

DM #2 PE3 DM Module I Introduction cont 20210810 090111 Meeting Recording

Tuesday, August 03, 2021 1:02 PM

Add Page

SQL  
OLAP  
levels  
grain

$I = f(D)$

$K = f(I)$

$D \rightarrow I \rightarrow \text{knowledge}$

DM #2 PE3 DM Module I Introduction cont 20210810 090111 Meeting Recording

## Data Mining

How can I analyze this data?

We have rich data, but poor information

Data mining-searching for knowledge (interesting patterns) in your data.

J. Han, M. Kamber, Data Mining: Concepts and Techniques  
Morgan Kaufmann, 2006 (Second Edition)

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg. Sangli, INDIA

## Definition

**The non-trivial process of analyzing large databases by applying specific algorithm to extract interesting**

**(hidden, implicit, previously unknown and potentially useful)**

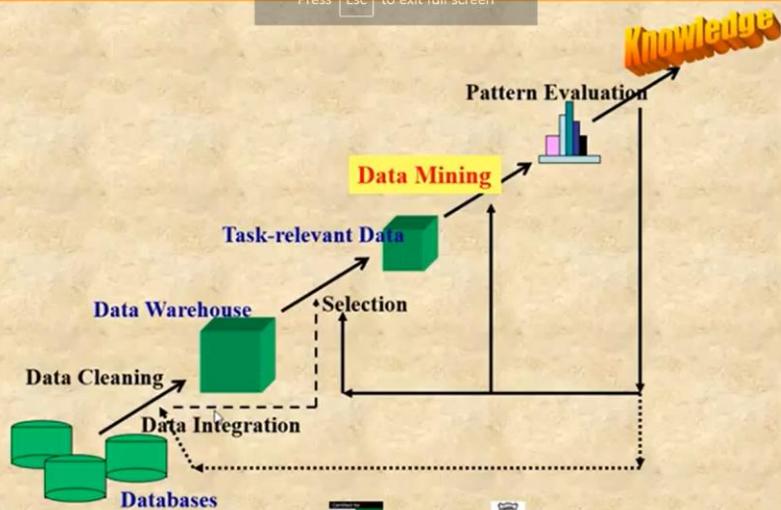
**information (knowledge) or patterns from data.**

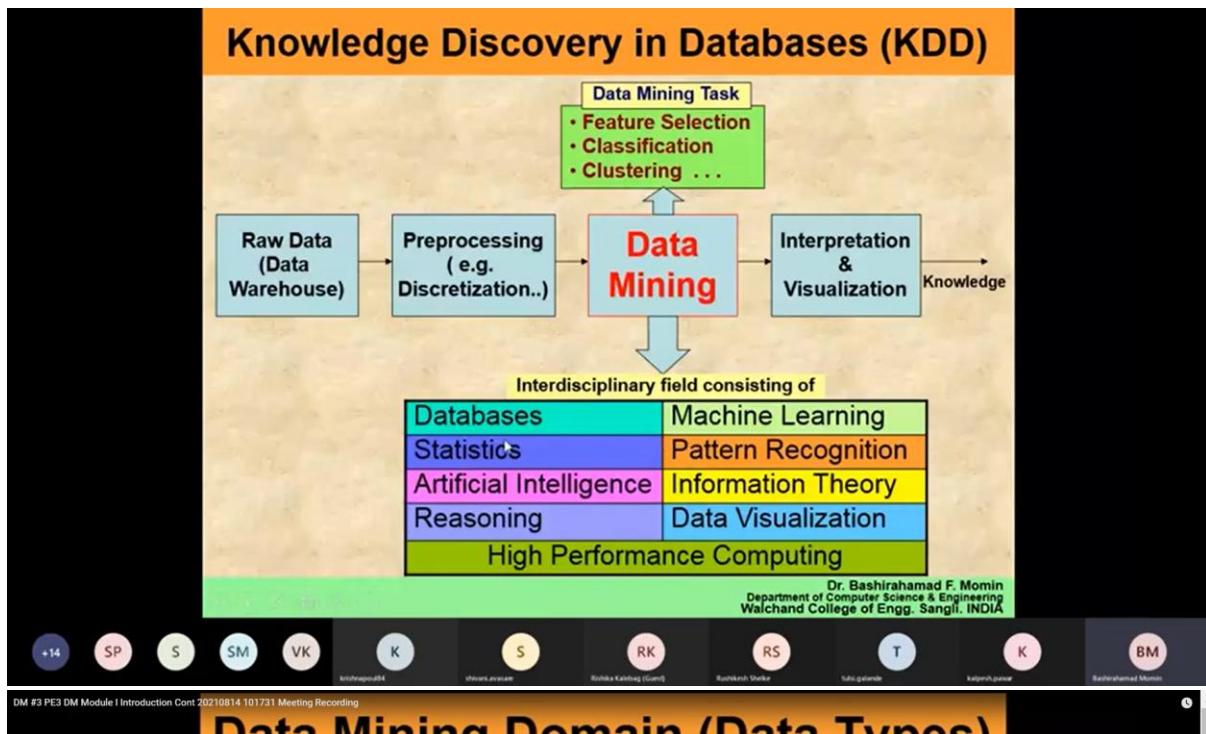
Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg. Sangli, INDIA



## Knowledge Discovery (KDD) Process

Press Esc to exit full screen





DM #3 PE3 DM Module I Introduction Cont 20210814 101731 Meeting Recording

## Data Mining Domain (Data Types)

- Database-oriented data sets and applications
  - Relational database, data warehouse, transactional database
- Advanced data sets and advanced applications
  - Data streams and sensor data
  - Time-series data, temporal data, sequence data (incl. bio-sequences)
  - Structure data, graphs, social networks and multi-linked data
  - Object-relational databases
  - Heterogeneous databases and legacy databases
  - Spatial data and spatiotemporal data
  - Multimedia database
  - Text databases
  - The World-Wide Web

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg. Sangli. INDIA

Participants (bottom bar): +7, A, A, RS, S, VK, S, K, AJ, T, SD, BM.



## Data Mining Functionalities

**Data mining functionalities are used to specify the kinds of patterns to be mined.**

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg., Sangli, INDIA

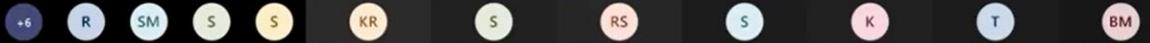


## Data Mining Functionalities

- Characterization and discrimination
- Mining of frequent patterns, associations, and correlations
- Classification and regression
- clustering analysis
- Outlier analysis

**Data mining functionalities are used to specify the kinds of patterns to be mined.**

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg., Sangli, INDIA



# Kinds of Patterns

- **Descriptive :**

- Describes the characteristics of the data in a target data set.
- It determines, what happened in the past by analyzing stored data.

- **Predictive :**

- Carry out the induction over the current and past data so that predictions can be made.
- It determines, what can happen in the future with the help past data analysis.

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg., Sangli, INDIA



## Examples

- **Descriptive patterns**

- Summarizing past events such as sales and operations data or marketing campaigns.
- Social media usage and engagement data such as Instagram or Facebook likes.
- Reporting general trends.
- Collating survey results.

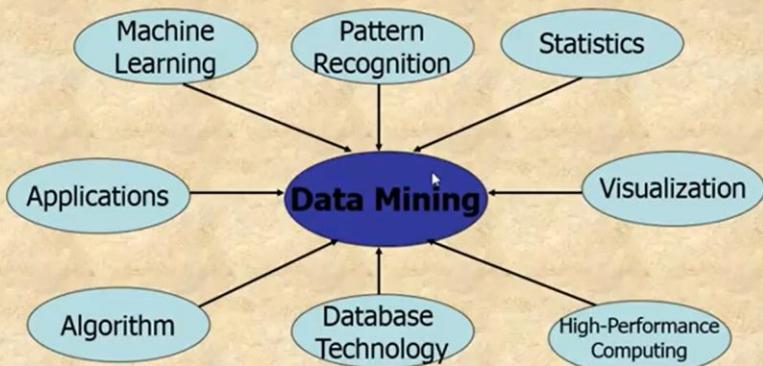
- **Predictive Patterns**

- Predicting buying behavior in retail. ...
- Detecting sickness in healthcare. ...
- Curating content in entertainment. ...
- Predicting maintenance in manufacturing. ...
- Detecting fraud in cyber security. ...
- Predicting employee growth in HR. ...
- Predicting performance in sports. ...
- Forecasting patterns in weather.

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg., Sangli, INDIA



## Data Mining: Inter-disciplinary



17

IBM DB2



+8 R SM Y A KR S RS S K T BM

## Why Multiple Disciplines?

- **Tremendous amount of data**
  - Algorithms must be highly scalable to handle such as tera-bytes of data
- **High-dimensionality of data**
  - Micro-array may have tens of thousands of dimensions
- **High complexity of data**
  - Data streams and sensor data
  - Time-series data, temporal data, sequence data
  - Structure data, graphs, social networks and multi-linked data
  - Heterogeneous databases and legacy databases
  - Spatial, spatiotemporal, multimedia, text and Web data
  - Software programs, scientific simulations
- **New and sophisticated applications**

18

IBM DB2



18

**Which kinds of applications are targeted ?**

*Where there are data,  
there are data mining applications.*

Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.

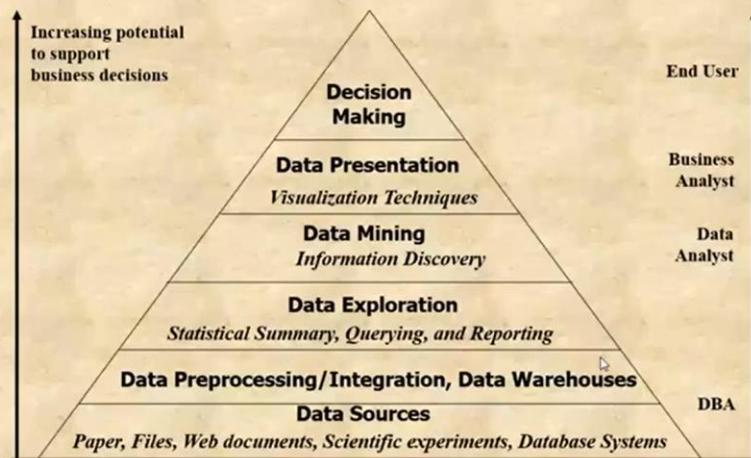
**Data Mining Applications**

- For analysis and decision support
  - Market analysis and management
    - target marketing, customer relation management, market basket analysis, cross selling, market segmentation
  - Risk analysis and management
    - Forecasting, customer retention, improved underwriting, quality control, competitive analysis
  - Fraud detection and management
- Business Intelligence (BI)
- Web page analysis: from web page classification, clustering to PageRank & HITS algorithms
- Biological and medical data analysis: classification, cluster analysis (microarray data analysis), biological sequence analysis, biological network analysis

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg. Sangli. INDIA

+2    GK    AJ    S    S    NM    S    Y    R    K    RS    BM

## Data Mining in Business Intelligence



21



## Major Issues in Data Mining

- **Mining Methodology**
  - Mining various and new kinds of knowledge
  - Mining knowledge in multi-dimensional space
  - Data mining: An interdisciplinary effort
  - Boosting the power of discovery in a networked environment
  - Handling noise, uncertainty, and incompleteness of data
  - Pattern evaluation and pattern- or constraint-guided mining
- **User Interaction**
  - Interactive mining
  - Incorporation of background knowledge
  - Presentation and visualization of data mining results

23



## cont... Major Issues in Data Mining

- **Efficiency and Scalability**
  - Efficiency and scalability of data mining algorithms
  - Parallel, distributed, stream, and incremental mining methods
- **Diversity of data types**
  - Handling complex types of data
  - Mining dynamic, networked, and global data repositories
- **Data mining and society**
  - Social impacts of data mining
  - Privacy-preserving data mining
  - Invisible data mining

24

24



## Lecture 5

**PE3:DATA MINING**  
**About data and its pre-processing**  
**Module-II**

Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.

**ER Model**

Set of Entity sets

Each entity set is set of tuples  
Each tuple is a set of attributes/fields  
Each field is set of permitted values called as domain

Sub1 sub2 sub3 sub4 Result

Feature

Sample

Conditional Attributes a1...a9  
Decision Attributes a10

Feature Vector

4:47 PM 11:38 AM 8/28/2021

+4 S S A R S KR RS K AP Y BM

The screenshot shows a Microsoft OneNote slide titled "Types of Attributes". The slide content is as follows:

**NOMINAL ATTRIBUTES**

Categorical data : finite set

Gender = { Male, Female, others }

Grade = { AA, AB BB CD FF }

Dept = { Civil, Mech, Elec, Eln, CSE, IT }

Dept = { 1,2,3,4,5,6 }

IFSC = { 1507, 1301 , 4589... }

The screenshot shows a Microsoft OneNote page titled "BINARY DATA - OneNote". The page contains the following text and notes:

Wednesday, August 18, 2021 11:45 AM

## BINARY DATA

Data type : boolean - true / false

Data type : logical true / false

Calibri - 11 -  $\frac{A}{a}$  -  $\frac{B}{B}$  -  $\frac{C}{C}$  -  $\frac{D}{D}$  -  $\frac{E}{E}$  -  $\frac{F}{F}$  -  $\frac{G}{G}$  -  $\frac{H}{H}$  -  $\frac{I}{I}$  -  $\frac{J}{J}$  -  $\frac{K}{K}$  -  $\frac{L}{L}$  -  $\frac{M}{M}$  -  $\frac{N}{N}$  -  $\frac{O}{O}$  -  $\frac{P}{P}$  -  $\frac{Q}{Q}$  -  $\frac{R}{R}$  -  $\frac{S}{S}$  -  $\frac{T}{T}$  -  $\frac{U}{U}$  -  $\frac{V}{V}$  -  $\frac{W}{W}$  -  $\frac{X}{X}$  -  $\frac{Y}{Y}$  -  $\frac{Z}{Z}$

Symmetric : both values have equal weightage.  
male/female

Asymmetric : positive / negative , Pass / Fail

The OneNote interface includes a ribbon bar with tabs like FILE, HOME, INSERT, DRAW, HISTORY, REVIEW, and VIEW. The DRAW tab is selected, showing tools for Type, Lasso, Panning, Eraser, Select, Hand, and Tools. The BINARY DATA section is highlighted in the sidebar.

Wednesday, August 18, 2021 11:51 AM

## Ordinal Attributes

Special case of categorical data with meaningful ordering / ranking / flow

**Shirt Size = { Small , Medium , Large }**

**Cadre = { Asst. prof , Associate Prof, Professor }**

I

Wednesday, August 18, 2021 11:55 AM

## Integer - OneNote

**Numeric Data**

- Discrete
- Finite set of values
  - e.g. Bonus rate = { 10, 15, 25 }
  - Number(2) = { 00..99 }
- Continuous
  - e.g. Salary = { 10000.. 250000 }
  - Sales = { 10.. 34234342.34 }

Loss of information

Discretization : bining , categorization - range1 10k to 50k , 2 : 51 k to 100k ....

Wednesday, August 18, 2021 12:11 PM

## Lecture 6

The screenshot shows a Microsoft OneNote page with a title slide and handwritten notes.

**Title Slide:**

# Statistical Description of Data

Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.

**Handwritten Notes:**

- Weighted Average mean:**  $\bar{x} = 76.75$
- Equation:**  $\bar{x} = \frac{\sum w_i f_i}{\sum w_i}$
- Trimmed Mean:**  $\hookrightarrow$  remove/clipped  $\rightarrow$  80 and 92
- Marks:** { 89, 90, 88, 40 }  
{ 89, 90, 88, 92 }
- 89.75**

**OneNote Interface:**

- Top ribbon: FILE, HOME, INSERT, DRAW, HISTORY, REVIEW, VIEW.
- DRAW tab selected.
- Tools: Type, Lasso, Panning, Eraser, Select, Hand.
- Shapes: Color & Thickness, Shapes, Insert Space, Delete, Arrange, Rotate, Link to Text, Link to Math, Convert.
- Bottom ribbon: demo1, New Section 1, Search (Ctrl+E), Add Page, etc.
- Taskbar: teams.microsoft.com is sharing your screen.
- Taskbar icons: +2, ST, SD, S, RS, 2H, K, R, NM, R, BM.
- System tray: 11:36 AM, 8/24/2023.

DM #6 PE3 DM Module II Statistical Description of Data 20210824 111525 Meeti...

Median : physical middle value - center of data set

It is the value that divide the data into two parts

Even  
↓  
two middle

Median = Average of two middle most values

$\frac{52+56}{2} = \text{median}$

Mode : value that occurs frequently / repeated

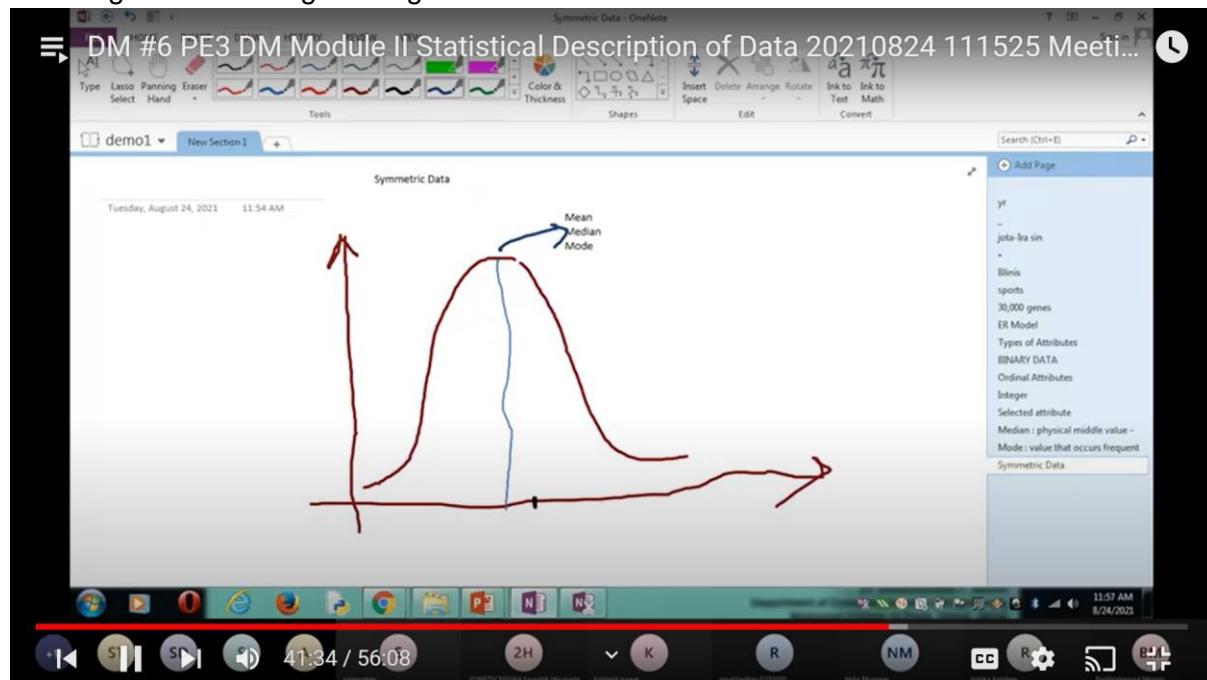
It should occur more than once  
e.g. 52, 70

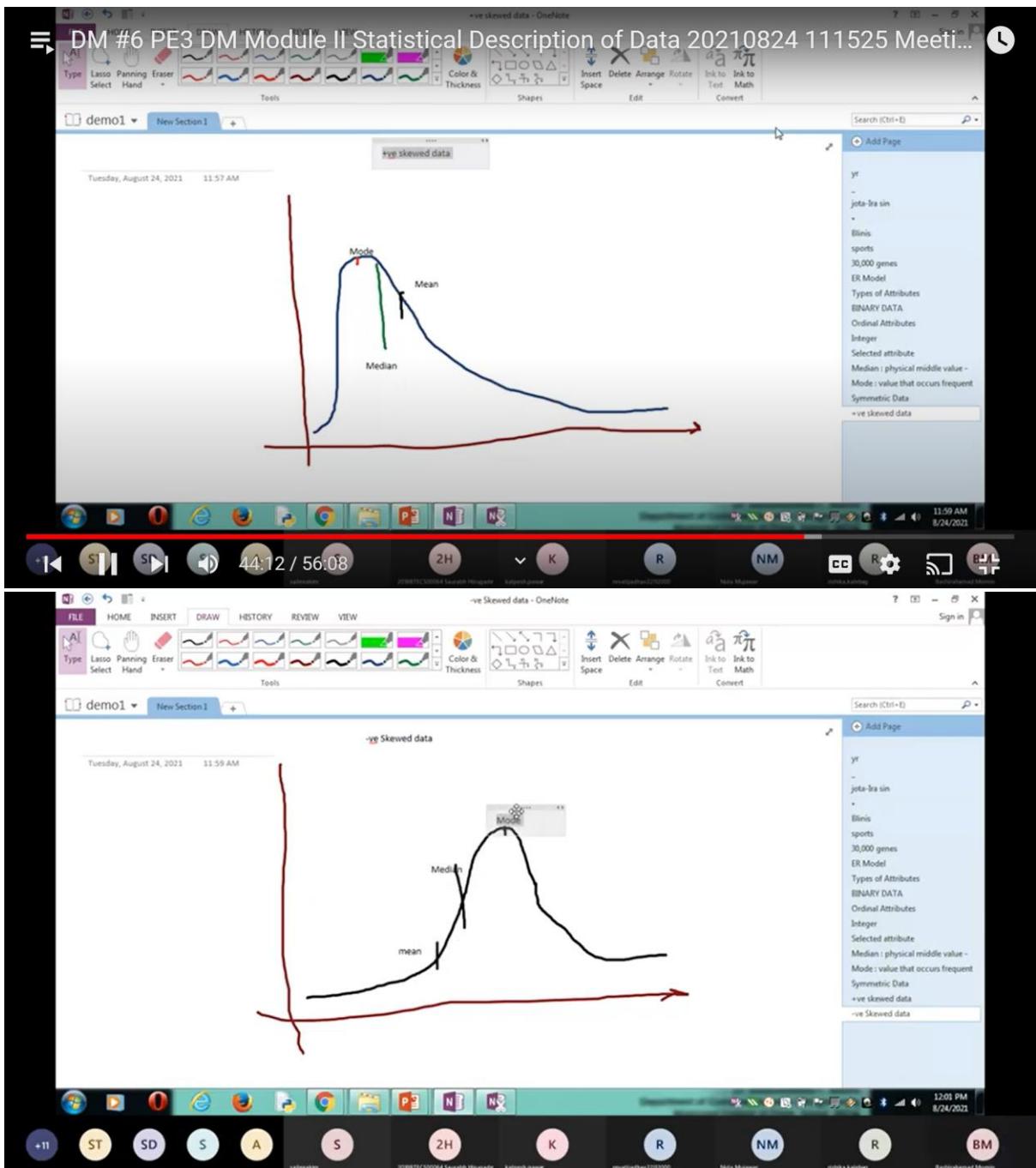
With one mode : unimodal  
Two modes : bimodal  
Three modes : trimodal  
More than three : multimodal

YR  
jota-Ira sin  
Bilnis  
sports  
30,000 genes  
ER Model  
Types of Attributes  
BINARY DATA  
Ordinal Attributes  
Integer  
Selected attribute  
Median : physical middle value -

Mode : value that occurs frequent

Mid range is the average of largest and smallest values





# Measuring the Dispersion of Data

## Spread of numeric data

- Range
- Quantiles
- Quartiles
- Percentile
- Inter-quartile range

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg., Sangli, INDIA



Dispersion : Spread of data - OneNote

**Range : diff. between largest and smallest**  
**Range :  $\max(X) - \min(X)$**

Quartiles : are the data points that split the given data in consecutive sets ... uniform sets

2-quantile : divide the given data sets into lower and upper halves  
 Median is the data point that divide into two part named as "quartiles"

4-quantile : have 3 data points that divide the data set into four quartiles

100-quantile : division of data into 100 consecutive sets : percentile

Inter-quartile range IQR =  $Q_3 - Q_1$

4-quantile : have 3 data points that divide the data set into four quartiles

100-quantile : division of data into 100 consecutive sets : percentile

Inter-quartile range IQR =  $Q_3 - Q_1$

FILE HOME INSERT DRAW HISTORY REVIEW VIEW

Dispersion : Spread of data - OneNote

Type Lasso Panning Eraser Select Hand Tools

Color & Thickness Shapes Insert Space Edit Ink to Text Convert

demo1 New Section 1

100-quantile : division of data into 100 consecutive sets : percentile

Min value      25th percentile      Median      75th percentile      Max value

Inter-quartile range IQR = Q3-Q1

Drawback : don't provide the information about the end points - trails .. extremes

Five number summary : consists of above quartiles + end points

Five number summary = { minimum value , Q1 , median (Q2) , Q3 , maximum value }

Search (Ctrl+E)

Add Page

yr  
jota-ira sin  
Ellipsis  
sports  
30,000 genes  
ER Model  
Types of Attributes  
BINARY DATA  
Ordinal Attributes  
Integer  
Selected attribute  
Median : physical middle value -  
Mode : value that occurs frequent  
Symmetric Data  
+ve skewed data  
-ve Skewed data  
Python code  
Dispersion : Spread of data

2H S NM A S A Y K S T K BM

abhishekmore710 yash\_loke kalpesh\_pawar shivanshavaram tubi\_gautam Krishnapudi44 Badiresham\_Morani

FILE HOME INSERT DRAW HISTORY REVIEW VIEW

Type Lasso Panning Eraser Select Hand Tools

Color & Thickness Shapes Insert Space Edit Ink to Text Math Convert

demo1 New Section 1

Drawback : don't provide the information about the end points - trails .. extremes

Five number summary : consists of above quartiles + end points

Five number summary = { minimum value , Q1 , median (Q2) , Q3 , maximum value }

Graphically represented by BOX PLOT

Whiskers

Minimum - Lower extreme      Q1      Median (Q2)      Q3      Maximum - upper extreme

Spread = IQR = Q3-Q1

Search (Ctrl+E)

Add Page

yr  
jota-ira sin  
Ellipsis  
sports  
30,000 genes  
ER Model  
Types of Attributes  
BINARY DATA  
Ordinal Attributes  
Integer  
Selected attribute  
Median : physical middle value -  
Mode : value that occurs frequent  
Symmetric Data  
+ve skewed data  
-ve Skewed data  
Python code  
Dispersion : Spread of data

Variance and standard deviation - OneNote

Wednesday, August 25, 2021 11:55 AM

**Variance and standard deviation**

It indicates How spread out of a data distribution is

Properties : lower Std. Dev -> data very close to mean  
High std. dev ->large spread over range of values

**Variance**

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

$\sigma$  = Standard deviation = square root of variance

DM #7 Module II Graphic Displays of Basic Statistical Descriptions of Data 2021...

Lower value of Std.Dev

Higher value of Std.Dev

12:08 PM 8/25/2021

12:11 PM 8/25/2021

## Lecture 8:

### DM #8 Module II Data Visualization 20210826 101923 Meeting Recording Graphic Displays of Basic Statistical Descriptions

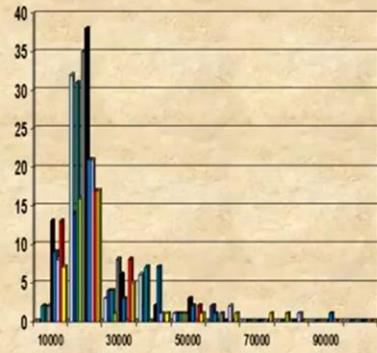
- **Boxplot:** graphic display of five-number summary
- **Histogram:** x-axis are values, y-axis repres. frequencies
- **Quantile plot:** each value  $x_i$  is paired with  $f_i$  indicating that approximately  $100 f_i\%$  of data are  $\leq x_i$
- **Quantile-quantile (q-q) plot:** graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- **Scatter plot:** each pair of values is a pair of coordinates and plotted as points in the plane

15



## Histogram Analysis

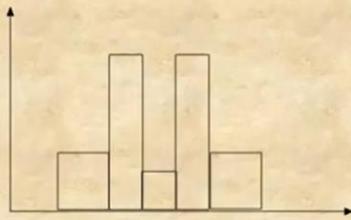
- Histogram: Graph display of tabulated frequencies, shown as bars
- It shows what proportion of cases fall into each of several categories
- Differs from a bar chart in that it is the area of the bar that denotes the value, not the height as in bar charts, a crucial distinction when the categories are not of uniform width
- The categories are usually specified as non-overlapping intervals of some variable. The categories (bars) must be adjacent



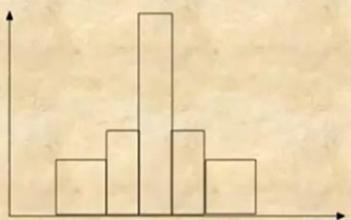
16



## Histograms Often Tell More than Boxplots

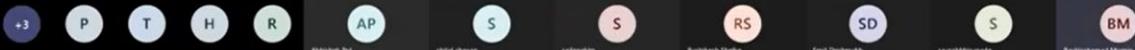


- The two histograms shown in the left may have the same boxplot representation
  - The same values for: min, Q1, median, Q3, max
- But they have rather different data distributions



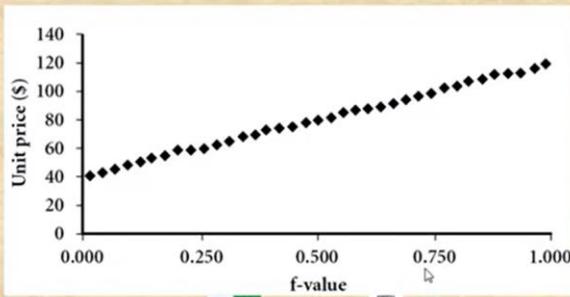
17

17



## Quantile Plot

- Displays all of the data (allowing the user to assess both the overall behavior and unusual occurrences)
- Plots **quantile** information
  - For a data  $x_i$  data sorted in increasing order,  $f_i$  indicates that approximately  $100 f_i\%$  of the data are below or equal to the value  $x_i$



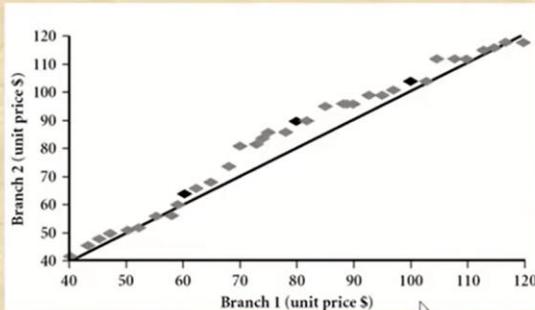
18

18



## Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- View: Is there is a shift in going from one distribution to another?



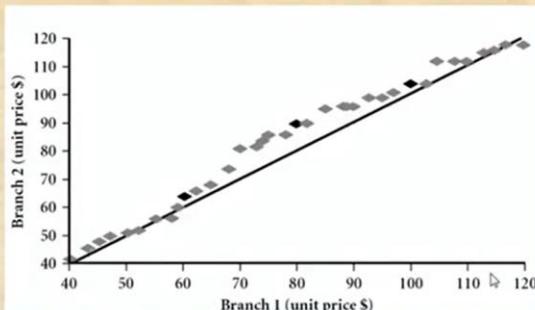
19

TIME 082



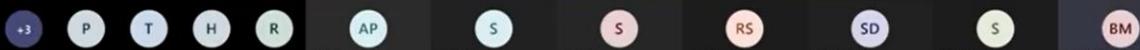
## Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- View: Is there is a shift in going from one distribution to another?



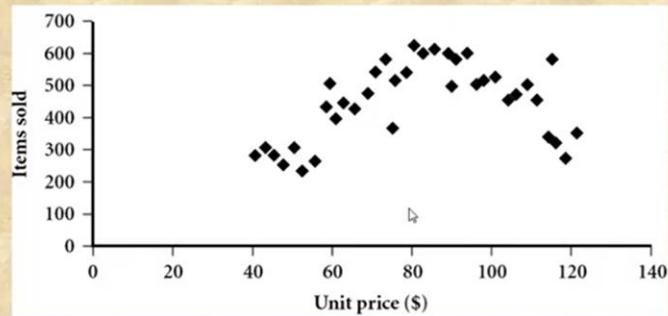
19

TIME 082



## Scatter plot

- Provides a first look at bivariate data to see clusters of points, outliers, etc
- Each pair of values is treated as a pair of coordinates and plotted as points in the plane



20

ITEM #82



20



K

T

H

VK

AP

S

S

RS

Bushra Shukla Smit Deshmukh

mohakakotwag2009

BushraShahid Momin

⇒ DM #8 Module II Data Visualization 20210826 101923 Meeting Recording



Watch later

21

- The left half fragment is positively correlated
- The right half is negative correlated



K

T

H

VK

AP

S

S

RS

Bushra Shukla Smit Deshmukh

mohakakotwag2009

BushraShahid Momin

21 21:53 / 59:08



S



S

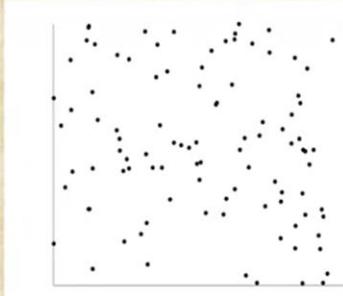
RS

SD

R



## Uncorrelated Data



22

+8 K T H VK AP S S RS SD R BM

# Data Visualization

Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.

+8 K T H VK AP S S RS SD R BM

# Data Visualization

- Why data visualization?
    - Gain insight into an information space by mapping data onto graphical primitives
    - Provide qualitative overview of large data sets
    - Search for patterns, trends, structure, irregularities, relationships among data
    - Help find interesting regions and suitable parameters for further quantitative analysis
    - Provide a visual proof of computer representations derived
  - Categorization of visualization methods:
    - Pixel-oriented visualization techniques
    - Geometric projection visualization techniques
    - Icon-based visualization techniques
    - Hierarchical visualization techniques
    - Visualizing complex data and relations



## Geometric Projection Visualization Techniques

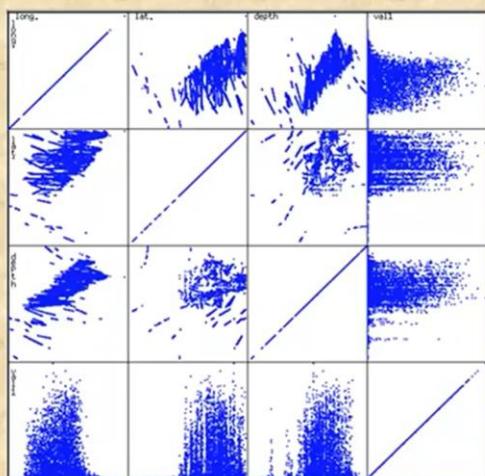
- **Visualization of geometric transformations and projections of the data**
- **Methods**
  - Direct visualization
  - Scatterplot and scatterplot matrices
  - Landscapes
  - Projection pursuit technique: Help users find meaningful projections of multidimensional data
  - Prosection views
  - Hyperslice
  - Parallel coordinates

data courtesy of NCSA, University of Illinois at Urbana-Champaign

27



## Scatterplot Matrices



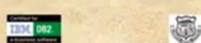
29



# Icon-Based Visualization Techniques

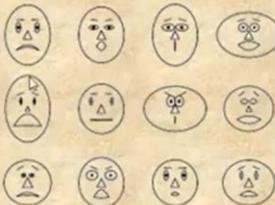
- Visualization of the data values as features of icons
- Typical visualization methods
  - Chernoff Faces
  - Stick Figures
- General techniques
  - Shape coding: Use shape to represent certain information encoding
  - Color icons: Use color icons to encode more information
  - Tile bars: Use small icons to represent the relevant feature vectors in document retrieval

33

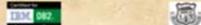


## Chernoff Faces

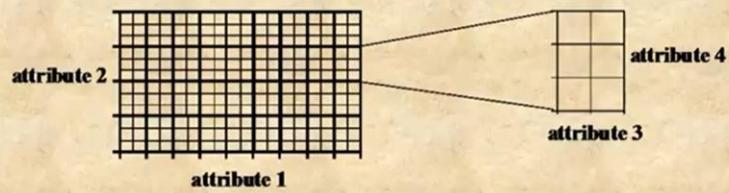
- A way to display variables on a two-dimensional surface, e.g., let  $x$  be eyebrow slant,  $y$  be eye size,  $z$  be nose length, etc.
- The figure shows faces produced using 10 characteristics--head eccentricity, eye size, eye spacing, eye eccentricity, pupil size, eyebrow slant, nose size, mouth shape, mouth size, and mouth opening): Each assigned one of 10 possible values, generated using [Mathematica](#) (S. Dickson)
- REFERENCE: Gonick, L. and Smith, W. [The Cartoon Guide to Statistics](#). New York: Harper Perennial, p. 212, 1993
- Weisstein, Eric W. "Chernoff Face." From [MathWorld](#)--A Wolfram Web Resource. [mathworld.wolfram.com/ChernoffFace.html](http://mathworld.wolfram.com/ChernoffFace.html)



34



## Dimensional Stacking



- Partitioning of the n-dimensional attribute space in 2-D subspaces, which are 'stacked' into each other
- Partitioning of the attribute value ranges into classes. The important attributes should be used on the outer levels.
- Adequate for data with ordinal attributes of low cardinality
- But, difficult to display more than nine dimensions

37



## Tree-Map

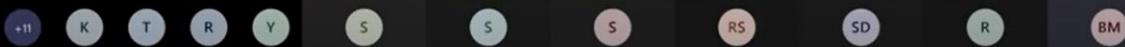
- Screen-filling method which uses a hierarchical partitioning of the screen into regions depending on the attribute values
- The x- and y-dimension of the screen are partitioned alternately according to the attribute values (classes)

MSR Netscan Image



Ack.: <http://www.cs.umd.edu/hdi/treemap-history/all102001.jpg>

40

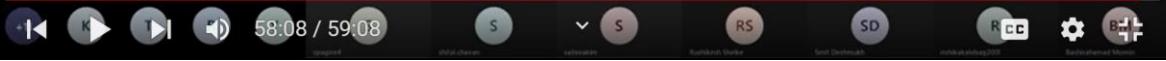




- A 3-D visualization technique where hierarchical information is displayed as nested semi-transparent cubes
- The outermost cubes correspond to the top level data, while the subnodes or the lower level data are represented as smaller cubes inside the outermost cubes, and so on



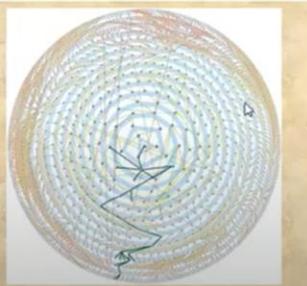
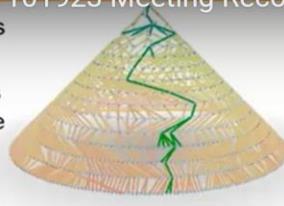
42



## DM #8 Module II Data Visualization 20210826 101923 Meeting Recording



- 3D cone tree visualization technique works well for up to a thousand nodes or so
- First build a 2D circle tree that arranges its nodes in concentric circles centered on the root node
- Cannot avoid overlaps when projected to 2D
- G. Robertson, J. Mackinlay, S. Card. "Cone Trees: Animated 3D Visualizations of Hierarchical Information", ACM SIGCHI'91 website: Visualize a social network data set that models the way an infection spreads from one person to the next



Ack.: <http://madeausoftware.com/articles/visualization>

43

58:13 / 59:08



## Visualizing Complex Data and Relations

- Visualizing non-numerical data: text and social networks

- The importance of tag is represented by font size/color
- Besides text data, there are also methods to visualize relationships, such as visualizing social networks



Newsmap: Google News Stories in 2005

44

spage4



## Recent Trend in Visualization

- Multimedia – animation
- Augmented reality AR
- Virtual reality VR

Dr. Bashirahamed F. Momin  
Department of Computer Science & Engineering  
Wachand College of Engg., Sion(El), INDIA



## Lecture 9 started

The screenshot shows a OneNote meeting recording interface. At the top, a slide titled "Data Pre-processing" is displayed, listing four main topics:

- Data Cleaning
- Data Integration
- Data Transformation
- and data discretization

Below the slide, a footer contains the text: "Dr. Bashirahamad F. Momin", "Department of Computer Science & Engineering", and "Walchand College of Engr., Sangli, INDIA". The OneNote ribbon is visible at the top, and the status bar at the bottom shows the date and time.

The main content area contains a handwritten note in red ink: "Quality of data mining is directly proportional to quality of data". Below this, a mathematical equation is written:  $Q(m) \propto Q(d)$ . A handwritten question "Why?" is followed by an arrow pointing to a list of six factors: Accuracy, Completeness, Consistency, Timeliness, Availability, and Interpretability.

A screenshot of the OneNote slide is also shown within the note area, illustrating the concept of data quality.

What is the source of data?

OLTP / IoT devices / machines / equipments

The screenshot shows a Microsoft OneNote page titled "Data Pro-processing - OneNote". The main content area contains a handwritten note: "Quality of data mining is directly proportional to quality of data" followed by a formula  $Q(m) \propto Q(D)$ . A bracket next to the formula lists six qualities: Accuracy, Completeness, Consistency, Timeliness, Belieavability, and Interpretability. Below the formula is a question "Why?" with an arrow pointing to the list. A sidebar on the right lists various data processing concepts. The taskbar at the bottom shows several open applications, and the status bar indicates the date as August 31, 2021.

DM #9 General 20210831 111611 Meeting Recording

Data Pro-processing - OneNote

demo1 New Section 1

Tuesday, August 31, 2021 11:22 AM

**Quality of data mining is directly proportional to quality of data**

$Q(m) \propto Q(D)$

Why?  $\rightarrow$

Accuracy  
Completeness  
Consistency  
Timeliness  
Belieavability  
Interpretability

Sources : OLTP / IoT devices / machines / Equipments

Data Entry form : care to be taken to avoid such errors.

Search (Ctrl+E)

Add Page

jota-ira sin  
Bilis  
sports  
30,000 genes  
ER Model  
Types of Attributes  
BINARY DATA  
Ordinal Attributes  
Integer  
Selected attribute  
Median : physical middle value  
Mode : value that occurs frequ  
Symmetric Data  
+ve skewed data  
-ve Skewed data  
Python code  
Dispersion : Spread of data  
Variance and standard deviation  
Histogram  
Data Pro-processing

11:29 AM  
8/31/2021

22:51 / 57:23

The screenshot shows a Microsoft OneNote page titled "Data Cleaning - OneNote". The main content area displays the following text:

**Data Cleaning**

Tuesday, August 31, 2021 11:41 AM

**Missing values !**

**Student table : PNR, NAME, ADDRESS, EMAIL, MOBILE, PARENT\_MOBILE**

{2018BTECS0010, 'Aishwarya', null, email@email.com, 0, 0}

A vertical search sidebar on the right lists various data concepts:

- jata-Ira sin
- Bilnis
- sports
- 30,000 genes
- ER Model
- Types of Attributes
- BINARY DATA
- Ordinal Attributes
- Integer
- Selected attribute
- Median : physical middle value
- Mode : value that occurs freqe
- Symmetric Data
- +ve skewed data
- ve Skewed data
- Python code
- Dispersion : Spread of data
- Variance and standard deviatio
- Histogram
- Data Pro-processing
- Data Cleaning

The taskbar at the bottom shows various application icons, and the system tray indicates the date and time as 11:47 AM / 8/31/2021.

**Data Cleaning**      **Missing values !**

Tuesday, August 31, 2021 11:41 AM

**Student table : PNR, NAME, ADDRESS, EMAIL, MOBILE, PARENT\_MOBILE  
{2018BTECS0010, 'Aishwarya', 'A/P K. Maha.', email@email.com, ?, ?}**

**Solution**

1. Ignore the tuples
2. Fill the missing values manually
3. Use a global constants to fill in this missing values
4. Use a statistical measures like Central Tendency of attribute
5. Use the attribute mean/median
6. Use the most probable value

Windows taskbar icons: Start, File Explorer, Edge, Firefox, File, Home, Insert, Draw, History, Review, View, Data Cleaning - OneNote, Sign in, Search (Ctrl+E), Add Page, jot-a-ira sin, Illinois, sports, 30,000 genes, ER Model, Types of Attributes, BINARY DATA, Ordinal Attributes, Integer, Selected attribute, Median: physical middle value, Mode: value that occurs frequ, Symmetric Data, +ve skewed data, -ve Skewed data, Python code, Dispersion: Spread of data, Variance and standard deviatio, Histogram, Data Pre-processing, Data Cleaning.

## Lecture 10

**DM #10 Module II cont Data Pre processing 20210901 111743 Meeting Recording**

**Date field : 03/05/2021 ... dd/mm/YYYY  
mm/dd/YYYY**

**Central Tendency of attributes : mean , median**

**For normal symmetric attribute , replace missing values with "mean"**  
**For skewed data : median is replace for missing values.**

**Compute the mean or median for group of tuples in that class and then replace that mean/median for missing values.**

**Regression , Bayesian theory**

Windows taskbar icons: Start, File Explorer, Edge, Firefox, File, Home, Insert, Draw, History, Review, View, Data Cleaning - OneNote, Sign in, Search (Ctrl+E), Add Page, jot-a-ira sin, Illinois, sports, 30,000 genes, ER Model, Types of Attributes, BINARY DATA, Ordinal Attributes, Integer, Selected attribute, Median: physical middle value, Mode: value that occurs frequ, Symmetric Data, +ve skewed data, -ve Skewed data, Python code, Dispersion: Spread of data, Variance and standard deviatio, Histogram, Data Pre-processing, Data Cleaning.

Noisy Data : Noise is a random error or variance in measured variables.

Wednesday, September 01, 2021 11:28 AM

**How to smooth it or clean ?**

**Using data smoothing techniques :**

**Binning : Equal frequency**

Smoothing by Bin "Mean"      Smoothing by Bin "Median"      Smoothing by Bin "boundaries"

1. Sort the given data in ascending order  
2. Form equal frequency/count bin (bucket/grouping)  
3. Compute "mean" or "median" for each bin  
4. Replace each bin values with these "mean" or "median"

1. Sort the given data in ascending order  
2. Form equal frequency/count bin (bucket/grouping)  
3. Replace each bin value with its neighbor/nearest boundaries

Example data set = { 4,8,15,21,21,24,25,28,34}

```

Bin 1 = {4,8,15} mean=9
Bin 2 = {21,21,24} mean = 22
Bin 3 = {25,28,34} mean = 29
Smoothed data Bin1={9,9,9}

```

11:29 AM 9/1/2021

Noisy Data - OneNote

**Using data smoothing techniques :**

**Binning : Equal frequency**

**Smoothing by Bin "Mean"**

1. Sort the given data in ascending order
2. Form equal frequency/count bin (bucket/grouping)
3. Compute "mean" or "median" for each bin
4. Replace each bin values with these "mean" or "median"

**Smoothing by Bin "Median"**

1. Sort the given data in ascending order
2. Form equal frequency/count bin (bucket/grouping)
3. Replace each bin value with its neighbor/nearest boundaries

**Smoothing by Bin "boundaries"**

Example data set = { 4,8,15,21,21,24,25,28,34}

```

Bin 1 = {4,8,15} mean=9
Bin 2 = {21,21,24} mean = 22
Bin 3 = 3 {25,28,34} mean = 29
Smoothed data
Bin1={9,9,9}
Bin2={22,22,22}
Bin3={29,29,29}

```

11:47 AM 9/1/2021

DM #10 Module II cont Data Pre processing 20210901 111743 Meeting Recording

Noisy Data - OneNote

**Smoothing by Bin "Mean"**

1. Sort the given data in ascending order
2. Form equal frequency/count bin (bucket/grouping)
3. Compute "mean" or "median" for each bin
4. Replace each bin values with these "mean" or "median"

**Smoothing by Bin "Median"**

1. Sort the given data in ascending order
2. Form equal frequency/count bin (bucket/grouping)
3. Replace each bin value with its neighbor/nearest boundaries

**Smoothing by Bin "boundaries"**

Example data set = { 4,8,15,21,21,24,25,28,34}

```

Bin 1 = {4,8,15} mean=9 median=8
Bin 2 = {21,21,24} mean = 22 median=21
Bin 3 = 3 {25,28,34} mean = 29 median=28
Smoothed data =>mean
Bin1={9,9,9}
Bin2={22,22,22}
Bin3={29,29,29}
Smoothed data =>median
Bin1={8,8,8}
Bin2={21,21,21}
Bin3={28,28,28}

```

11:49 AM 9/1/2021

**Smoothing by Bin "Mean"**

- Sort the given data in ascending order
- Form equal frequency/count bin (bucket/grouping)
- Compute "mean" or "median" for each bin
- Replace each bin values with these "mean" or "median"

**Example data set = { 4,8,15,21,21,24,25,28,34}**

```

Bin 1 = {4,8,15} mean=9 median=8
Bin 2 = {21,21,24} mean = 22 median=21
Bin 3 = {25,28,34} mean = 29 median=28
Smoothed data =>mean
Bin1={9,9,9}
Bin2={22,22,22}
Bin3={29,29,29}
Smoothed data =>median
Bin1={8,8,8}
Bin2={21,21,21}
Bin3={28,28,28}
  
```

**Smoothing by Bin "Median"**

- Sort the given data in ascending order
- Form equal frequency/count bin (bucket/grouping)
- Replace each bin values with its neighbor/nearest boundaries

**Smoothed data => boundaries**

```

Bin1={4,4,15}
Bin2={21,21,24}
Bin3={25,25,34}
  
```

**Smoothing by Bin "boundaries"**

- Sort the given data in ascending order
- Form equal frequency/count bin (bucket/grouping)
- Replace each bin values with its neighbor/nearest boundaries

**Smoothed data => boundaries**

```

Bin1={4,4,15}
Bin2={21,21,24}
Bin3={25,25,34}
  
```

How to smooth the data?

Binning

Regression: Linear/ Multi linear

Data Cleaning as a “process”

Detection of Discrepancy

Whenever there is an human intervention, there is a possibility of an error

**Data Cleaning as a "Process"** - OneNote

Wednesday, September 01, 2021 11:57 AM

**Detection of Discrepancy**

1. Poorly feed/input data by human beings
2. Human errors
3. Deliberate errors
4. Data decay - obsolete / outdated / old - expired
5. ...

**How to proceed with detection ?**

1. Use domain knowledge - recall ER model - values/contents of attributes => domain i.e. set of permitted values  
e.g. DOB <= today's date 23/10/2022 ? **31/04/2020**
2. Date field dd/mm/yyyy , mm/dd/yyyy
3. Primary key .... Unique PRN
4. Field overloading : common error - attribute defined for one purpose, actual used for other

12:03 PM 9/1/2021

+9 NM S K SM P S Y SS R S BM

prachi.messore sagnik.yash.hoke Supriya.Senj mithakakshag2009 salilakshi Bashirahmed.Momin

FILE HOME INSERT DRAW HISTORY REVIEW VIEW

Type Lasso Panning Eraser Select Hand Tools

Color & Thickness Shapes Insert Space Delete Arrange Rotate Ink to Text Convert

Search (Ctrl+E)

Add Page

Ellipsis  
sports  
30,000 genes  
ER Model  
Types of Attributes  
BINNARY DATA  
Ordinal Attributes  
Integer  
Selected attribute  
Median : physical middle value  
Mode : value that occurs frequ  
Symmetric Data  
+ve skewed data  
-ve Skewed data  
Python code  
Dispersion : Spread of data  
Variance and standard deviatio  
Histogram  
Data Pro-processing  
Data Cleaning  
Noisy Data :  
Data Cleaning as a "Process"

## Lecture 12

**Merging of data from multiple sources**

Matching of schemas / objects ~ Entity Identifications .... ER model  
Attribute correlation  
Detection and removal of conflict data .... Duplication , inconsistency

**Entity Identification problem :**

Branch\_code = { CV , ME , EE , EN , CS , IT }  
Dept\_code = { 01,02,03,04,05 }

Discount field %

**Redundancy and Correlation analysis : Detection of such attributes**

Chi-square test : for nominal attributes  
Correlation coefficient & Co-variance for numerical attributes

**Pearson statistics**

$\chi^2$

**Chi-Square Test**

Steps

- Find the distinct values in each attribute i.e. domain {}
- Let there are 'c' distinct values in first attribute and 'r' distinct values in second attribute
- Build contingency matrix of r X c
- Entry in this table is refer as "join event" .... Occurrences /count of event
- Then compute chi-square value

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

Sample Data Set - { gender , perfect\_reading , salary }  
I  
Gender = { M, F, O }  
Perfect\_reading = {fiction, non\_fiction}  
Salary = { 10K to 200K }  
No of samples/tuples/rows = 1500  
Distinct values <= {domain values}

8:29 AM 9/3/2023

8:31 AM 9/3/2023

## Lecture 13

DM #13 cont Chi Square Test, Correlation Coefficient 20210907 111635 Meeting Recording

Test parameters : Significance level , degree of freedom  
 Degree of freedom =  $(r-1)(c-1) = (2-1)(2-1) = 1$

How to choose significance level alpha  $\alpha$   
 Risk factor or accuracy in % 0.1% = 0.1% 0.001

df	0.200	0.100	0.075	0.050	0.025	0.010	0.005	0.001	0.0005
1	1.642	2.706	3.170	3.841	5.024	6.635	7.879	10.828	12.116

**Conclusion :**  
**If calculated chi-test value > std, then two attributes are strongly correlated**

Correlation coefficient for numeric data (Pearson coefficient)

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{\sqrt{n} \sigma_A \sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - n \bar{A} \bar{B}}{\sqrt{n} \sigma_A \sigma_B},$$

Where  
 A1 AND B1 are the attribute values  
 A bar & B bar = "MEAN" center of attribute  
 Sigma A and Sigma B = Std. deviation  
 N is the no. of tuples/rows/records in data set  
 A1B1 are dot product

$$-1 \leq r_{A,B} \leq +1$$

If  $r_{A,B}$  is greater than 0 then A & B are positively correlated.

DM #13 cont Chi Square Test, Correlation Coefficient 20210907 111635 Meeting Recording

coefficient)

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n\sigma_A\sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - n\bar{A}\bar{B}}{n\sigma_A\sigma_B},$$

Where  
 A1 AND B1 are the attribute values  
 A bar & B bar = "MEAN" center of attribute  
 Sigma A and Sigma B = Std. deviation  
 N is the no. of tuples/rows/records in data set  
 AiBi are dot product

$$-1 \leq r_{A,B} \leq +1$$

If Rab is greater than 0 then A & B are positively correlated.  
 Else negative

11:45 AM 9/7/2021

$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$

Covariance of Numeric data - OneNote

DM #13 cont Chi Square Test, Correlation Coefficient 20210907 111635 Meeting Recording

Tuesday, September 07, 2021 11:51 AM

Consider attributes A & B with n number of rows/tuples

$$E(A) = \bar{A} = \frac{\sum_{i=1}^n a_i}{n}$$

Calculate "Mean" Center of each attribute

$$E(B) = \bar{B} = \frac{\sum_{i=1}^n b_i}{n}$$

$$\text{Cov}(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

Covariance of Numeric data - OneNote

DM #13 cont Chi Square Test, Correlation Coefficient 20210907 111635 Meeting Recording

Tuesday, September 07, 2021 11:51 AM

Consider attributes A & B with n number of rows/tuples

$$E(A) = \bar{A} = \frac{\sum_{i=1}^n a_i}{n}$$

Calculate "Mean" Center of each attribute

$$E(B) = \bar{B} = \frac{\sum_{i=1}^n b_i}{n}$$

$$\text{Cov}(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

**Relationship between Covariance and Pearson Coefficient**

$$\text{Cov}(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

*6A6*

$$r_{A,B} = \frac{\text{Cov}(A, B)}{\sigma_A \sigma_B},$$

**Covariance of Numeric data - OneNote**

**Relationship between Covariance and Pearson Coefficient**

$$\text{Cov}(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

**Interpretation :**  
+ values show growth of both attributes i.e. positively correlated

**DM #13 cont Chi Square Test, Correlation Coefficient 20210907 111635 Meeting Recording**

**Data from heterogeneous sources**

**Data Value conflict detection and resolution (solution)**

- Examples : Grade AA BB CC FF , CPI in autonomous colleges
- University : First class 65%
- Dist. 75%

## Lecture 14

Data transformation and Discretization - OneNote

Wednesday, September 08, 2021 11:21 AM

### Data Transformation :

**It involves :**  
Summary : Applying aggregation function  
ETL in Data warehouse : Extract Transform Load

### Discretization

### Normalization

### By Binning - equal width/equal frequency

### Histogram analysis

I|

11:34 AM 9/8/2021

Normalization - OneNote

Wednesday, September 08, 2021 11:36 AM

### Mapping between -1 to +1 or 0 to 1

#### Types

1. Min-Max normalization
2. Z-Score normalization
3. Normalization by decimal scaling

Let "A" is attribute with values from  $v_1, v_2, \dots, v_n$   
Find  $\text{Min}_A, \text{Max}_A$

$$v'_i = \frac{v_i - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{new}_{\text{max}} - \text{new}_{\text{min}}) + \text{new}_{\text{min}}$$

0 to 1

11:44 AM 9/8/2021

DM #14 Data transformation and Discretization 20210908 112213 Meeting Recording

$v'_i = \frac{v_i - \bar{A}}{\sigma_A}$

Z-Score Normalization / Zero Mean Normalization

$$s_A = \frac{1}{n}(|v_1 - \bar{A}| + |v_2 - \bar{A}| + \dots + |v_n - \bar{A}|)$$

Mean Absolute Deviation

$$v'_i = \frac{v_i - \bar{A}}{s_A}$$

Mean Std. Deviation

Normalizing by decimal scaling:

$$v'_i = \frac{v_i}{10^j}$$

where  $j$  is the smallest integer such that  $\max(|v'_i|) < 1$ .

DM #14 Data transformation and Discretization 20210908 112213 Meeting Recording

$s_A = \frac{1}{n}(|v_1 - \bar{A}| + |v_2 - \bar{A}| + \dots + |v_n - \bar{A}|)$

Mean Absolute Deviation

$$v'_i = \frac{v_i - \bar{A}}{s_A}$$

Mean Std. Deviation

Normalizing by decimal scaling:

$$v'_i = \frac{v_i}{10^j}$$

where  $j$  is the smallest integer such that  $\max(|v'_i|) < 1$ .

Normalization - OneNote

DM #14 Data transformation and Discretization 20210908 112213 Meeting Recording

Normalization by decimal scaling :

$$v'_i = \frac{v_i}{10^j}, \quad \text{Decimal scaling factor}$$

where  $j$  is the smallest integer such that  $\max(|v'_i|) < 1$ .

$A = \{-986 \dots 917\}$

min max

Absolute max(A) = 986 --- > 0.986  
So decimal scaling should be 1000  
Hence value of  $j=3$

12:04 PM 9/8/2021

Discretization by Binning - OneNote

DM #14 Data transformation and Discretization 20210908 112213 Meeting Recording

Discretization by Binning

Wednesday, September 08, 2021 12:04 PM

It is unsupervised discretization

Equal width / equal frequency

Equal frequency binning : Divide the data into bins having equal frequency/count of values

Total values are 120 and frequency/count=4  
No of bins = 120/4 30 bins

Median : physical middle value  
Mode : value that occurs frequ  
Symmetric Data  
+ve skewed data  
-ve Skewed data  
Python code  
Dispersion : Spread of data  
Variance and standard deviatio  
Histogram  
Data Pro-processing  
Data Cleaning  
Noisy Data :  
Data Cleaning as a "Process"  
Data Integration  
Chi-Square Test  
Correlation coefficient for nu  
Covariance of Nu  
Covariance of Numeric data  
Tuple Duplication  
Data transformation and Discr  
Normalization

12:12 PM 9/8/2021

DM #15 Meeting in General 20210909 101913 Meeting Recording

Discretization by Binning - OneNote

It is unsupervised discretization

Equal width / equal frequency

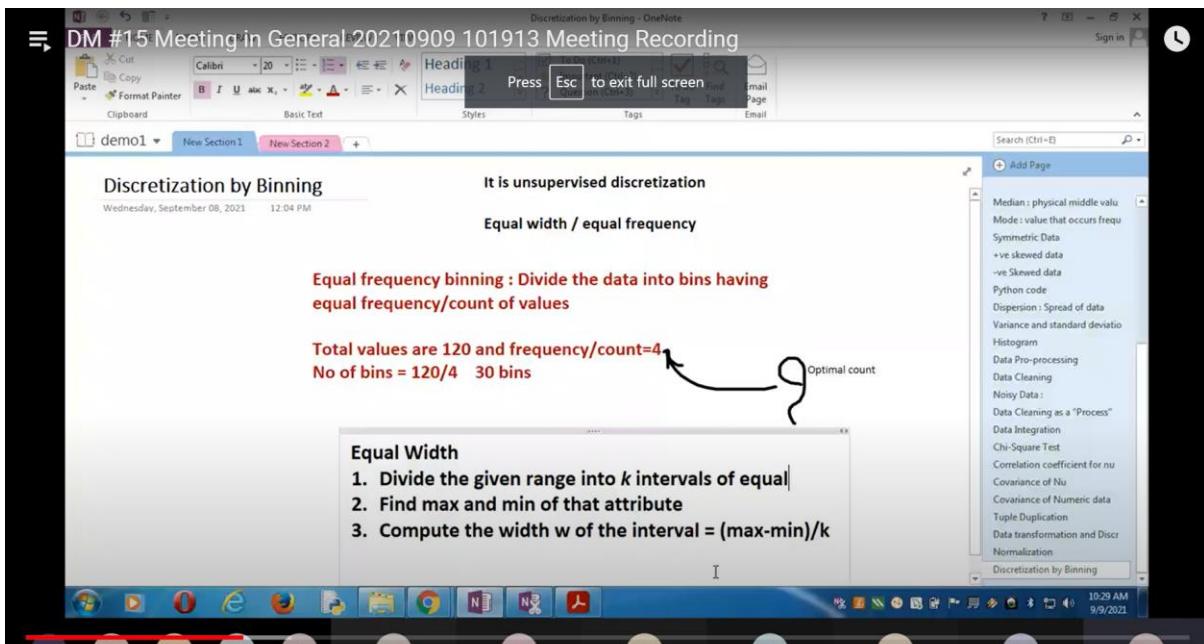
Equal frequency binning : Divide the data into bins having equal frequency/count of values

Total values are 120 and frequency/count=4  
No of bins =  $120/4 = 30$  bins

Optimal count

Equal Width

1. Divide the given range into  $k$  intervals of equal width  
2. Find max and min of that attribute  
3. Compute the width  $w$  of the interval =  $(\text{max-min})/k$



DM #15 Meeting in General 20210909 101913 Meeting Recording

Discretization by Binning - OneNote

Interval =  $\text{min} + \{i-1\} * w$  where  $i=1$  to  $k$

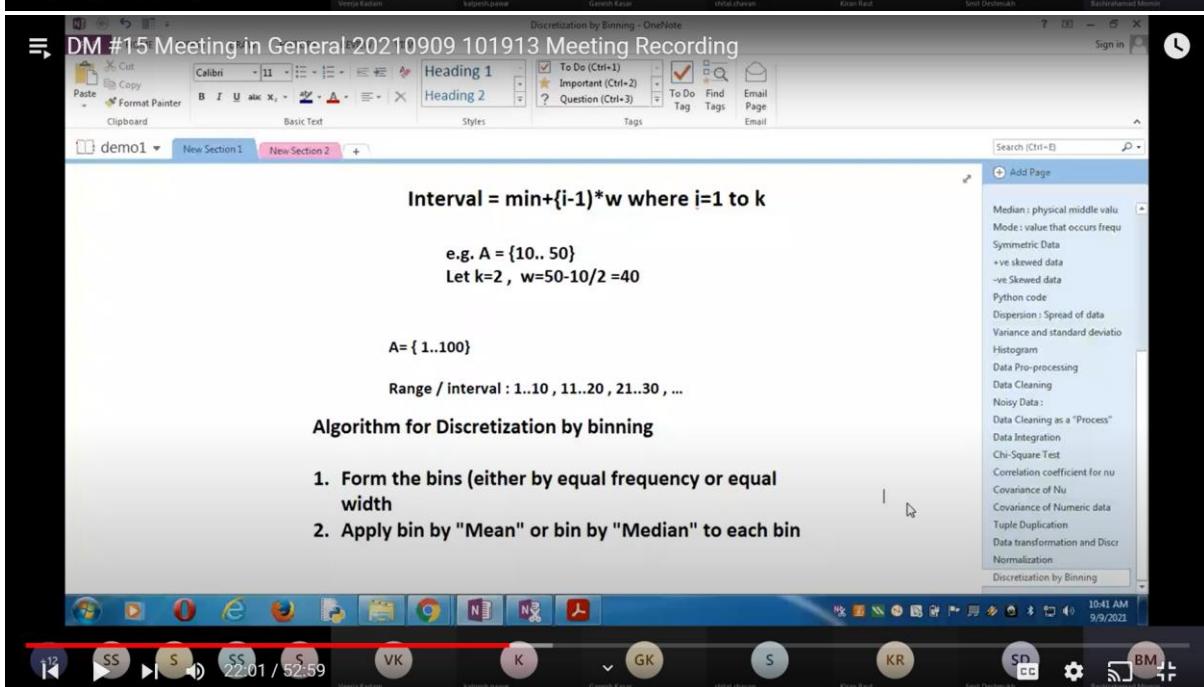
e.g.  $A = \{10.. 50\}$   
Let  $k=2$ ,  $w=50-10/2 = 40$

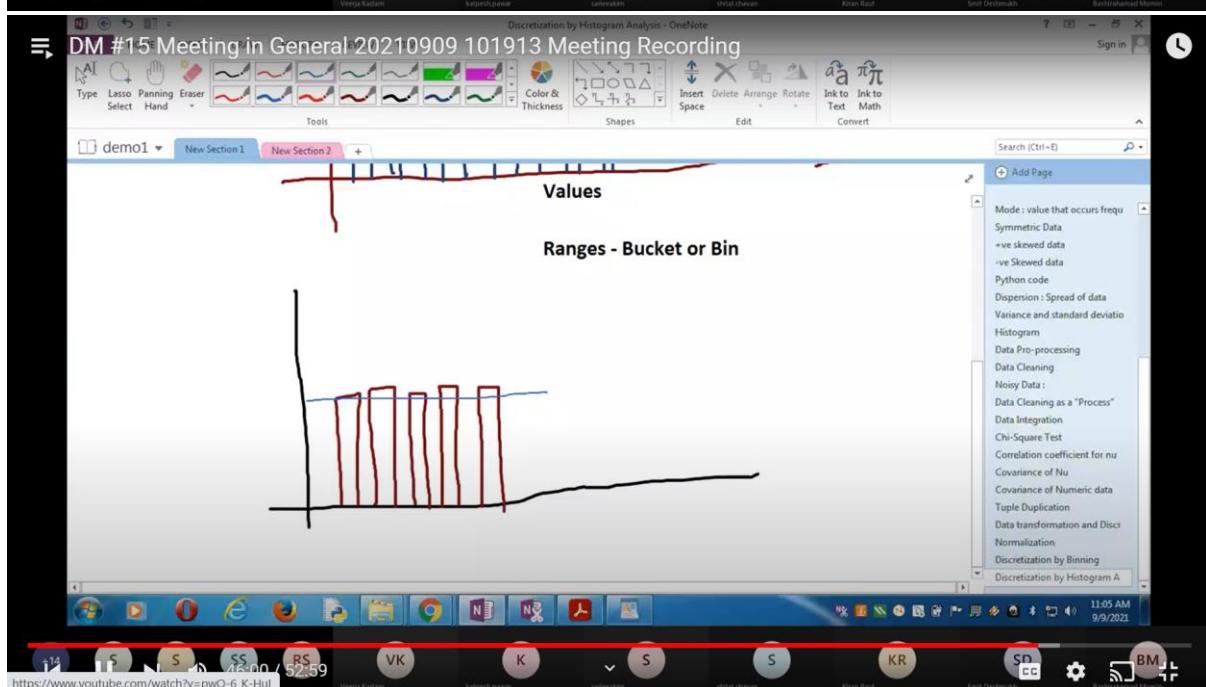
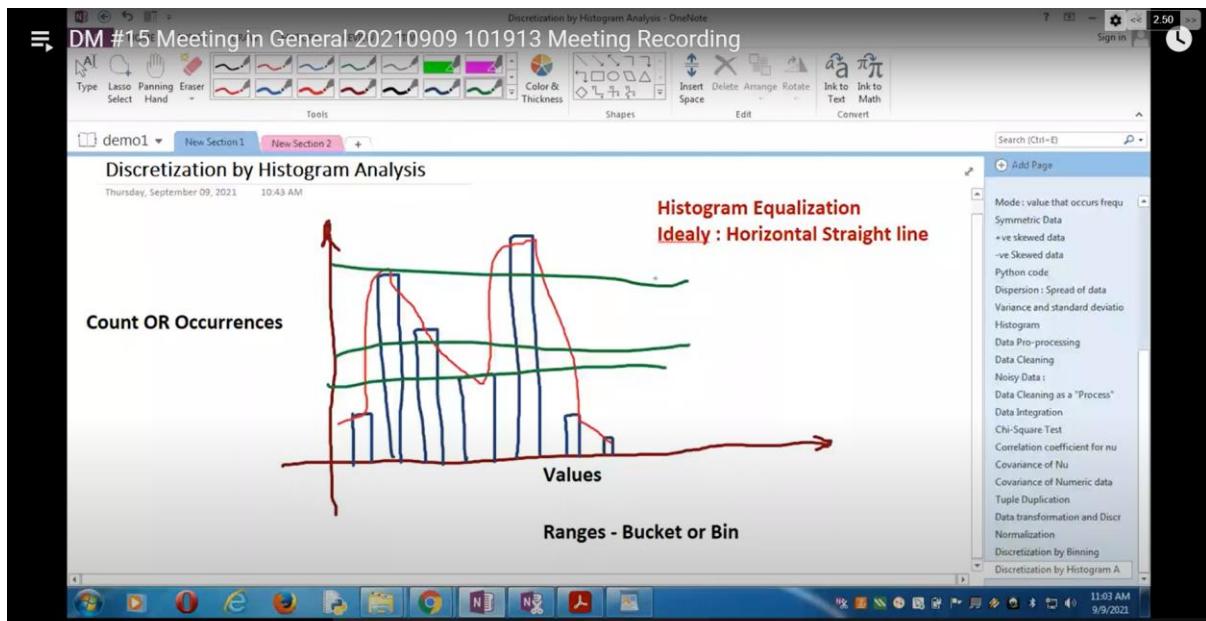
$A = \{1..100\}$

Range / Interval :  $1..10, 11..20, 21..30, \dots$

Algorithm for Discretization by binning

1. Form the bins (either by equal frequency or equal width)
2. Apply bin by "Mean" or bin by "Median" to each bin





## Lecture 16

# Supervised vs. Unsupervised Learning

- **Supervised learning (classification)**

- Supervision: The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations

- New data is classified based on the training set

- **Unsupervised learning (clustering)**

- The class labels of training data is unknown
- Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

3

Certified by  
IBM DB2  
a database software



## Prediction Problems: Classification vs. Numeric Prediction

- **Classification**

- predicts categorical class labels (discrete or nominal)
- classifies data (constructs a model) based on the training set and the values (**class labels**) in a classifying attribute and uses it in classifying new data

- **Numeric Prediction**

- models continuous-valued functions, i.e., predicts unknown or missing values

- **Typical applications**

- Credit/loan approval:
- Medical diagnosis: if a tumor is cancerous or benign
- Fraud detection: if a transaction is fraudulent
- Web page categorization: which category it is

4

Certified by  
IBM DB2  
a database software



## Classification - Definition

- Method of categorizing or assigning class labels to a pattern set under supervision of teacher : **Supervised Learning.**
- The decision boundaries are generated to discriminate between patterns of different classes.
- Prediction of class from set of samples
- Different methods :
  - Decision Trees
  - Probabilistic or generative models
  - Nearest Neighbor classifier
  - Artificial Neural Network (ANN)

5

Created by  
IBM DB2  
A business software

A

KP

S

K

R

revatpadhanav2212000

Rushikesh Shelke

Suthanshu Puraskar

SP

NM

S

salleekan

SD

Smit Deshmukh

Bashrahmad Monin

BM

DM #16 Module III Classification Basic concept 20210921 111551 Meeting Recor...

- Given a database  $D=\{t_1, t_2, \dots, t_n\}$  and a set of classes  $C=\{C_1, \dots, C_m\}$ , the **Classification Problem** is to define a mapping  $f:D \rightarrow C$  where each  $t_i$  is assigned to one class.
- Actually divides  $D$  into **equivalence** classes.
- **Prediction** is similar, but may be viewed as having infinite number of classes.

6

Created by  
IBM DB2  
A business software

A

KP

S

R

31'45 / 59'53

revatpadhanav2212000

Rushikesh Shelke

Suthanshu Puraskar

SP

NM

S

salleekan

SD

Smit Deshmukh

Bashrahmad Monin

BM

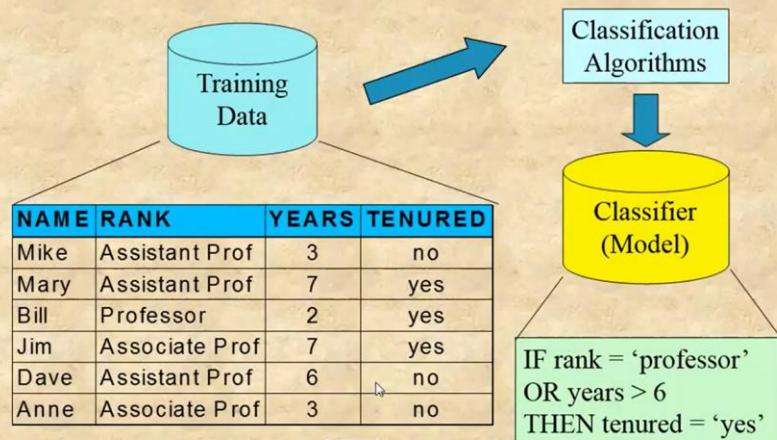
## Classification—A Two-Step Process

- **Model construction:** describing a set of predetermined classes
  - Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
  - The set of tuples used for model construction is **training set**
  - The model is represented as classification rules, decision trees, or mathematical formulae
- **Model usage:** for classifying future or unknown objects
  - **Estimate accuracy** of the model
    - The known label of test sample is compared with the classified result from the model
    - Accuracy rate is the percentage of test set samples that are correctly classified by the model
    - Test set is independent of training set, otherwise over-fitting will occur
  - If the accuracy is acceptable, use the model to **classify data** tuples whose class labels are not known

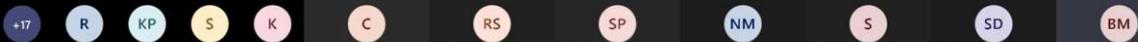
7

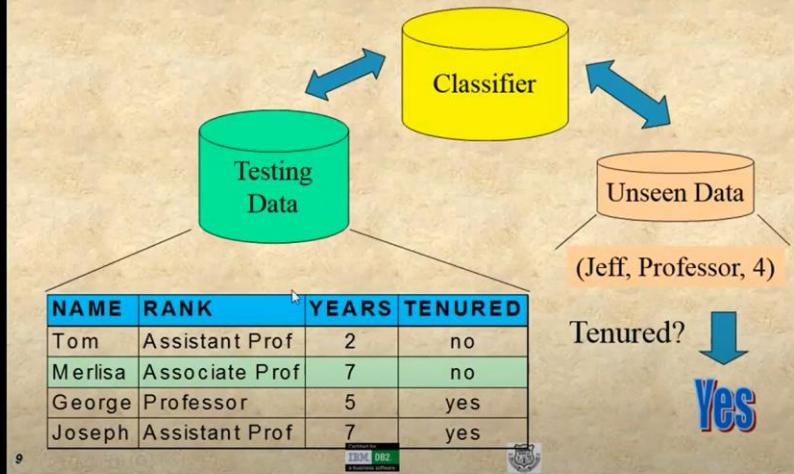
Created by  
IBM DB2  
A database software

## Process (1): Model Construction



8

Created by  
IBM DB2  
A database software



9

◀ R || K ► 🔊 52:52 / 59:53 RS SP NM S SD CC Smit Debnath Bachirahmed Monim

## Classification Examples

- Teachers classify students' grades as A, B, C, D, or F.
- Identify mushrooms as poisonous or edible.
- Predict when a river will flood.
- Identify individuals with credit risks.
- Speech recognition
- Pattern recognition

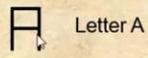
10

+17 R KP S K C RS SP NM S SD BM chavanhiralika Bushra Shekhar Sultanzha Purakar Nida Mujawir salirekha Smit Debnath Bachirahmed Monim

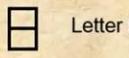
# Classification Ex: Letter Recognition

View letters as constructed from 5 components:

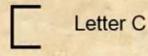
Features : Strokes , Tees , joint , end points etc



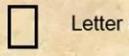
Letter A



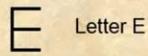
Letter B



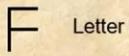
Letter C



Letter D



Letter E



Letter F

11



Content by  
TERM 062  
A Custom Edition



+17

R

KP

S

K

C

chavanshil04

RS

Rushikesh Sheke

SP

Suthanshu Pusadkar

NM

Nida Mujawar

S

salleekan

SD

Smit Deshmukh

BM

Bushrahamd Moroni

# Decision Tree Induction

Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.

S S SD RS BM

Salinakim Shivansavant Smit Deshmukh Rushikesh Sherkar Bashirahamad Momin

## Introduction

- learning of decision trees from class-labeled training tuples.
- A **decision tree** is a flowchart-like tree structure, where each **internal node** (nonleaf node) denotes a test on an attribute, each **branch** represents an outcome of the test, and each **leaf node** (or *terminal node*) holds a class label.

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg., Sangli, INDIA

+2 VK R SP 2H KP S S S SD RS BM

Kristina Pod Seemantidatt Salinakim Shivansavant Smit Deshmukh Rushikesh Sherkar Bashirahamad Momin

## ...Cont Introduction

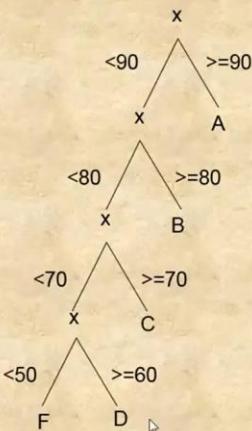
- Internal nodes are denoted by **rectangles**, and leaf nodes are denoted by **ovals**.
- Some decision tree algorithms produce only **binary** trees (where each internal node branches to exactly two other nodes), whereas others can produce **nonbinary** trees.

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Wничанд College of Engg., Sancili, INDIA



## Decision Tree for Grading

- If  $x \geq 90$  then grade =A.
- If  $80 \leq x < 90$  then grade =B.
- If  $70 \leq x < 80$  then grade =C.
- If  $60 \leq x < 70$  then grade =D.
- If  $x < 50$  then grade =F.



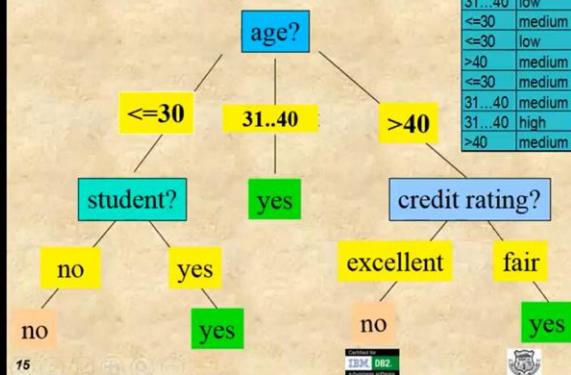
16

Created by  
**IBM DB2**  
A business software



## Decision Tree Induction: An Example

- Training data set: Buys\_computer
- Resulting tree:



age	income	student	credit rating	buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31..40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31..40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31..40	medium	no	excellent	yes
31..40	high	yes	fair	yes
>40	medium	no	excellent	no

## How are the decision trees used for classification ?

- Given a tuple,  $X$ , for which the associated class label is unknown, the attribute values of the tuple are tested against the decision tree.
- A path is traced from the root to a leaf node, which holds the class prediction for that tuple – **tree traversal**.
- Decision trees can easily be converted to **classification rules**
- #rules = #path from root to leaf node

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg. Sangli, INDIA

## **Why are decision tree classifiers so popular?**

- The construction of decision tree classifiers does not require any domain knowledge or parameter setting.
- appropriate for exploratory knowledge discovery.
- can handle multidimensional data.
- representation of acquired knowledge in tree form is intuitive and easy to assimilate/understand by humans
- decision tree classifiers have good accuracy.

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg., Sangli, INDIA



## **How to build Decision Tree ?**

Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.



## History

- During the late 1970s and early 1980s, J. Ross Quinlan, a researcher in machine learning, developed a decision tree algorithm known as ID3 (Iterative Dichotomiser).
- Quinlan later presented C4.5 (a successor of ID3), which became a benchmark to which newer supervised learning algorithms are often compared.
- In 1984, a group of statisticians (L. Breiman, J. Friedman, R. Olshen, and C. Stone) published the book *Classification and Regression Trees (CART)*, which described the generation of binary decision trees.
- ID3 and CART were invented independently of one another at around the same time, yet follow a similar approach for learning decision trees from training tuples.
- ID3 , C4.5 & CART** algorithms becomes de-facto standards for decision tree induction

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg., Sangli, INDIA

The screenshot shows a Microsoft OneNote interface during a meeting. The title bar reads "DM #17 Module III Decision Tree Induction 20210922 111519 Meeting Recording". The main content area displays a hand-drawn decision tree diagram with red and green lines and circles. Below the diagram, a box contains the following text:

1] Constructing/building the model : Needs to train by training set  
2] Evaluation / testing : test data set  
Division : 80 20 or 70 30

The OneNote ribbon shows tabs like Tools, Shapes, Edit, and Convert. A search bar is visible on the right. A sidebar on the right lists various data science and machine learning topics. The taskbar at the bottom shows icons for various applications, and the system tray indicates the date and time as 11:45 AM, 9/22/2021.

## Greedy Algorithm

- **ID3, C4.5, and CART** adopt a greedy (i.e., nonbacktracking) approach.
- decision trees are constructed in a top-down recursive divide-and-conquer manner.
- It starts with a training set of tuples and their associated class labels.
- The training set is recursively partitioned into smaller subsets as the tree is being built.

Dr. Bashirahamed F. Momin  
Department of Computer Science & Engineering  
Wничанд College of Engg., Savitri, INDIA

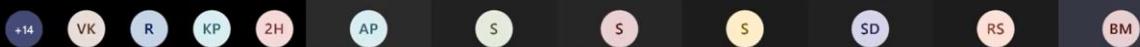


## Algorithm for Decision Tree Induction

- Basic algorithm (a greedy algorithm)
  - Tree is constructed in a **top-down recursive divide-and-conquer manner**
  - At start, all the training examples are at the root
  - Attributes are categorical (if continuous-valued, they are discretized in advance)
  - Examples are partitioned recursively based on selected attributes
  - Test attributes are selected on the basis of a heuristic or statistical measure – **Entropy , Gain Ratio & Gini Index**)

22

Powered by  
**IBM DB2**  
A business software



## Algorithm for Decision Tree Induction

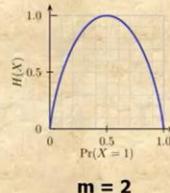
- Basic algorithm (a greedy algorithm)
  - Tree is constructed in a **top-down recursive divide-and-conquer manner**
  - At start, all the training examples are at the root
  - Attributes are categorical (if continuous-valued, they are discretized in advance)
  - Examples are partitioned recursively based on selected attributes
  - Test attributes are selected on the basis of a heuristic or statistical measure – **Entropy , Gain Ratio & Gini Index**)
- Conditions for stopping partitioning
  - All samples for a given node belong to the same class
  - There are no remaining attributes for further partitioning – **majority voting** is employed for classifying the leaf
  - There are no samples left

22

Generated by  
IBM DB2  
A database software

## Entropy

- Entropy (Information Theory)
  - A measure of uncertainty associated with a random variable
  - Calculation: For a discrete random variable  $Y$  taking  $m$  distinct values  $\{y_1, \dots, y_m\}$ ,
    - $H(Y) = -\sum_{i=1}^m p_i \log(p_i)$ , where  $p_i = P(Y = y_i)$
  - Interpretation:
    - Higher entropy => higher uncertainty
    - Lower entropy => lower uncertainty
- Conditional Entropy
  - $H(Y|X) = \sum_x p(x)H(Y|X = x)$



23

Generated by  
IBM DB2  
A database software

## Gini Index (CART, IBM IntelligentMiner)

- If a data set  $D$  contains examples from  $n$  classes, gini index,  $gini(D)$  is defined as

$$gini(D) = 1 - \sum_{j=1}^n p_j^2$$

where  $p_j$  is the relative frequency of class  $j$  in  $D$

- If a data set  $D$  is split on A into two subsets  $D_1$  and  $D_2$ , the gini index  $gini(D)$  is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

- Reduction in Impurity:

- The attribute provides the smallest  $gini_{\text{split}}(D)$  (or the largest reduction in impurity) is chosen to split the node (*need to enumerate all the possible splitting points for each attribute*)

26

Created by  
IBM DB2  
A business software



## Attribute Selection using Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- Let  $p_i$  be the probability that an arbitrary tuple in  $D$  belongs to class  $C_i$ , estimated by  $|C_i \cap D| / |D|$
- Expected information** (entropy) needed to classify a tuple in  $D$ :

$$Info(D) = - \sum_{i=1}^m p_i \log_2(p_i)$$

- Information** needed (after using A to split D into v partitions) to classify D:

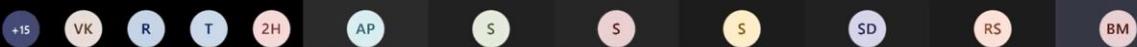
$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

- Information gained** by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

24

Created by  
IBM DB2  
A business software



# ENTROPY AND INFORMATION GAIN

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

$Values(Wind) = Weak, Strong$

$$S = [9+, 5-]$$

$$S_{Weak} \leftarrow [6+, 2-]$$

$$S_{Strong} \leftarrow [3+, 3-]$$

$$\begin{aligned} Gain(S, Wind) &= Entropy(S) - \sum_{v \in \{Weak, Strong\}} \frac{|S_v|}{|S|} Entropy(S_v) \\ &= Entropy(S) - (8/14)Entropy(S_{Weak}) \\ &\quad - (6/14)Entropy(S_{Strong}) \\ &= 0.940 - (8/14)0.811 - (6/14)1.00 \\ &= 0.048 \end{aligned}$$

Subscribe to Mahesh Huddar

Visit: vtupulse.com

## Gain Ratio for Attribute Selection (C4.5)

- Information gain measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$

$$- \text{GainRatio}(A) = \text{Gain}(A)/\text{SplitInfo}(A)$$

- Ex:  $SplitInfo_{income}(D) = -\frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) - \frac{6}{14} \times \log_2 \left( \frac{6}{14} \right) - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) = 1.557$
- $- \text{gain\_ratio}(income) = 0.029/1.557 = 0.019$
- The attribute with the maximum gain ratio is selected as the splitting attribute

25

Powered by  
IBM DB2  
A business software



⇒ DM #17 Module III Decision Tree Induction 20210922 111519 Meeting Recording 

## Gain Ratios for Attribute Selection (C4.5)

- Information gain measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)
$$SplitInfo_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$
- GainRatio(A) = Gain(A)/SplitInfo(A)
- Ex.
$$SplitInfo_{income}(D) = -\frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) - \frac{6}{14} \times \log_2 \left( \frac{6}{14} \right) - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) = 1.557$$
- gain\_ratio(income) = 0.029/1.557 = 0.019
- The attribute with the maximum gain ratio is selected as the splitting attribute

25

45:12 / 1:00:22

Abhishek Pati seminarparallelität

S

S

S

SD

RS

CC

Settings

BM

Rushikesh Shekde

Bashirahmed Munim

## Comparing Attribute Selection Measures

- The three measures, in general, return good results but
  - **Information gain:**
    - biased towards multivalued attributes
  - **Gain ratio:**
    - tends to prefer unbalanced splits in which one partition is much smaller than the others
  - **Gini index:**
    - biased to multivalued attributes
    - has difficulty when # of classes is large
    - tends to favor tests that result in equal-sized partitions and purity in both partitions

27

45:12 / 1:00:22

Abhishek Pati seminarparallelität

S

S

S

SD

RS

CC

Settings

BM

Rushikesh Shekde

Bashirahmed Munim

Decision Tree - OneNote

FILE HOME INSERT DRAW HISTORY REVIEW VIEW

Type Lasso Panning Eraser Select Hand Tools

Shapes Edit

Press Esc to exit full screen

Ink to Text Math Convert

demo1 New Section 1 New Section 2 +

+ Add Page

+ve skewed data  
-ve Skewed data  
Python code  
Dispersion : Spread of data  
Variance and standard deviation  
Histogram  
Data Pro-processing  
Data Cleaning  
Noisy Data :  
Data Cleaning as a "Process"  
Data Integration  
Chi-Square Test  
Correlation coefficient for Nu  
Covariance of Nu  
Covariance of Numeric data  
Tuple Duplication  
Data transformation and Discr  
Normalization  
Discretization by Binning  
Discretization by Histogram A  
Classification - Basic Concept  
Decision Tree

1] Constructing/building the model : Needs to train by training set  
2] Evaluation / testing : test data set  
Division : 80 20 or 70 30

1. Start at node N and present whole training data set  
2. Apply the attribute selection strategy (Entropy, Gain Ratio or Gini Index) to attribute split  
3. Divide the data set based on selected attribute.  
4. Repeat step 2 , 3 till all data set covered or there is no attribute remain

12:15 PM  
9/22/2021

+13 VK R S 2H AP S S S S SD RS BM

Abhishek Patil sadiqdarw13102000 Saitheekan shivansavant Smit Deshmukh Rushikesh Sheke Bushrahamad Morin

## Lecture 18

⇒ DM #18 101520 Meeting Recording

**Training Data Set**

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Dr. Bashirahamad F. Momin  
Department of Computer Science & Engineering  
Walchand College of Engg., Sangli, INDIA

◀ ▶ 4:14 / 59:26

RS NM AP SE CC G S B

⇒ DM #18 101520 Meeting Recording

$|D| = 14$

$$Info(D) = \sum_{i=1}^m p_i \log_2 p_i$$

$$Info(D_{youth}) = -\{p_1 \log_2 p_1 + p_2 \log_2 p_2\}$$

$$p_1 = \frac{|C_1, \text{Dage}|}{|D|} = \frac{\text{No. of tuples for class 1 having youth}}{\text{Total No. of tuples for youth value}}$$

$$= \frac{2}{5}$$

X 1.6

\*◀ S ▶ □ 🔍 AF II ● ○ ■ 🔍 \*

27:59 / 59:26

K RS NM AP SE CC G S B

⇒ DM #18 101520 Meeting Recording

$$= 0.140 - 0.093 = \underline{0.246}$$

(compute similarity for other attributes)

$$\text{Gain}(\text{income}) = 0.029$$

$$\text{Gain}(\text{student}) = 0.151$$

$$\text{Gain}(\text{creditrating}) = 0.048$$

Choose among the attributes having higher gain  
∴ "Age" is split attribute

X 1.6

38:10 / 59:26

Veetja Kadam

Kishorekumar

Ruturaj Shete

Nida Majeed

Ashish Patel

Smit Deshmukh

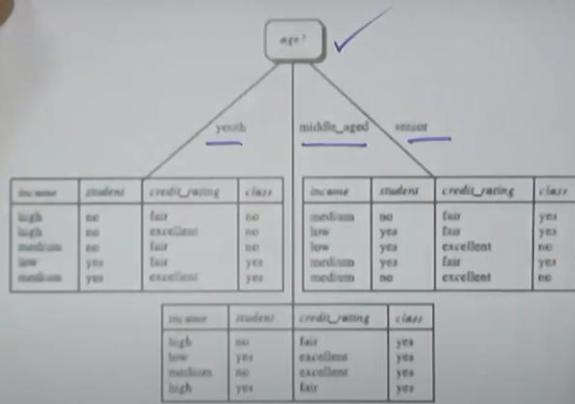
SE CC

Settings



⇒ DM #18 101520 Meeting Recording

#### Splitting of Dataset at attribute "age"



X 1.6

Settings

40:04 / 59:26

Veetja Kadam

Kishorekumar

Ruturaj Shete

Nida Majeed

Ashish Patel

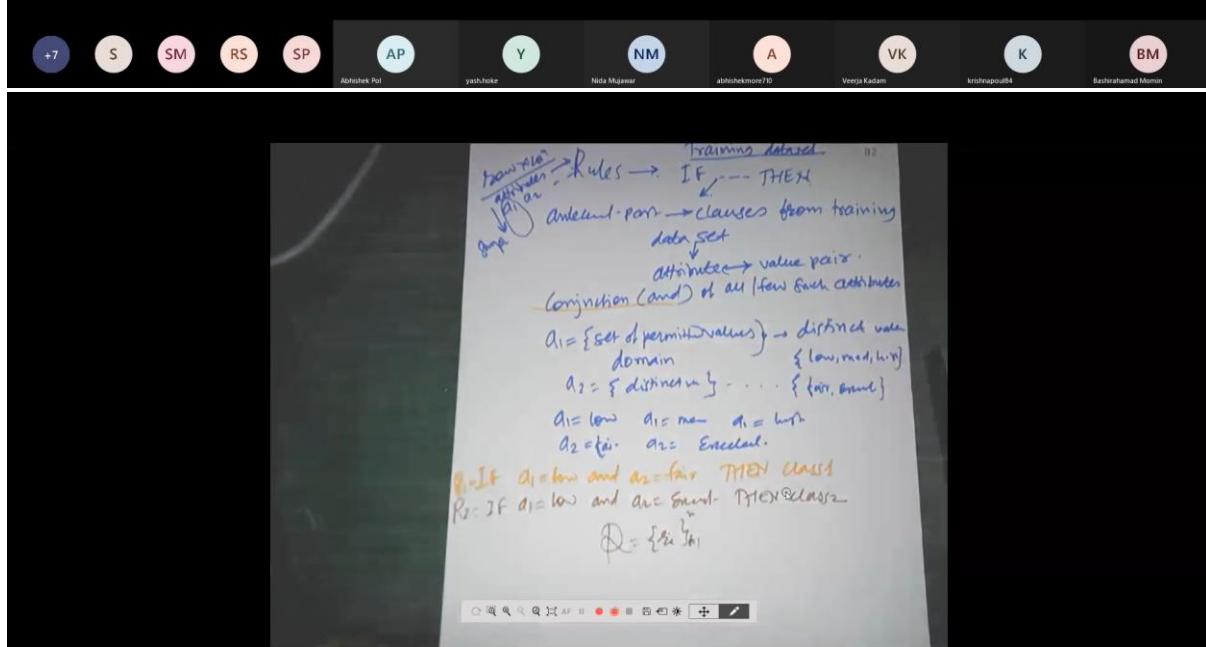
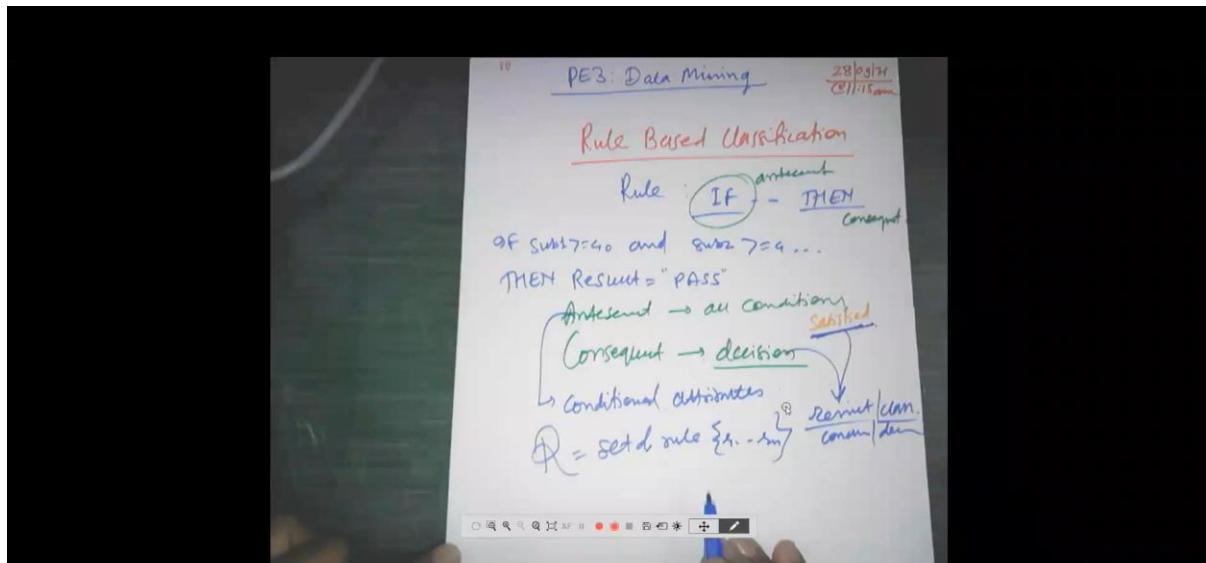
Smit Deshmukh

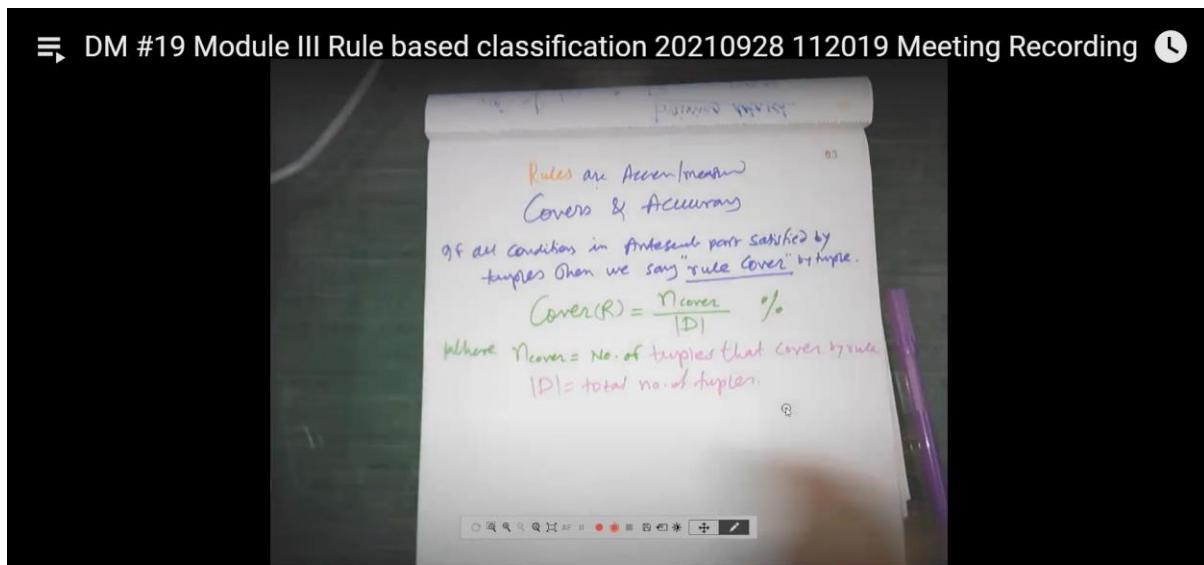
SE CC

HD



## Lecture 19





◀ S ▶ 🔍 28:45 / 55:34
Y NM A VK K CC HD BM

Ashishk Pali
yash.hiske
Nida Mujawir
abhishekmore710
Veenja Kadam
krishnapo04
Bashirahmed Monim

Rules are Accen/measur  
Covers & Accuracy

If all conditions in Antecedent part satisfied by tuples Then we say "rule Cover" by tuple.

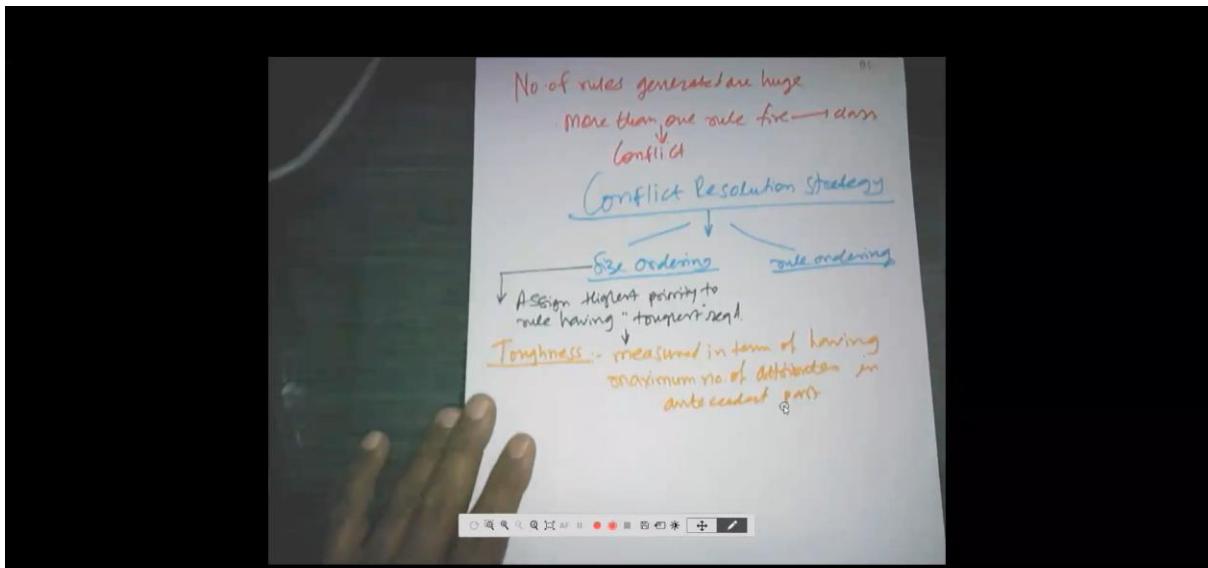
$$\text{Cover}(R) = \frac{n_{\text{cover}}}{|D|} \% \quad \text{---} \odot$$

Where  $n_{\text{cover}}$  = No. of tuples that cover by rule.  
 $|D|$  = total no. of tuples.

$$\text{Accuracy} = \frac{n_{\text{correct}}}{n_{\text{cover}}} \% \quad \text{---} \odot$$

$n_{\text{correct}}$  = No. of tuples that correctly classified by rule R.

+7 S SM RS SP AP Y NM A VK K BM
Abhishek Pali yash.hiske Nida Mujawir abhishekmore710 Veenja Kadam krishnapo04 Bashirahmed Monim



DM #19 Module III Rule based classification 20210928 112019 Meeting Recording

No of rules generated are huge  
more than one rule fire down  
Conflict

Conflict Resolution Strategy

↓  
Size ordering      rule ordering

↓ Assign highest priority to rule having "toughest test".

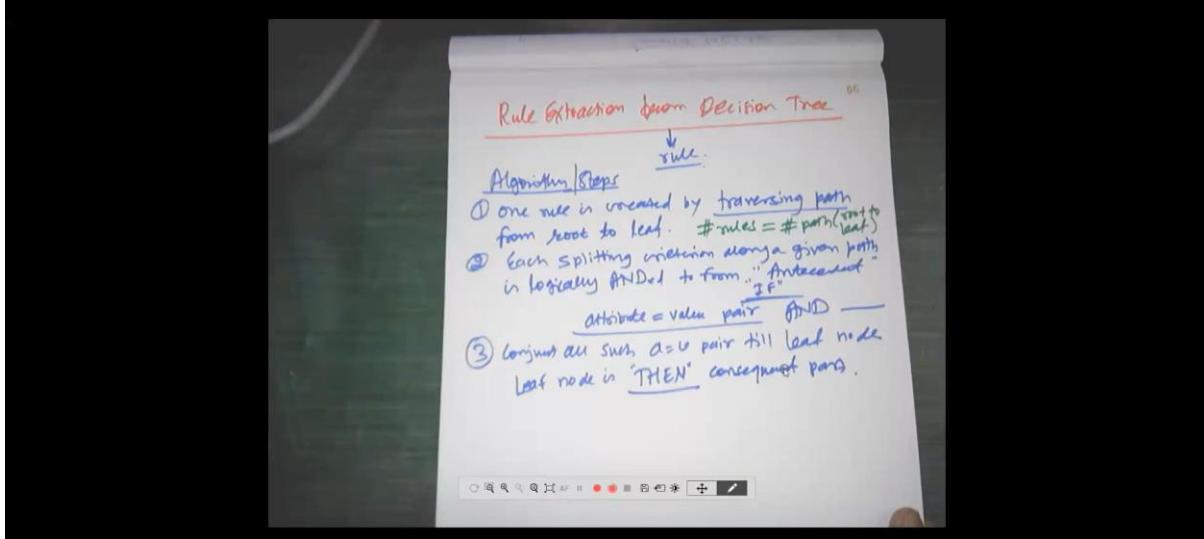
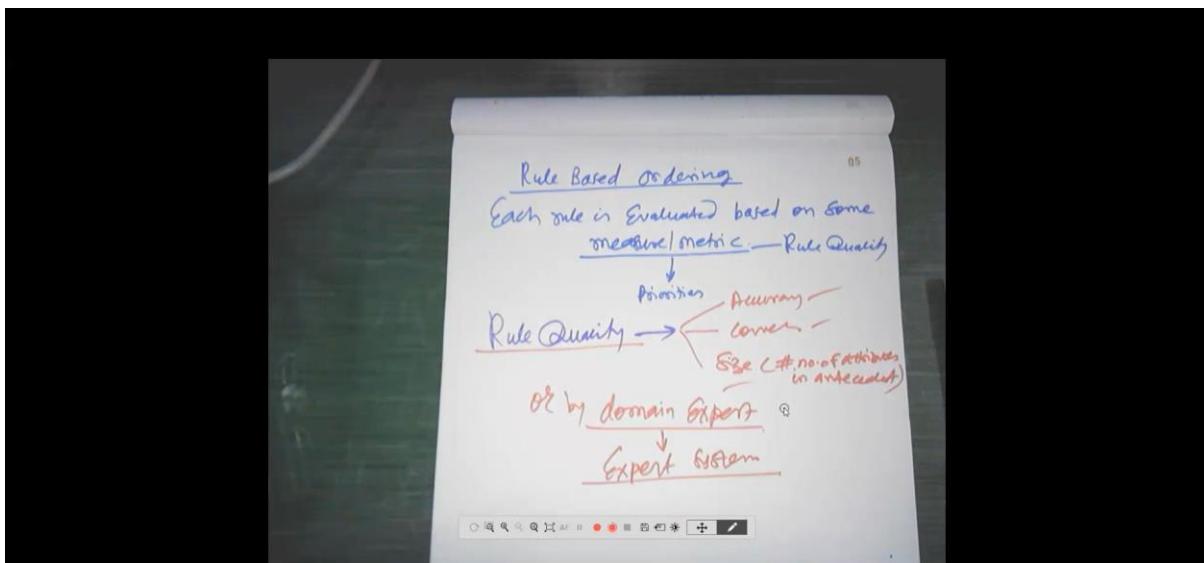
Toughness - measured in term of having maximum no. of attributes in antecedent part

↓ Class-Based ruleBased  
Soci on classes eg. a rule with C1  
Importance

Participants shown in the video player interface:

- +7
- S
- SM
- RS
- SP
- AP
- Y
- NM
- A
- VK
- K
- BM

Abhishek Pati, yashchokhe, Nida Mujawar, abhishekmore710, Veerja Kadam, krishnapaul84, Bashrahmad Momin



⇒ DM #19 Module III Rule based classification 20210928 112019 Meeting Recording ⏱

The whiteboard displays a handwritten decision tree diagram titled "Decision Tree To Buy Computer". The root node is labeled "Age". The tree branches into three categories: "Young", "middle", and "senior". The "Young" branch leads to two nodes: "Student" (Yes) and "Student" (No). The "Yes" node leads to a leaf node labeled "BuyComputer = No". The "No" node leads to another node labeled "CreditRating". From "CreditRating", the tree further branches into "fair" (Yes) and "Excellent" (Yes), both leading to a final leaf node labeled "BuyComputer = Yes".

Handwritten text below the tree:

- If age = "Young" and Student = "No" THEN BuyComputer = No
- If age = "Young" and CreditRating = "Yes" THEN BuyComputer = Yes
- If age = "middle"
- THEN BuyComputer = "Yes"

Below the whiteboard, there is a video player interface with various controls and participant names.

Video player controls: +, ▶, S, ||, ⏪, 🔍, 54:40 / 55:34, RS, A, VK, CC, HD, B4+, etc.

Participants: Rohitkesh Pati, yashikesh, Rukhsar Shekar, gRohitkeshwar710, Veena Kadam, kritiagupta14, Bachrahamad Munim

## Lecture 20

DM #20 Module III Metrics for Evaluating classifier performance Statistical base... 

Accuracy (Good)

① accuracy / Recognition Rate =  $\frac{TP+TN}{P+N}$

where TP: True Positive  
 TN: True Negative  
 P: #positive tuples  
 N: #negative tuples.

e.g. Person Buys Computer = "Yes" & "No"

Correctly classified.  
 No. of tuples having positive class → positive class  
 No. of tuples having negative class → negative class

CLASS Imbalance

with Data set having less no. of tuples/rare Samples of "main class of interest".

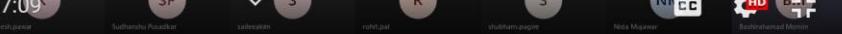
e.g. XYZ dataset

100 Samples Two class problem

Class Imbalance →  
 C1: +ve class - class of interest  
 = 3 sample tuple  
 C2: -ve class  
 = 97 tuples.

Accur % =  $\frac{TP+TN}{P+N} = 94\% \rightarrow$  class inter.

→ Contribution is less/rare.


For Class Imbalance problem 12

Metrics      Sensitivity & Specificity

true faceted picture of Quality classifier

$\frac{TP}{P}$        $\frac{TN}{N}$

Covid-19       $\Leftrightarrow$  true Posit.

✓  $TP: 3 \quad \frac{3}{5} \cdot 100\% = 60\%$

Uniform distribution in data set: No. of samples  
Balanced samples  $\Rightarrow$  Scatter plot

Position mean Bd.

abhishekmore710 Sudhanshu Pusadkar salilekum sawanpandita21 Hritik Belari Veeja Kadam Bachrahmad Momin

Performance of Classifier is given by

- ① Speed ✓ Computational Cost
- ② Robustness → handle noisy data  
missing value dataset
- ③ Scalability → To handle large Data Set
- ④ Interpretability → Understanding of Knowledge Extracted  
Rule → Human being



## Lecture 21

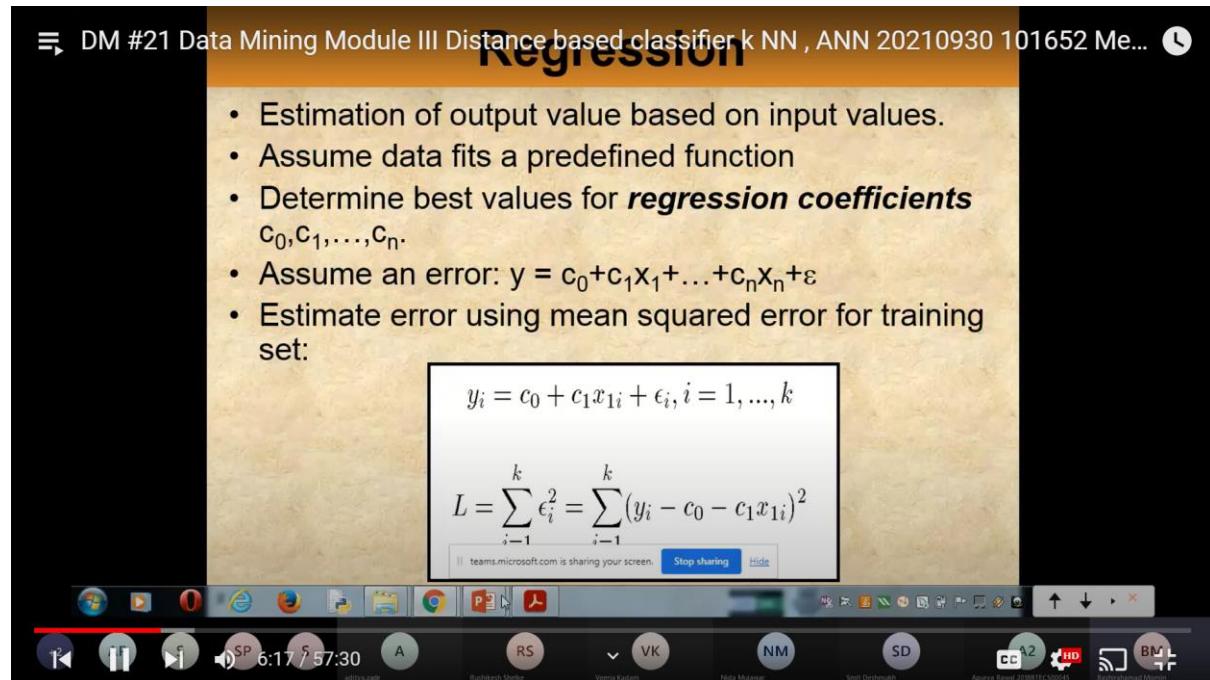
DM #21 Data Mining Module III Distance based classifier k NN , ANN 20210930 101652 Me... ⏱

# Regression

- Estimation of output value based on input values.
- Assume data fits a predefined function
- Determine best values for **regression coefficients**  $c_0, c_1, \dots, c_n$ .
- Assume an error:  $y = c_0 + c_1 x_1 + \dots + c_n x_n + \epsilon$
- Estimate error using mean squared error for training set:

$$y_i = c_0 + c_1 x_{1i} + \epsilon_i, i = 1, \dots, k$$
$$L = \sum_{i=1}^k \epsilon_i^2 = \sum_{i=1}^k (y_i - c_0 - c_1 x_{1i})^2$$

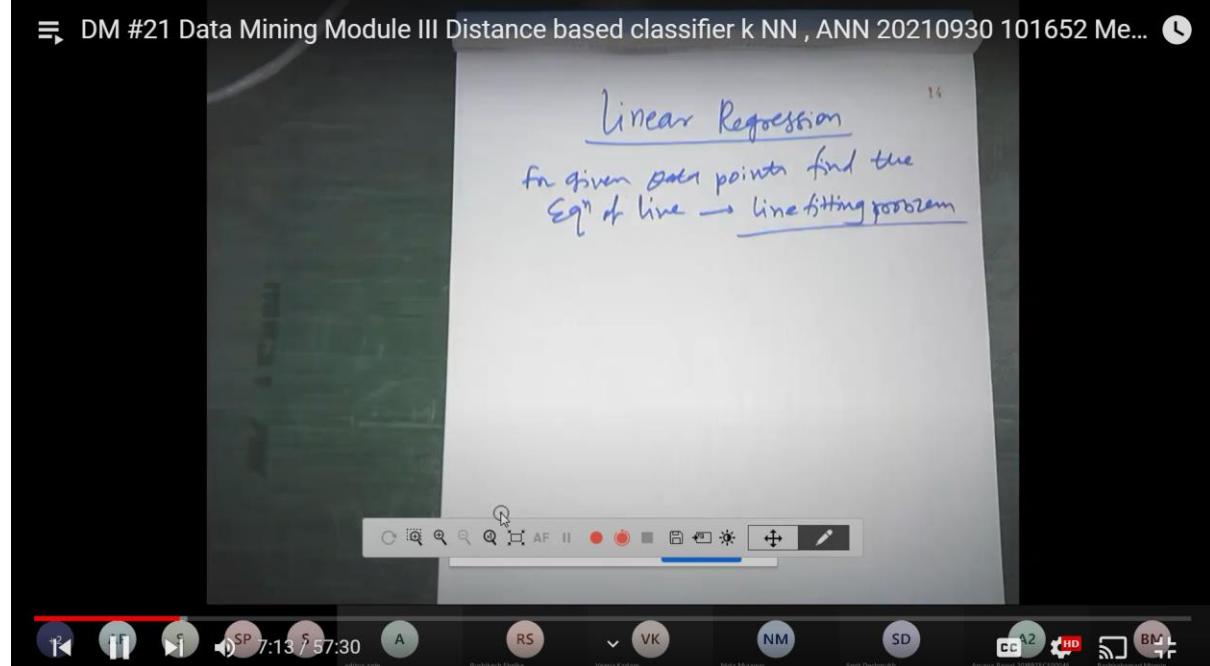
|| teams.microsoft.com is sharing your screen. Stop sharing Hide

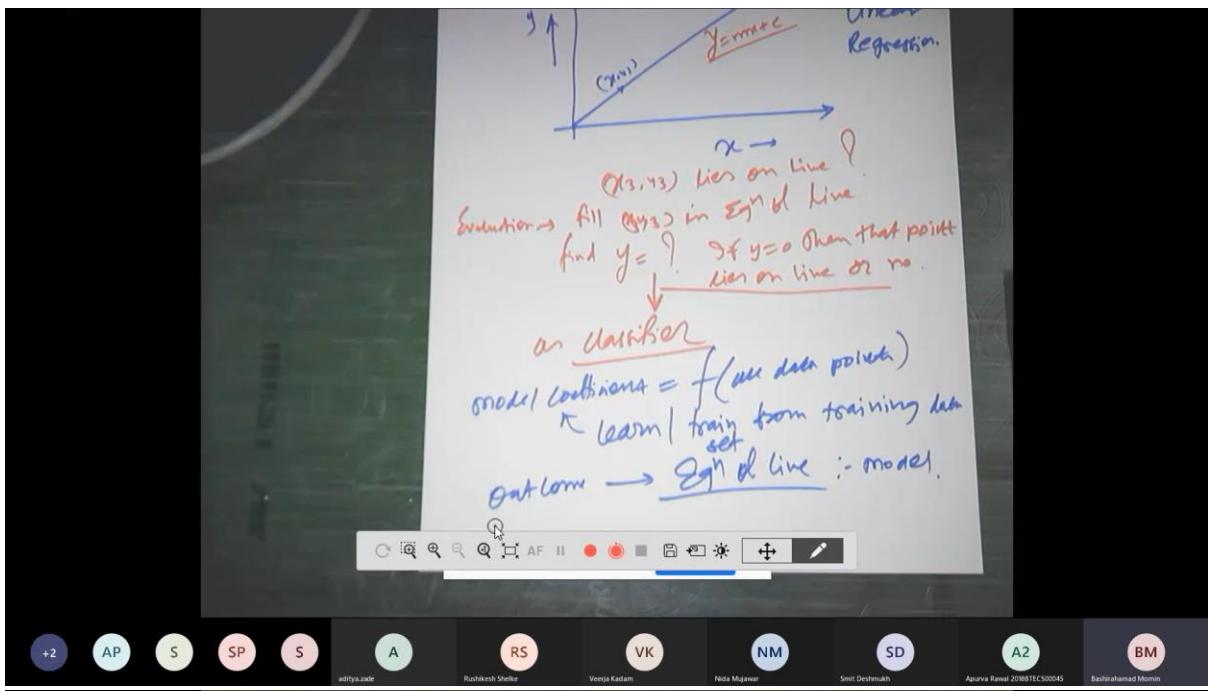


DM #21 Data Mining Module III Distance based classifier k NN , ANN 20210930 101652 Me... ⏱

# Linear Regression

for given data points find the Eq<sup>n</sup> of line → line fitting problem





DM #21 Data Mining Module III Distance based classifier k NN , ANN 20210930 101652 Me... (1)

# Bayesian Classifier

5 teams.microsoft.com is sharing your screen. Stop sharing Hide

TSM 082

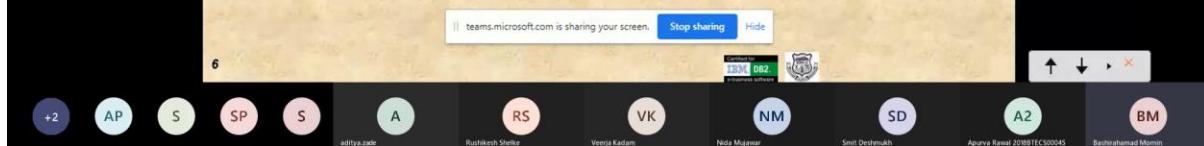
18:37 / 57:30

aditya.zade Rushikesh Shinde Veerja Kadam Nida Mujawar Smit Deshmukh Apurva Rawal 20BITEC50045 Basitrahmad Momin

AP SP S A RS VK NM SD A2 BM

## Bayesian Classification: Why?

- A statistical classifier: performs *probabilistic prediction*, i.e., predicts class membership probabilities
- Foundation: Based on Bayes' Theorem.
- Performance: A simple Bayesian classifier, *naive Bayesian classifier*, has comparable performance with decision tree and selected neural network classifiers
- Incremental: Each training example can incrementally increase/decrease the probability that a hypothesis is correct — prior knowledge can be combined with observed data
- Standard: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured



## Classification Using Distance

- Place items in class to which they are “closest”.
- Must determine distance between an item and a class.
- Classes represented by
  - **Centroid:** Central value.
  - **Medoid:** Representative point.
  - Individual points
- Algorithm:  $k$ -NN

## $k$ -NN Algorithm

```

Input:
  D      //Training data
  K      //Number of neighbors
  t      //Input tuple to classify
Output:
  c      //Class to which t is assigned
KNN Algorithm:
    //Algorithm to classify tuple using KNN
    N = ∅;
    //Find set of neighbors, N, for t
    foreach d ∈ D do
      if |N| ≤ K then
        N = N ∪ d;
      else
        if ∃ u ∈ N such that sim(t, u) ≥ sim(t, d) then
          begin
            N = N - u;
            N = N ∪ d;
          end
        //Find class for classification
    c = class to which the most u ∈ N are classified;
  
```

# Nearest Neighbor (NN) Classifier

- Based on learning by analogy
- Better for data set which are :
  - Non-parametric in nature
  - Multivariate dependencies
- Closeness measure
  - Euclidean Distance (ED)
  - Cosine Coefficient (CC)
  - Pearson Coefficient etc. ...

18



## Nearest Neighbor Classifier (1-NN)

Unknown pattern  
 $U_1, \dots, U_n$

- Find the nearest\* known sample
- Assign class label of nearest sample to unknown pattern

\*Distance metric : Euclidean Distance

$$ED(U, X) = \sqrt{\sum_{i=1}^n (U_i - X_i)^2}$$

Samples	X <sub>1</sub> , ..., X <sub>n</sub>	Class Label
V <sub>1</sub>	X <sub>1</sub> , ..., X <sub>n</sub>	C <sub>1</sub>
V <sub>m</sub>		C <sub>k</sub>



19



↑

↓

→

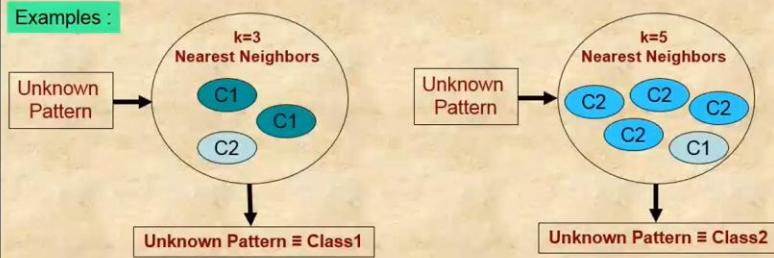
←

+2

## k-Nearest Neighbor (k-NN) Classifier

- Choose “ $k$ ”  $\geq 3$  (odd number)
- Find “ $k$ ” nearest neighbors using Euclidean Distance.
- Majority voting : select class label assign to maximum of “ $k$ ”-neighbors.
- Assign selected class label to unknown pattern.

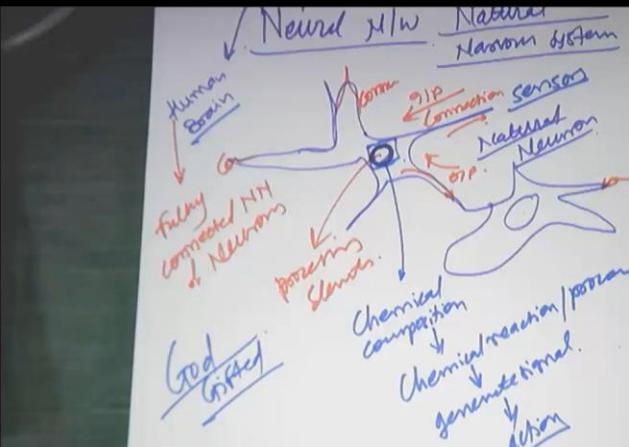
Examples :



20



rohit.pai +2



DM #23 Meeting in General 20211005 112139 Meeting Recording

The whiteboard contains the following handwritten notes:

- Machine learning
- NNN  $\rightarrow$  Artificial NN (ANN)
- Neuron  $\rightarrow$  mathematical model
- NNN  $\rightarrow$  Set of Neurons (ANN)
  - fully connected
  - (topology)

Below the whiteboard is a video player interface showing the timestamp 39:08 / 53:54.

DM #23 Meeting in General 20211005 112139 Meeting Recording

## Perceptron

A diagram of a single-layer neural network (Perceptron) with three input nodes ( $x_1$ ,  $x_2$ ,  $B$ ) and one output node ( $f_4$ ). The connections and their weights are:

- $x_1$  to  $f_4$  with weight 3
- $x_2$  to  $f_4$  with weight 2
- $B$  to  $f_4$  with weight -6

Below the diagram is a bulleted list:

- Perceptron is one of the simplest NNs.
- No hidden layers.

At the bottom of the slide is the number 39.

Below the slide is a video player interface showing the timestamp 39:22 / 53:54.

**Types of NNs**

- Different NN structures used for different problems.
- Perceptron
- Self Organizing Feature Map
- Radial Basis Function Network

38

SD Smit.Deshmukh abhishekmore710  
 NM Nida.Mujawar VK Veerja.Kadam  
 RS Rushikesh.Shetke prachi.wanware  
 S saileesam kalpresh.pansar  
 Y yash.hoke srujanabhirugade  
 S sanyarpandita21 shivani.avasthi  
 SS Suyash.Sajji KP Krishna.Poul  
 SM Shreyas.Mandale S shubham.page  
 R rohit.pai +2

(4) Meeting | Microsoft Teams

DM\_Module\_III\_2\_of\_2 [Compatibility Mode] - PowerPoint

**Perceptron Example**

- Suppose:
  - Summation:  $S = 3x_1 + 2x_2 - 6$
  - Activation: if  $S > 0$  then 1 else 0

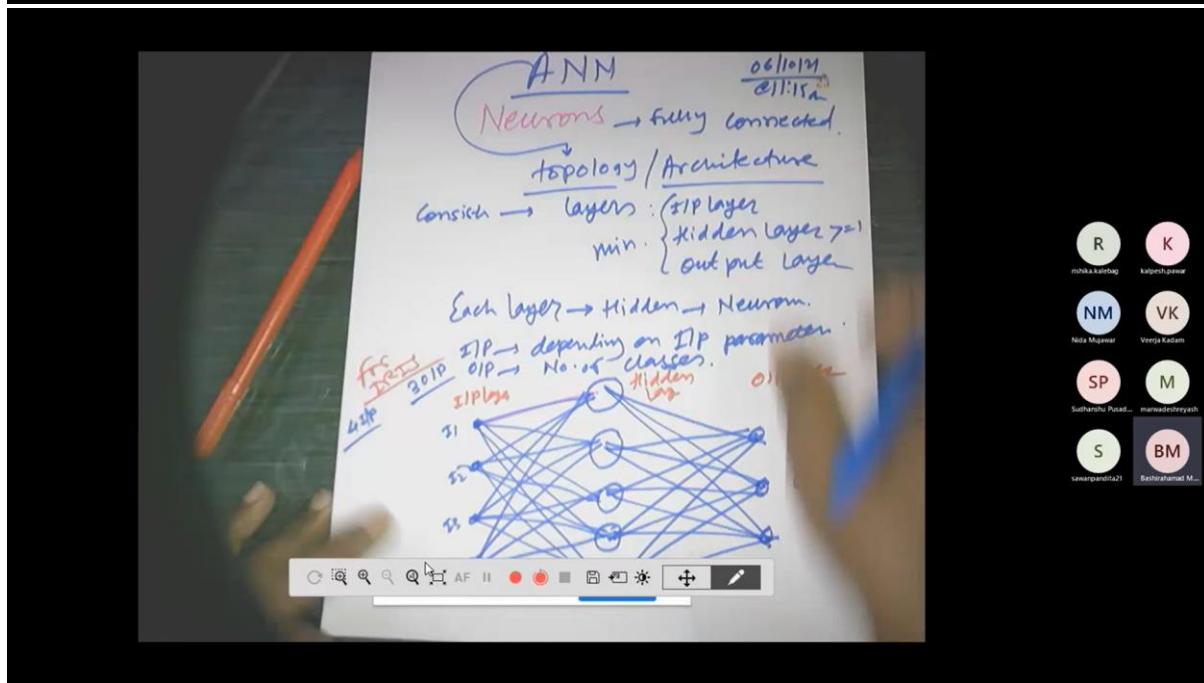
a) Classification Perceptron b) Classification Problem

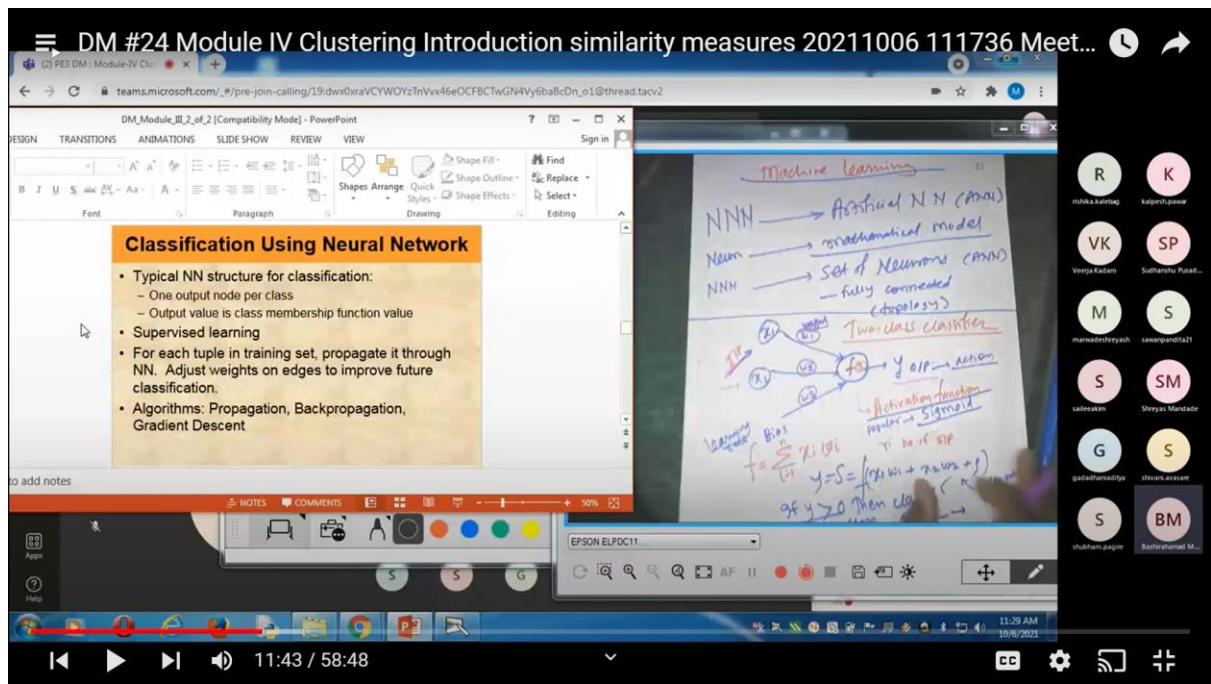
NN → Artificial NN (ANN)  
 Neuron → mathematical model  
 ANN → Set of Neurons (ANN) – fully connected (topology)  
 $f(x) = \sum_i w_i x_i + b$  Activation function  
 Bias  
 $S = x_1 w_1 + x_2 w_2 + b$

+11 A2 meetingAttendance.csv 12:07 PM 10/5/2021

## DM #23 Meeting in General 20211005 112139 Meeting Recording

The screenshot shows a Microsoft Teams meeting interface. On the left, a PowerPoint slide titled "Perceptron Example" is displayed. It contains a diagram of a perceptron with two inputs and one output, labeled "Classification Perceptron". To its right is a graph of a linear classification boundary in a 2D space with axes  $x_1$  and  $x_2$ . Below the slide is a note: "Suppose: - Summation:  $S = 3x_1 + 2x_2 - 6$  - Activation: if  $S > 0$  then 1 else 0". On the right side of the screen, a whiteboard is visible with handwritten notes. The notes include: "Natural N → Remembrance Learning", "How in Artificial / machine", "Pattern recognition → ST.", "Weights → Associative memory", "Training → learning in weights update", "Iteration loop", "Velocity", and "MIN in trained". A red arrow points from the text "Weights" to the word "Weights" on the whiteboard. The Teams interface shows a list of participants on the right and a control bar at the bottom.





## NN Issues

- Number of source nodes
- Number of hidden layers
- Training data
- Number of sinks
- Interconnections
- Weights
- Activation Functions
- Learning Technique
- When to stop learning (convergence)

## NN Learning

- Adjust weights to perform better with the associated test data.
- **Supervised:** Use feedback from knowledge of correct classification.
- **Unsupervised:** No knowledge of correct classification needed.

## NN Supervised Learning

Input:

```
N    //Starting Neural Network  
X    //Input tuple from Training Set  
D    //Output tuple desired
```

Output:

```
N    //Improved Neural Network
```

SupLearn Algorithm:

```
    //Simplistic algorithm to illustrate approach to NN learning
```

```
    Propagate X through N producing output Y;
```

```
    Calculate error by comparing D to Y;
```

```
    Update weights on arcs in N to reduce error;
```

## Supervised Learning

- Possible error values assuming output from node  $i$  is  $y_i$  but should be  $d_i$ :

$$\begin{aligned} & |y_i - d_i| \\ & \frac{(y_i - d_i)^2}{2} \\ & \sum_{i=1}^m \frac{(y_i - d_i)^2}{m} \end{aligned}$$

- Change weights on arcs based on estimated error

## NN Backpropagation

- Propagate changes to weights backward from output layer to input layer.
- Delta Rule:**  $\Delta w_{ij} = c x_{ij} (d_j - y_j)$
- Gradient Descent:** technique to modify the weights in the graph.

DM #24 Module IV Clustering Introduction similarity measures 20211006 111736 Meet...

**Backpropagation Algorithm**

```

Input:
N           //Starting Neural Network
X ==< x1, ..., xk > //Input tuple from Training Set
D ==< d1, ..., dm > //Output tuple desired
Output:
X           //Improved Neural Network
Backpropagation Algorithm:
Propagation(N, X);
E = 1/2 ∑i=1m (di - yi)2;
Gradient(N, E);
  
```

**Backpropagation Algorithm**  
Advanced version of feed forward training  
Hidden layer (Output layer)  
Error = desired/actual value - computed value.  
Iteration I in progress. Reckon once.  
Now weight = old weight + error  
W<sub>ij</sub> = W<sub>ij</sub> + D<sub>ij</sub>;  
After adjust (Iteration) & repeat till  
Error is in acceptable range.

**Types of NNs**

- Different NN structures used for different problems.
- Perceptron
- Self Organizing Feature Map
- Radial Basis Function Network

**Diagram Annotations:**

- Supervised learning
- Input layer (4 nodes)
- Hidden layer (3 nodes)
- Output layer (2 nodes)
- Connections between layers
- Activation function
- Iterations → Convergence criterion
- Ideally → Error = 0
- Permissible error
- Weights: Initially random values

## Clustering - Introduction

- Unlike classification, the class label of each training sample is not known : ***Unsupervised Learning.***
- It is the process of grouping a set of training samples or examples into similar classes.
- It groups similar objects based on their distance, connectivity, relative density or some specific characteristics.
- A ***cluster*** is therefore a collection of objects which are “similar” between them and are “dissimilar” to the objects belonging to other clusters
- Different algorithms / techniques :
  - Hierarchical Algorithms
  - Partitional Algorithms
  - Clustering Large Databases

### Clustering Problem

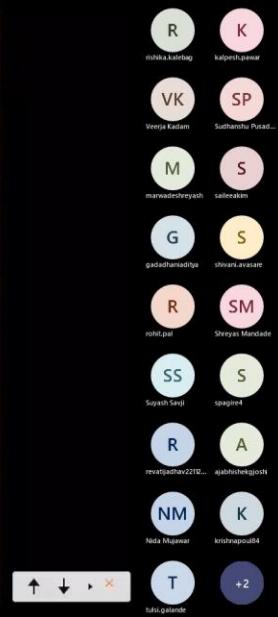
- Given a database  $D=\{t_1, t_2, \dots, t_n\}$  of tuples and an integer value  $k$ , the ***Clustering Problem*** is to define a mapping  $f:D \rightarrow \{1, \dots, k\}$  where each  $t_i$  is assigned to one cluster  $K_j$ ,  $1 \leq j \leq k$ .
- A ***Cluster***,  $K_j$ , contains precisely those tuples mapped to it.
- Unlike classification problem, clusters are not known a priori.



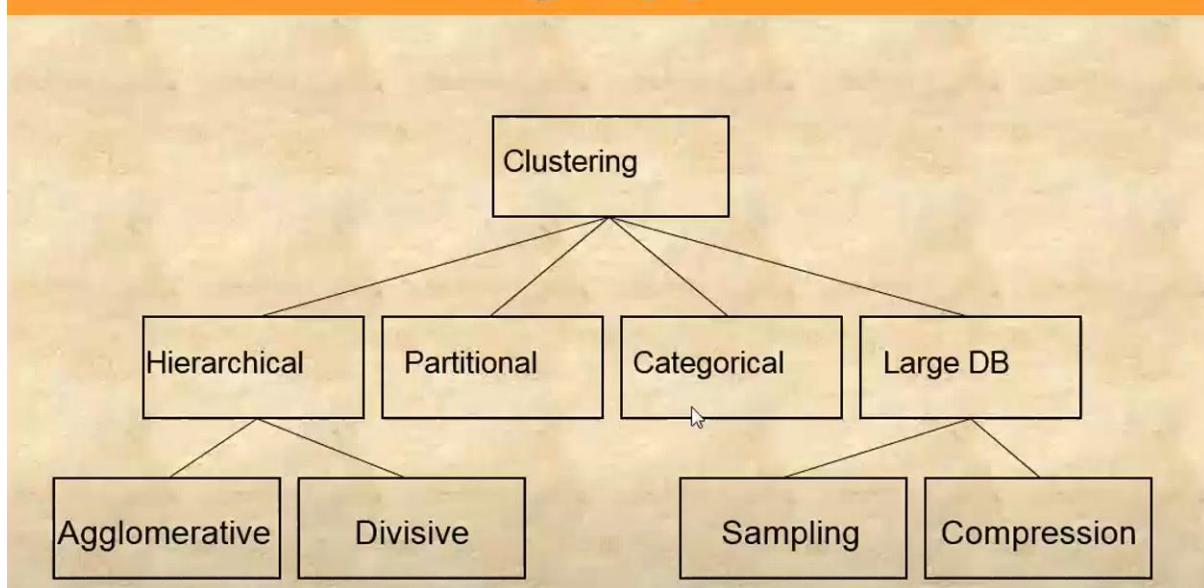
## Types of Clustering

- **Hierarchical** – Nested set of clusters created.
- **Partitional** – One set of clusters created.
- **Incremental** – Each element handled one at a time.
- **Simultaneous** – All elements handled together.
- **Overlapping/Non-overlapping**

5



## Clustering Approaches



## Clustering Examples

- **Segment** customer database based on similar buying patterns.
- Group houses in a town into neighborhoods based on similar features.
- Identify new plant species
- Identify similar Web usage patterns



## Clustering Issues

- Outlier handling
- Dynamic data
- Interpreting results
- Evaluating results
- Number of clusters
- Data to be used
- Scalability



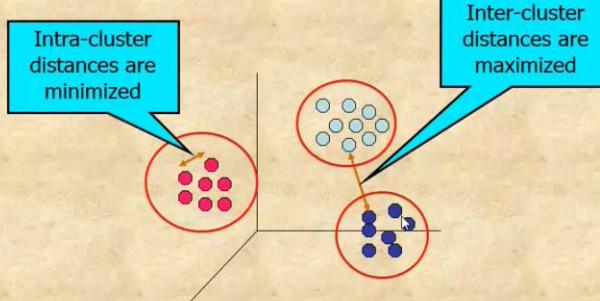
# Clustering Example

Income	Age	Children	Marital Status	Education
\$25,000	35	3	Single	High School
\$15,000	25	1	Married	High School
\$20,000	40	0	Single	High School
\$30,000	20	0	Divorced	High School
\$20,000	25	3	Divorced	College
\$70,000	60	0	Married	College
\$90,000	30	0	Married	Graduate School
\$200,000	45	5	Married	Graduate School
\$100,000	50	2	Divorced	College

9

## Similarity and Distance Measures

Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups



11



# Dissimilarity / Similarity metric

Similarity is expressed in terms of a distance function,  
typically metric :  $d(i, j)$

- Distances are normally used to measure the similarity or dissimilarity between two data objects
- Some popular ones include: *Minkowski distance*:

$$d(i, j) = \sqrt[q]{(|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{ip} - x_{jp}|^q)}$$

where  $i = (x_{i1}, x_{i2}, \dots, x_{ip})$  and  $j = (x_{j1}, x_{j2}, \dots, x_{jp})$  are two  $p$ -dimensional data objects, and  $q$  is a positive integer

- If  $q = 1$ ,  $d$  is Manhattan distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

## Cont... Dissimilarity / Similarity metric

- If  $q = 2$ ,  $d$  is Euclidean distance:

$$d(i, j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

– Properties

- $d(i, j) \geq 0$

- $d(i, i) = 0$

- $d(i, j) = d(j, i)$

- $d(i, j) \leq d(i, k) + d(k, j)$

- Also, one can use weighted distance, parametric Pearson product moment correlation, or other dissimilarity measures



## The Distance between Clusters

- **Single link**: smallest distance between an element in one cluster and an element in the other, i.e.,  $\text{dis}(K_i, K_j) = \min(t_{ip}, t_{jq})$
- **Complete link**: largest distance between an element in one cluster and an element in the other, i.e.,  $\text{dis}(K_i, K_j) = \max(t_{ip}, t_{jq})$
- **Average**: avg distance between an element in one cluster and an element in the other, i.e.,  $\text{dis}(K_i, K_j) = \text{avg}(t_{ip}, t_{jq})$
- **Centroid**: distance between the centroids of two clusters, i.e.,  $\text{dis}(K_i, K_j) = \text{dis}(C_i, C_j)$
- **Medoid**: distance between the medoids of two clusters, i.e.,  $\text{dis}(K_i, K_j) = \text{dis}(M_i, M_j)$ 
  - **Medoid**: one chosen, centrally located object in the cluster

14



## Centroid, Radius and Diameter of a Cluster

- **Centroid**: the “middle” of a cluster

$$C_m = \frac{\sum_{i=1}^N (t_{ip})}{N}$$

- **Radius**: square root of average distance from any point of the cluster to its centroid

$$R_m = \sqrt{\frac{\sum_{i=1}^N (t_{ip} - c_m)^2}{N}}$$

- **Diameter**: square root of average mean squared distance between all pairs of points in the cluster

$$D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (t_{ip} - t_{jq})^2}{N(N-1)}}$$

15

## Clustering

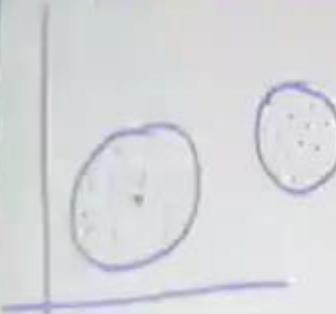
Cluster → Centroid → Medoid

mathematically

arbitrary  
computed  
scalar  
values

center  
middle  
Need not  
be physical point

most centrally  
located  
PHYSICAL  
POINT



## Partitioning Algorithms: Basic Concept

- **Partitioning method:** Construct a partition of a database  $D$  of  $n$  objects into a set of  $k$  clusters, s.t., min sum of squared distance

$$\sum_{m=1}^k \sum_{t_{mi} \in K_m} (C_m - t_{mi})^2$$

- Given a  $k$ , find a partition of  $k$  clusters that optimizes the chosen partitioning criterion
  - Global optimal: exhaustively enumerate all partitions
  - Heuristic methods:  $k$ -means and  $k$ -medoids algorithms
  - **$k$ -means** (MacQueen'67): Each cluster is represented by the center of the cluster
  - **$k$ -medoids** or **PAM (Partition Around Medoids)** (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

32

## **$k$ -Means**

- Initial set of clusters randomly chosen.
- Iteratively, items are moved among sets of clusters until the desired set is reached.
- High degree of similarity among elements in a cluster is obtained.
- Given a cluster  $K_i = \{t_{i1}, t_{i2}, \dots, t_{im}\}$ , the **cluster mean** is  $m_i = (1/m)(t_{i1} + \dots + t_{im})$

33

## Partitioning Clustering

Distance metric

↓  
Similarity / Closeness / nearest.

K-mean

center

Centroid

cluster value  
(avg)  $\sum_{i=1}^n z_{di}$

K-medoid

center

medoidcenter

physically  
centrally  
located  
point.

K-Means

Converging  
till no change in  
centroid

### K-means

- ① Randomly choose k-points  
as k-centres  
where k given by user
- ② Merge/assign new points closer  
to k
- ③ Recompute Centroid
- ④ Repeat step 2 & 3 till convergence

## **k-Means Example**

- Given: {2,4,10,12,3,20,30,11,25}, k=2
- Randomly assign means:  $m_1=2, m_2=4$ 
  - $K_1=\{2,3\}, K_2=\{4,10,12,20,30,11,25\}, m_1=2.5, m_2=16$
  - $K_1=\{2,3,4\}, K_2=\{10,12,20,30,11,25\}, m_1=3, m_2=18$
  - $K_1=\{2,3,4,10\}, K_2=\{12,20,30,11,25\}, m_1=4.75, m_2=19.6$
  - $K_1=\{2,3,4,10,11,12\}, K_2=\{20,30,25\}, m_1=7, m_2=25$
- Stop as the clusters when these means are the same.

$\text{Centroid} = \frac{1}{n} \sum_{i=1}^n x_i$

K-means Example. k=Input

$D = \{2, 4, 10, 12, 3, 20, 30, 11, 25\}$   $k=2$

Step ① : Choose randomly k-clusters from D

	①	②		
Centroid	$m_1$	$m_2$	$(m_1 + m_2)/2$	$m_1 = 2$
	2	4	3	$m_2 = 4$
2	0	2	1	$k_1 = \{2, 3\}$
4	2	0	1.5	$k_2 = \{4\}$
10	8	2	5	$m_1 = 2.5$
12	10	8	9	$k_1 = \{4, 10, 12, 20, 30, 11, 25\}$
3	1	1	1.5	$m_2 = \frac{4+10+12+20+30+11+25}{7}$
20	18	10	14	$= 16$
30	28	26	27	$k_1 = \{2, 4, 3\}, G=3$
11	5	7	6	$k_2 = \{10, 12, 20, 30, 11, 25\}$
25	23	21	22.5	$G=18$

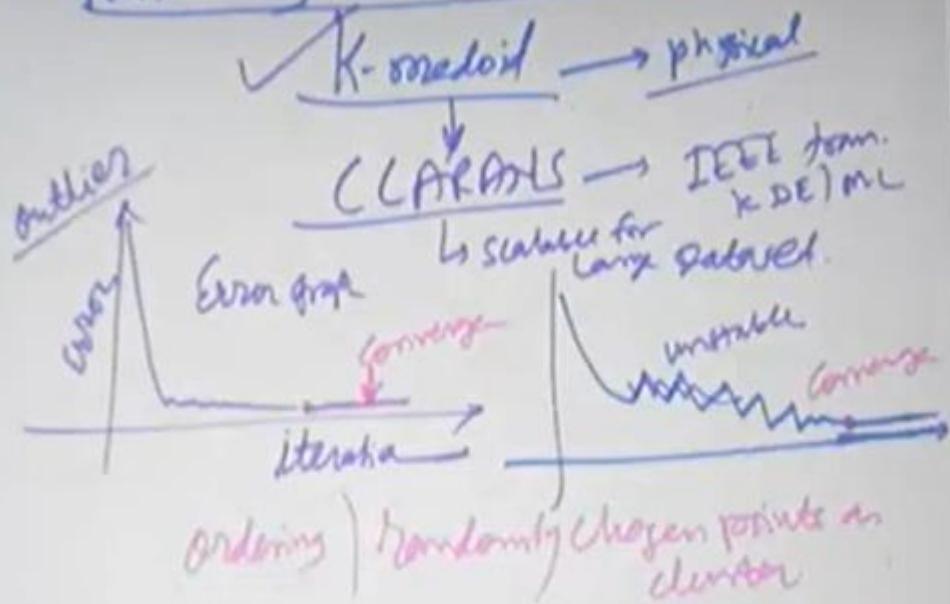
Repeat till Convergence  
till no change in centroid.

## **Partitioning Around Medoids (PAM) ( $k$ -Medoids)**

- Handles outliers well.
- Ordering of input does not impact results.
- Does not scale well.
- Each cluster represented by one item, called the **medoid**.
- Initial set of  $k$  medoids randomly chosen.

# PAM

## Partitioning Around medoid

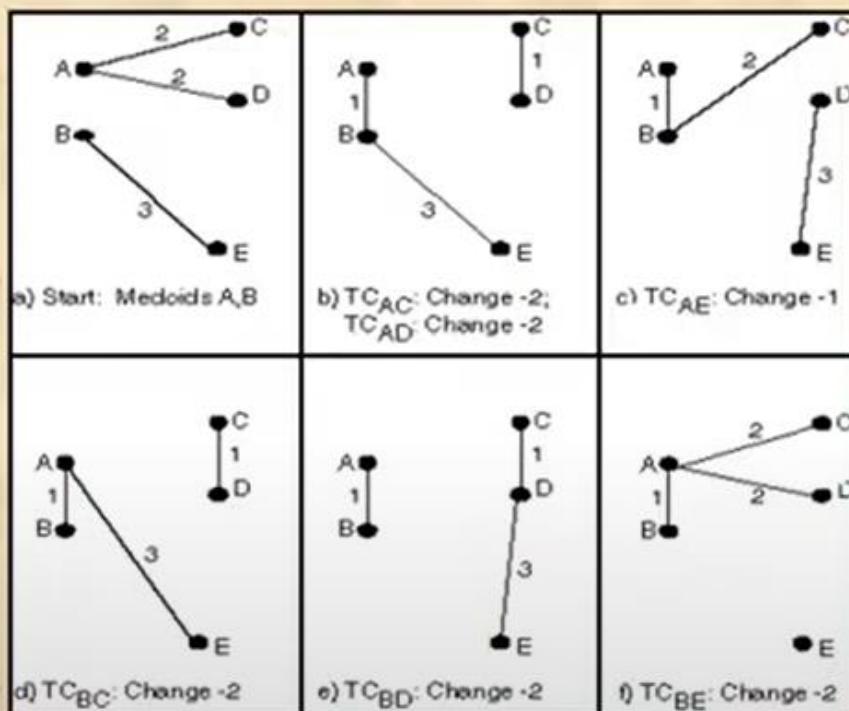


# PAM Algorithm

- Use real object to represent the cluster
  1. Select  $k$  representative objects arbitrarily
  2. For each pair of non-selected object  $h$  and selected object  $i$ , calculate the total swapping cost  $TC_{ih}$
  3. For each pair of  $i$  and  $h$ ,
    - a) If  $TC_{ih} < 0$ ,  $i$  is replaced by  $h$
    - b) Then assign each non-selected object to the most similar representative object
  4. repeat steps 2-3 until there is no change

37

## PAM



38

## PAM Cost Calculation

- At each step in algorithm, medoids are changed if the overall cost is improved.
- $C_{jih}$  – cost change for an item  $t_j$  associated with swapping medoid  $t_i$  with non-medoid  $t_h$ .

1.  $t_j \in K_i$ , but  $\exists$  another medoid  $t_m$  where  $dis(t_j, t_m) \leq dis(t_j, t_h)$
2.  $t_j \in K_i$ , but  $dis(t_j, t_h) \leq dis(t_j, t_m) \forall$  other medoids  $t_m$ ;
3.  $t_j \in K_m$ ,  $\notin K_i$ , and  $dis(t_j, t_m) \leq dis(t_j, t_h)$ ; and
4.  $t_j \in K_m$ ,  $\notin K_i$ , but  $dis(t_j, t_h) \leq dis(t_j, t_m)$ .

39

## PAM Algorithm Listing

### Input:

```
 $D = \{t_1, t_2, \dots, t_n\}$  // Set of elements  
 $A$  // Adjacency matrix showing distance between elements.  
 $k$  // Number of desired clusters.
```

### Output:

```
 $K$  // Set of clusters.
```

### PAM Algorithm:

```
arbitrarily select  $k$  medoids from  $D$ ;
```

```
repeat
```

```
    for each  $t_h$  not a medoid do
```

```
        for each medoid  $t_i$  do
```

```
            calculate  $TC_{ih}$ ;
```

```
        find  $i, h$  where  $TC_{ih}$  is the smallest;
```

```
        if  $TC_{ih} < 0$  then
```

```
            replace medoid  $t_i$  with  $t_h$ ;
```

```
until  $TC_{ih} \geq 0$ ;
```

```
for each  $t_i \in D$  do
```

```
    assign  $t_i$  to  $K_j$  where  $dis(t_i, t_j)$  is the smallest over all medoids;
```

40

Swapping cost = Sum of distance of each  
non-selected point to  
selected (Medoid)

Example

$$D = \{2, 4, 10, 12, 3, 20, 30, 11, 25\} = \sum \text{dist.}$$

Medoid:

$$k_1 = 2$$

$$k_2 = 4$$

Choose unselected point 3

Pair with unselected  $\{2, 3\}$

pair with remaining  $\{8, 4\} (2, 10, 12, 25)$

$$nCr = \frac{n!}{r!(n-r)!}$$

It continues

Compute cost of new pair & old pairs.

Keep pair having smallest cost | discard else.

Repeat above swapping till no change  
in Medoid



## Clustering Large Databases

- Most clustering algorithms assume a large data structure which is memory resident.
- Clustering may be performed first on a sample of the database then applied to the entire database.
- Algorithms
  - **BIRCH**
  - **DBSCAN**
  - CURE

## **Desired Features for Large Databases**

- One scan (or less) of DB
- Online
- Suspendable, stoppable, resumable
- Incremental
- Work with limited main memory
- Different techniques to scan (e.g. sampling)
- Process each tuple once

## **BIRCH : Balanced Iterative Reducing and Clustering using Hierarchies**

- Incremental, hierarchical, one scan
- Save clustering information in a tree
- Each entry in the tree contains information about one cluster
- New nodes inserted in closest entry in tree

# Clustering Feature

- CT Triple:  $(N, \overrightarrow{LS}, SS)$ 
  - N: Number of points in cluster
  - $\overrightarrow{LS}$ : Sum of points in the cluster
  - SS: Sum of squares of points in the cluster
- CF Tree
  - Balanced search tree
  - Node has CF triple for each child
  - Leaf node represents cluster and has CF value for each subcluster in it.
  - Subcluster has maximum diameter

## Clustering Large Databases

BIRCH

DBSCAN

BIRCH

Incremental & hierarchical / Single scan

→ Representing TREE → CF-tree

↓  
Tree structure

Clustering Feature (CF): triplet

$$CF = (N, LS, SS)$$

Where N: No. of data points in D.

LS: Linear sum of points

SS: Sum of Squared points

CF-tree consists nodes represented by CF

Statistical form CF

$$CF = (N, LS, SS)$$

Let D consist of N points

$$LS = \sum_{i=1}^N x_i \quad SS = \sum_{i=1}^N (x_i)^2$$

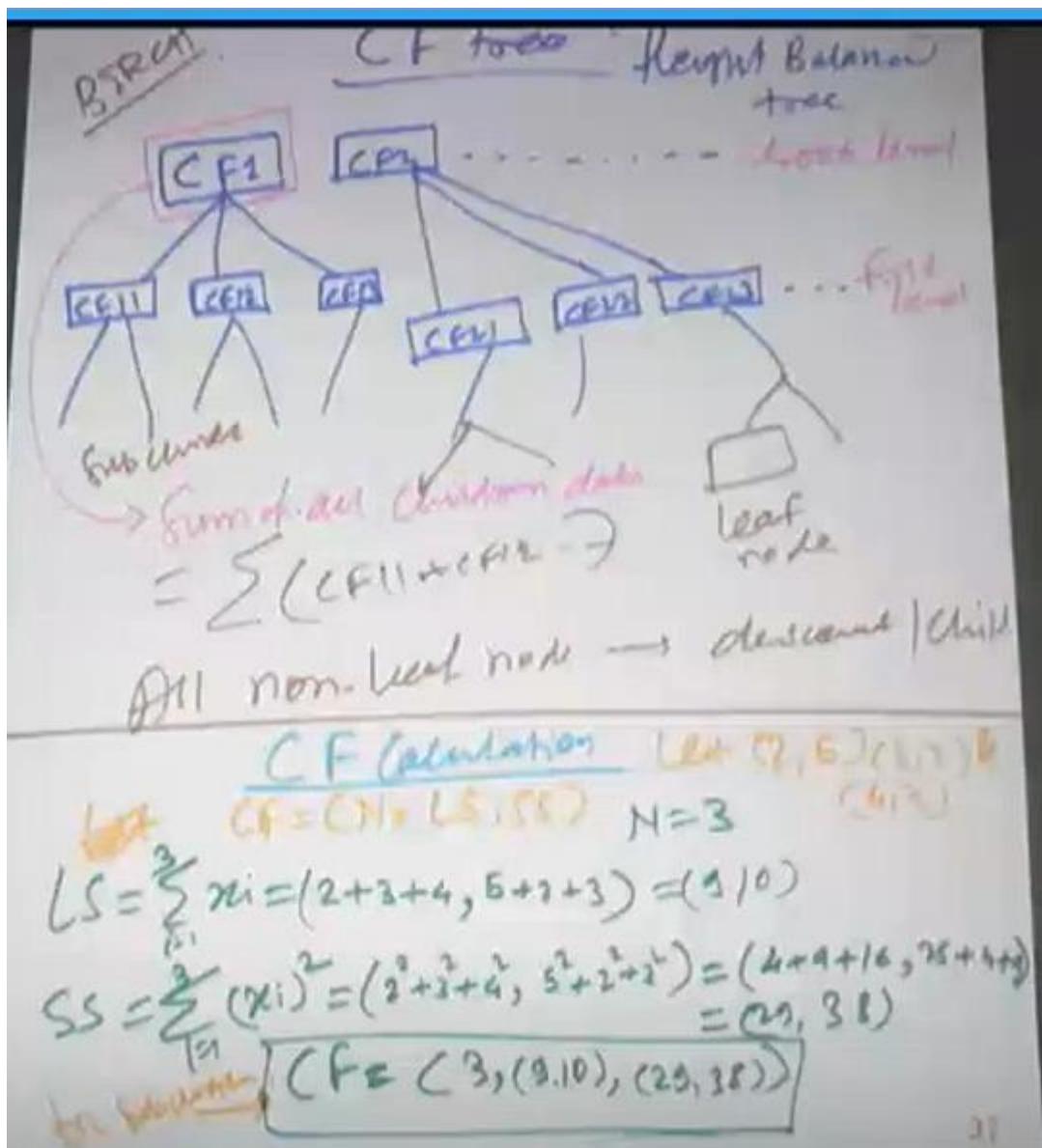
$$\text{Cluster Centroid } x_0 = \frac{\sum x_i}{n} = \frac{LS}{n}$$

$$\text{Radius} = \sqrt{\frac{\sum_{i=1}^N (x_i - x_0)^2}{n}} = \sqrt{\frac{nSS - 2LS^2 + nLS}{n^2}}$$

$$\text{Diameter } D = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (x_i - x_j)^2}{n(n-1)}} = \sqrt{\frac{2nSS - 2LS^2}{n(n-1)}}$$

R: Average distance from member objects to centroid.

D: Average pairwise distance within cluster



## BIRCH Algorithm $\rightarrow$ BT tree

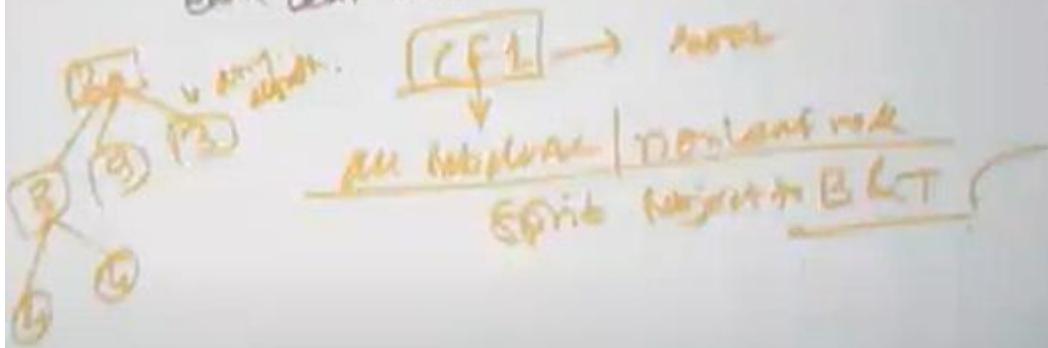
Two parameters: i) Branching factor  $B$   
ii) Threshold  $T$

$B$  = max. no. of children per non-leaf node

$T$  = max. diameter of subcluster stored  
at leaf node of tree

Two phases CF tree BT tree

- ① Scan the database to build initial tree.
- ② Apply any other second algorithm to cluster each leaf node.



# BIRCH Algorithm

**Input:**

$D = \{t_1, t_2, \dots, t_n\}$  // Set of elements  
 $T$  // Threshold for CF tree construction.

**Output:**

$K$  // Set of clusters.

**BIRCH Clustering Algorithm:**

```
for each  $t_i \in D$  do
    determine correct leaf node for  $t_i$  insertion;
    if threshold condition is not violated then
        add  $t_i$  to cluster and update CF triples;
    else
        if room to insert  $t_i$  then
            insert  $t_i$  as single cluster and update CF triples;
        else
```

add notes

▲ NOTES □ COMMENTS

# Improve Clusters

1. Create initial CF tree using Algorithm. If there is insufficient memory to construct the CF tree with a given threshold, the threshold value is increased, and a new smaller CF tree is constructed.
2. Apply another global clustering approach applied to the leaf nodes in the CF tree. Here each leaf node is treated as a single point for clustering.
3. The last phase (which is optional) reclusters all points by placing them in the cluster which has the closest centroid.

4 notes

# A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise

## DBSCAN

Martin Ester, Hans-Peter Kriegel, Jiirg Sander, Xiaowei Xu

Institute for Computer Science, University of Munich  
Oettingenstr. 67, D-80538 Miinchen, Germany  
{ester | kriegel | sander | xwxu } @informatik.uni-muenchen

KDD-96 Proceedings. pp 226-231  
Copyright © 1996, AAAI ([www.aaai.org](http://www.aaai.org)).

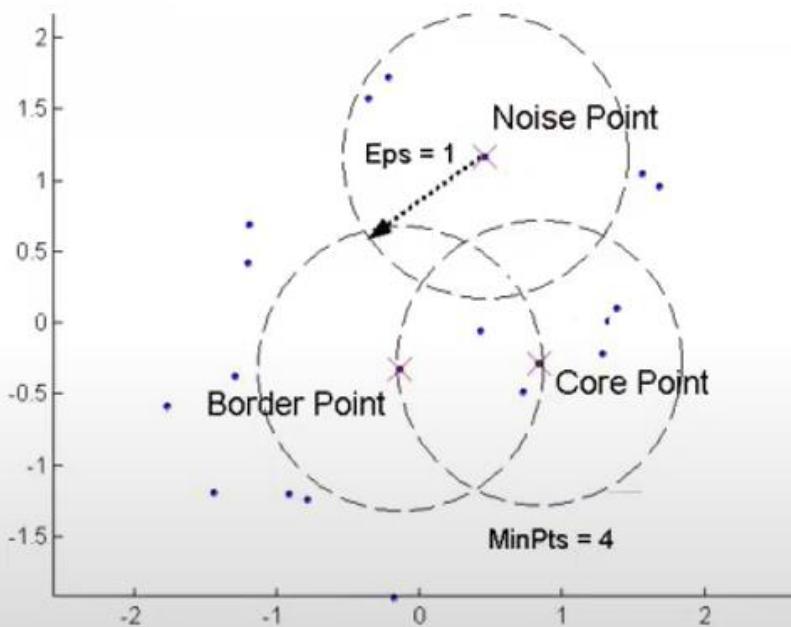
## DBscan Algorithm

Density-Based Spatial Clustering of Applications with Noise

DBSCAN is a density-based algorithm.

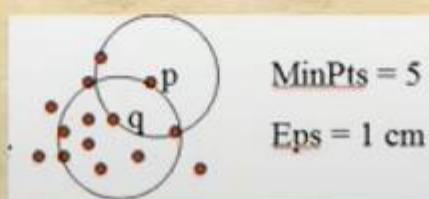
- **Eps** : Maximum radius of the neighbourhood
- **MinPts** : Minimum number of points in an Eps-neighbourhood of that point
- Density = number of points within a specified radius (**Eps**)
- A point is a *core point* if it has more than a specified number of points (**MinPts**) within **Eps**
  - These are points that are at the interior of a cluster
- A *border point* has fewer than **MinPts** within **Eps**, but is in the neighborhood of a core point
- A *noise point* is any point that is not a *core point* or a *border point*.

## DBSCAN: *Core*, *Border*, and *Noise Points*



## DBscan : other parameters

- **Directly density-reachable**: A point  $p$  is directly density-reachable from a point  $q$  w.r.t.  $Eps$ ,  $MinPts$  if
  - $p$  belongs to  $N_{Eps}(q)$
  - $N_{Eps}(q) : \{ p \text{ belongs to } D \mid dist(p,q) \leq Eps \}$
  - core point condition:  
 $|N_{Eps}(q)| \geq MinPts$

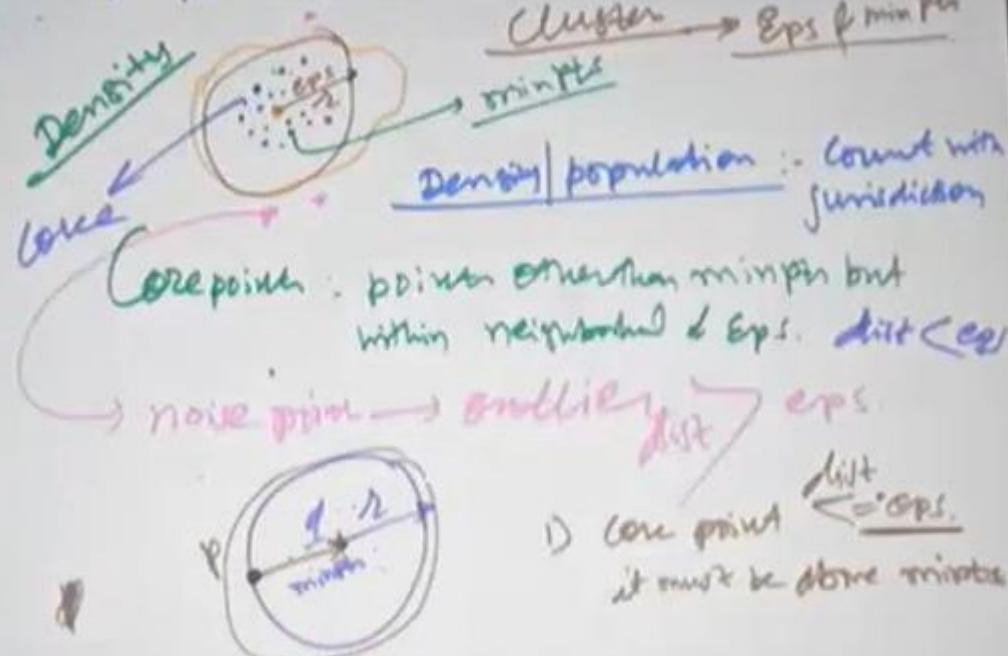


## DBSCAN

### Density Based Algorithm

Eps: neighborhood radius

minpt: minimum no. of points around Eps



35

## DBscan Algorithm : Generic Steps

- Arbitrary select a point  $p$
- Retrieve all points **density-reachable** from  $p$  w.r.t. **Eps** and **MinPts**.
- If  $p$  is a **core point**, a cluster is formed.
- If  $p$  is a **border point**, no points are density-reachable from  $p$  and DBSCAN visits the next point of the database.
- Continue the process until all of the points have been processed.

# DBSCAN Algorithm

**Algorithm:** DBSCAN: a density-based clustering algorithm.

**Input:**

$D$ : a data set containing  $n$  objects,  
 $\epsilon$ : the radius parameter, and  
 $MinPts$ : the neighborhood density threshold.

**Output:** A set of density-based clusters.

Dr. Bashirahamad F. Momin

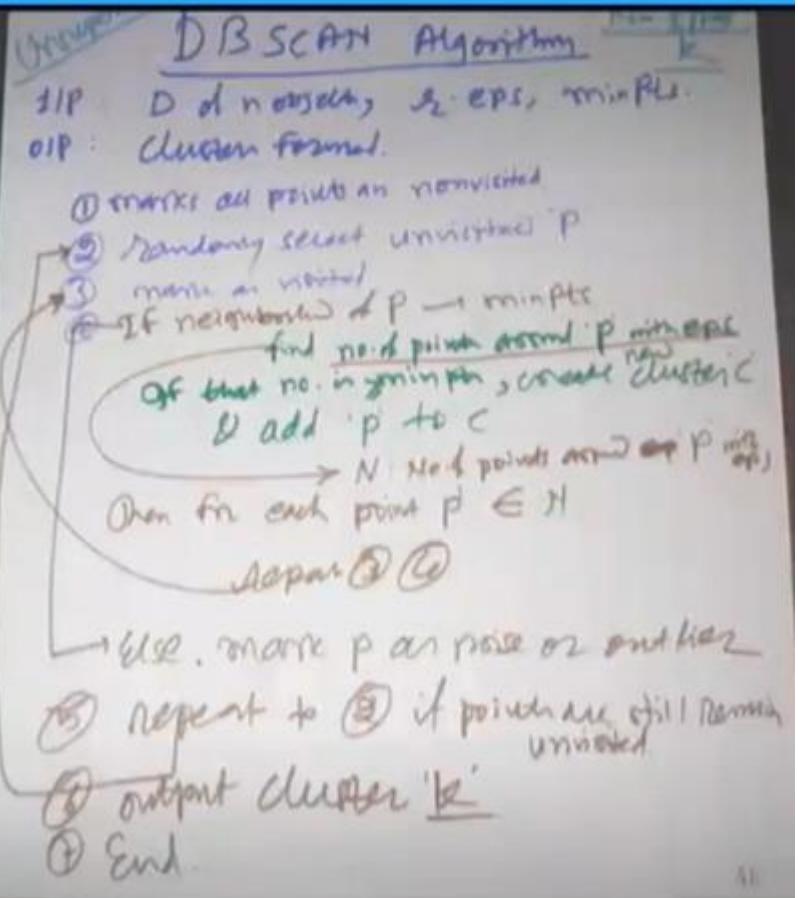
## Cont..

**Method:**

- (1) mark all objects as unvisited;
- (2) **do**
- (3) randomly select an unvisited object  $p$ ;
- (4) mark  $p$  as visited;
- (5) **if** the  $\epsilon$ -neighborhood of  $p$  has at least  $MinPts$  objects
- (6) create a new cluster  $C$ , and add  $p$  to  $C$ ;
- (7) let  $N$  be the set of objects in the  $\epsilon$ -neighborhood of  $p$ ;
- (8) **for** each point  $p_0$  in  $N$
- (9) if  $p_0$  is unvisited
- (10) mark  $p_0$  as visited;
- (11) if the  $\epsilon$ -neighborhood of  $p_0$  has at least  $MinPts$  points,  
add those points to  $N$ ;
- (12) if  $p_0$  is not yet a member of any cluster, add  $p_0$  to  $C$ ;
- (13) **end for**
- (14) output  $C$ ;
- (15) **else** mark  $p$  as noise;
- (16) **until** no object is unvisited;



Dr. Bashirahamad F. Momin



How to select

Selection of  $\epsilon$  &  $\text{minPts}$  → affects DBSCAN

tips

$\text{minPts} = \text{Dimensional Data}$

e.g. for IRIS: 4  
 $\therefore \text{minPts} = 4 + 1 = 5$

with  $\epsilon$ : maximum distance among pairs of data points.

Whole DBSCAN cluster depends on selection (optimal) of  $\epsilon$  &  $\text{minPts}$ .

## Lecture 27

## Evaluation of clustering

To determine how clusters form  $\downarrow$   
Assessment

- ① Assessing cluster tendency
  - ② Determine the no. of good optimal clusters  $K$
  - ③ measuring the cluster quality

## ① Clustering Tendencies

for given Data set  $\Rightarrow$  formulates  
↳ eligible to form clusters  
↳ characters ..

For Good clustering  $\Rightarrow$  Dataset  $\Rightarrow$  non-uniform



AF II

At your disposal  
John Clinton

Conington

1.6 non-uniform + good class

### Uniform pressure



As per Agreed  
it from Clinton

it from direct  
written meanin

If non-uniform  $\rightarrow$  Good chm

To determine the tendency of clustering

### Hopkin's Statistics

Test non-uniformity of dataset

- ① Let  $n$  points  $D = \{P_1, \dots, P_n\}$  for each point  $P_i$   
find the nearest neighbor of  $P_i \rightarrow x_i$

$$x_i = \min_{v \in D} \{ \text{dist}(P_i, v) \}$$

- ② Sample  $n$  points  $\{q_1, \dots, q_n\} \subset D$  find nearest  
neighbor of  $q_i$  from  $D - \{q_i\}$

$$y_i = \min_{\substack{v \in D, \\ v \neq q_i}} \{ \text{dist}(q_i, v) \}$$

- ③ Hopkins statistic  $H$

$$H = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i + \sum_{i=1}^n y_i}$$

$H=0.5$  for uniform dataset and

H=0 for non-uniform dataset

Determining No of Clusters (K)  $k > 0$   
optimal  
Good Balance b/w  
Compressibility & granularity

Depends on :

- ① Distribution Shape
- ② Scale in Data set
- ③ Cluster Resolution.

### ④ Empirical method :-

$$\text{No. of clusters} = \sqrt{\frac{n}{2}} \quad \text{where } n \text{ is no. of data points}$$

e.g. IRIS 150 Samples.

$$k = \sqrt{\frac{150}{2}} = \sqrt{75} = 8.6 \approx 9$$

Physical  $\rightarrow$  There 3 classes

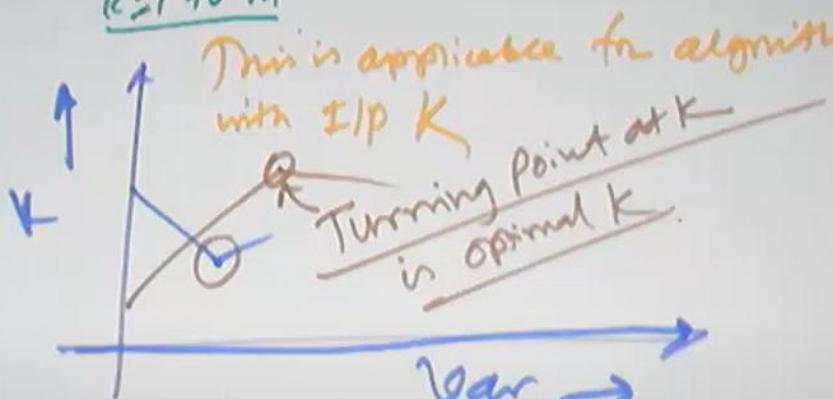
## ⑪ Elbow method

Increasing no. of cluster  $\rightarrow$  Reduce the sum of "within-cluster" variance.

"within cluster" - Variance

Find the sum of within cluster variance  $\forall k$   $\text{Var}(k)$

$k=1 \text{ to } m$



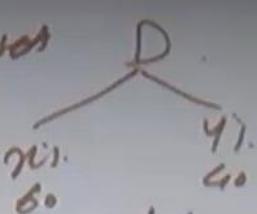
- ① Compute variance for each cluster.
- ② plot  $\text{Var} \rightarrow k \rightarrow$  Check turning point in graph  $\rightarrow$  Turning Point @  $k$  will be min



### IV on-Validation / m-partitioning

Given Dataset

①



① Build clusters with  $X\%$  data set.

② Use  $Y\%$  data sets for cluster evaluation.

for each data point  $p \in Y\%$ ,  
 $\sum$  compute the dist to centroid.  
 $\sum$  sum of dist

Do this for all clusters with smallest sum will be optimal.

③ Repeat ① with different  $X\%$  &  $Y\%$ .

avg. Average out Results with all sets.

### III Measuring Clustering Quality

Clustering & validation for Very large Database (VLDB) 2006 IEEE proceeding

all/major/popular

cluster validation techniques

Sir has published an IEEE paper for the above mentioned topic.

### III. Measuring clustering Quality

- Clustering & validation for Very Large Database (VLDB)  
2006 IEEE proceeding  
↓  
all/major/popular  
Cluster validation techniques

dist high Inter Cluster Similarity → low  
dist < Emax intra cluster Similarity → high  
Inter cluster dist → high  
Intra cluster dist → small

Sonelors Distance → most similar  
NN → Sonelors distance  
→ Silhouettes Statistics

## Lecture 29

The screenshot shows a video player interface with a presentation slide. The slide has a teal header with white text: "PE3:DATA MINING", "Association Rule", "Mining", and "Module - V". Below the header is a large blue rectangular area containing the text "Decision Rules vs Association Rules". The slide is divided into two sections: "Decision Rules :" and "Association Rules :". Under "Decision Rules :", it says "IF < > part contain all conditional attributes" and "THEN < > part contain only Decision attribute". Under "Association Rules :", it says "IF < > part contain any attributes" and "THEN < > part contain any attributes". The video player interface includes a progress bar at 4:26 / 57:22, a list of participants on the right, and standard video controls.

DM #29 Module V Association Rule Mining Introduction 20211020 111834 ...

SD  
Smit Deshmukh

R  
rohit.pal

NM  
Nida Mujawar

VK  
Veerja Kadam

BM  
Bashirahamed M...

Dr. Bashirahamed F. Momin  
CSE Dept., Walchand COE, Sangli.

CC HD

Decision Rules vs Association Rules

**Decision Rules :**

**IF** < > part contain all conditional attributes

**THEN** < > part contain only Decision attribute

**Association Rules :**

**IF** < > part contain **any** attributes

**THEN** < > part contain **any** attributes

2

## Association Rules

Antecedent  
IF  $\in D$

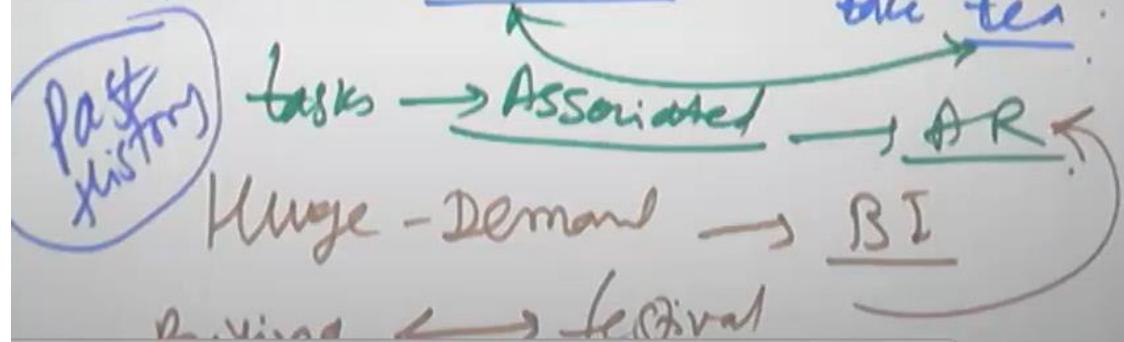
consequent  
THEN  
 $\in D$ .

e.g.  $\{sub1, sub2, sub3 \dots\} \quad \wedge$

IF  $sub1 >= 45$  and  $sub2 >= 60$  Then  
 $sub4 >= 55$

## Market Basket Data

of person take Break fast Then he/she also take tea:



## Definitions of various terms

- **Set of items:**  $I = \{I_1, I_2, \dots, I_m\}$
- **Transactions:**  $D = \{t_1, t_2, \dots, t_n\}, t_j \subseteq I$
- **Itemset:**  $\{I_{i1}, I_{i2}, \dots, I_{ik}\} \subseteq I$
- **Support of an itemset:** Percentage of transactions which contain that itemset.
- **Large (Frequent) itemset:** Itemset whose number of occurrences is above a threshold.

## Itemset

For 'n' items in data set,  
there will be  
 **$2^n - 1$**  subset

An **itemset** is any subset of the set of all items

## Association Rules Example

Transaction	Items
$t_1$	Bread, Jelly, PeanutButter
$t_2$	Bread, PeanutButter
$t_3$	Bread, Milk, PeanutButter
$t_4$	Beer, Bread
$t_5$	Beer, Milk

$$I = \{ \text{Beer, Bread, Jelly, Milk, PeanutButter} \}$$

Support of {Bread, PeanutButter} is 60%

# Example: Market Basket Data

- Items frequently purchased together:  
*Bread  $\Rightarrow$  PeanutButter*
- Uses:
  - Placement
  - Advertising
  - Sales
  - Coupons
- Objective: increase sales and reduce costs

5

The screenshot shows a Microsoft Teams meeting interface. On the left, a slide titled "Association Rules Example" displays a transaction table:

Transaction	Items
$t_1$	Bread, Jelly, PeanutButter
$t_2$	Bread, PeanutButter
$t_3$	Bread, Milk, PeanutButter
$t_4$	Beer, Bread
$t_5$	Beer, Milk

Below the table, the text  $I = \{ \text{Beer, Bread, Jelly, Milk, PeanutButter} \}$  and "Support of {Bread, PeanutButter} is 60%" are shown.

On the right, a video feed shows a whiteboard with handwritten notes on association rules:

Basic Terminology

$D = \{\dots\}$        $I = \text{distinct/unique values in } D$        $|I| = 5$

$\text{Itemset} = \text{Subset of } I$        $n=5$

$\text{Itemset} = 2^5 - 1 = 32 - 1 = 31$

$\hookrightarrow \text{Subset of set of all itemsets}$

$I = \{\text{Beer, Bread, Jelly, Milk, PeanutButter}\}$

$\{\text{Beer, Bread}\} \subset I$        $\{\text{Beer, Milk}\} \subset I$        $\{\text{Beer, PeanutButter}\} \subset I$

$\{\text{Beer, Bread, Jelly}\} \subset I$        $\{\text{Beer, Bread, Milk}\} \subset I$        $\{\text{Beer, Bread, PeanutButter}\} \subset I$

$\{\text{Beer, Bread, Jelly, Milk}\} \subset I$        $\{\text{Beer, Bread, Jelly, PeanutButter}\} \subset I$

$\{\text{Beer, Bread, Milk, PeanutButter}\} \subset I$

$\{\text{Beer, Milk, PeanutButter}\} \subset I$

$\{\text{Beer, Bread, Milk, PeanutButter}\} \subset I$

$\{\text{Beer, Bread, Jelly, PeanutButter}\} \subset I$

$\{\text{Beer, Bread, Jelly, Milk}\} \subset I$

$\{\text{Beer, Bread, Jelly, Milk, PeanutButter}\} \subset I$

$\text{Itemset} \Rightarrow \text{Itemset}$

$\text{Itemset} \nRightarrow \text{Itemset}$

A list of participant icons is visible on the right side of the video feed.

DM #29 Module V Association Rule Mining Introduction 20211020 111834...

DESIGN TRANSITIONS ANIMATIONS SLIDE SHOW REVIEW VIEW

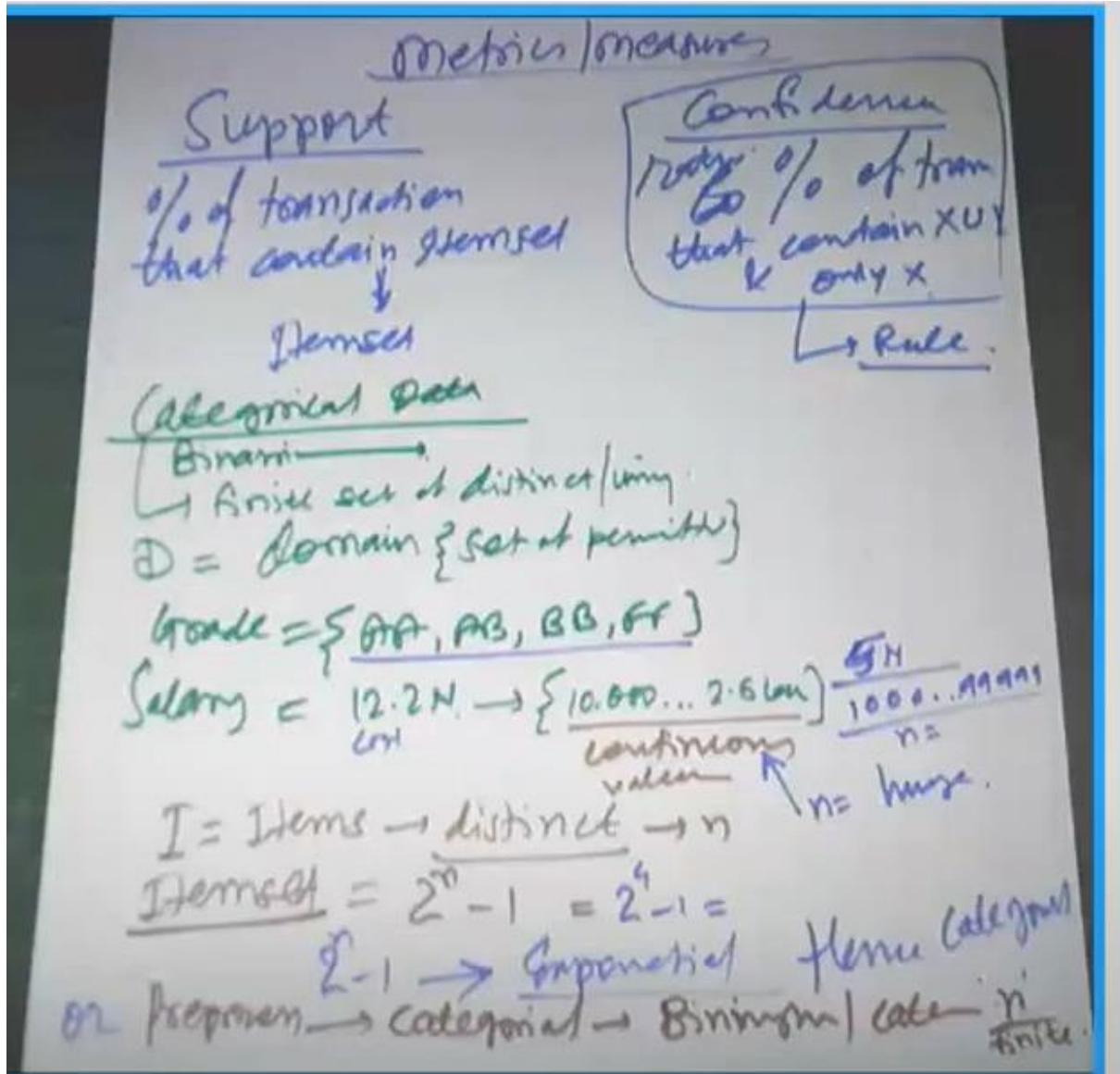
## Association Rule Definitions

- Association Rule (AR):** implication  $X \Rightarrow Y$  where  $X, Y \subseteq I$  and  $X \cap Y = \text{null}$ ;
- Support of AR (s)  $X \Rightarrow Y$ :** Percentage of transactions that contain  $X \cup Y$
- Confidence of AR ( $\alpha$ )  $X \Rightarrow Y$ :** Ratio of number of transactions that contain  $X \cup Y$  to the number that contain  $X$

Metrics / measures

Support  
% of transaction that contain itemset  
↳ Dense

Confidence  
Ratio % of transaction that contain  $X \cup Y$  & only  $X$   
↳ Rule.



Lecture 30

(2) PE3 DM Module-V : Assoc

teams.microsoft.com/\_#/pre-join-calling/19:dwx0xraVCYWOYzTnVx46eOCFBTwGN4Vy6baBcDn\_o1@thread.tacv2

DESIGN TRANSITIONS ANIMATIONS SLIDE SHOW REVIEW VIEW

## Association Rules Example

Transaction	Items
$t_1$	Bread, Jelly, PeanutButter
$t_2$	Bread, PeanutButter
$t_3$	Bread, Milk, PeanutButter
$t_4$	Beer, Bread
$t_5$	Beer, Milk

$I = \{ \text{Beer, Bread, Jelly, Milk, PeanutButter} \}$

Support of {Bread, PeanutButter} is 60%

Support of Demand  
 $\frac{\% \text{ of trans that } I \text{ is}}{\% \text{ of tuples that contain item } I} = \frac{\text{No. of tuples that contain item } I}{\text{Total No. of Transactions}}$

Support(Bread, PB) =  $\frac{3}{5} = 60\%$

Why support right? If in 5 trans there will be [2] -> 2/5 = 40% probable

Are all interesting? How to prove / remove

Should there be some threshold?

User specified threshold

Support =  $\frac{1}{5} = 20\%$

Support(Bread, PB) =  $\frac{1}{5} = 20\%$ . Below it goes low confidence

80% - 90% - 90%

The screenshot shows a Microsoft Teams meeting interface. On the left, a PowerPoint slide titled "Association Rules Example" displays a transaction table and a support statement. On the right, a whiteboard contains handwritten notes explaining the concept of association rules, including the formula  $R = \frac{\text{Support}}{\text{Confidence}}$ , an example rule  $X \rightarrow Y$  with support 60% and confidence 75%, and a note about user-supervised parameter tuning.

DESIGN TRANSITIONS ANIMATIONS SLIDE SHOW REVIEW VIEW

## Association Rules Example

Transaction	Items
$t_1$	Bread, Jelly, PeanutButter
$t_2$	Bread, PeanutButter
$t_3$	Bread, Milk, PeanutButter
$t_4$	Beer, Bread
$t_5$	Beer, Milk

$I = \{ \text{Beer, Bread, Jelly, Milk, PeanutButter} \}$

Support of {Bread, PeanutButter} is 60%

to add notes

SD Smit Deshmukh  
NM Nida Majeed  
VK Veerja Kadam  
S shivaniavasani  
GK Ganesh Kesar  
K krishnapo04  
T tulsi.galande  
AP Abhishek Patil  
M marwadeenreyaah  
SP Sutharsan Prasad...  
T Tanmay\_Padale  
P prachi.maniare  
SS Sayali.Sayali  
R rohit.pai  
S saileshakan  
S spagie4  
S seempanchal21  
+2

Sign in

Association Rule

$X \rightarrow Y$  antecedent consequent  $X, Y \in I$

$X \rightarrow Y$  is called common

$\text{IF } X \text{ THEN } Y$

Homogeneous item I say  $X, Y$  are true antecedent & consequent of rule R

Support

Confidence

Ratio

Ratio

$R = \frac{\text{Support}}{\text{Confidence}}$

$\text{Support} = 60\%$

$\text{Confidence} = \frac{3}{4} = 75\%$

User Supervised Parameter tuning with DM, S & R

**DM #30 Module V Apriori ARM algorithm 20211021 101734 Meeting Recording**

The screenshot shows a Microsoft Teams meeting interface. On the left, a PowerPoint slide titled "Association Rule Mining" is displayed, featuring a "Two Step Process". The right side shows a whiteboard with handwritten notes about the Apriori algorithm, including formulas for association rules and support/confidence calculations.

**Handwritten Notes on Whiteboard:**

- ASSOCIATIVE
- If  $A \Rightarrow B \wedge B \Rightarrow C$   
Then  $A \Rightarrow C$  OR  $B \Rightarrow C$
- OR  $A \Rightarrow B$  Then  $B \Rightarrow A$
- $X \Rightarrow Y \neq Y \Rightarrow X$
- Support =  $\frac{\text{Number of FB}}{\text{Total number of FB}}$  (60%)  
 $FB \Rightarrow FB \rightarrow 60\% \text{ (Confidence 100%)}$
- AR are not associative
- Support & confidence =  $|X \cup Y| / |U|$
- Support & confidence =  $|X \cup Y| / |U|$
- HR Mining
- Two steps
- I Generate all Large Itemset Frequent Itemset
- II Generate Rules from frequent itemsets

**Handwritten Notes on Whiteboard (Continued):**

- Large Itemset Property: If an itemset is not large, any subset of a large itemset is large.
- Contrapositive: If an itemset is not large, any subset of a large itemset is large.
- Interestingness: Itemsets → subset of itemset. Combination of itemset → K-itemset. Occurrences. Kth combination e.g. 3-itemset containing only one item.
- Pruning Rule: Needs some Pruning! Remove items that don't meet the support & confidence rule.
- Optimal Itemset generation → Rules: Apriori → Large Itemset frequent itemset → Occurrences.

**Participants (Visible on the right):**

- K (krishnapoud4)
- SD (Smit Deshmukh)
- GK (Ganesh Katar)
- ST (siddharth.tandale)
- SP (Sudhamshi.Pusad)
- S (sudhamshi.pusad21)
- P (prachi.waware)
- R (rishiika.karligade)
- S (saurabhirahude)
- R (rohit.pai)
- NM (Nida.Mujawar)
- BM (Bastinahamed.M...)
- K (krishnapoud4)
- SD (Smit Deshmukh)
- ST (siddharth.tandale)
- SP (Sudhamshi.Pusad)
- S (sudhamshi.pusad21)
- P (prachi.waware)
- R (rishiika.karligade)
- S (saurabhirahude)
- R (rohit.pai)
- NM (Nida.Mujawar)
- VK (Veenja.Kadam)
- S (shivani.avasthi)
- Y (yash.hoke)
- RS (Ruchikesh.Shelke)
- K (kaplesh.pansari)
- SS (Shreya.Singh)

**1<sup>st</sup> Step : Apriori – Large Itemsets generation**

- An itemset that contains  $k$  items is called as  **$k$ -itemset**.
- $k$ -itemset** that satisfies **minimum support**, is called as (**strong**) **frequent itemsets**.
- Large Itemset Property:**  
Any subset of a large itemset is large.
- Contrapositive:**  
*If an itemset is not large, none of its supersets are large.*

## Apriori Algorithm

- $C_1$  = Itemsets of size one in  $I$ ;
- Determine all large itemsets of size 1,  $L_1$ ;
- $i = 1$ ;
- Repeat
- $i = i + 1$ ;
- $C_i$  = Apriori-Gen( $L_{i-1}$ );
- Count  $C_i$  to determine  $L_i$ ;
- until no more large itemsets found;

## Apriori-Gen Example (cont'd)

$s=30\%$        $\alpha = 50\%$

Scan	Candidates	Large Itemsets
1	{Blouse}, {Jeans}, {Shoes}, {Shorts}, {Skirt}, {Tshirt}	{Jeans}, {Shoes}, {Shorts} {Skirt}, {Tshirt}
2	{Jeans, Shoes}, {Jeans, Shorts}, {Jeans, Skirt}, {Jeans, Tshirt}, {Shoes, Shorts}, {Shoes, Skirt}, {Shoes, Tshirt}, {Shorts, Skirt}, {Shorts, Tshirt}, {Skirt, Tshirt}	{Jeans, Shoes}, {Jeans, Shorts}, {Jeans, Skirt}, {Jeans, Tshirt}, {Shoes, Shorts}, {Shoes, Skirt}, {Shoes, Tshirt}, {Shorts, Tshirt}, {Skirt, Tshirt}
3	{Jeans, Shoes, Shorts}, {Jeans, Shoes, Tshirt}, {Jeans, Shorts, Tshirt}, {Jeans, Skirt, Tshirt}, {Shoes, Shorts, Tshirt}, {Shoes, Skirt, Tshirt}, {Shorts, Skirt, Tshirt}	{Jeans, Shoes, Shorts}, {Jeans, Shoes, Tshirt}, {Jeans, Skirt, Tshirt}, {Jeans, Tshirt}, {Shoes, Shorts}, {Shoes, Tshirt}
4	{Jeans, Shoes, Shorts, Tshirt}	{Jeans, Shoes, Shorts, Tshirt}
5		

## 2<sup>nd</sup> Step : Algorithm to Generate Association Rules

**Input:**

$D$  //Database of transactions  
 $I$  //Items  
 $L$  //Large itemsets  
 $s$  //Support  
 $\alpha$  //Confidence

**Output:**

$R$  //Association Rules satisfying  $s$  and  $\alpha$

**ARGen Algorithm:**

```
 $R = \emptyset;$ 
for each  $l \in L$  do
    for each  $x \subset l$  such that  $x \neq \emptyset$  and  $x \neq l$  do
        if  $\frac{\text{support}(l)}{\text{support}(x)} \geq \alpha$  then
             $R = R \cup \{x \Rightarrow (l - x)\};$ 
```

18

## Association Rule Generation Apriori Algorithm

If there are  $n$  items

Then maximum no. of rules generated

$$R = 3^n - 2^{n+1} + 1$$

Huge no. of rules generated, hence built-in rule pruning in Apriori

Generate the strong rules from strong frequent / large item set, satisfying minimum support and confidence

## Lecture 33:

DM #33 Module V Interesting patterns and its evaluation methods 20211110 111514 Meeting Recor...

The whiteboard contains the following handwritten notes:

- (S, (30.2, 18.2), (17, 18.2), 11.3, 13.1, 10.3, 11.6)
- Interest Pattern & Evaluation methods
- Support & Confidence → given
- Interestness (highlighted in red) with arrows pointing to:
  - Subjective (with arrows to person dependent and domain expert)
  - Objectivity (with arrows to Statistical theory, Statistical Description of Data, and Data mining framework)
- Intention Yes/No

Below the whiteboard is a video player interface showing the time 12:10 / 58:40.

DM #33 Module V Interesting patterns and its evaluation methods 20211110 111514 Meeting Recor...

The whiteboard contains the following handwritten notes:

- Association Analysis to Correlation Analysis
- Correlation measures
- A  $\leftrightarrow$  B [support, confidence, correlation]
- Interest or not.
- ① Lift: Simple Correlation measure
- Rule:  $A \leftrightarrow B$  where A, B are items
- Occurrence of A is Independent of B if  $P(A \cup B) = P(A) \cdot P(B)$
- otherwise A & B are dependent & strongly correlated events.
- $$\text{Lift}(A|B) = \frac{P(A \cup B)}{P(A) \cdot P(B)}$$
- $$\text{Lift}(A|B) = \frac{\text{Support}(A, B)}{\text{Supp}(A) \cdot \text{Supp}(B)}$$

Below the whiteboard is a video player interface showing the time 17:25 / 58:40.

DM #33 Module V Interesting patterns and its evaluation methods 20211110 111514 Meeting Recor...

Association Analysis to Correlation Analysis

Correlation measures

A  $\leftrightarrow$  B [support, confidence, correlation]

Interesting or not:

① Lift: Simple Correlation measure

Rule:  $A \leftrightarrow B$  where A, B are items

Outcome of A is Independent of B  
 $\text{if } P(A \cup B) = P(A) \cdot P(B)$

otherwise A & B are dependent & strongly correlated events.

$\text{Lift}(A, B) = \frac{P(A \cup B)}{P(A) \cdot P(B)}$

$\text{Lift}(A, B) = \frac{\text{Support}(A, B)}{\text{Supp}(A) \cdot \text{Supp}(B)}$

$\begin{cases} < 1 & \text{-ve correlation} \\ = 1 & \text{no correlation} \\ > 1 & \text{+ve correlation} \end{cases}$

$\begin{cases} < 1 & \text{-ve correlation} \\ = 1 & \text{no correlation} \\ > 1 & \text{+ve correlation} \end{cases}$

$= \frac{\text{Conf}(A \rightarrow B)}{\text{Supp}(B)}$

SD Smit Deshmukh  
K krishnapo04  
S shivansavare  
A2 sambhavita21  
VK Veerja Kadam  
A abhishek.pal  
AP Abhishek.Pal  
T tulsi.gandhe  
R revatijathav2210...  
NM Nida Majeed  
S spagireel  
BM Bashirahamed M...

DM #33 Module V Interesting patterns and its evaluation methods 20211110 111514 Meeting Recor...

Interesting or not:

① Lift: Simple Correlation measure

Rule:  $A \leftrightarrow B$  where A, B are items

Outcome of A is Independent of B  
 $\text{if } P(A \cup B) = P(A) \cdot P(B)$

otherwise A & B are dependent & strongly correlated events.

$\text{Lift}(A, B) = \frac{P(A \cup B)}{P(A) \cdot P(B)}$

$\text{Lift}(A, B) = \frac{\text{Support}(A, B)}{\text{Supp}(A) \cdot \text{Supp}(B)}$

$\begin{cases} < 1 & \text{-ve correlation} \\ = 1 & \text{no correlation} \\ > 1 & \text{+ve correlation} \end{cases}$

$\begin{cases} < 1 & \text{-ve correlation} \\ = 1 & \text{no correlation} \\ > 1 & \text{+ve correlation} \end{cases}$

$= \frac{\text{Conf}(A \rightarrow B)}{\text{Supp}(B)}$

SD Smit Deshmukh  
K krishnapo04  
S shivansavare  
A2 sambhavita21  
VK Veerja Kadam  
A abhishek.pal  
AP Abhishek.Pal  
T tulsi.gandhe  
R revatijathav2210...  
NM Nida Majeed  
S spagireel  
BM Bashirahamed M...

DM #33 Module V Interesting patterns and its evaluation methods 20211110 111514 Meeting Recor...

Sample : T : 10,000

6000 Buy Computer Games  
7500 Buy Videos  
4000 Buy both Video & Computer Games

$S: 401. d = 68\%$

$\text{Buy}(X; \text{Game}) = \text{Buy}(X; \text{Video})$

Build Contingency table :

	game	not game	
Video	4010	3580	6590
Not video	2000	500	2500
	6000	4080	10,000

$\text{lift} = \frac{P(\{\text{game}, \text{video}\})}{P(\{\text{game}\}) \cdot P(\{\text{video}\})}$

$P(\{\text{gm}\}) = \frac{6000}{10,000} = 0.6 \quad \text{prob(Bom)} = \frac{4010}{10,000} = 0.4$

$P(\{\text{vid}\}) = \frac{3580}{10,000} = 0.35$

$\text{lift} = \frac{0.4}{0.6 \times 0.35} = 0.89 < 1$

(lift identity)  
up to conf.

L +ve cor → junkets

29:47 / 58:40

Microsoft Teams

Untitled - File - Edit

Camera

③ Chi Square Test  $\chi^2$  measure

$$\chi^2 = \sum \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

Where Expected Count  $E_{ij} = \frac{\text{Count}(A=i) \cdot \text{Count}(B=j)}{n}$

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

Expected (game, video) =  $\frac{\text{Count(game)} \cdot \text{Count(video)}}{n}$

$$= \frac{6000 \cdot 3580}{10,000} = 4570$$

$$\chi^2 = \frac{(4010 - 4570)^2}{4570} + \frac{(3580 - 3570)^2}{3570} + \frac{(2000 - 1500)^2}{1500} + \frac{(500 - 500)^2}{500} = 555.6$$

as  $\chi^2 > 1$ . It is negatively correlated

EPSON ELPDC11

meetingAttendance.csv 27901542

Share invite

In this meeting (14) Mute all

hirahamad Momin (Owner)

Nishik Pol (Guest)

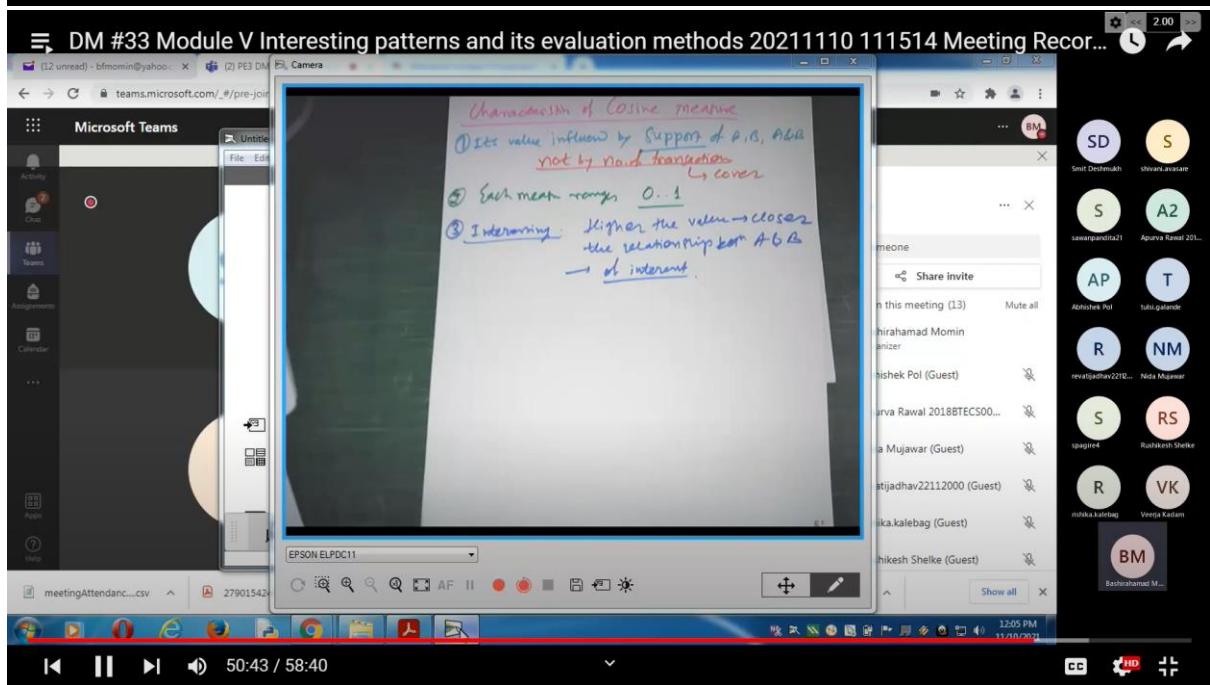
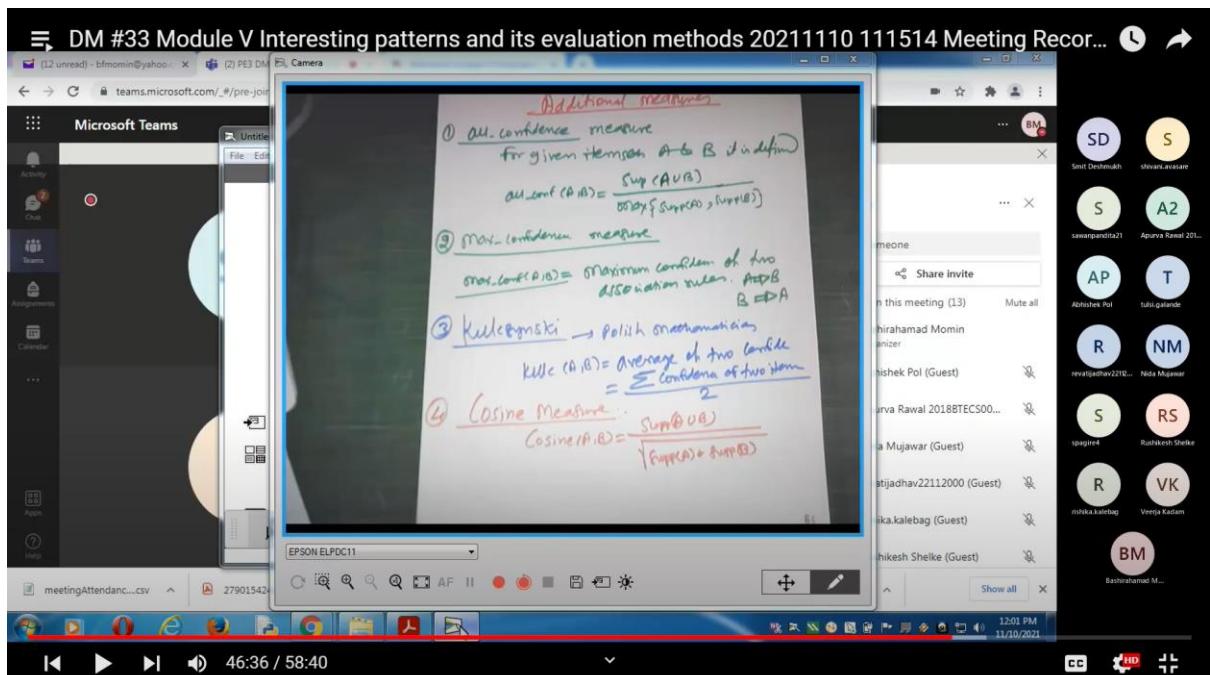
Arpura Rawal 2018BTTECS00...

revatjadav22112000 (Guest)

Nida Mujawar (Guest)

isha.kalebag (Guest)

11:56 AM 11/10/2021



## Web Data Mining

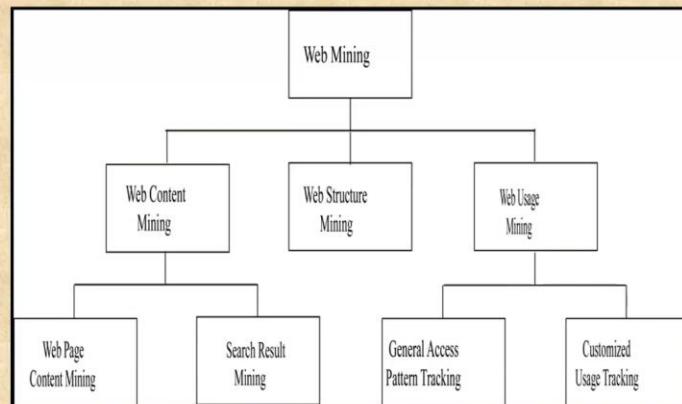
**Use of data mining techniques to automatically discover interesting and potentially useful information from Web documents and services.**

Web mining may be divided into three categories:

1. Web content mining
2. Web structure mining
3. Web usage mining



## Web Mining Taxonomy



## Web Data

- Web pages
- Intra-page structures
- Inter-page structures
- Usage data
- Supplemental data
  - Profiles
  - Registration information
  - Cookies



5

## Web details

- More than 20 billion pages in 2008
- Many more documents in databases accessible from the Web
- More than 4m servers
- A total of perhaps 100 terabytes
- More than a million pages are added daily
- Several hundred gigabytes change every month
- Hyperlinks for navigation, endorsement, citation, criticism or plain whim

**Web data is *BIG Data***



6

## Web Index Size in Pages

Total Web is estimated to be about 56 B pages.

Search Engine	Number of Pages in Billions
Google	16
MSN Search	7
Yahoo	50
Ask	4.2

Source

<http://www.worldwidewebsize.com/>



## Web Representation by Graph

- Web is a graph – vertices and edges (V,E)
- Directed graph – directed edges (p,q)
- Undirected graph - undirected edges (p,q)
- Strongly connected component - a set of nodes such that for any (u,v) there is a path from u to v



8

## Graph terminology

- **Breadth first search** - layer 1 consists of all nodes that are pointed by the root, layer k consists of all nodes that are pointed by nodes on level k-1
- **Diameter of a graph** - maximum over all ordered pairs  $(u,v)$  of the shortest path from u to v

9



## Web size

- **In-degree** is number of links to a node
- **Out-degree** is the number of links from a node
- **Fraction of pages** with i in-links is proportional to  $1/i^{2.1}$
- With i out-links, it is  $1/i^{2.72}$

10

i	In-links	Out-links
2	23%	15%
3	10%	5%
4	5%	2%
5	3%	1%



## Challenges

- There are several major challenges for Web mining research:
  - Most Web documents are in HTML format and contain many **markup** tags, mainly used for formatting.
  - While traditional IR systems often contain structured and well-written documents, this is **NOT** the case on the Web.
  - While most documents in traditional IR systems tend to remain static over time, Web pages are much more **dynamic**.
  - Web pages are **hyperlinked** to each other, and it is through hyperlink that a Web page author “**cites**” other Web pages.
  - Lastly, the size of the Web is **larger** than traditional data sources or document collections by several orders of magnitude.

11



DM #34 Module VI Web Mining Introduction 20211111 101726 Meeting Recording

# Web Content Mining

12



◀ ▶ 🔍 32:47 / 55:53

## Web Content mining

- Discovering useful information from contents of Web pages.
- Web content is very rich consisting of textual, image, audio, video etc and metadata as well as hyperlinks.
- The data may be unstructured (free text) or structured (data from a database) or semi-structured (html) although much of the Web is unstructured.

13



## Web Content Mining

- Extends work of basic search engines
- Search Engines
  - IR application
  - Keyword based
  - Similarity between query and document
  - Crawlers
  - Indexing
  - Profiles
  - Link analysis
- **Two approaches :**
  - agent-based
    - » software agents perform the content mining
  - database oriented
    - » view the Web data as belonging to a database

14



## Problems with Search Engine

- Use of the search engines to find content in most cases does not work well, posing an abundance problem. Searching phrase “data mining” gives 2.6m documents.
- It provides no information about structure of content that we are searching for and no information about various categories of documents that are found.
- *Need more sophisticated tools for searching or discovering Web content.*

15



DM #34 Module VI Web Mining Introduction 20211111 101726 Meeting Recording

## Text Mining for Web Documents

- Text mining for Web documents can be considered a sub-field of **Web content mining**.
- **Information extraction techniques** have been applied to Web HTML documents
  - E.g., *Chang and Lui (2001)* used a PAT tree to construct automatically a set of rules for information extraction.
- **Text clustering algorithms** also have been applied to Web applications.
  - E.g., *Chen et al. (2001; 2002)* used a combination of noun phrasing and SOM to cluster the search results of search agents that collect Web pages by meta-searching popular search engines.

16



## Crawler : Introduction

- **Robot (spider)** traverses the hypertext structure in the Web.
- Collect information from visited pages
- Used to construct indexes for search engines
- **Traditional Crawler** – visits entire Web (?) and replaces index
- **Periodic Crawler** – visits portions of the Web and updates subset of index
- **Incremental Crawler** – selectively searches the Web and incrementally modifies index
- **Focused Crawler** – visits pages related to a particular subject

18

↑ ↓ ← → ×

## How it works ?

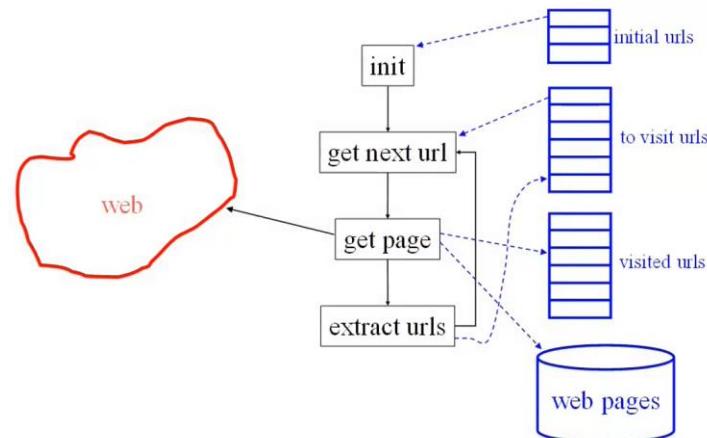
- It starts with a list of URLs to visit, called the **seeds..** As the crawler visits these URLs, it identifies all the hyperlinks in the page and adds them to the list of visited URLs, called the **crawl frontier**.
- URLs from the frontier are recursively visited according to a set of policies.

Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.



SP	SD
Sudharshu Pusad...	Smit Deshmukh
S	S
sawargandita21	shital.chavhan
NM	VK
Nida Mujawar	Veenja Kadam
P	S
prachi.wosware	Shivani.avasare
K	ST
kriti.shrivastava4	siddharth.lendale
A2	S
Apurva Raval 20...	salmehakim
Y	BM
yash.hoke	Bashirahamad M...

## What is a Crawler?



2



SP	SD
Sudharshu Pusad...	Smit Deshmukh
S	S
sawargandita21	shital.chavhan
NM	VK
Nida Mujawar	Veenja Kadam
P	S
prachi.wosware	Shivani.avasare
K	ST
kriti.shrivastava4	siddharth.lendale
A2	S
Apurva Raval 20...	salmehakim
Y	BM
yash.hoke	Bashirahamad M...

# Algorithms

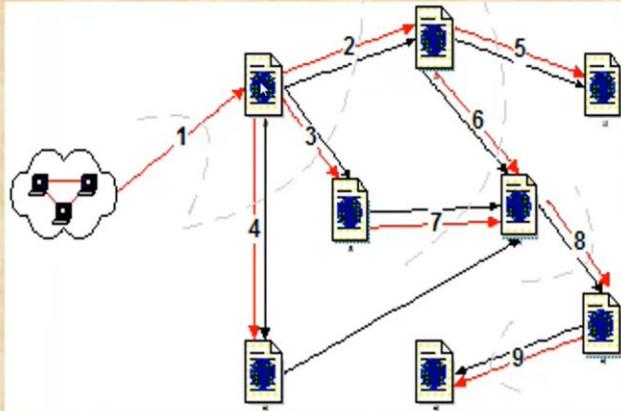
- Build the web graph
- Apply breadth first / depth first search algorithm

Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.



SP	SD
Sushanthu Pusad...	Smit Deshmukh
S	S
sawarpanita21	shital.chavan
NM	VK
Nida Majeer	Venja Kadam
P	S
prachi.wanire	shivani.avasare
K	ST
kmskrupali04	siddharth.tandale
A2	S
Apurva Raval 20...	saihilekam...
Y	S
yash.hoke	shubham.pagine
AP	R
Abhishek.Pol	revati.jathav2210...
A	BM
abhishekmore70	Bashirahamad M...

## Breadth First Crawler

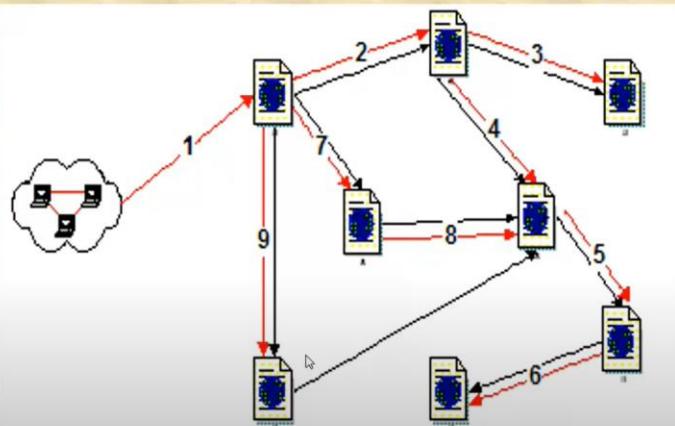


Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.



SP	SD
Sushanthu Pusad...	Smit Deshmukh
S	S
sawarpanita21	shital.chavan
NM	VK
Nida Majeer	Venja Kadam
P	S
prachi.wanire	shivani.avasare
K	ST
kmskrupali04	siddharth.tandale
A2	S
Apurva Raval 20...	saihilekam...
Y	S
yash.hoke	shubham.pagine
AP	R
Abhishek.Pol	revati.jathav2210...
A	BM
abhishekmore70	Bashirahamad M...

## Depth First Crawler



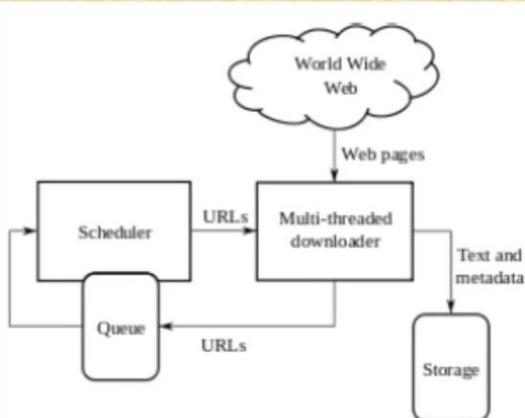
Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.



◀ ▶ ⏪ ⏩ 18:21 / 49:29

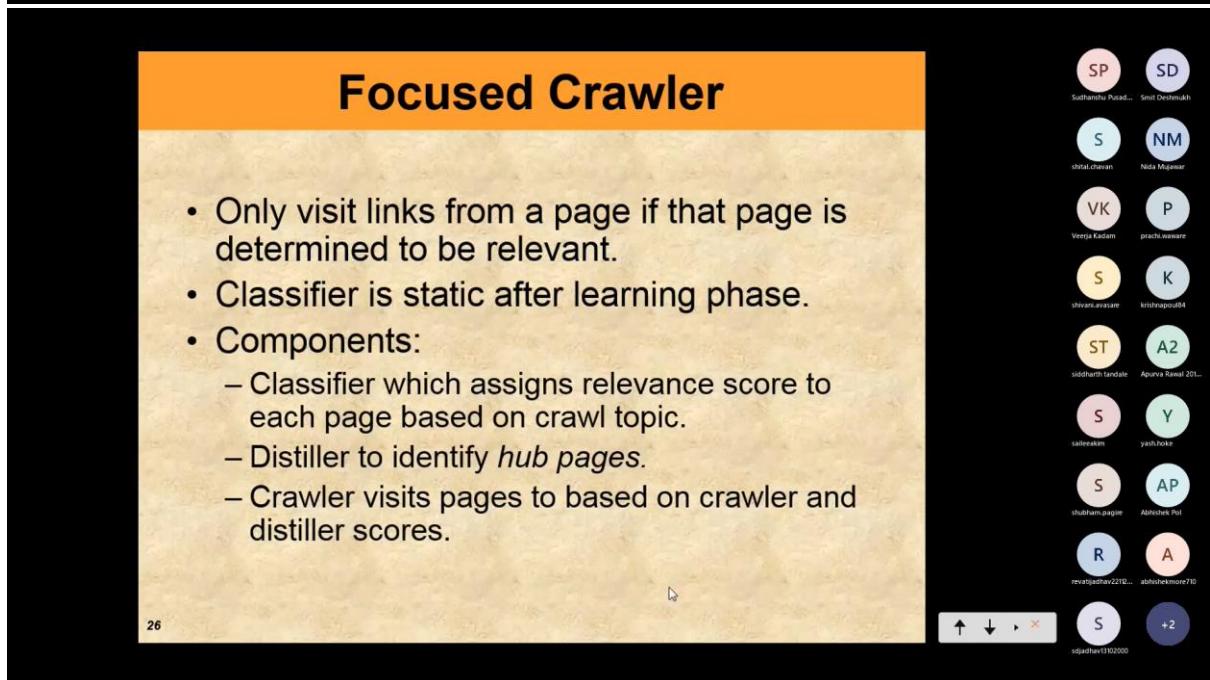
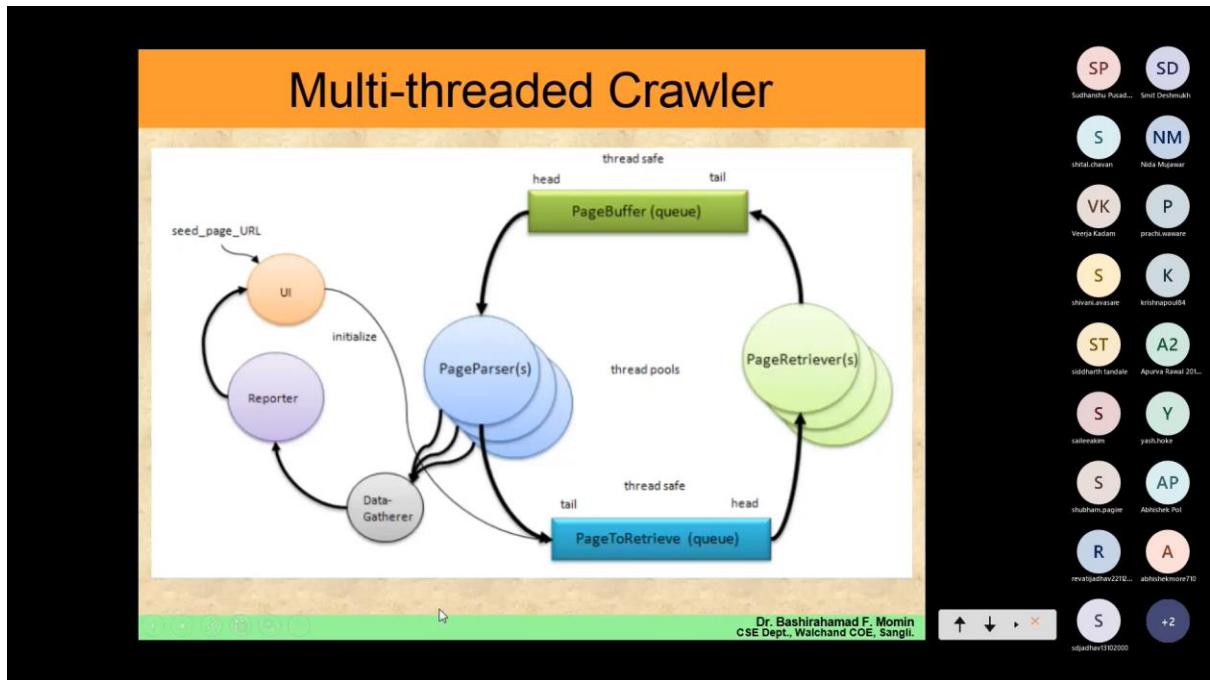


## Multi-threaded Crawler



Dr. Bashirahamad F. Momin  
CSE Dept., Walchand COE, Sangli.





## Focused Crawler

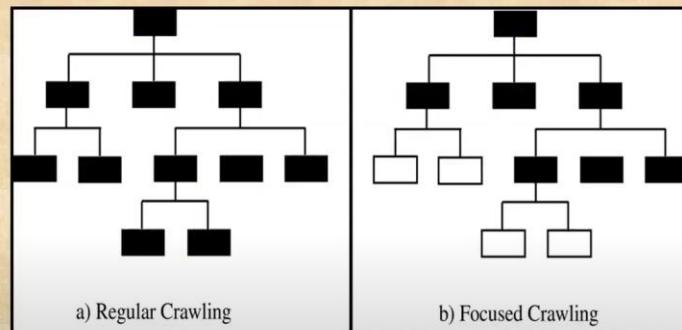
- Classifier to related documents to topics
- Classifier also determines how useful outgoing links are
- **Hub Pages** contain links to many relevant pages. Must be visited even if not high relevance score.

27

24:44 / 49:29

↑ ↓ × CC HD

## Focused Crawler



28

24:46 / 49:29

↑ ↓ × CC HD

## Context Focused Crawler

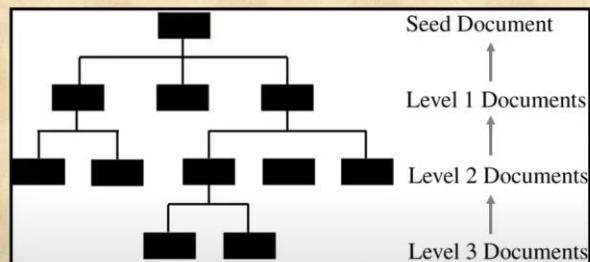
- Context Graph:
  - Context graph created for each seed document .
  - Root is the seed document.
  - Nodes at each level show documents with links to documents at next higher level.
  - Updated during crawl itself .
- Approach:
  1. Construct context graph and classifies using seed documents as training data.
  2. Perform crawling using classifiers and context graph created.

29

24:48 / 49:29



## Context Graph



30

24:50 / 49:29



## Virtual Web View

- **Multiple Layered DataBase (MLDB)** built on top of the Web.
- Each layer of the database is more generalized (and smaller) and centralized than the one beneath it.
- Upper layers of MLDB are structured and can be accessed with SQL type queries.
- Translation tools convert Web documents to XML.
- Extraction tools extract desired information to place in first layer of MLDB.
- Higher levels contain more summarized data obtained through generalizations of the lower levels.

A screenshot of a Microsoft Teams meeting interface. At the top, there's a header with the meeting title 'DM #35 Meeting in General 20211115 102328 Meeting Recording'. Below the header is a presentation slide with the title 'Virtual Web View'. The slide contains the list of points about MLDB. On the right side of the slide, there's a grid of user icons and names. At the bottom of the slide, there's a navigation bar with icons for back, forward, and search, and a timestamp '24:52 / 49:29'. To the right of the slide, there's a video player showing a person speaking. The video player has controls for volume, brightness, and a camera icon. The overall interface is dark-themed.

A screenshot of a Microsoft Teams meeting interface. It shows a presentation slide titled 'Virtual Web View' which lists the points about MLDB. To the right of the slide, there's a whiteboard with a hand-drawn diagram. The diagram illustrates the 'Virtual web view' as a layer above 'modular web as DBMS'. It shows 'HTML page' being converted to 'XML' via 'Extract temp information URLs'. There are also boxes for 'Page-to-page extraction' and 'parallel processing'. The whiteboard is signed off by 'Nidhi'. The Teams interface includes a sidebar with user profiles and a bottom navigation bar with various icons. The timestamp at the bottom right is '10:56 AM 11/15/2021'.

## Personalization

- Web access or contents tuned to better fit the desires of each user.
- Manual techniques identify user's preferences based on profiles or demographics.
- **Collaborative filtering** identifies preferences based on ratings from similar users.
- **Content based filtering** retrieves pages based on similarity between pages and user profiles.

# Web Structure Mining

- Mine structure (links, graph) of the Web
- Web as graph
  - Pages = nodes, hyperlinks = edges
- Techniques
  - PageRank
  - CLEVER (*modified as HITS*)
- Create a model of the Web organization.
- May be combined with content mining to more effectively retrieve important pages.

`<html>  
class="in">  
</html>`

Web structure  
↓  
Pages  
HTML

HTTP stack

X → fast but · Analytics

HTML → Can't use for analysis

Can't apply any algorithm

→ my algorithm to operate  
model it mathematically

ANN → Perceptron

Web → Pages → Graph  
↓  
Nodes  
Link ↓ Edges

Web page | website → modeled as graphs,  
once modeled as Graph → All graph theory is applicable

## CLEVER

- Identify authoritative and hub pages.
- **Authoritative Pages :**
  - Highly important pages.
  - Best source for requested information.
- **Hub Pages :**
  - Contain links to highly important pages.

37

## Hyperlink-Induced Topic Search (HITS)

- Based on a set of keywords, find set of relevant pages – R.
- Identify hub and authority pages for these.
  - Expand R to a base set, B, of pages linked to or from R.
  - Calculate weights for authorities and hubs.
- Pages with highest ranks in R are returned.

38





≡ DM #37 Module VI Web usage mining 20211117 111633 Meeting Recording



## Introduction

- **Web usage mining:** automatic discovery of patterns in clickstreams and associated data collected or generated as a result of user interactions with one or more Web sites.  
*Discovering user 'navigation patterns' from web data.*
- **Goal:** analyze the behavioral patterns and profiles of users interacting with a Web site.  
*Prediction of user behavior while the user interacts with the web.*
- The discovered patterns are usually represented as collections of pages, objects, or resources that are frequently accessed by groups of users with common interests.

3



5:14 / 58:47

## Web Usage Mining Applications

- Personalization
- Improve structure of a site's Web pages
- Aid in caching and prediction of future page references
- Improve design of individual pages
- Improve effectiveness of e-commerce (sales and advertising)

4



## Sources of Data for Usage Mining

- Automatically generated data stored in server access logs, *referrer* logs, *agent* logs, and client-side *cookies*.
- user profiles.
- metadata: page attributes, content attributes, usage data.

5

