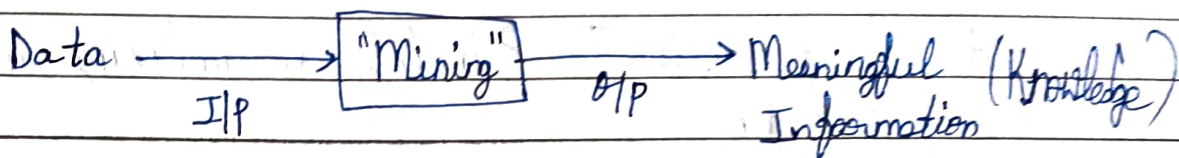


Data Mining Overview

Data vs Information?

Data Mining is the automated/semi-automated process of discovering meaningful patterns and knowledge from large volumes of data and involves retrieval, querying, inference, pattern discovery, and predictive modelling.



Eg: Using aggregated Google search data of particularly frequent/trending topics to identify informative and even predictive patterns from the huge database containing several queries' convergent outcomes.

Eg. 1. $\text{buys}(X, \text{"computer"}) \Rightarrow \text{buys}(X, \text{"webcam"})$ [support = 1%, confidence = 50%]
 where X is a variable representing a customer.

Eg. 2. $\text{age}(X, \text{"20..29"}) \wedge \text{income}(X, \text{"40K...49K"}) \Rightarrow \text{buys}(X, \text{"laptop"})$
 [support = 0.5%, confidence = 60%]

This rule indicates that of all customers under study, 0.5% are 20 to 29 years old with annual income of \$40,000 to \$49,000 and have purchased a laptop computer. There is a 60% probability that a customer in this age and income group will purchase a laptop.

"Multidimensional association Rule"

Knowledge Discovery in Databases (KDD):-

KDD - [Data cleaning, data integration, data selection, data transformation, data mining, data pattern evaluation, and knowledge presentation.

Data transformation → includes normalizing and/or homogenizing the data.

Data Categories:

1. Structured data: Tables, databases, spreadsheets
2. Semi-structured data: XML, JSON, emails
3. Unstructured data: Text, images, audio, video
4. Big Data: Volume, velocity, variety and veracity; ^{high} complexity and high-scaling. _{verifiable}
5. Streaming Data: Social media feeds, real-time sensor data (Netflix/Amazon Prime, sensor measurements, etc.)
6. Spatial and Temporal Data: GPS, time-series, geolocation, etc.

Multidimensional Data Summarization:

Summarizing data along dimensions (eg: sales by region, time-period, product) enables understanding; often used

OLAP (Online Analytical Processing) and aggregation for data summarization

Mining frequent patterns, associations and correlations

1. Uncovering relationships (eg: Market-based analysis) (metrics used: confidence, support, etc.)
2. Apriori Algorithm (works on Union and Intersection)
3. Classification (predicting categorical labels), Regression (predicting continuous numerical values).

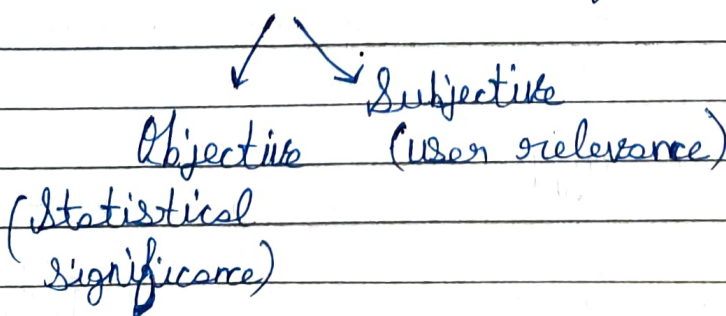
Eg: If-then, Decision Trees, Random Forest, Neural networks, etc.

Cluster Analysis:- Discovering structure in data

Deep learning: Using multi-layered neural networks; better for text, audio, images, video etc.

Outlier Analysis:- Detecting anomalies or unusual patterns (eg: fraud detection, network intrusion) medical imaging, etc.

Patterns must be novel, valid, useful and understandable



08/01/2026

Data, Measurements and Data Preprocessing

Classifying data into different attribute types: binary, ordinal, nominal and others (numeric attributes).

↓
Distinguishing the data "types"

1. Nominal Attributes:

Categories without order (no inherent order)

Operations:- Equality ($=$) and Inequality (\neq); no arithmetic operations.

Eg: Apples vs Oranges

We can calculate MODE to count here what is most occurring.

2. Binary Attributes:-

There are two categories.

Eg:- Tea \Rightarrow Yes/No.

3. Ordinal Attributes:-

Meaningful Order/Rank

Eg:- Assist. Prof. \rightarrow Associate \rightarrow Prof. \rightarrow Prof. (HAG)

Covariance and Correlation Analysis:

Covariance measures how numeric values/variables change together.

Normalized covariance ranges from -1 to +1
 Pearson's correlation coefficient (ρ)

$$\text{Cov}(A, B) = E(A \cdot B) - \bar{A}\bar{B}$$

Table 2.1

Q ₁	Time point	All Electronics	High Tech
	t1	6	20
	t2	5	10
	t3	4	14
	t4	3	5
	t5	2	5

$$E(A) = \frac{20}{5} = 4$$

$$E(B) = \frac{54}{5} = 10.8$$

FORMULAS:-

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

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

$$\text{Cov}(A, B) = E((A - \bar{A})(B - \bar{B}))$$

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

$$\therefore \text{Cov}(A, B) = \frac{\sum_{i=1}^n (a_i - 4)(b_i - 10.8)}{5}$$

$$= \textcircled{7} \left(\frac{as + 18.4 - 0.8 + 0 + 5.8 + 11.6}{5} = 7. \right)$$

Answer

Chi-square coefficient:

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

$$e_{ij} = \frac{\text{count}(A=a_i) \times \text{count}(B=b_j)}{n}$$

Table 2.2

a_i	Male	Female	Total
fiction	250 (90)	200 (360)	450
non fiction	50 (210)	1000 (840)	1050
Total	300	1200	1500

$e_{21} = \frac{1000 \times 450}{1500} = 300$
 $e_{22} = \frac{1000 \times 1050}{1500} = 700$
 $e_{11} = \frac{250 \times 450}{1500} = 75$
 $e_{12} = \frac{250 \times 1050}{1500} = 175$

$e_{12} = \frac{250 \times 1050}{1500} = 175$
 $e_{21} = 300$
 $e_{22} = 700$

$$\chi^2 = \frac{(250 - 75)^2}{75} + \frac{(50 - 210)^2}{210} + \frac{(200 - 360)^2}{360} + \frac{(1000 - 840)^2}{840}$$

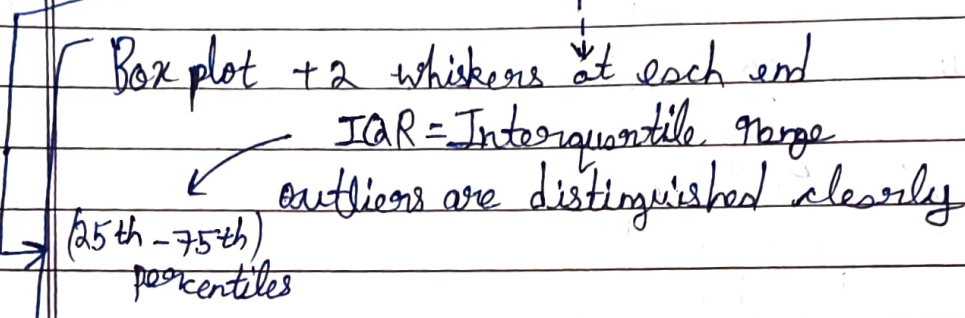
$$= 284.44 + 121.90 + 71.11 + 30.48 = 507.93$$

Answer

13/01/2026

(1 hr. LAB#)

- Data plots: -
- Quartile plot
 - Quartile - Quartile plot
 - Scatterplot
 - Histogram plot
 - Box and whisker plot



Range: $(Q_1 - 1.5(IQR), Q_3 + 1.5(IQR))$

Bar Graph, Grouped Bar Charts

Stacked Bar Chart, 100% - Stacked Bar Chart

Pie Chart (Circular)

- ↳ Donut Chart
- ↳ Nested Donut Chart

(इसमें सब Bars की length same है summing up to 100%)

13/01/2026

Main Class

Data Similarity & Quality Assessment: -

Recommendation Problem

Data Matrix: contains row attribute values for each object. $(n \times p)$

Square MATRIX

Dissimilarity Matrix $(n \times n)$: contains pairwise distances for symmetric & diagonal is 0. (d)

Proximity Measures: Categorical

1. Nominal: Simple matching $d = \frac{p-m}{p}$ a. Customer 1: {Gender: M, City: NY, Product: A}

2. Symmetric Binary: Simple matching $\frac{a+d}{\text{Total}}$ Customer 2: {Gender: M, City: LA, Product: B}

3. Asymmetric Binary: $d = \frac{(3-1)}{3} = \frac{2}{3} = 0.667$

$$J = \frac{a}{a+b+c}$$

Jaccard's coefficient

M	NY	A
M	LA	B

2x3 Data Matrix

2x3

→ Nominal → Used for distinct categories like City or Color

0	0.667
0.667	0

2x2 Data dissimilarity Matrix (nxn)

→ Symmetric Binary → Both states (0 and 1) like Yes/No, Gender, etc.

→ Asymmetric Binary (Ignores negative matches) purchases, for rare occurrences like disease, award-winning, selection, etc.

Object j

1	0
---	---

Numerical Distance: Minkowski

Object i	1	a	b
0	c	d	

$$d(i, j) = \left(\sum_{k=1}^p |x_{ik} - x_{jk}|^{\alpha} \right)^{1/\alpha}$$

- $q_1 = 1$ (Manhattan) (City Block)
- $q_2 = 2$ (Euclidean) (Straight Line)
- $q_\infty = \infty$ (Chebyshev) (Max dimension difference) (Supremum)

Normalization: $\rightarrow x_{new} = \frac{x - \min}{\max - \min}$
 (Prerequisite) for Numeric distance calculations

अर्थ? Without normalization, large values can bias / skew the distance estimates.

Q1 A: (Age = 25, Income = 50K)
 B: (Age = 35, Income = 80K)

A_{norm} = (0.125, 0.286) After normalizing (Age: 20-60, Income: 30-100K)
 B_{norm} = (0.375, 0.714)

Euclidean: $\sqrt{[(0.25)^2 + (0.428)^2]}$
 $= 0.495$ Ans.

$x_{new} = \frac{25 - 25}{35 - 25} = 0?$

\therefore Euclidean distance = $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$
 $= \sqrt{(0.375 - 0.125)^2 + (0.714 - 0.286)^2}$
 $= \sqrt{0.0625 + 0.1831} = \sqrt{0.2456}$
 $= 0.495$ (approx.)
Ans.

Ordinal & Mixed Attributes

order is important

$$d(i, j) = \frac{\sum \delta_{ij} d_{ij}}{\sum \delta_{ij}}$$

Maps values to ranks (1 to M)
Normalize ranks to [0.0, 1.0]
Uses numeric measures

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

a) Education $M_f - 1$

Levels:-

High School: 1 Bachelor's: 2
Master's: 3 PhD: 4

Normalized: HS = 0, Bachelor = 0.33
Master's = 0.67, PhD = 1.0

$$z_{HS} = \frac{1-1}{4-1} = 0$$

$$z_B = \frac{2-1}{4-1} = \frac{1}{3} = 0.33$$

$$z_M = \frac{3-1}{4-1} = 0.667$$

$$z_P = \frac{4-1}{4-1} = 1.0$$

a) Customer 1: { Age: 30 (numeric), Gender: M (nominal), Income: 70K (numeric) }

Customer 2: { Age: 25, Gender: M, Income: 50K }

$$d_{\text{income}} = \frac{|70-50|}{\text{max_income_diff}} = \frac{20}{60} = 0.3332$$

$$d_{\text{age}} = \frac{|30-25|}{\text{max_age_diff}} = \frac{5}{40} = 0.125$$

$$d_{\text{gender}} = 0 \text{ (match)}$$

∴ Data Matrix: $\begin{bmatrix} 30 & M & 70K \\ 25 & M & 50K \end{bmatrix}$

$$\text{Income}_1 = \frac{70-30}{70} = 0.57$$

$$\text{Income}_2 = \frac{50-30}{70} = 0.286$$

$$\text{Overall correct distance} = \frac{(0.125 + 0.284 + 0)}{3}$$

$$d_{\text{income}} \text{ (CORRECT value)} = 0.137 \text{ Answer}$$

$$= 0.284 \in 0.570 - 0.286$$

Cosine Similarity: - Used for text and document analysis

$$\cos \theta = \frac{A \cdot B}{\|A\| \|B\|}$$

Range: -1 to 1

(opposite) (identical)

Q1 Doc 1 = [2, 0, 1, 3]

Doc 2 = [1, 2, 0, 2]

(counts of: data, mining, science, analysis)

$$\cos \theta = \frac{(2*1 + 0*2 + 1*0 + 3*2)}{\sqrt{14} * \sqrt{9}}$$

$$= \frac{8}{(3.74 * 3)} = 0.713 \text{ Ans}$$

KL Divergence: - $D_{KL}(P||Q) = \sum P(i) \log \frac{P(i)}{Q(i)}$

⇒ How much P diverges from Q?

a) $P = [0.5, 0.3, 0.2]$
 $Q = [0.4, 0.4, 0.2]$

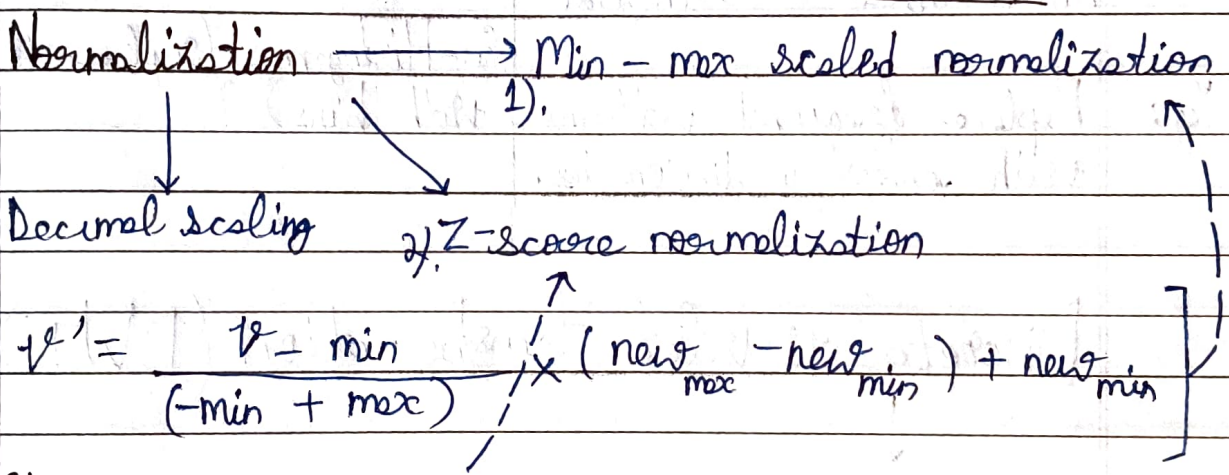
$$D_{KL} = 0.5 \times \log_e \left(\frac{0.5}{0.4} \right) + 0.3 \log_e \left(\frac{0.3}{0.4} \right) + 0.2 \log_e \left(\frac{0.2}{0.2} \right)$$

$$= 0.5 \times 0.223 + 0.3 \times (-0.288) + 0$$

$$= 0.1115 - 0.0864 = 0.0251 \text{ Ans.}$$

16/01/2026

Data Cleaning & Preprocessing:-



Formula:

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

$$v' = \frac{v}{10^j} \text{ where } j = \text{smallest integer s.t. } \max(v') < 1$$

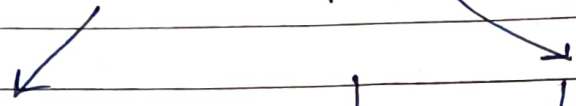
Discretization:-

1. Equal width: Divide range into N equal intervals.
2. Equal frequency: Each bin has same count.

3. Clustering Based: Use K-means to find natural groups

4. Decision tree: Use entropy for optimal splits.

Data Compression



Lossless:

Eg: Run-length encoding:

AAAA BBCC → 4A3B2C

Eg: Replace frequent patterns with codes in dictionaries.

Lossy

Eg: PCA (dimensionality reduction)

Eg: Histograms (replace values with bins)

Discrete wavelet transformation (DWT)



Used in linear signal processing for data compression.

20/10/12025

PCA: Find the angle that shows the most information (variance).

"Kaiser Criterion"

Used in Face Recognition etc.

For PCA, we need to calculate Eigen values for the components.

Attribute Subset Selection

Filter

Variance threshold,
 correlation,
 Chi-square

(Low-variance/highly
 correlated features may
 be preprocessed & removed)

Wrapper

We have forward
 selection, backward
 elimination.

(Evaluate feature sets
 iteratively)

Embedded

LASSO Regression,
 Random Forest

(In-built feature
 selection)

Kernel - PCA: Maps data to higher dimensions to find linear
 separations

t-SNE: Preserve local structure but not for feature
 extraction

Autoencoders: NNs that compress (encode) and reconstruct
 (decode) and learn automatically, feature selection.

Normal PCA v/s Kernel PCA

dimensionality to
 get reduced

dimensionality increases

20/01/2026

Data Warehousing

Warehouse

Operational
(OLTP)

Analytical
(OLAP)

(Online Transaction Processing with only INSERT, UPDATE, DELETE)

(Online Analytical Processing with SELECT and AGGREGATE commands)

→ Simpler queries

→ Complex queries
→ Denormalized for speed

A data warehouse is a subject-oriented, integrated, time-variant and non-volatile collection of data in support of management's decision-making process.

Father of "Data Warehousing"

"Prof. Bill Inmon's Definition"

ETL: Extract Transform Load

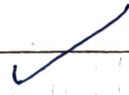
Data Lakes में पहले Load करते हैं तदोपरांत आवश्यकता पड़ने पर उसे Transform किया जाता है।

EDW vs Data Marts:-

Enterprise Data
Warehousing (EDW)



Data Marts
(DM)

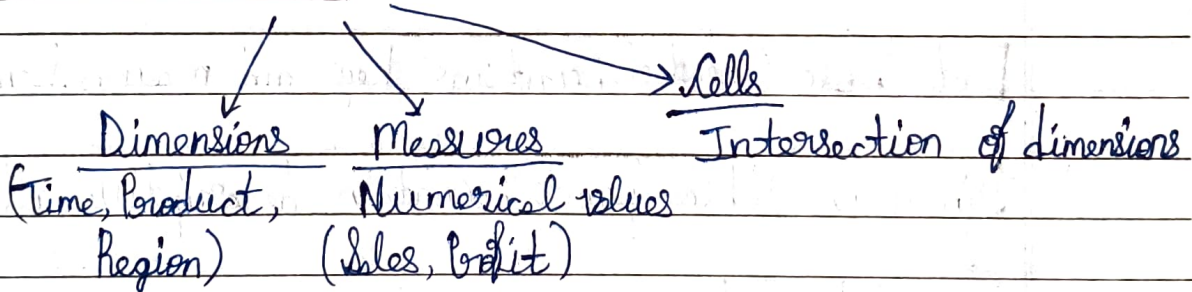


(Slide में
34th & 35th
differences!!)

EDA (Exploratory Data Analysis)

Data lakes → have a lot of unstructured data usually

Data Cubes:



Lattice of Cuboids (Power of Multidimensionality)

(0D, 1D, 2D, 3D)

apex
cuboid

base
cuboid

For n dimensional cuboid, there are n features
 for representing the base cuboid (Check again once!!)

n -dimensions

→ approx. 2^n dimensions cuboids will be present.

23/01/2026

Data Cube:-

Dimensions: Perspectives (Time, Product, Region)

Measures: Numerical values

Cells: Intersection of dimensions

Note:- "Base" Cuboid is the largest and "0" Cuboid is the smallest cuboid in the dataset.

→ Star Schema:-

Fact table (center): contains keys and measures. Denormalized

Dimension tables (points): contain descriptive attributes

→ Snowflake Schema:

Normalized dimensions are there in the architecture,

→ Fact Constellation (Galaxy) Schema:-

Multiple Fact Tables: a complex scheme.

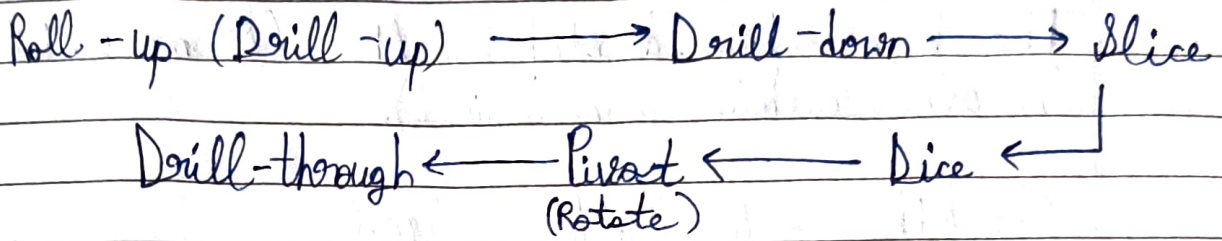
Types of Measures:-

→ Additive

→ Non-Additive

→ Semi-Additive

OLAP (Online Analytical Processing) Operations:-



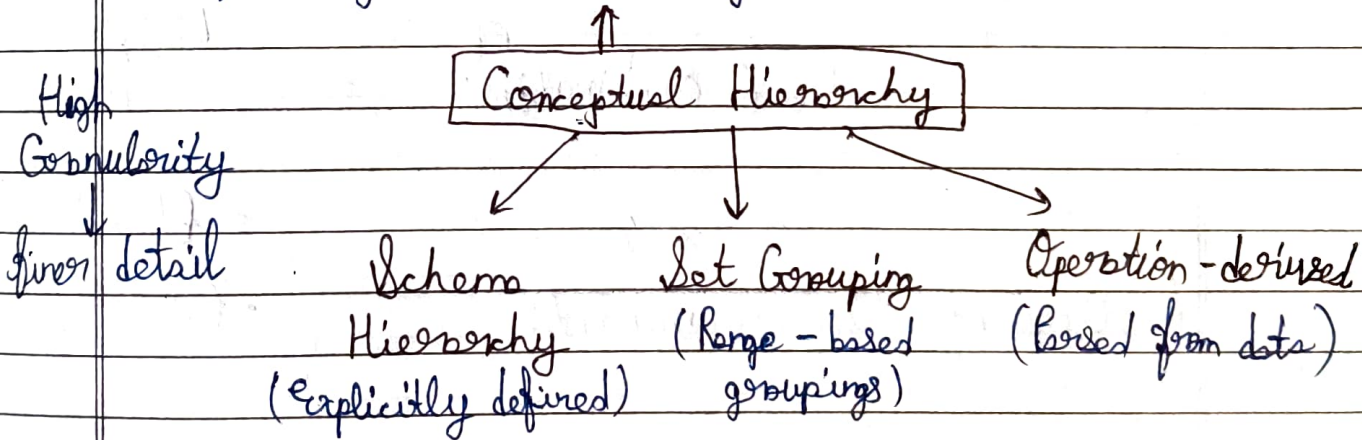
Advanced Data Warehousing

'Amazon' Scenario

5-dimensional data cube (Product × Time × Region × Customer Segment × Sales channel).

We have to understand the concept hierarchies.

(*) (*) A sequence of mapping from lower-level (highly detailed) to ~~low~~ higher-level concepts within a dimension.



Total order: strict hierarchy where every child has exactly one parent.

Partial order: Employees might belong to multiple committees and departments.

Multisway array cube computation:-

The base cuboid, denoted by ABC (from which all the other cuboids are derived).

Chunk No.	Eqn. A	Eqn. B	Eqn. C
1	a_0	b_0	C_0
2	a_0	b_0	C_1
3	a_0	b_0	C_2
4	a_0	b_0	C_3
5	a_0	b_1	C_0
6	a_0	b_1	C_1
7	a_0	b_1	C_2
8	a_0	b_1	C_3
⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮
16	a_1	b_0	C_0

≈ 6 combinations

$$40 \times 400 \times 4000$$

$$\overline{3} \times \overline{2} \times \overline{1} = 6$$

$$= 64 \times 10^6$$

Another possible arrangement:

$$= 6.4 \text{ crore}$$

$a_0 \ b_0 \ C_0$
 $a_1 \ b_0 \ C_0$
 $a_2 \ b_0 \ C_0$
 $a_0 \ b_1 \ C_0$
 $a_1 \ b_1 \ C_0$

and similarly,
 more arrangements
 possible

(As explained in Book,
 (Before Ex 3.13) ← so careful!!)

for b:

for a:

Cube	Aggregate over
AB	C
AC	B
BC	A

- $a_0 \ b_0 \ c_0$
- $a_0 \ b_0 \ c_1 \ // \ b_0 \ c_2$
- $a_0 \ b_0 \ c_3 \ // \ b_1 \ c_0$
- $a_0 \ b_1 \ c_1$
- $a_0 \ b_1 \ c_2$
- $a_0 \ b_1 \ c_3$
- $a_0 \ b_2 \ c_0$
- $a_0 \ b_2 \ c_1$
- $a_0 \ b_2 \ c_2$
- $a_0 \ b_2 \ c_3$
- $a_0 \ b_3 \ c_1$
- $a_0 \ b_3 \ c_2$
- $a_0 \ b_3 \ c_3$
- $a_0 \ b_3 \ c_4$

$$\begin{aligned}
 &10 \times 1000 + 10 \times 4000 \\
 &\quad + 400 \times 4000 \\
 \hline
 &16,41,000
 \end{aligned}$$

$$\begin{aligned}
 &40 \times 400 + 100 \times 1000 \\
 &\textcircled{A} \quad \textcircled{B} \quad \textcircled{C}
 \end{aligned}$$

$$+ 40 \times 1000$$

Similarly,
other arrangement

$$\begin{aligned}
 &= 16,000 + 100,000 \\
 &\quad + 40,000
 \end{aligned}$$

$$= 156,000$$

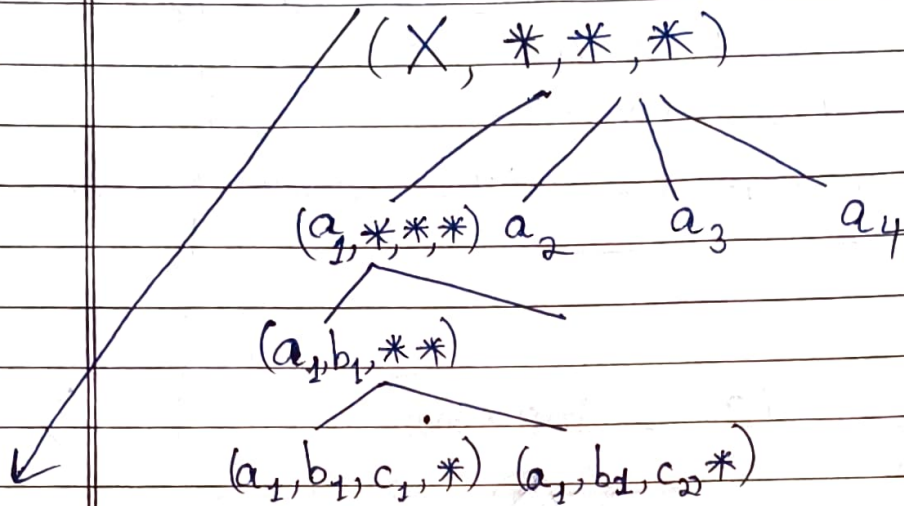
(श्री दर्ष जी class में explain करिये ****)

- $a_0 \ b_0 \ c_0$
- $a_1 \ b_0 \ c_0$
- $a_2 \ b_0 \ c_0$
- $a_3 \ b_0 \ c_0$
- $a_0 \ b_1 \ c_0$
- $a_1 \ b_1 \ c_0$
- $a_2 \ b_1 \ c_0$
- $a_3 \ b_1 \ c_0$
- $a_0 \ b_2 \ c_0$
- $a_1 \ b_2 \ c_0$
- $a_2 \ b_2 \ c_0$
- $a_3 \ b_2 \ c_0$
- $a_0 \ b_3 \ c_0$
- $a_1 \ b_3 \ c_0$
- $a_2 \ b_3 \ c_0$
- $a_3 \ b_3 \ c_0$

"Bottom Up"
"Computation"

BUC construction of an iceberg cube

"Checking antimonotonicity"
 ("Pruning") } check these terms to get BUC construction (Please!!)



03/02/2025 Table 3.4 Original Database (in Slides)

TID	A	B	C	D	E	inserted index table
1	a ₁	b ₁	c ₁	d ₁	e ₁	बताए!!
2	a ₁	b ₂	c ₁	d ₂	e ₁	
3	a ₁	b ₂	c ₁	d ₁	e ₂	
4	a ₂	b ₁	c ₁	d ₁	e ₂	
5	a ₂	b ₁	c ₁	d ₁	e ₃	

a₁ → 1, 2, 3

a₂ → 4, 5

b₁ → 1, 4, 5

b₂ → 2, 3

c₁ → 1, 2, 3, 4, 5

$d_1 \rightarrow 1, 3, 4, 5$

$d_2 \rightarrow 2$

$e_1 \rightarrow 1, 2$

$e_3 \rightarrow 5$

$e_2 \rightarrow 3, 4$

वाराणसी

अनंदवन

काशी

बनारस

नाथद्वारा

वृंकावन

Apriori Algorithm:-

- Candidate set
- Check whether they follow apriori
- Calculate Support
- Finalize

(Sup. Count)

Frequent Item Set $\Rightarrow \{11, 12\} \Rightarrow 4$

$\{11, 13\} \Rightarrow 4$ "Rule Mining"

$\{11, 15\} \Rightarrow 2$

$\{12, 13\} \Rightarrow 4$

$\{12, 14\} \Rightarrow 2$

$\{12, 15\} \Rightarrow 2$

Transaction Table

TID List of item IDs

T100 11, 12, 15

T200 12, 14

T300 12, 13

T400 11, 12, 14

T500 11, 13

T600 12, 13

T700 11, 13

T800 11, 12, 13, 15

T900 11, 12, 13

I_1 determines $(I_1, I_2) = ?$

$I_1 \rightarrow (I_1, I_2)$

$= \text{Support}(I_1 \cup I_2)$

$\text{Support}(I_1)$

$$I_2 \rightarrow I_5 = \frac{\text{Support}(I_2 \cup I_5)}{\text{Support}(I_2)}$$

$I_5 \rightarrow I_1 = 1$ I_1 is always purchased if purchase I_5 as such

$$(I_5 \rightarrow I_{(1,2)}) \text{ or } I \rightarrow \{I_1, I_2\}$$

$$\frac{\text{Support}(I_1, I_2, I_5)}{\text{Support}(I_1 \cup I_2)}$$

$\{11, 12\} \Rightarrow 15$ confidence = $2/4 = 50\%$

Bucket constraints:- Hash-based Technique
 Reduce candidate set

Partitioning the data (mining):

Sampling: Speed vs Accuracy "Trade-off"



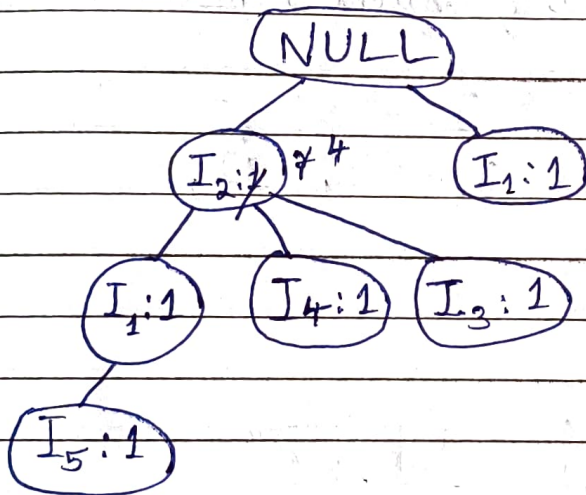
"False Negatives" Risk

Digital Itemset Counting (DIC):

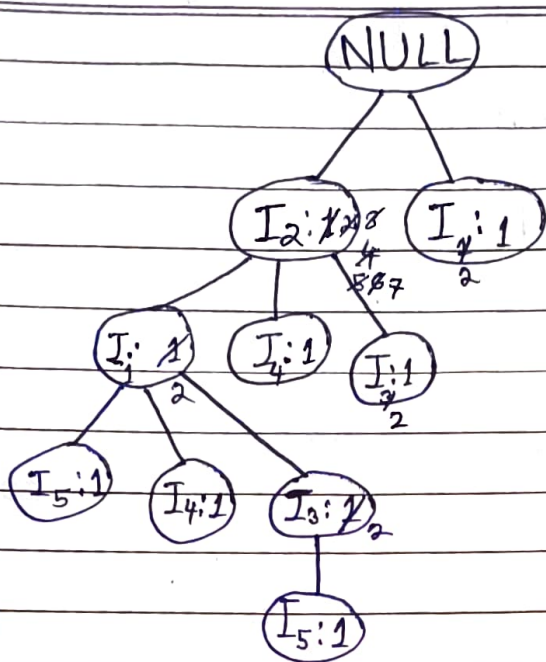
Data Mining (continued)

FP - Growth Algorithm:-

T ID	List of Item-IDs
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I1, I3
T700	I2, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3;



(कुछ next page)



Mining for the least frequent item (I_5):

path to reach $I_5 \Rightarrow I_2 \rightarrow I_1 \rightarrow I_5 \text{ (1)}$
 $I_2 \rightarrow I_1 \rightarrow I_3 \rightarrow I_5 \text{ (1)}$

Prefix path:

$\langle I_2, I_1 \rangle$
 $\langle I_2, I_1, I_3 \rangle$
 $I_2: 2 \quad I_1: 2 \quad \text{"} I_3: 1 \text{"}$

(less than minimum support count)

Apriori suffers from generating huge numbers of candidate items. FP-Growth adopts a Divide-and-Conquer strategy.

Compress DB into an FP-tree.
 Mine the tree recursively.

* Path to reach $I_4 \Rightarrow$ $I_2 \rightarrow I_4$ (1)
 $I_2 \rightarrow I_1 \rightarrow I_4$ (1)

Prefix Path \Rightarrow

$\langle I_2 \rangle$
 $\langle I_2, I_1 \rangle$
 $I_2:2 \quad I_1:1$

-----> (less than minimum support count = 1)

freq. set $\Rightarrow \{I_1, I_2\}$

* Path to reach $I_3 \Rightarrow$ $I \rightarrow I_3$ (2)
 $I_2 \rightarrow I_1 \rightarrow I_3$ (2)
 $I_1 \rightarrow I_3 \rightarrow (2)$

Prefix Path \Rightarrow $\langle I_3 \rangle$
 $\langle I_2, I_4 \rangle$

$I_2:4 \quad I_1:4$

freq. set: $\{\{I_2, I_3\}, \{I_2, I_3\}, \{I_2, I_1, I_3\}\}$

* Path to reach $I_2 \Rightarrow I_2$ (7)

10/02/2020

Frequency Pattern Tree

FP-Growth adopts a Divide-and-Conquer strategy.

- Compress the DB into FP-tree.
- Mine the tree recursively.

Avoiding Candidate Generation

Mining using Vertical Data format:

Transactional data set

TID	List of item IDs
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

I1 → T100, T400, T500, T700, T800, T900

I2 → T100, T200, T300, T400, T600, T800, T900

I3 → T300, T500, T600, T700, T800, T900

I4 → T200, T400,

I5 → T100, T800

1.
$$I_1 \cup I_3 = T500, T700, T800, T900$$

2.
$$I_1 \cup I_2 = T100, T400, T800, T900$$

$$3. \quad I_1 \overset{U}{\cap} I_4 = T400$$

$$4. \quad I_1 \overset{U}{\cap} I_5 = T100, T800$$

$$5. \quad I_1 \cup I_2 \cup I_3 = T800, T900 \quad \checkmark$$

$$6. \quad I_1 \cup I_2 \cup I_5 = T100, T800 \quad \checkmark$$

ECLAT, mine closed/~~minimal~~ maximal

→ Maximal itemsets have "lossy" compression, lose support information.

→ 100 closed itemsets: lossless compression. Can derive all support counts (lossless compression).

Strong vs Interesting ASSOCIATION rules:

buys (X, "computer games") \Rightarrow buys (X, "videos")
[support = 40%, confidence = 66.7%]

→ a misleading "strong" association rule!!

$$\text{lift}(A, B) = \frac{P(A \cap B)}{P(A) \cdot P(B)} \quad \neq$$

$$\text{lift}(A \rightarrow B) = \frac{P(A \cup B)}{[P(A) * P(B)]}$$

U \equiv intersection here, not union ("∩")

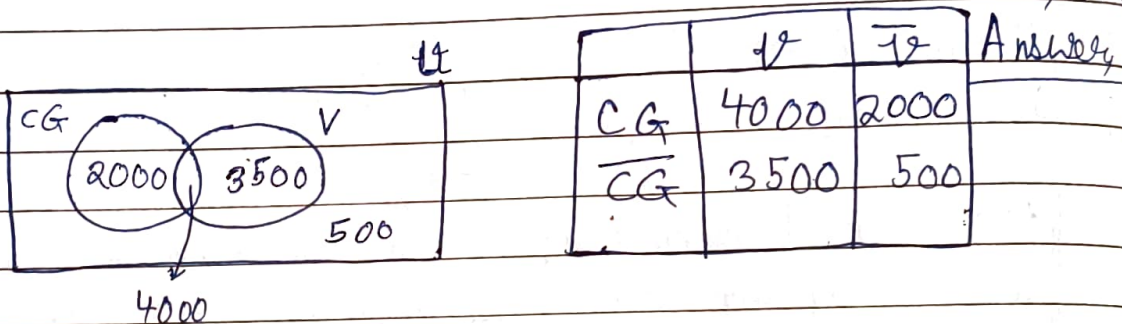


$\text{lift}(A \rightarrow B) = 1 \rightarrow \text{independent}$ ✓

buys (X, "computer games") \Rightarrow buys (X, "videos")

[support = 40%, confidence = 66%]

Support Value = $\frac{40,000}{10,000.0} = 40\%$ and Confidence Value = $\frac{4000.0}{6000.0} = \frac{4}{6} = 2/3$



$$e_{11} = \frac{6000 \times 7500}{10,000} = 4500$$
 ✓

$$e_{12} = \frac{6000 \times 2500}{10,000} = 1500$$

$$e_{21} = \frac{4000 \times 7500}{10,000} = 3000$$

$$e_{22} = \frac{2500 \times 4000}{10,000} = 1000$$

$$\chi^2 = \frac{(4000 - 4500)^2}{4500} = \frac{500 \times 500}{4500} = \frac{500}{9}$$

∴ negative correlation actually at final.

Answer

Krishna

Jaccard Similarity :-

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

→ Cluster similar patterns together

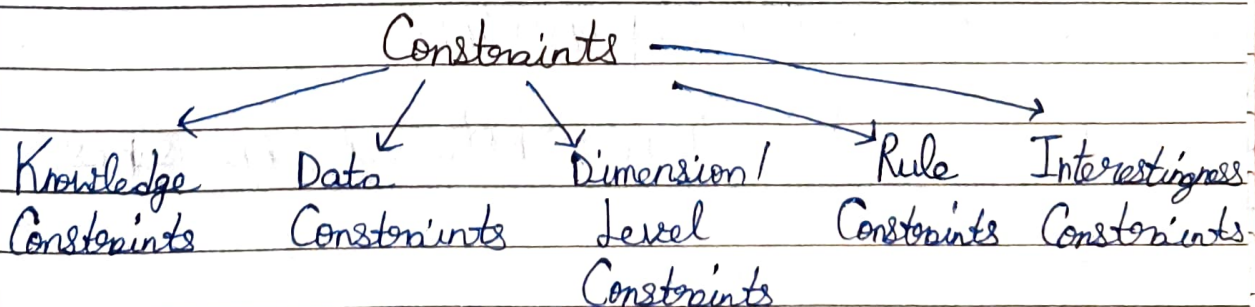
→ Balancing Score = $\alpha \times \text{Support} + \beta \times (1 - \text{max similarity to be deleted})$

→ For Geometrically "FAR" patterns,

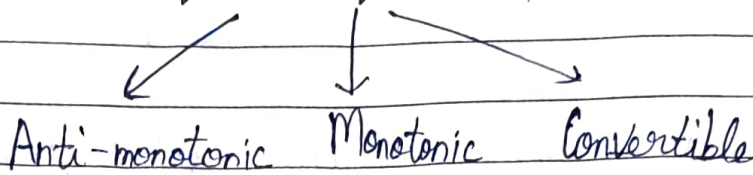
$$\text{Distance}(\{A, B\}, \{A, C\}) = 1 - \text{Jaccard Similarity}$$

Larger size = higher support
Greater distance = more diversity

"Constraint-based Mining"



Classification of Constraints



Order ID Customer ID

Constraint-based SQL Query

Note:

given in page 60
in slide sent on 3rd Year Advanced Algo
(17/02/2026)

- Rare Events
- Compressed Results
- Focused Mining

20/02/2026

Graph Mining:-

Graph $G: (V, E)$

Subgraph:- G' 's subgraph of G if $V' \subseteq V$
Graph Isomorphism and $E' \subseteq E$

Downward closure property:-

FSG Miner → better for sparseness