

SF-417

## Data Warehousing &

### Data Mining

Unit 1 → warehousing

Unit 2-3-4 → 3 different techniques of Data Mining

#### # Database Warehousing

DBMS  
SQL  
SYSTEM & SOFTWARE

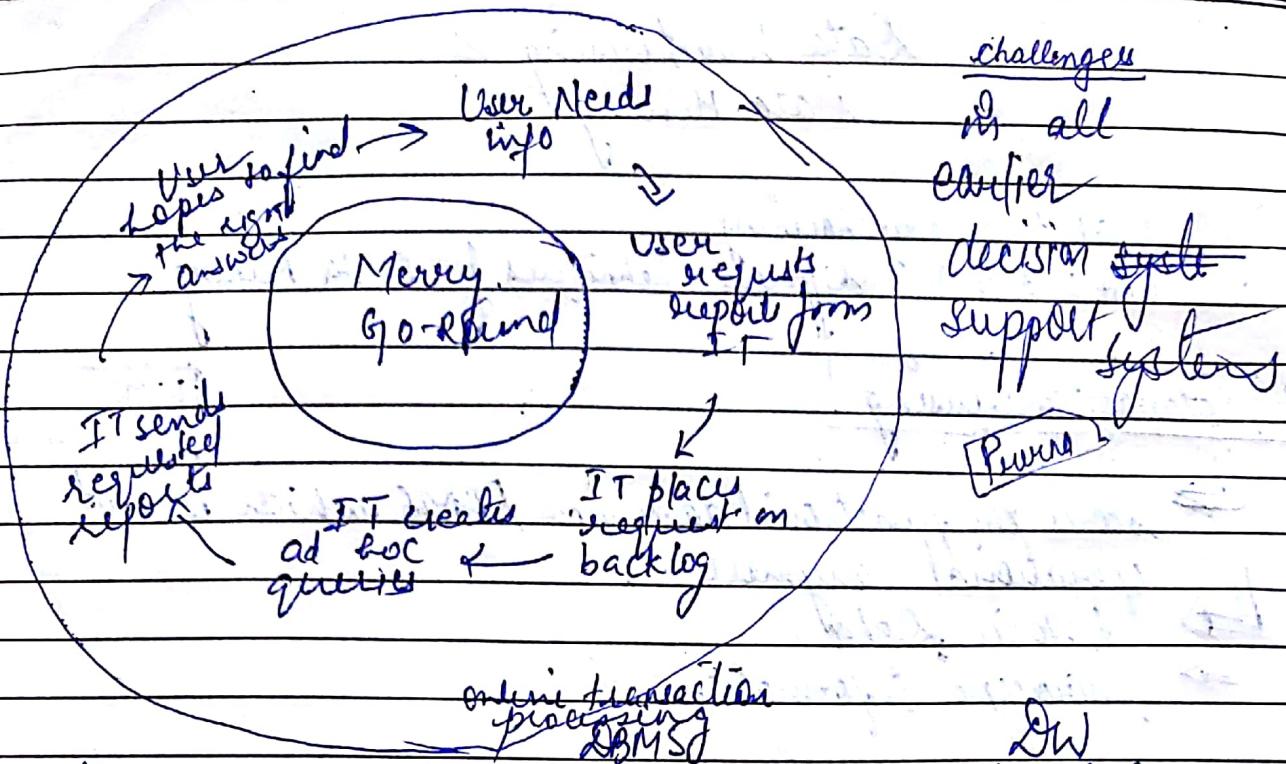
- ↳ stores the functional data from DBMS which is not operational anymore.
- ↳ Historic Data
- ↳ strategic information

#### # Data Warehouse

- It is an informational environment that provides
- (i) an integrated and total view of an Enterprise
- (ii) makes the enterprise's current and historic information easily available for strategic decision making
- (iii) makes decision support transactions possible without hindering the operational systems
- (iv) encloses the organizations' information consider
- (v) presents a flexible and intuitive source of strategic information

#### # Earlier proposed Decision Support Systems.

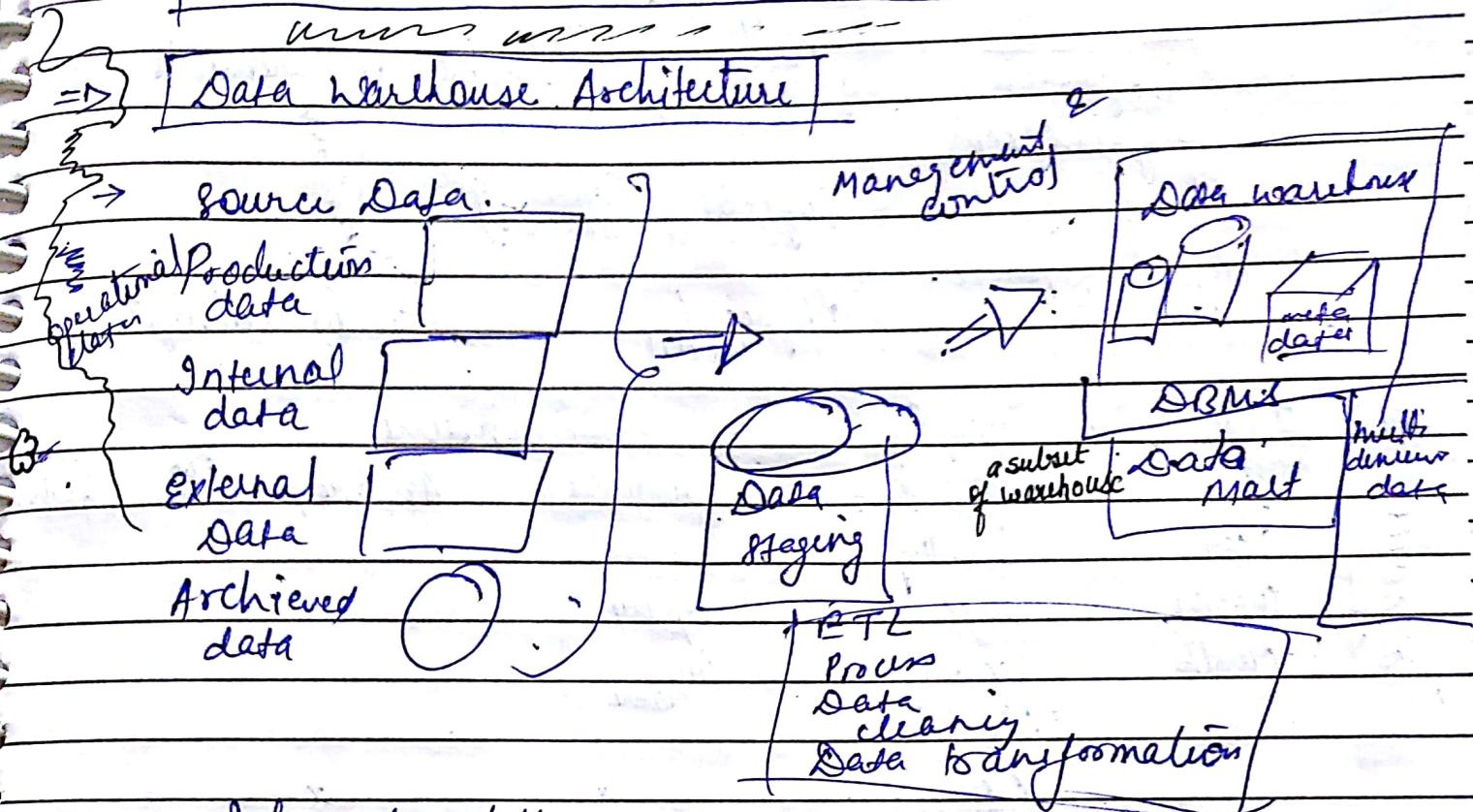
- 1) ad hoc Reports
- 2) special Extract reports
- 3) Information Centers
- 4) small applications
- 5) Decision Support Systems → Digitized
- 6) Executive Information System



## # OTLP systems      Informational Systems

	Operation System	Decision Support System
Data content	Current Values	Historic, achieved & aggregated
Access frequency	Very high	Very low.
No of users	Very large Number	Small number (Executives)
Typing access	Read, write & update	only read
Usage	Predictable, Repetitive	Random, Ad-hoc & heuristic
Response time	Subseconds	several seconds to few minutes
Data structure	optimized for transactions	optimized for complex Queries
Design goal	designed to put the data designed to get strategic	information out of database
	into the database	wheeling the business turn
	making the wheels	wheeling the business turn
	of business turn	

# # Data Warehouse Architecture



Information delivery

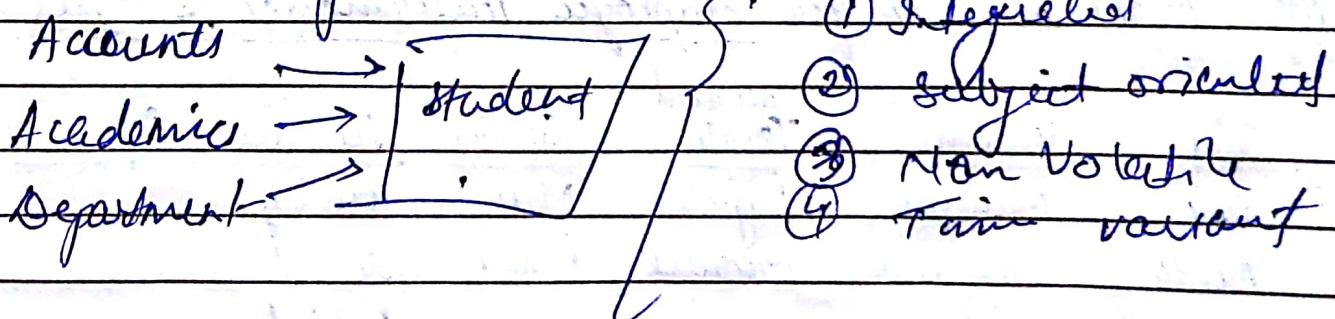
Data mining

OLAP

Query / report

Pearson

## # features of Data in Data Warehouse .

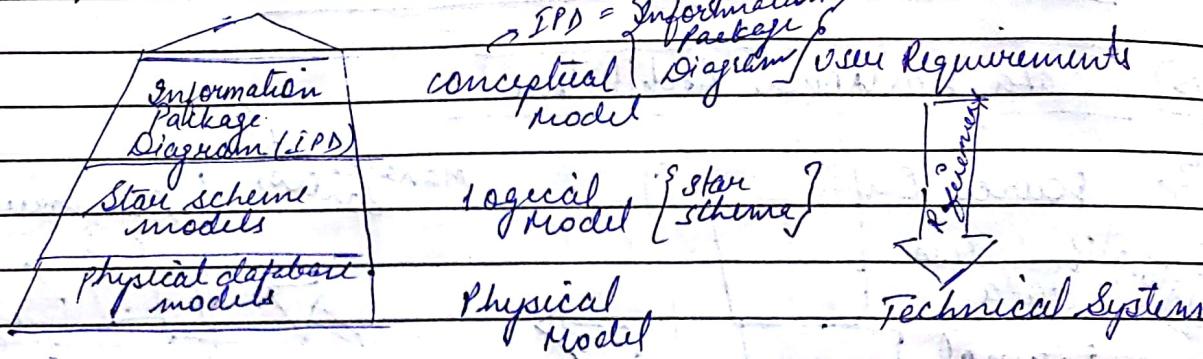


Best tool for Data Warehousing  
Oracle is it used 80%

5<sup>th</sup> Sep 2018

Page No. \_\_\_\_\_  
Date / /

## # Data Warehouse Modeling



### # IPD | Sales Analysis → Dimensions

Time	Locations	Products	Age Groups/Eco class	Gender
Year	Country	class		
Quarter	Area	Group		
Month	Region	Product		
	District	Name		
	store		x	x

Measures : Forecast Sales, Budget Sales, Actual Sales, Forecast Variance, Budget variance

# Warehouse is based on business organisation. Based on that an IPD is created. For example → Sales analysis

// categories are hierarchical //

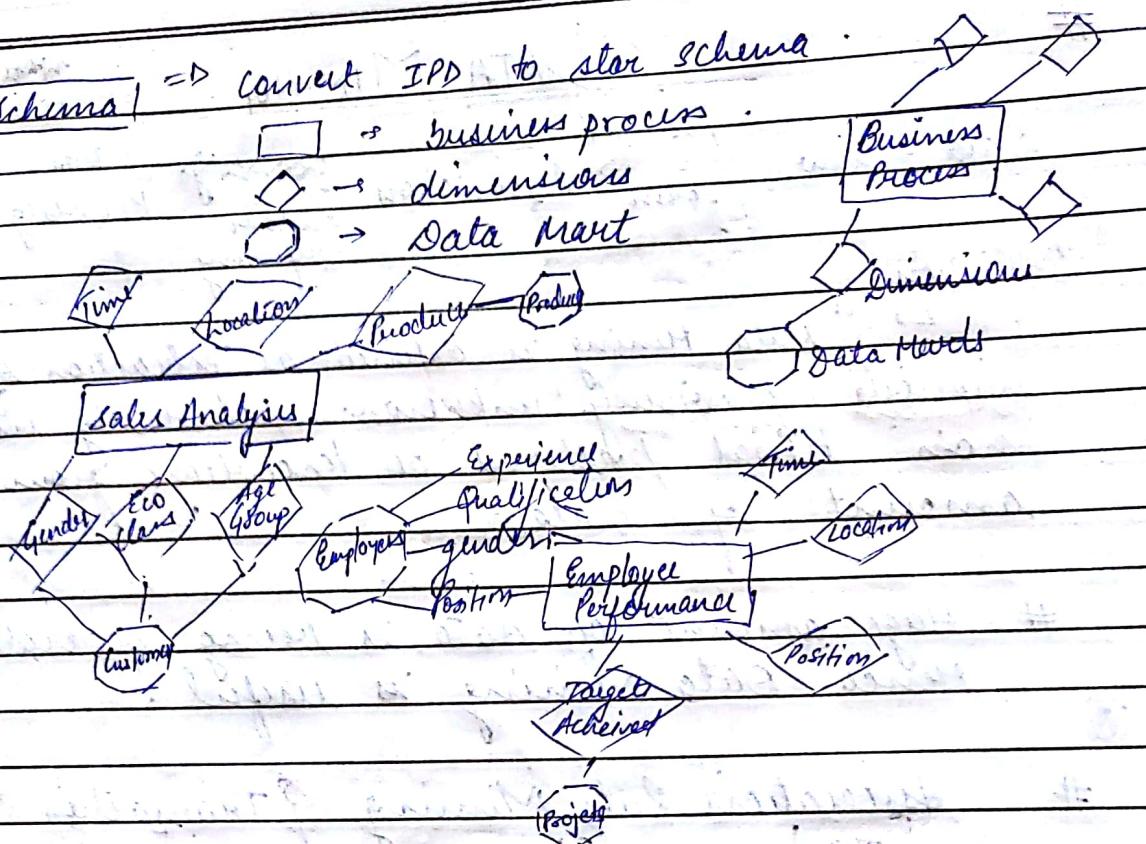
// time will always be there //

11/10/11 Make an IPD for Employee Performance Analysis

Time	Position	Achieved Targets	Gender	Location	Experience	Qualification
Year	Manager	Projects	Male	Country		
Month	Team Leader	Modules	Female	State		
Week	Employee			District		
				Branch		

Measures : Daily Working hours, Avg working hours, Expected working hours

# [star schema]  $\Rightarrow$  convert IPD to star schema



7<sup>th</sup> Sep 2018

# Data

7 Sep 2018



Definition

→ "Data Mining is defined as extraction of interesting, implicit, previously unknown, potentially useful and non trivial patterns or knowledge from huge amount of data"

# Huge amount of data is being generated everyday.  
Hence Data Mining is useful.

8

# Association Rule Mining of Transaction Data Set

TID	Items	D : Transaction database
T0	Beer, Nuts, Diaper	I : Itemset
T1	Beer, Coffee, Diaper	T : event $T \subseteq I$
T2	Nuts, Eggs, Milk	$T$ contains itemset
T3	Beer, diaper, Eggs	$x \iff y \subseteq T$
T4	Nuts, coffee, Diaper, Eggs, Milk	

Association Rule :  $X \rightarrow Y$

$X \subseteq I$ ,  $Y \subseteq I$  and  $X \cap Y = \emptyset$

$X \rightarrow Y$  : support =  $\frac{|X \cup Y|}{|D|}$  can be calculated for individual items & itemsets as well.

$$\text{confidence} = \frac{|X \cup Y|}{|X|}$$

# Apriori Algorithm

→ gives frequent itemsets  
Support = 50% (0.5), Confidence (0.5)

Iteration :- Candidates

F1	B1	Frequent
1-itemset	Beer, Nuts, diaper Coffee, Eggs, Milk	Beer, Nuts, Diaper, Egg

Iteration 2 : Candidates

$F_2$  2-itemset :  
 (Beer, Nuts), (Beer, Diaper)  
 (Beer, Coffee), (Beer, Eggs) (Beer, Milk)  
 (Nuts, Diaper), (Nuts, Coffee), (Nuts, Eggs) (Nuts, Milk)  
 (Diaper, Coffee), (Diaper, Eggs), (Diaper, Milk)  
 (Coffee, Eggs), (Coffee, Milk), (Eggs, Milk)

2222222222

Frequent : - (Beer, Diaper) // (Nuts, Eggs)

$F_3 = \phi$  Association Rule :- Beer  $\rightarrow$  Diaper  
 Diaper  $\rightarrow$  Beer

# Making Association Rules = RHS contains only 1 item

$$\text{confidence } (\text{Beer} \rightarrow \text{Diaper}) = \frac{3}{3} = 1 \quad \left. \begin{array}{l} \text{Both} \\ \text{are valid} \end{array} \right.$$

$$(\text{Diaper} \rightarrow \text{Beer}) = \frac{3}{4} = 0.75$$

## # APRIORI ALGORITHM

$C_k$  : Candidate itemsets of size K

$L_k$  : frequent itemsets of size K

$L_1$  : { frequent items }

for ( $K=1$  ;  $L_K = \phi$  ;  $K++$ )  
 do begin

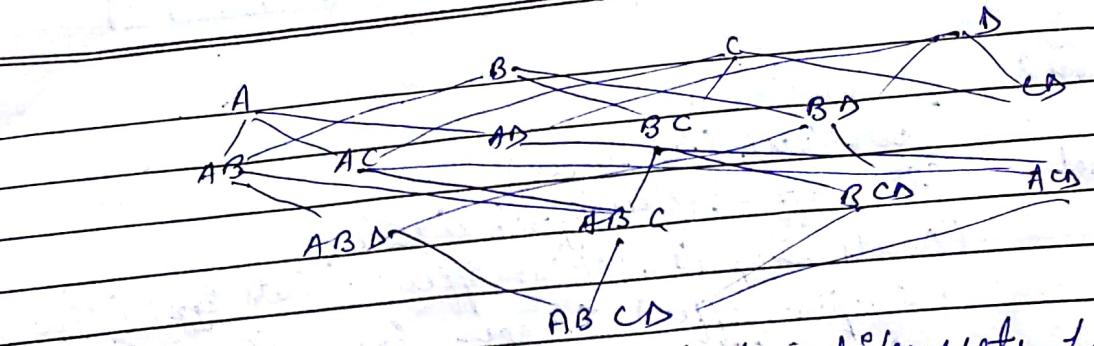
$C_{K+1}$  = Candidates generated from  $L_k$  (cartesian product)  
 for each transaction  $t$  in  $D$ .

{ increment the count of all  
 candidates that are contained in t  
 $C_{K+1}$

$L_{K+1}$  = candidates in  $C_{K+1}$  with min-support

return  $L_k$

3

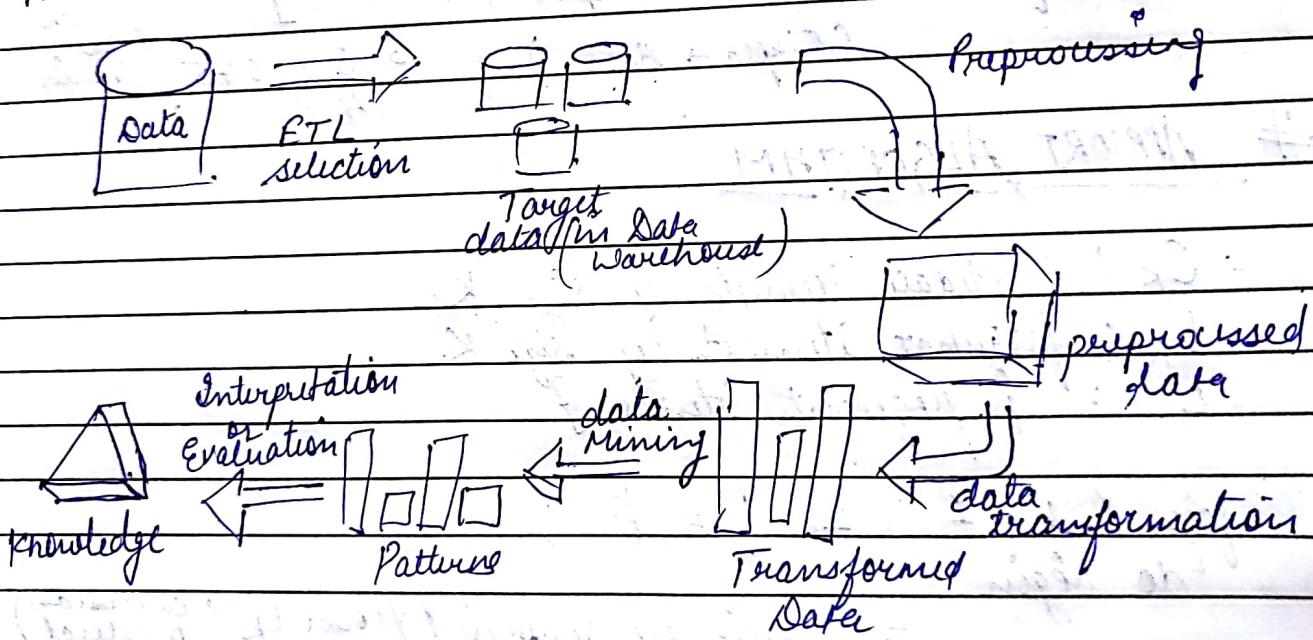


Spiritus Principle : combinations of itemsets form a lattice structure.

If an itemset doesn't satisfy minimum support threshold, then its superset will also not satisfy  $\Rightarrow$  they can be discarded.

$\equiv$  PRUNING

# KDD Knowledge discovery Data  
or knowledge discovery in DataBase



# Architecture : Data Mining System.

GUI (graphical user interface)

knowledge base  $\xrightarrow{\text{Pattern}} \text{Evaluation}$

Data Mining Engine

Data Base or Data Warehouse Server



12/9/18

Page No. classroom  
Date / classroom  
cl

## # Types of Attributes

### Example

1. Nominal : Id, number, names. (salt, short, correlation,  $\chi^2$  test, mode, entropy)
2. Ordinal : Ranking, grade, Height (median, percentiles, rank correlation, run tests, t-test, kruskall wallis test)
3. Interval : calendar date, temperature ( $\mu$ ,  $\sigma$ ,  $\text{cov}$ ,  $r$ ) (mean,  $<$ ,  $>$ ,  $\neq$  test)
4. Ratio marks in (0 to 100) temperature in (Celsius, time length) ( $<$ ,  $>$ ,  $=$ ,  $\neq$ ,  $\text{cm}, \text{mm}$ , percent variation)

## # Why Data Preprocessing :-

# Data Quality problems → (1) Noise & Outliers

(2) missing values

(3) Duplicate values.

# Data Preprocessing :- to improve data mining analysis with respect to time, cost & quality.

# (1) Aggregation (2) Sampling (3) Dimensionality Reduction.

why data Preprocessing ? → redundant data

missing data, (1) to improve data quality -

(2) after preprocessing data may become more valuable

(3) reduce computational load as it may lead unnecessary things thus we remove it

# Outliers # elements do not belong to that data.

# Aggregation → many aggregation along n dimension

~~disadvantages~~ Sales Analysis

informed by	Time	Location	Product	→ computation
may be lost	year (1)	Country (10)	Category (50)	cost reduces.
	Quarter (4)	State (20)	Sub Category (500)	higher level view of data
	month (12)	District (300)	Product (1-5000)	gives a better insight
	Day (365)	Town (3000)		→ time save