

Data - Raw and unprocessed fact.
Information - Processed and structured data.

Ex- The marks of all the students can be considered as data while the average calculated is information.

Data mining is important because not all data is relevant to us and can be discarded.

Extracting or mining knowledge from a large amount of data is called data mining.

knowledge discovery in database (KDD)

- It is a synonym of data mining.
- Database consists of a collection of inter-related data.
- Datawarehouse is used to store large amount of data.
- It is a repository of information collected from multiple sources.

Data mart - Subset of datawarehouse used at dept. level

Q. Difference between database and data warehouse

DATA WAREHOUSE is a step-by-step process:

- i) Cleaning - Removing inconsistent and noisy data.

Table A: $AB \rightarrow C$; Table B: $AB \rightarrow E$

→ Inconsistent

- Missing value - Solution:

(1) Take average of data

(2) Manually insert data (in effective)

(3) Use global variable

[Ex - unknown or -∞]

(4) Ignore the tuple (may lead to loss of important info)

ii) Integration - Combining data from multiple sources.

iii) Data Selection - where data relevant to the analysis task are retrieved from the database.

iv) Data Transformation - where data are transformed into appropriate forms for mining (representing available data according to need)

v) Data mining

vi) Pattern analysis

vii) Knowledge presentation - displaying items based on the pattern recognized in the previous step

Data Mining → Patterns → Evaluation & Knowledge Representation

Databasewarehouse

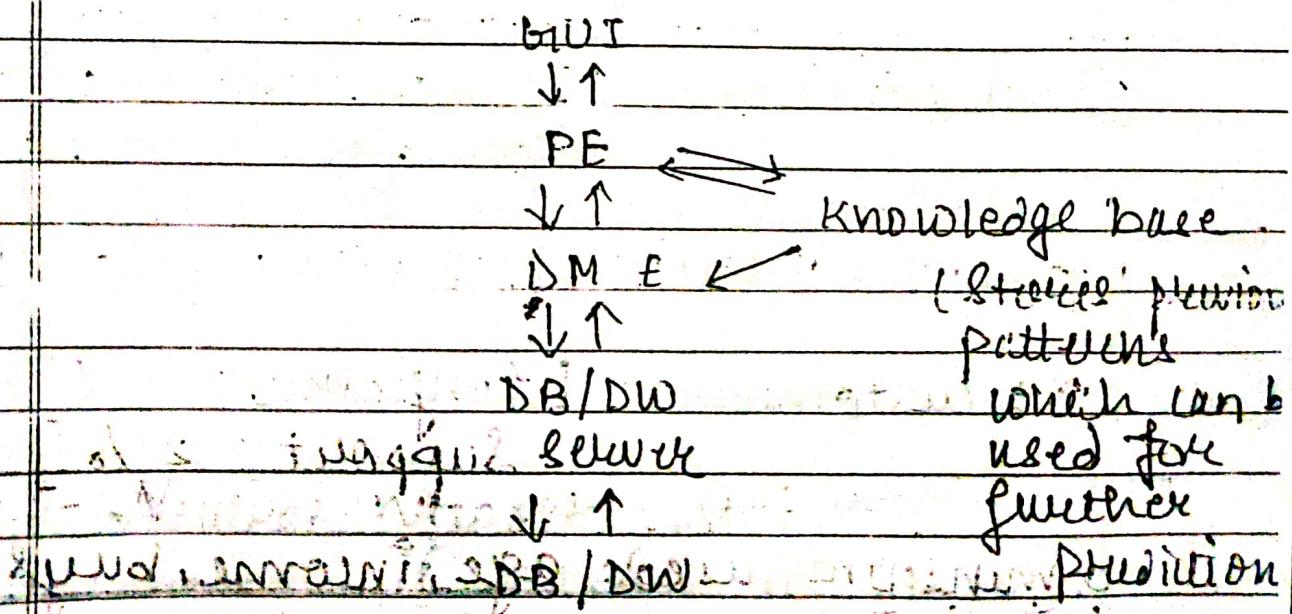
↑ cleaning and integration

+ Database / Flat Files

01/7/19

Datamining components:

- i) DB / DW (Database or Datawarehouse)
Data can be mined from either of these.
- ii) DB / DW Server
When user fires a query, this server is responsible for fetching the data from DB or DW.
- iii) Datamining Engine
It determines the technique by which data is to be mined.
- iv) Pattern Evaluation
Ex- To find the probability of an user buying certain items after buying an item.
- v) GUR



① Datamining to obtain kind of data

in i) operational (DB) ii) Data warehouse

Link -

- iii) Transactional DB. iv) Advanced DB system.
 v) Advanced DB applications. (Object Oriented DB)
 i) Spatial DB:
 stores info. relating
 to geographical
 resources.
- vi) Temporal and time-series database
 (Ex- Transactions in a bank) stroke
 To store data which changes based with
 respect to time.
- vii) Multi-media DB | Text based DB
 viii) World wide Web database. (Search history)

Data Mining Functionalities:

i) Association:

$$\text{age}(x, "20\dots 29") \wedge \text{income}(x, "20k\dots 30k") \Rightarrow \\ \text{buys}(x, "CD Player")$$

Interpretation - If a customer within the range of 20 to 29 years and income of 20k to 30k, it is probable that he/she will buy a CD Player.

$x \rightarrow \text{customer}$ [confidence = 50%]

[rule support = 2%]

History

dimensions used: age, income, buys
 (attribute)

∴ Multi-dimensional association rule
 means [50%] means that CD Player is bought and

20% entries in the database which abide by this rule].

$\Rightarrow \text{contains}(T, \text{"Computer"}) \rightarrow \text{contains}(T, \text{"S/w")}$
 $P(\text{confidence}) = 60\% \quad S = 1\%$

single dimensional association rule.

If transaction T contains comp, there is a 60% chance that it contains software.

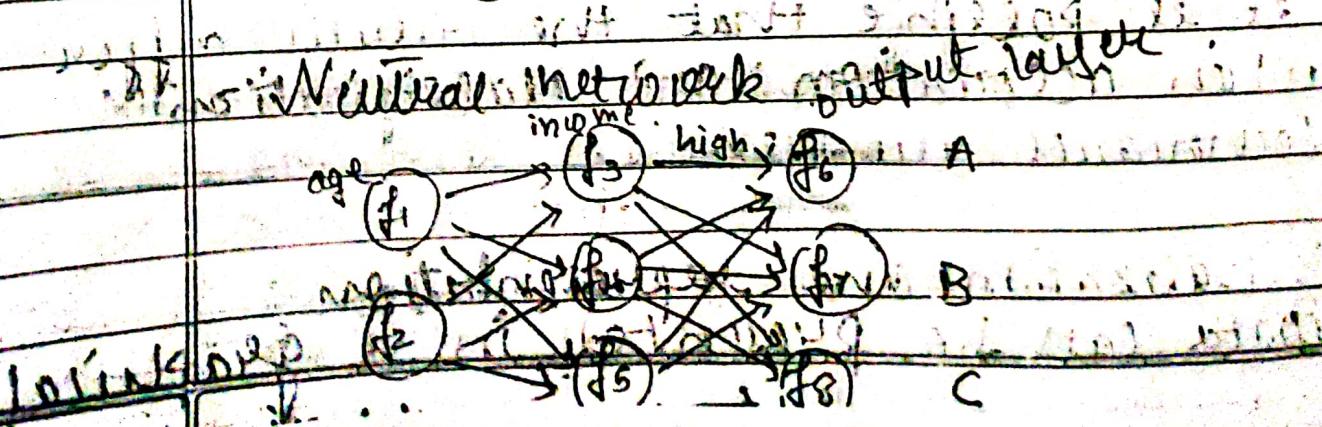
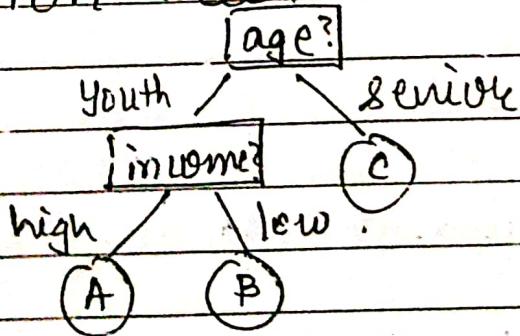
ii) Classification and Prediction

- IF-THEN rules

Ex - $\text{age}(x, \text{"youth"}) \wedge \text{income}(x, \text{"high"}) \rightarrow \text{class}(x, \text{"A"})$.
 If a person can be classified under age group of 'youth' and has a high income, then the response is put under class A.

- Decision Tree

Ex:



- iii) Cluster Analysis. (involves finding similarities & then grouping them)
- Inner-class similarity higher.
Outer-class similarity lower.
- similar to classification but, classification has certain pre-defined ^{labels} classes which are not available for clustering.
- iv) Outlier analysis
Used to find certain anomaly in the pattern. Ex - Fraud.
- v) Evolution analysis
useful for analyzing data which keeps changing with respect to time (undergo continuous evolution)

12/1/19 Parameters to check whether pattern is interesting:

- i) understandable
- ii) Valid
- iii)
- iv) Novel.

Issues of DM

- i) Based on human interaction and data mining techniques.
It is possible that the results differ when alligation and classification techniques are used.
- ii) Visualization and representation
Data can be presented in a graphical

format or tabular format.

Ex - Most sold item - table, more helpful
Profit/Loss - graph more helpful

b) iii) Pattern detection (evaluation)

Threshold value needs to be correctly determined for detecting outliers and minimize the variation to evaluate a pattern.

c) iv) Scalability and efficiency of DM technique

The size of the dataset shouldn't have an effect on the results.

The algorithm should produce quicker results in optimal time (efficient).

d) v) Parallel processing & distributed env.

Given the vastness of the dataset, the DM technique should be able to perform parallel processing in a distributed env.

e) vi) Incremental approach

When new data flows in, the DM technique should be able to pick up the process where it left or paused. Should start from scratch.

(iii) Variety: Population database

Ex - A DM tech. should be global and be able to mine data based on the parameters of a particular database. Ex - Coordinates for long distance flight.

geographical data and pixel info. for images, etc.

17/2/19

Characteristics of Data warehouse

i) Subject - Oriented

It is not transaction oriented. Rather, it is based on the specific topic for which we intend to find patterns.

ii) Integrated

It ensures that the naming / dimensions, structural conventions are same (i.e. taken care of) in all the databases used to create the warehouse.

iii) Time variant.

Databases may have a wide range of time. While they are combined to form a warehouse, a specific range is selected.

iv) Non-volatile

It deals with historic data, mostly. And thus, does not need any concurrency or reconciling measures.

OLAP & OLTP :

OLAP - Online Analytical Processing - DW

OLTP - Online Transaction Processing - DB

OLTP is high operational

OLAP

Informational

Orientations:	Transaction	Analysis
User:	DBA / Clerk / DB Professional	Knowledge workers / analysts
Storage:	Day-to-day	Long-term informational storage
Design:	ER-based designing schemes	Star or Snowflake schema
Data:	Up-to-date (current)	Historical
Dimensionality:	2-D	Multi-dimensional
Operations:	Read / Write	Read
No. of users:	In terms of thousands	In terms of hundreds
Memory size:	100 MB to GB	100 GB to TB

→ When the requirements are clear, top-down approach can be used to design the structure of data warehouse. (Waterfall model)
 Bottom-up approach is preferred when the req. are unclear or may change. (Iterative model).

Q. Advantage & disadvantage of top-down and bottom-up approach.

3-tier architecture of DW. (Diagram in TE)
 External data source is not a 'tier' in the model but it serves as a base over which the entire architecture depends.

Meta-data acts as an index for the actual data stored in DW.

Virtual Warehouse

When we represent the information residing in a DW, it can be said as a view which is termed as virtual warehouse.

Type of OLAP Server

- Relational OLAP (ROLAP)
- Multi-dimensional OLAP (MOLAP)
- Hybrid OLAP (HOLAP)

Q. Difference between ROLAP, MOLAP, HOLAP

ROLAP: Has large memory but, limited speed. Data stored in RDBMS.

HOLAP = (ROLAP + MOLAP).

MOLAP: Stores summarized data and has comparatively higher speed than ROLAP.

19/1/19

location "Chicago"

item 10151, w4 1, wellington 114

time taken for entertainment computer phone security

Q ₁	1854	882	89	623
Q ₂	943	890	64	698
Q ₃	1032	924	59	789
Q ₄	1129	992	63	877

location "New York"

item

time	HF	Comp.	Ph.	Sec.
Q ₁	1087	968	38	872
Q ₂	1130	1024	41	925
Q ₃	1034	1048	45	1002
Q ₄	1142	1091	54	989

location "Toronto"

item

time	HF	Comp.	Ph.	Sec.
Q ₁	818	746	43	591
Q ₂				682
Q ₃				728
Q ₄				784

location "Vancouver"

item

time	HF	Comp.	Ph.	Sec.
Q ₁	605	825	14	400
Q ₂				512
Q ₃				501
Q ₄				580

	C	854	882	89	623	
	NY	1027	968	38	872	
	T	818	746	43	591	
	V					
Q ₁	605	825	74	400		
Q ₂	682	952	31	512		
Q ₃	812	1023	30	501		
Q ₄	927	1038	38	580		
	HF	C	P	S		

keep same

Q₂ fix

city taken

item 825

Dice:
 location = "V" OR "T"
 time = "Q₁" OR "Q₂"
 item = "H" OR "COMP"

T	618	146
Q ₁	605	825
Q ₂	680	952

ME C

Slice: Perform select operation for single attribute.

Ex: time = "Q₁"

C	854	882	84	623
NY	1087	968	38	272
T	818	746	43	591
V	605	825	14	400

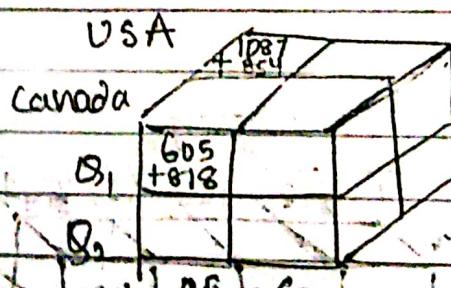
ME C P S

Pivot: Rotating the axis.

ME	605	818	1087	854
C
P
S

V NY C

Roll Up:
 (Generalization)



Country

State

City

Roll-up

Drill-down

Drill down: Q1 = Jan
 (Specification) Feb
 Mar
 Apr.

- i) Schemas for multi-dimensional DB:
 Star schema (faster than snowflake)
 (generally used to obtain fast).
 - Time:
 - Fact Table. (larger)
 The analysis carried out is stored.
 - Dimension Table.
 Attributes used for analysis are stored.

Ex: Dimension Table

{item, time, branch, location}

Time	Branch	Item
(Dimension Table)	(DT)	(DT)
- time key	- b-key	- i-key
- day	- b-name	- i-name
- day of the week	- b-type	- brand
- month		- type
- quarter		- supplier
- year		

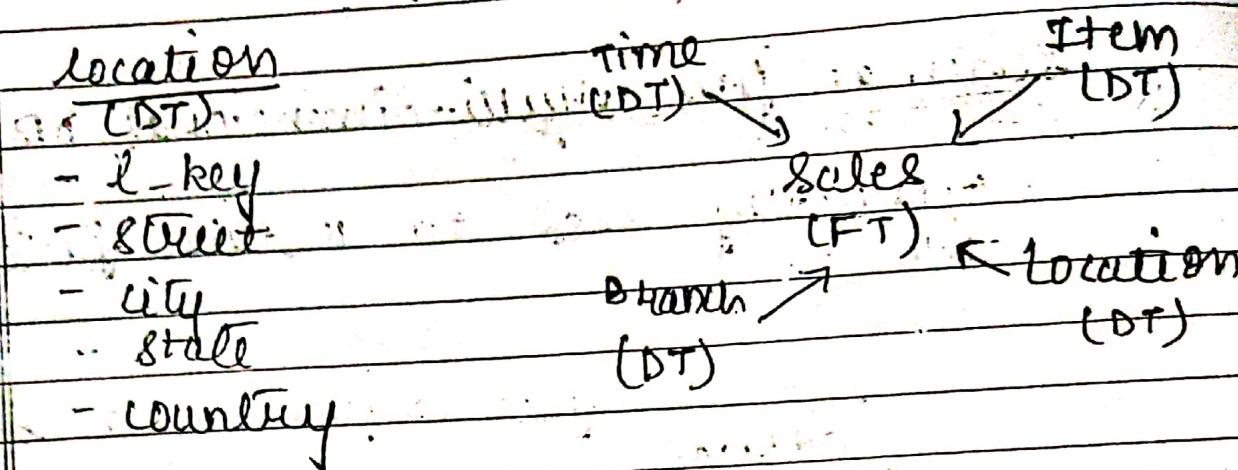
Sales (Fact Table) Every primary key stored as foreign key

- time-key

b+branch+item
 ...

- dollar sold
- sold unit

] generated item which stores the analysis of the entire DW



Dimension Table

Dimension Table — Fact Table — Dimension Table

Dimension Table

D.T are not normalized in star schema
Redundancy present

To remove this redundancy, we use the snow-flake schema. It has more than one fact table.

Course

Student — Student Enroll — Time

Data: student, course

fact: student-enrollment, course-sections

We need to find how many students

are enrolled in a particular course.

Health care:

Patient (DT)

- Patient-id
- P-name
- Address
- Age

Payer (DT)

- Payer-id
- Name
- Address
- Phone-no

Claim (FT)

- Physician-id
- Patient-id
- Service-id
- Payer-id

Physician (DT)

- Physician-id
- Physician-name
- Speciality-id

{ Foreign
key }

Service (DT)

- Claim-no
- Date-of-service
- Amount-of-charge
- Unit-of-service

{ Fact table }
table's :- service-id
attribute :-
or
measures

Disadvantage of Star Schema:

i) Data integrity.

Tables are denormalized. Hence, redundant data persist.

ii) Not flexible.

We cannot ^{easily} determine the attribute to be considered in case of multi-valued attribute.

Snowflake Schema

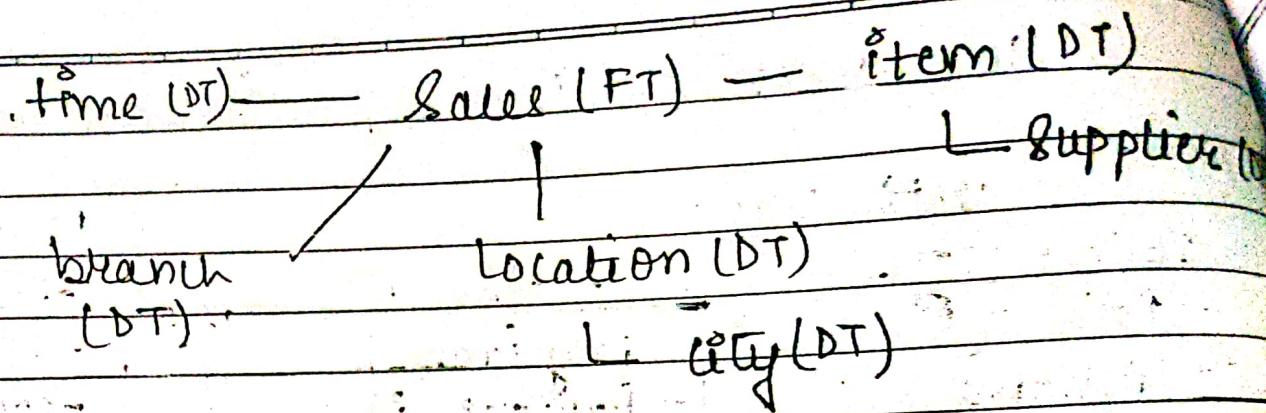
More than one fact table possible.

Every dimension table is normalized.

More no. of joins need to be performed.

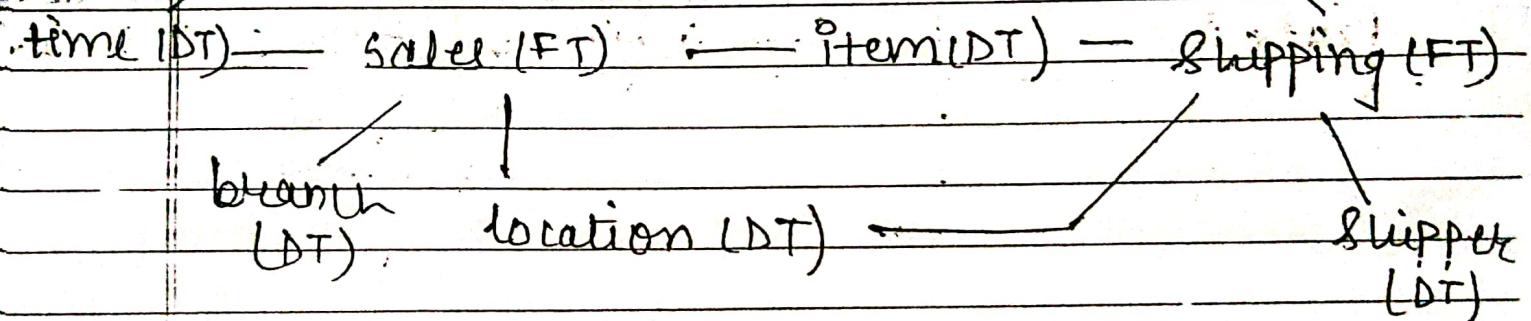
Take more query response time.

More reliable.

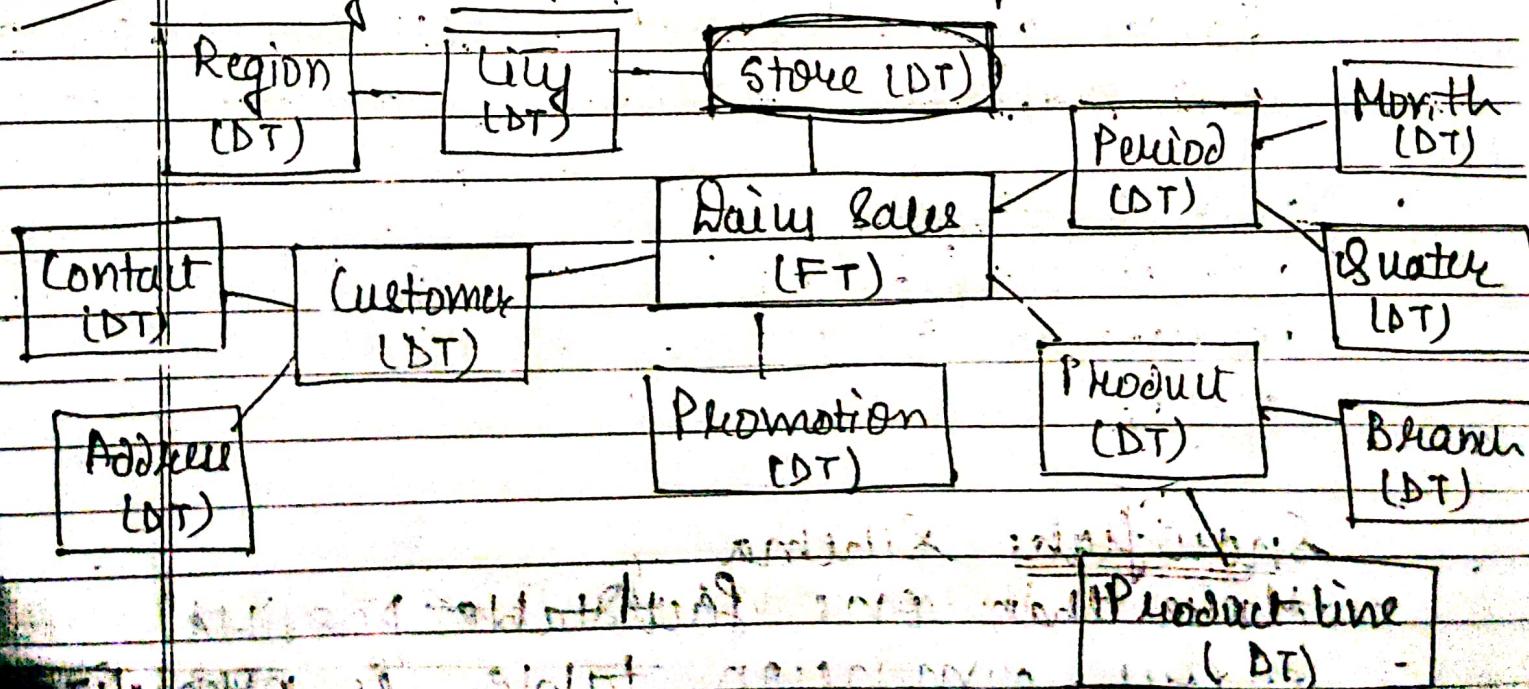


Fact Constellation Schema:

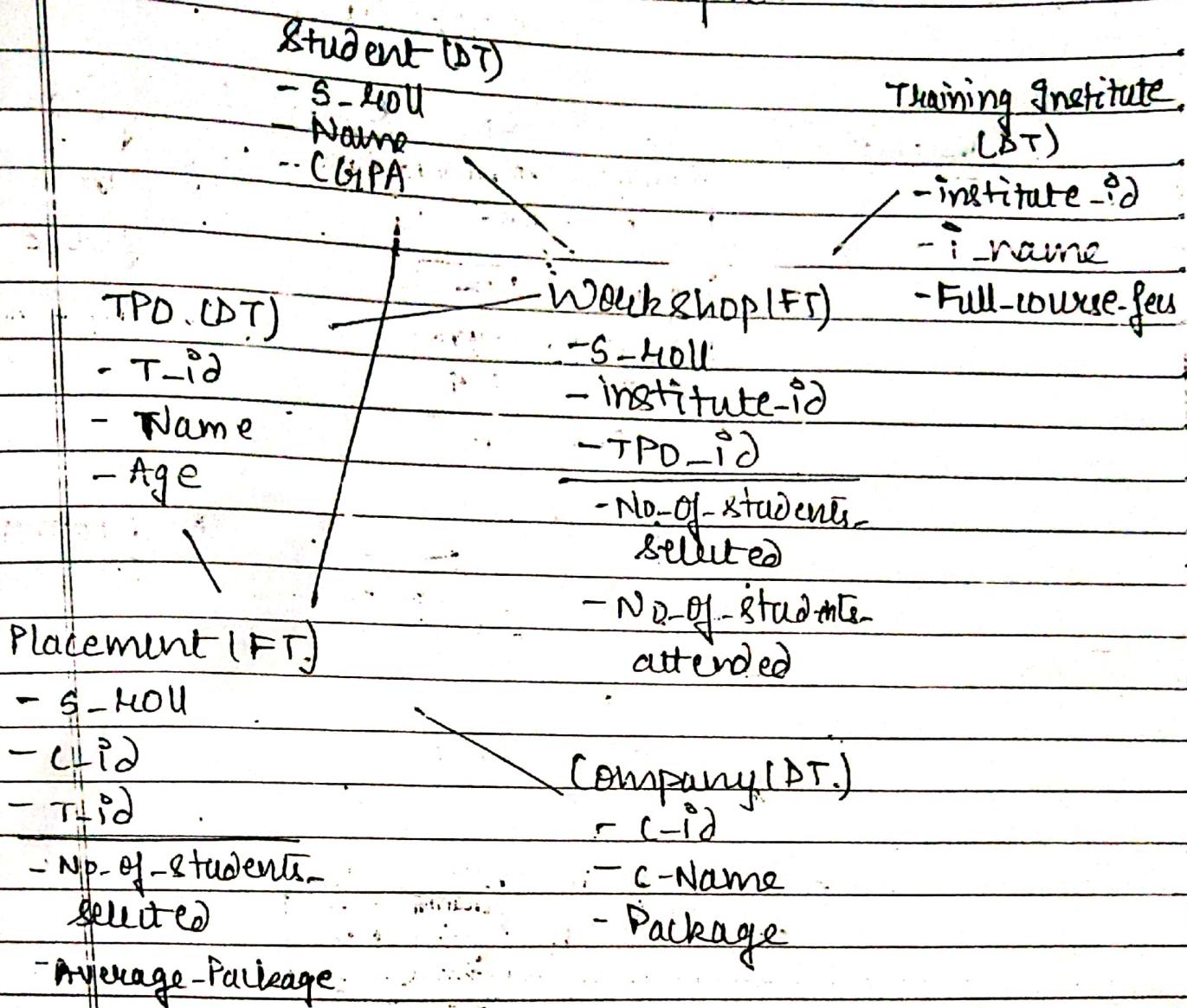
- collection of two or more star schema
- tables may exist in denormalized form.



26/7/19 Snowflake Schema Example:



Fact constellation example:



From OLAP to OLAM

↑
just summarizing
top over available
data. Not
extracting any new
info.)

finding out
patterns and
predicting future
results.

Applications of Data warehouse:

- Info. processing
 - Analytical processing
- (Querying / Mining)

iii) Data Mining (Classification / Clustering / Predicting)

Architecture of OLAM (Online Analytical Mining)

