predefined number of clusters clusters is reached.

* A threshold distance or dissimilarity is met.

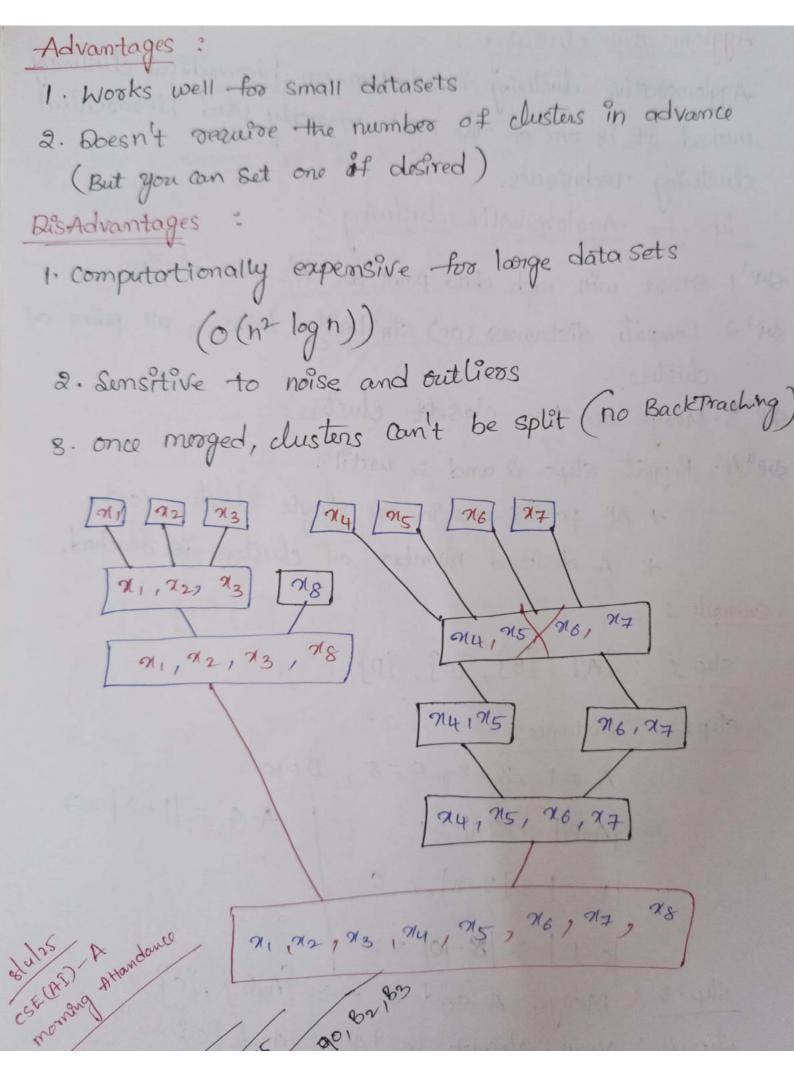


Step I: Top-Down approach Start with all in one cluster -> {A,B,C,D} step. 2 : split into 2 groups (or) clusters Board on distance, a good split is ZAIB and Ecip? A=5, B=2, C=4, D=2 C-D = 4-2 = 2 Examples Imagine you have 10 data points 1 Initially chister -1 = { 71, 12, 12 -- no } 2 After 1st split cluster_I = {01, 22 23 24 cluster-2 = { no xx no x9 no? (3) After next split: duster-Ia = {11, 72} duster-16 = { m3, mu, n5 } Advantages of Divisive clustering 1. May produce better results if the global structure is more important than local partierns. 2. Good When a natural top-down structure exists in the data

Scanned with OKEN Scanner



Agglomosative clustering Agglomorative clustering is a bottom-up hierarchical clustering method. It is one of the most commonly used hierarchical clustering techniques. steps for Agglomerative clustering: step 1. Start with each data point on it's own cluster. stop 2. Compute distances (or) similarity, between all pairs of clusters ste 3. Merge the two closeste, clusters. step 4. Repeat steps 2 and 3 until * All points one 9n one single cluster (or) * A desired number of clusters is reached. Example Step 1: {A}, {B}, {C}, {D} step 2 : Distances. A=1,8=2, C=8, D=10 A-C = |1-8|=7 · | A-B | = (2-1) = 1 |B-c| = (2-8) = 6 |c-D| = |8-10| = 2 | step-3: Merge A and B: -> {A,B}, {c} step. 4: Next closest is EA,B3 And EC8 merge -> {A1B,c}



2	portitional d	ustering :	Euclida	an distance	Loomula
9 -	Data points	f, f2	\ (n,-g,	32 -1 (912.	-y_)2
3	N,	1 1	where		V
9	72	9, 3		() 30 "	1 to point
9	23	3 3			, data point
9	214	C2 3 3	(3)	y) 98 th	ie cluster conter
0	75	8 2	Initial.	clusters :	
0	n 6			- (1,1)	
9	717	C3 7 9			
-	718	8 9	C2	- (712)	1 30
0	Example of do	ta patterns with		- (719)	2000
0	Example: 071	to patterns with	two tealinges	e treat	
0	0°01 0	1 (11) to	$c_1 = (1,1)$		
9	oustance - from	m \ (1-1)2+	(1-1)2 = 1	0+0 = 0	
9 6	Distance -	from $x_2 = (1/3)$	14.	•	
3	112.	12 = 013) to t ₂ =	(7,2)	
•	(1-+) +	(3-2)2 = \(\frac{36+}{}	-1 = √37		
0	3) Distance	from 12 = (3	(3) to C2	=(7,9)	
-	(2-7)2	+(3-9)2 =	1/10		
-3	70 13	70 7	16+36 =√	52	
0	taclidean dis	stances with clus	ster assignm	ents.	
2	Data points	Eudidean	V		
2	X1	distance from	C2	C3	Assigned cluste
2	72	6	. \37	10	2
2	23	2	87	V72	
2	24	18	(F	V52	C,
2	75	√37	0	7	C2
2	76	√68 √50	15	150	C2 C2
2	27	10	7	0	C ₃
		1115	167	O Sca	nned with OKEN Scanr

Euclideam distance Petrinulo

$$d = \sqrt{(f_{1}x - f_{1}c)^{2} + (f_{2}x - f_{3}c)^{2}}$$

$$\pi_{1} = (f_{1}) \quad \text{To } c_{1}$$

$$\sqrt{(f_{1} - f_{1}c)^{2} + (f_{2}x - f_{3}c)^{2}}$$

$$= \sqrt{0.4489 + 1.7689}$$

$$= \sqrt{2.2178}$$

$$\approx 1.49$$

$$To - c_{2}$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

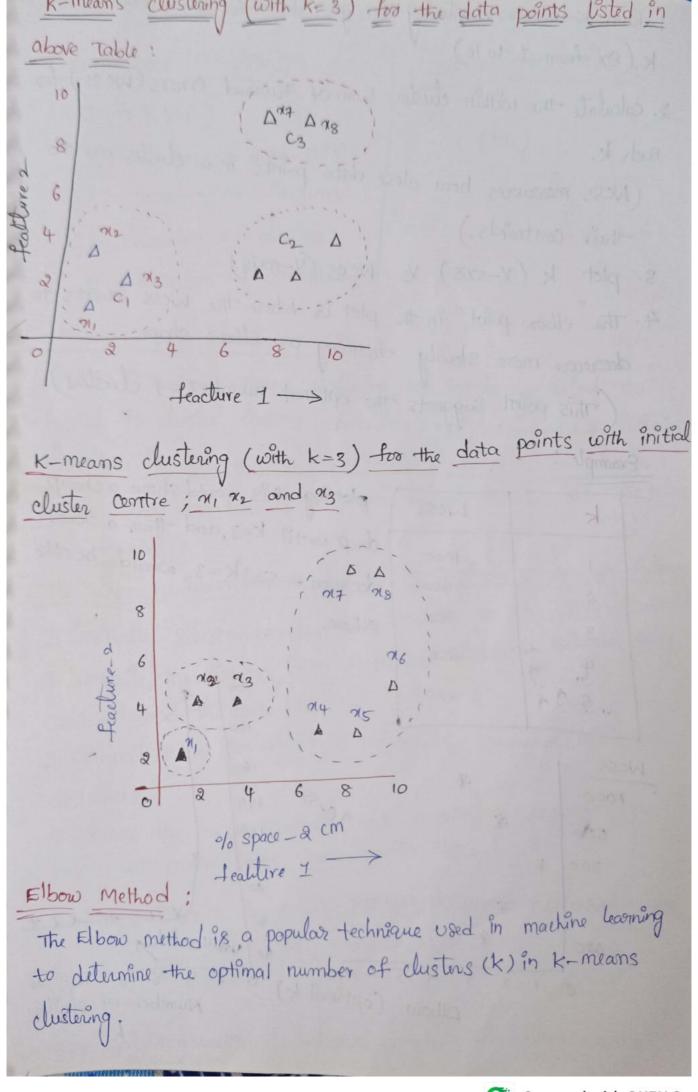
$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689}$$

$$\approx 7.13$$

$$\sqrt{(f_{1}-e)^{2} + (f_{2}-e)^{2}} = \sqrt{49 + 1.7689$$



1. Run K-means clustering on your data too different values of K. (En from I to 10) 2. calculate the within-cluster sum of squared errors (wess) for each k. (Wess measures how close data points in a cluster one to their controids.) 3. plot k (X-axis) vs wess (y-axis) 4. The elbow point" in the plot is where the wess starts to decreare more slowly forming an elbow shape. (This point Suggests the optimal number of clusters) Example: plotting this would show a sharp WCSS drop until k=3, and then a slow. k decrease - so k=3 would be the 1000 500 300 elbow. 250 200 220 5 WCSS 160 * 1000 140 * 500 120 100 300 80 250 60 value of K 220 8 Elbow (optimal k) Number of clusters

Elbow method for Elbow the k- value.

Elbow mother

K-MERCHS + A chastering Time & space Complexity 0 (n.1.k.p) n = number of data points 1 = number of feactures k = number of clusters p = number of Iterations K-Means + + clustering

space complexity

o(kn)

In the case of k-means clustering, we randomly select the initial k-cluster centres (controids) The k-means ++ clustering algorithm is mainly used for identifying the initial cluster Contres. The idea here is to good representatives of the clus The steps to identify the initial k-centroids one on follows ters fromed. I. Initially select a random data point's as the first controid 2. Compute the distance (say Euclidean distance) between any Controld and the data point. 3. compute the probabilities among the computed minimum 4. choose the next controld with maximum probability among all probabilities computed in step3. 5. Repeat steps 2 to 4 till we get k-number of controids P/ 213,50

Example: Co	onsides the	data set with -	following da	ita points	
Prost Rand 1/2 = Distance from	(8,15), $x_3 =$ lomly chasen contine (9,14), $x_5 =$	(6,10), x4 = 000 000 000 000 000 000 000 000 000	(4,4), 26,10 (2,14) (4,14) (6,10)	6 = (11,2) (8,5) (8,5) (9,4) (1,2) (10 Colitroid uve I	
Formula:		= (21-21)2+		Set	
step-by-step		privateula eruna			
Data point	coordinate	s Distance from 11,	Distance	Distance 1	
71,	date	2-10-42	in publicati	· 8-410	
χ2	(914)	$(8-9)^{2}+(5-9)^{2}$ $\sqrt{(1)^{2}+(1)^{2}}$	1.41	2 122	
7/3	(6,10)	$\frac{(8-6)^2+(5-10)^4}{(5-4)^2+(5-4)^2}$	5.39	29	
714	(414)	$\sqrt{(4)^2+(-11)^2}$ $(8-2)^2+(5-4)^2$	4.12	137	
75	(214)	[(6)2+(41)2	6.08	37	
26	(11,2)	$(8-11)^{2}+(5-2)^{2}$ $\sqrt{3^{2}+(13)^{2}}$	4.24	18	
	१ हे वृत्रीय स्त्री		13.34	178	
Sum of distance samore = 122+29+17+37+18 = 103					
Now compute probability for each point: : 623					
probability (i) - distance;					
Total Sum of distance 2					

Data point	Distance ²	probability	
92	122	2/1031 ≈ 0.019	122 = 0.195 623
Ng	29	29 ≈ 0.282	1
X4	137	17 20.165	137 = 0.219
95	157	37 × 0.359	$\frac{157}{623} = 0.25$
716	178	18 & 0.175	
Example:	25 Dist	tance (approx)	
Total		stance (approx)	12.53
	6.08	(12.53)	.08 = 0.020
31.0 20 21.3	1.41+5.39+4.12		1.24
			.286
After Selecting	9 95 or Second	Centroid:	Jandgreet.
Now yo	u have two	Centooids:	: towy year
21 = ((8,5)		The same of the sa
25=1	(214)		probabog
Now fo	s each remainin	g data point, yo	ula
		both centroids	
2 for each p		minimum distance	
centroid.	minimum dist	ances (or their squ	vares) to
compute p	probabilities.	3 60. (100 0)0	. 5 60

1 1	- 1-	-1	-	this:
Le	15	0	0	1.110

	Tenen se			•
Data	d(x1,x;)	d(15,10)	Min. Distance	Probability (a min dist)
×2	1-41	17.0 331	1-41	1.41
23	5.39	7-2	5.39	5·39 15·16 ≈ 0·36
24	4.12	2.0	2.0	2.0 ~ 0.13 15.16
	CX.	stance (appr	9 -8	(but table says
26	4.24	9.22	4.24	0.27, maybe typo or miscal C $\frac{4.24}{15.16} \sim 0.28$
	E-0 6			

"Highest probability = 23 - Chosen as third centroid

Key point:

* K-mean s++ usually uses squared distances for probability.

y But your example uses raw distances, so the probabilities are calculated accordingly. 5/04/2015

Soft portitioning In Machine learning soft portitioning refers to assigning data points to muttiple clusters (or) categories with degree of membership rather than assigning each point to dust one cluster (on in hoord partitioning) Muttiple membership approach: This approvach allows data points to belong to muttiple clusters. rather than aust one (as in traditional hand clustering like kmeans) Methods: I. Fuzzy clustering 2. Rough clustering 3. Newral network-bared clustering 50 stochastic Approach: 1. simulated annealing 2. Evolutionary Algorithms I Fuzzy clustering: Each data point is assigned to Muttiple clusters, typically more than one boned on a membership value. The value is comput a dusing the data point and the corresponding cluster controld. 2. Rough clustering: Each cluster is assumed to have both a non-over lapping port and an overlaping part. Data points in the non-overlaping postion belong enclusively to that cluster, while data points in the overlaping port may belong to multiple clusters. 3. Newsal Network-board clustering: In this method, varying weights CSE AR TO 3.15 PM 2U 1
2.15 PM PM 2 associated with the data points are bared to obtain a soft partition.

2) Stochastic Approach:

In this category, the focus is on achieving the global optimal Solution for cluster formation through probabilistic methods, some existing methods for this approach include.

- I. Simulated annealing
- 2. Tabu Seasch
- 3. Evolutionary algorithms.

I simulted annealing: In this care, the current solution is randomly updated, and the resulting solution is accepted with a certain probability. If the resulting solution is better than the current fromation, it is accepted; otherwise, it is accepted with a probability ranging from o to I.

2 Tabu Search : unlike simulated annealing, muttiple solutions are stored, and the current solution is modified in various ways to determine the next configuration.

3. Evolutionary Algorithms:

This method maintains population of solutions. In Addition to the fitness values of individuals, a random search board on the interaction among solutions which mutation is employed to generate the next population.

Fuzzy C-Means Clustering Fuzzy c-Means clustering (FCM) is an advanced version of k-means clustering used in ML, Especially when clusters over lap and a dapa points can belong to more than one cluster. steps for Fuzzy c-Means clustering I. Initialize: choose number of clusters c. Randomly initialize the membership matrix (values indicating how much a point belongs to each cluster) a update cluster conters: Each cluster center is calculated on a weighted average of all points, weighted by their membeoship values. 3. Update Membership Values Membership of a point a cluster depends on its distance to the cluster center. closer = higher membership. 4. Repeat steps 2 and 3 until convergence (small changes in values) Mathematical formulation: The objective function to minimize is I = Zi=1 Zi=1 21ij · ||ni-cj||2 where: nij - Membership of data point ni in chuster i m - Fuzziness coefficient (usually 2) c; - cluster center of cluster; 12:-59 | - Distance between point and cluster center.

Example: Consider the data points with two feactures,

Example of a data Set with two feactures.

Data points	f,	f2
N	1	1
912	å	2
23	4	3
214	5	3

let the number of clusters c=2 and M=2. To start with assume the following membership values between the data point and the clusters (40)

	911	2	N3	74
c,	110	the second	10	0
C2	0	0	- Clark	1 1

The cluster centres are calculated as follows

$$c_1 f_1 = u_{11} * \alpha_1(f_1) + u_{12} * \alpha_2(f_1) + u_{13} * \alpha_3(f_1) + u_{14} * \alpha_4(f_1)$$

or and red many thingu is

$$c_{1}(f_{1}) = \frac{1^{2} * 1 + 1^{2} * 2 + 0^{2} * 4 + 0^{2} * 5}{1^{2} + 1^{2} + 0^{2} + 0^{2}}$$

C1(f2) = U11 * X1 f2 + U12 * X2 f2 + U13 * X3 f2+ Mi4 * X4(f2). M11 + M12 + M13 + M142 CI(f2) = 12 * 1+12 * 2+02+3 * 02 * 2 c1(f2) = 3 C1(f2) = 1.5 So, the updated cluster centre for c, = (1.5, 1.5) Similarly c2(f1) = M21 * 91(f1) + M22 * 92(f1) + M23 * 93(f1) + M24 * 94 (fi) M2, + M22 + M23 + M24 C2(f1) = 02 *1+02 *2+12 *4+12 *5 $c_2(f_1) = \frac{9}{2}$ c2(f1) = 4.5 ca(f2) = 121 * X1 +2 + 122 * X2 f2 + 123 * 7(3(f2) + M24 * 7/4 (B2) M2 + M22 + M24 + M24 C2(f2) = 02 *1 + 02 *2 + 12 *3 + 12 *3 $c_2(f_2) = 6/2 = 3$ $c_2(f_2) = 6/2 = 3$

So, the updated duster centre for C2 = (4.5,3) member Recalculation of mombership values: d (c1, x1) = d((1.5, 1.5), (1.1)) 1= (1.5-1)2+ (1.5-1)2 d (cx, x2) = d (1.5, 1.5) (2,2) = \((1.5-2)^2 + (1.5-2)^2 = 0.71 d (C1,23 = d((1.5,1.5),(4,3)) $=\sqrt{(1.5-4)^2+(1.5-3)^2}=2.92$ d (C1,94) = d (1.5,1.5) (513) $= (1.5 - 5)^{2} + (1.5 - 3)^{2} = 3.81$ d(c2, n) = d((4.5,3)(1,1)) $= (4.5-1)^{2} + (3-1)^{2} = 4.03$ d(c2, 1/2) = d (4.5,3) (2,2) $= \left((4.5 - 2)^2 + (3 - 2)^2 = 2.69 \right)$ d (c2,93) = d (4.5,3) (413). $= \sqrt{(4.5-4)^2 + (3-3)^2}$ d (C2, 74) = d(4.6,3) (5,3) $= \sqrt{(4.5-5)^2 + (3-3)^2} = 0.25$

number Ship Values one
$$\frac{d(x_{3}, c_{1})^{-1}/(M-1)}{d(x_{3}, c_{1})^{-1}/(M-1)}$$

$$\frac{d(x_{3}, c_{1})^{-1}/(M-1)}{d(x_{3}, c_{1})^{-1}/(M-1)}$$

$$\frac{1}{d(x_{3}, c_{1})^{-1}/(M-1)}$$

$$\frac{1}{d(x_{3}, c_{1})^{-1}/(M-1)}$$

$$\frac{1}{d(c_{1}, x_{2})}$$

$$\frac{1}{d(c_{1}, x_{2})}$$

$$\frac{1}{d(c_{1}, x_{3})}$$

$$\frac{1}{d(c_{1}, x_{3})}$$

$$\frac{1}{d(c_{1}, x_{3})}$$

$$\frac{1}{d(c_{1}, x_{3})}$$

$$\frac{1}{d(c_{1}, x_{4})}$$

$$\frac{1}{d($$

simil	only w	121, 1122,	M23, M24	can b	e computed and the	Sr Sr
value	s ore	0.15, 0.2	1,0.92 an	nd 0.94	respectively.	33
		2(1	72	X3	74	3
	C,	0.85	0.79	0-08	0.06	55
	C2	0.15	0.21	0.92	0.94	3
	CI(fi) = 08			+ 0.08 *4 + 0.06 *	
		(f1) =.	2.04	=1.5		
	C1(f1	L) = 0.	0.852	+ 0.792	2+0.08 ² +3+0.06 ² +0.08 ² +0.06 ²	*3
T .		(f2) = _		.48.		B
C1 =((1.5, 1.4)	1 - 0.15	1.36	0.21 *2+	0.92 *4 + 0.94 *	25
	C2 (+1)		0.152+0	.2 +0.	92 +0.942	一里
		= 4.	4	r.pJE		平
	Co (f2) = 0.	15 *1+	2 . 21 *2	+0-92 * 3+0.943	13
		4	0.152+	0.212+0	-92 + 0.94	T
	C2	(f2) =	2.95.	3) 6	1172.3	79
		updated	C2 clu	ster vali	e c2 = (4.4,2.95	1
				100 +	- puls	T.

$$d(C_1, \mathcal{H}_1) = d(1.5, 1.48)(1.1)$$

$$= \sqrt{(1.5-1)^2 + (1.48-1)^2}$$

$$d(C_{11}\mathcal{H}_2) = 0.693$$

$$d(C_{11}\mathcal{H}_2) = d(1.5, 1.48)(2.12)$$

$$= \sqrt{(1.5-2)^2 + (1.48-2)^2}$$

$$= 0.721$$

$$d(C_{11}\mathcal{H}_3) = d(1.5, 1.48)(4.3)$$

$$= \sqrt{(1.5-4)^2 + (1.48-3)^2}$$

$$= 2.925$$

$$d(C_1, \mathcal{H}_4) = d(1.5, 1.48)(5.3)$$

$$= \sqrt{(1.5-5)^2 + (1.48-3)^2}$$

$$= 3.815$$

$$d(C_2, \mathcal{H}_1) = d(4.4, 2.95)(1.1)$$

$$= \sqrt{(4.4-1)^2 + (2.95-2)^2}$$

$$= 3.919$$

$$d(C_2, \mathcal{H}_3) = d(4.4, 2.95)(4.3)$$

$$= \sqrt{(4.4-4)^2 + (2.95-2)^2}$$

$$= 2.58$$

$$d(C_2, \mathcal{H}_3) = d(4.4, 2.95)(4.3)$$

$$= \sqrt{(4.4-4)^2 + (2.95-3)^2}$$

$$= 0.40$$

updalld membership values are: 0-693 Un = 0.85 0.721 1112 = = 0.79 0.721 + 2.58 2,925 113 3.815 114 = 0-14 6 Similarly 1/21, 1/22, 1/23, 1/24 can be computed and their Values are 0.15, 0.21, 0.88 and 0.86 respectively. 6 Now the updated membership values 11(2) one given below: 214 2/3 24 72 0.14 0.12 0.79 0.85 0-86 0.88 0-21 0.15 Since we do not end up with the same cluster point (or) Some membership values, we iterate for the next step. The updated cluster conters for the updated membership values (1(2)) are ar follows: c, (f1) = (0.85) * *1+(0.79) *2+(0.12) *44 + (0.14)2 *5 (0.85)2+(0.79)2+(0.12)2+(0.14)2

$$c_{1}(f_{1}) = \frac{a \cdot o_{1}}{1.365} = 15$$

$$c_{1}(f_{1}) = 1.5$$

$$c_{1}(f_{2}) = (0.85)^{2} \times 1 + (0.74)^{2} \times 2 + (0.12)^{2} \times 3 + (0.14)^{2} \times 3 + (0.15)^{2} \times 4 + (0.12)^{2} \times 4 + (0.12)^{2}$$

Since there is no change in the cluster centre, the algorithm Stops. The final membership values and cluster allocation for the data points is given below Table. chista formation using fuzzy c-Means algorithm Allocated clusters Data points 462 Ma 0.15 0.85 9/1 0.21 0.79 9/2 0.88 0.12 913 C2 0.86 0.14 214 15)2+(0.21)+(0.92)246 246 Je (100) 1 / 210 92 7.

Rough clustering:

In traditional clustering, each data point belongs to exactly one cluster. But in real-world data, this may not be ideal-some data points clearly belong to one cluster, while others are ambiguous and could belong to multiple clusters. This is where rough set theory helps.

This is abstracted by a rough set which is represented using two sets, there sets are called lower approximation and upper approximation.

1. Lower Approximation:

set of data points that definitely belong to clusters.

2. Upper Approximation:

Set of data points that possibly belong to cluster S.

Ex: Overlap with s but not fully contained.

In Essence:

* Lower - Approx - Certain members

* Upper Approx - possible members

* The difference R(S) - R(S) captures the "roughness

En: Un cortainity

RS = U; Gi Where Gi CS

R(s) = Ui Gi Where Gins = 0

Rough . K-Means clustering Algorithm:

This is a vociant of K-Means that integrate the rough set idea. It assigns data points to clusters with cortainty (lower) (or).

Possibility (upper), and update clusters accordingly.

Inputs :

* n : Number of data points

* X: 21, 1/2 --- 2/n Data points

* k: Number of clusters.

* k i w1, we : Weights for lower and upper Approximations.

* E: convergence threshold.

step-by-step Algorithm:

I Initialization :

* Randomly assign each data point to the lower approximation of exactly one cluster

* Each point also belongs to that cluster's upper approximation by definition.

2. Controid calculation:

the cluster center C3 depends on the nature of points in the

case 1: if only lower approx points exist (clean cluster)

*
$$C_j = \omega_l \cdot \frac{1}{|R(k)|} \sum_{n \in R(k)} n + \omega_l \cdot \frac{1}{|R(k) - R(k)|}$$

 $\sum_{n \in R(k)} R(k) - R(k) n = \sum_{n \in R(k)} R(k) - R(k) = \sum_{n \in R(k)} R(k) - R(k) = \sum_{n \in R(k)} R(k) = \sum_{n \in R$

1. Distance calculation

For each point xi, calculate the Euclidean distance to all cluster controids Ci.

2. Reassignment :

Assign each data point to the lower approximation of the cluster where it is closest (if it's clearly port of that cluster). Or to the upper approximation of multiple clusters it it's ambiguous.

3 Repeat :

Iterate steps 2-4 until changes in assignments are less than the threshold E.

Intuition

Rough K-means is more flexible than regular K-Means.

* It handles uncertainty

* porevents forcing ambiguous points into dust one cluster

* Gives more realistic clustering, especially in noisy (or) overlapping data.

Example: Consider the data set shown in Table with k=2, $w_1=0.7$ Wu = 0-3 and € = 2

Example of data set with two feature

Examp.			1 and each data points to
Data points	1 41	f2	1. Randomly assign each data points to
	-		exactly one lower approximation (two cluster)
χ,	1		5, 1 (2.22 (1.1)?
	1 37 500	100	$R(k_1) = \{(1,1), (2,2), (4,1)\}$
72	2		
73	4	1	R(K2) = {(5,2), (4,0), (8,0)}
200	01143	No log	a. Since (i) R(ki) + O and R(ki)-R(ki)
74	5	2	TO TO TO TO THE RICH TO THE TOTAL
25	1	0	Day 1 Det 1 Rike = B
			\mathbb{O} $R(k_2) \neq 0$ and $R(k_2) - R(k_2) = 0$
26	8	0	the centroid is calculated using
			$C_j = \sum_{x_i} \in R(k) \frac{x_i}{ R(k) }$

$$C_1 = \left(\frac{1+2+4}{3}, \frac{1+2+1}{3}\right) = (2.33, 1.33)$$

$$c_2 = \left(\frac{5+7+8}{3}, \frac{2=0+0}{3}\right) = \left(6.67, 0.67\right)$$

3. Find the Euclidean distance blw each data point and the cluster

with reference to c,

① d (2,2), (2,33,1.33) =
$$\sqrt{(2-2.33)^2 + (2-1.33)^2} = 0.75$$

(3) d (411), (2.33, 1.33) =
$$\sqrt{(4-2.33)^2+(1-1.33)^2}=1.70$$

(A)
$$d(512) / (2.33, 1-33) = \sqrt{(5-2.33)^2 + (2-1.33)^2} = 2.75$$

(a)
$$d(710) \cdot (2.33,11.33) = \sqrt{(7-2.33)^2 + (0-1.33)^2} = 4.86$$

(6) d (8,0), (2.33, 1.33) =
$$\sqrt{(8-2.33)^2 + (0-1.33)^2} = 5.82$$

With Reference to C2

d (111),
$$(6.67, 0.67) = \sqrt{(1-6.67)^2 + (1-0.67)^2} = 5.68$$

$$d(212), (6.67), (0.67) = \sqrt{(2-6.67)^2 + (2-0.67)^2} = 4.86$$

$$d(411),(6.67),(0.67) = \sqrt{(4-6.67)^2 + (1-0.67)^2} = 2.69$$

1 pries talebales of blooms of

4. Use the vortio d(n, c), where 1<i, P<k, to determine the membership.

$$\frac{d(111)}{d(111)}, (6.67, 0.67) = \frac{5.68}{1.37} = 4.14 \text{ 2.50 x, will be posit}$$

$$\frac{d(111)}{d(111)}, (2.33, 1.33) = \frac{5.68}{1.37} = 4.14 \text{ 2.50 x, will be posit}$$

$$(8.2) \rightarrow \frac{d(2.12),(6.67,0.67)}{d(2.12),(2.33,1.33)} = \frac{4.86}{0.75} = 6.48 \stackrel{1}{\cancel{2}} 2$$

so or will be post of R(ki)

$$\frac{d(411)}{d(411)}, (6.67, 0.67) = \frac{3.69}{1.7} = 1.58 \angle 2$$

so is will not be post of R(k) and R(k2)

$$5(12) \Rightarrow \frac{d(5(12),(2.33,1.33)}{d(5(12),(6.67,0.67)} = \frac{2.75}{2.14} = 1.28 \angle 2$$

so my will not be part of R(k1) and R(k2)

$$4(70) \Rightarrow \frac{d(70),(2.33,1.33)}{d(70),(6.64,0.67)} = \frac{4.86}{0.75} = 6.48 \frac{1}{2}$$

so 95 will be post of R(k2)

$$8(0) \Rightarrow \frac{d(8(0),(2.33,1.33))}{d(8(0),(6.67,0.67))} = \frac{5.82}{1.49} = 3.91 \stackrel{\cancel{4}}{\cancel{2}}$$

So 914 will be poort of R (k2)

Now we have clusters

We have clusters:

$$R(k_1) = \{(1,1), (2,2)\}$$
 $R(k_1) = \{(1,1), (2,2)\}$
 $R(k_2) = \{(4,1), (5,2), (4,0), (8,0)\}$

Here (1)
$$R(k_1) \neq 0$$
 and $R(k_1) - R(k_1) \neq 0$
(11) $R(k_2) \neq 0$ and $R(k_2) - R(k_2) \neq 0$

so the centroid is calculated using. $C_j = w_l \times \sum \frac{\eta_i}{\chi \in R(k)} + w_u \times \sum \frac{\eta_i}{\chi \in R(k) - R(k)} \frac{\eta_i}{R(k) - R(k)}$ $C_1 = 0.7 \times \left(\frac{1+2}{2}, \frac{1+2}{2}\right) + 0.3 \times \left(\frac{4+5}{2}, \frac{1+2}{2}\right)$ C1 = (2.4,1.5) $C_2 = 0.7 \times \left(\frac{7+8}{2}, \frac{0+0}{2}\right) + 0.3 \times \left(\frac{4+5}{2}, \frac{1+2}{2}\right)$ C2 = (6.6, 0.45) 5. Repeat from step 3 until Convergence 3:15 TO 4:15 19/04/2025 (1/42)-A 1:15 AM

Expectation Maximization-Based clustering D computing the probability that a data point belongs to cluster prob (n: e c;) = Zii exp[- = 202 (n: -ui)2 = exp [-1 (4:-41)2] Data points 01 = 1, x2 = 6, x3 = 6, x4 = 7 Initial cluster means cluster I (mu_1) =1 cluster 2 (mu-1) =7 Distance to $mu_{-1} = (1-1)^2 = 0$ $=(1-7)^2=36 \rightarrow enp(-18) 20$ $Z_{11} = \frac{1}{1 + e^{-18}} \int_{0}^{1} |0.9999|$ 212 = 0.0001 Similarly 92 = 6, mu-2=7 $M_1 = \frac{211 \cdot 91 + 221 \cdot 912 + 231 \cdot 93 + 241 \cdot 94}{211 \cdot 912 + 231 \cdot 913 + 241 \cdot 94}$ Dotta point (i) Zi2 Zil 0.0001 0.9999 21 =1 0.0001 0.9999 0.0001 0.9999 0.0001 0.999 94 = 7

