

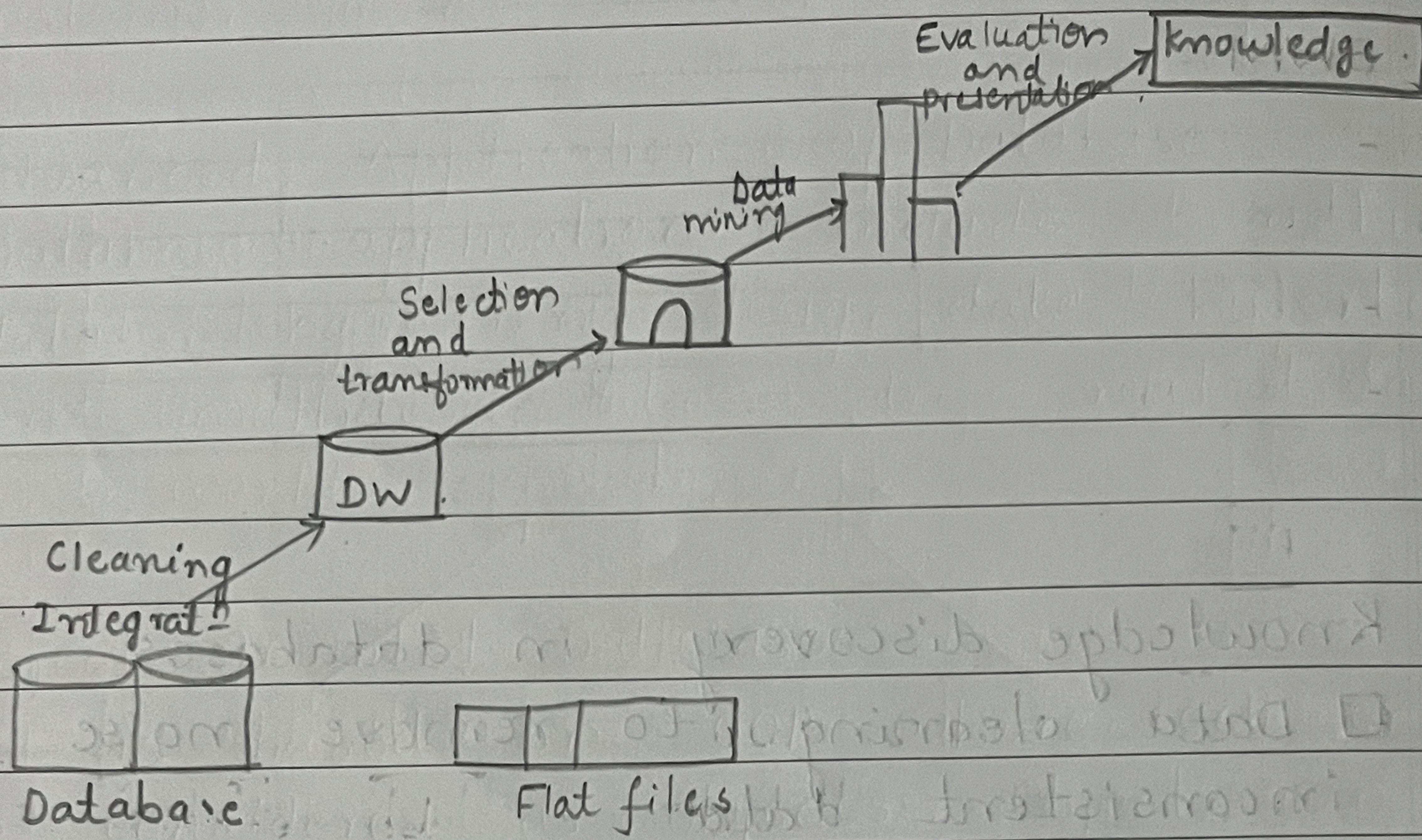
DWM:-

-
-
-

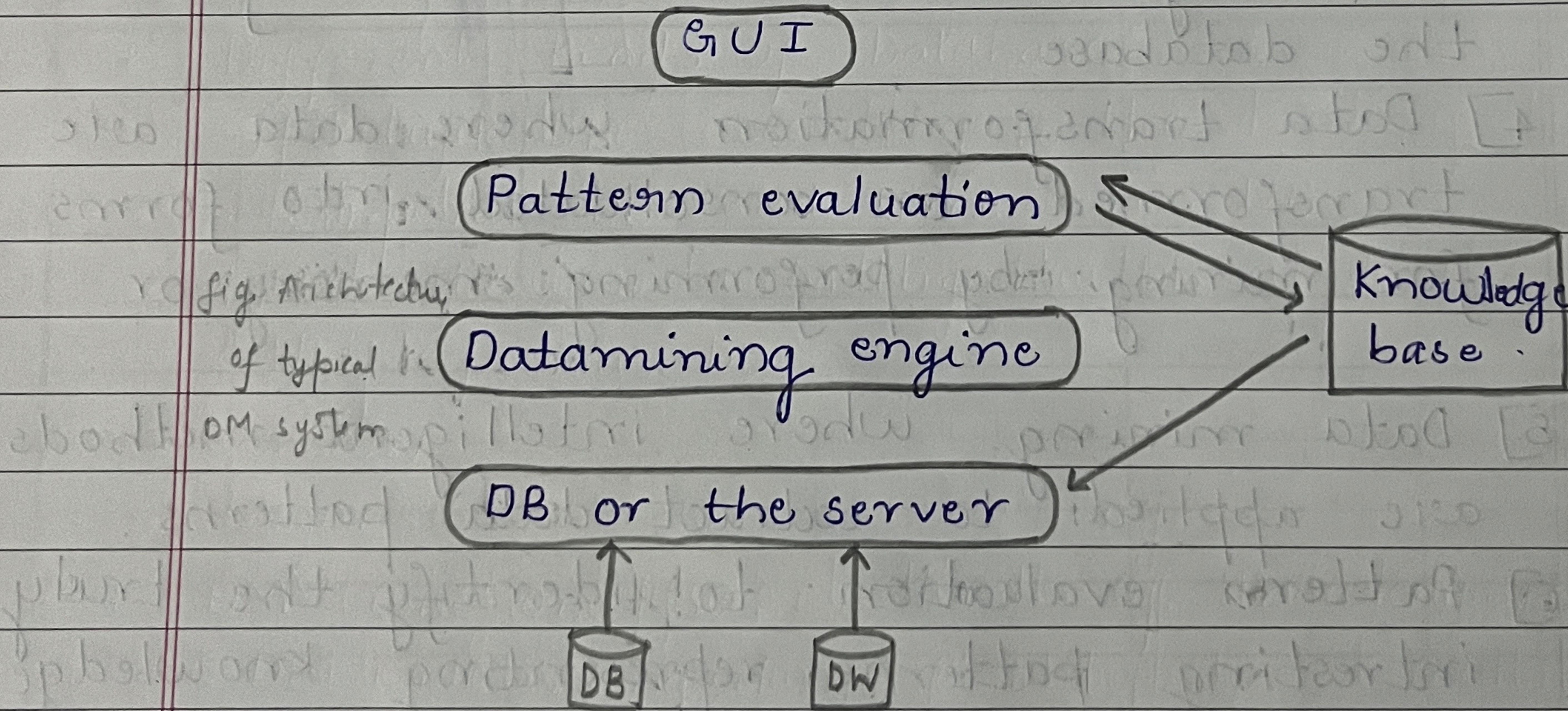
Knowledge discovery in databases.

- 1] Data cleaning to remove noise and inconsistent data.
- 2] Data integration where multiple data sources maybe combined
- 3] Data selection where data relevant to analysis task are retrieved from the database.
- 4] Data transformation where data are transformed or consolidated into forms for mining by performing memory or space.
- 5] Data mining where intelligent methods are applied to extract data patterns.
- 6] Pattern evaluation to identify the truly interesting patterns representing knowledge.
- 7] Knowledge presentation where visualization and knowledge representation techniques are used to present mind knowledge to the user.

6-8marks Process of KDD :-



* Architecture of typical DM system:-



* Types of Data:-

1) Relational Databases:-

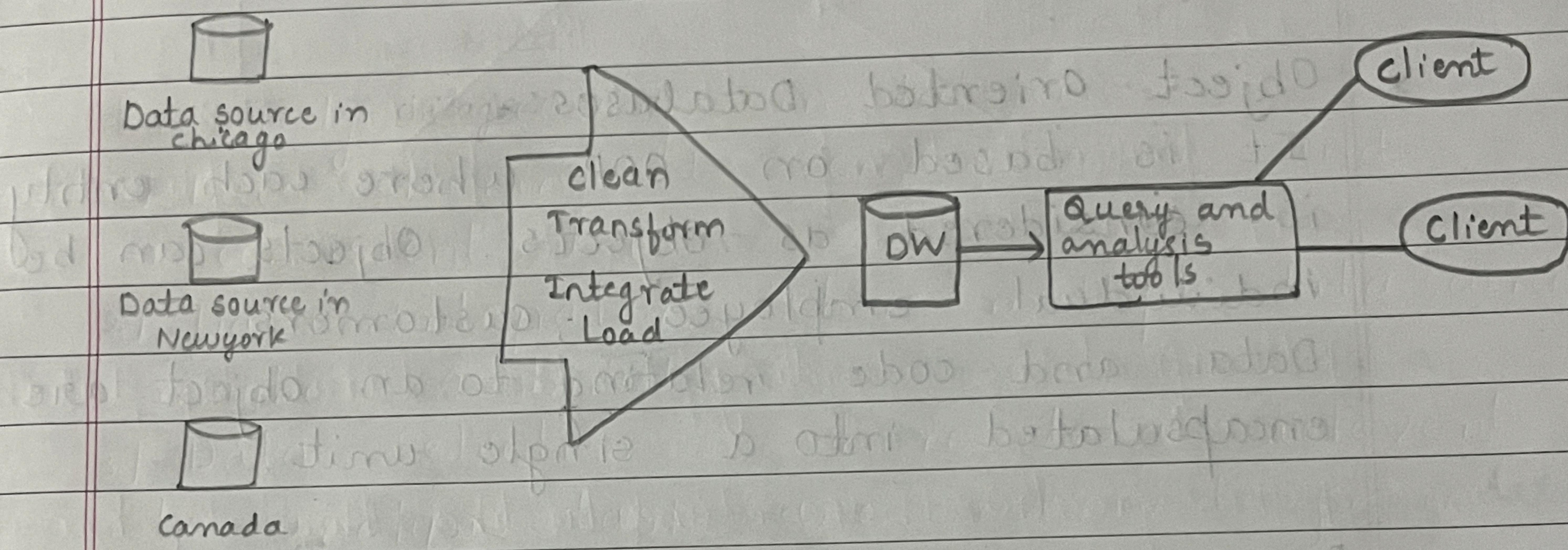
Database Management System

collection of interrelated data called as database. A relational Database is a collection of tables each of which is assigned a unique name. Each table consists of set of

attributes and stores large set of tuples.

2) Data Warehouses:

Modelled by multidimensional DB structure where each dimension corresponds to an attribute in the schema.



address(city)

Quarter
(Quarters)

Q1 605 825 14 400

Q2

Q3

Q4

Home
entertainment Computer Phone

item(types)

Fig: multidimensional data cube.

3. Transactional Databases:-

TDB consist of a file where each record represents a transaction. A transaction includes a unique transaction identity number (transaction), and list of items

Date _____
Page _____

making up the transaction, it contains other information such as date of transaction, customer id no., branch, id.no. of sales.

* Advanced Application of Database:-

Database applications include handling special data, engineering design data, hypertext and multidesign data, time related data and worldwide web.

* Object oriented Databases:-

It is based on OOP where each entity is considered as objects. Objects can be individual employee, customers.

Data and code relating to an object are encapsulated into a single unit.

* Different kinds of Database:-

- Spatial DB.
- Temporal DB and Time Series DB.
- Text DB and Multimedia DB.
- Heterogeneous DB.
- WWW.

* DM Functionalities - Kinds of Patterns can be Mined

- Data Mining Functionality.

(i) Descriptive

(ii). Predictive.

1] concept / class Description : characterizations and Discrimination

2] Association Analysis.

Association Analysis is the discovery of association rules showing attribute value

conditions that occur frequently together in a given set of data.

1] age ($x, 20 \dots 29$) \wedge income ($x, 20k \dots 29k$) \Rightarrow buys (x, CD)
 (support = 2%, confidence = 60%).

2] contains ($T, computer$) \Rightarrow contains ($T, software$)
 (support = 1+, confidence = 50%).

3] Classification and Prediction

- Distinguish data classes / concepts.
- Training data.
- (IF - THEN) rules.

Eg - All electronics company
 good response, mild response and no response.

4] Cluster Analysis

- Analyse data objects without consulting a known class label.
- Class labels are not present in training data.

Eg:-

S No.	Question	Label
1.	What is cluster analysis	Prescriptive
2.	Water content in banner.	Numeric
3.	Where GIT is located	Location

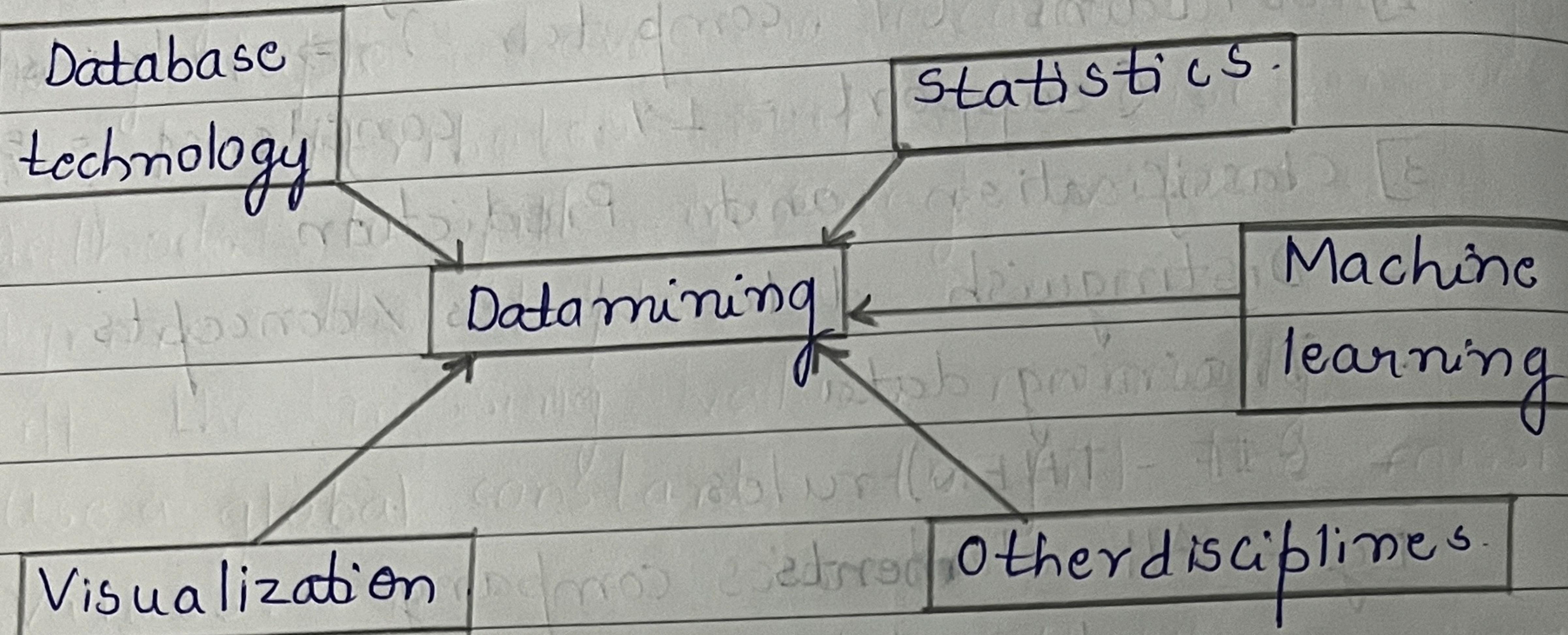
The objects are clustered or grouped based on the principle of maximizing the intra-class similarity and minimizing the inter-class similarity.

The cluster of objects are formed so that objects within a cluster have high similarity in comparison to one another and very dissimilar to objects in other clusters.

5] Outer Analysis

- Do not comply with the general behaviour

* Classification of Datamining systems:-



1] Classification according to the kinds of databases mined.

2] Classification according to the kinds of knowledge mined.

3] Classification according to the kinds of techniques utilized.

4] Classification according to the applications adapted.

* Major issues in Datamining:-

- Mining different kinds of knowledge in DB
- Interactive mining of knowledge at multiple level of abstraction.
- Data mining query languages.
- Presentation and visualisation of DM results.
- Handling noisy or incomplete data.

* Data Processing :-

- Need of Data Preprocessing
- What is Data Preprocessing?

- To improve quality of data.
- Data preprocessing techniques.

Data pre processing techniques:

Data cleaning, data integration, data transformation, Data reduction.

- * Data cleaning can be applied to remove noise and correct inconsistency in the data.
- * Data integration merges data from multiple sources into a coherent datasource such as data warehouse.
- * Data transformation normalization maybe applied to improve efficiency and accuracy of mining algorithm.
- * Data reduction can reduce the datasize, elimination and aggregating redundant features or clustering.

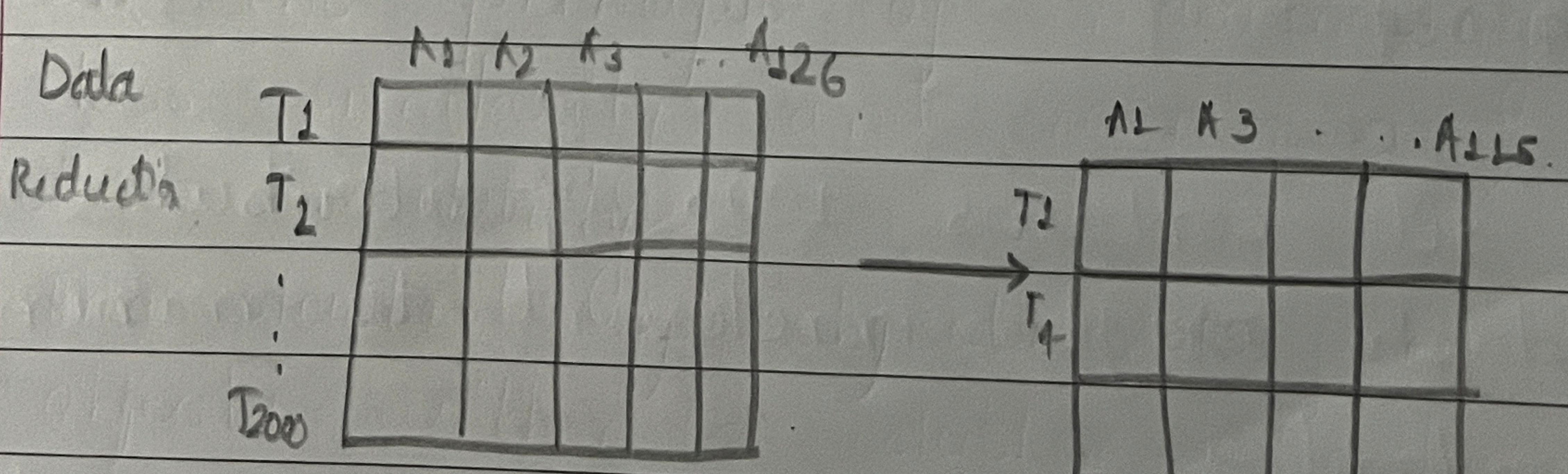
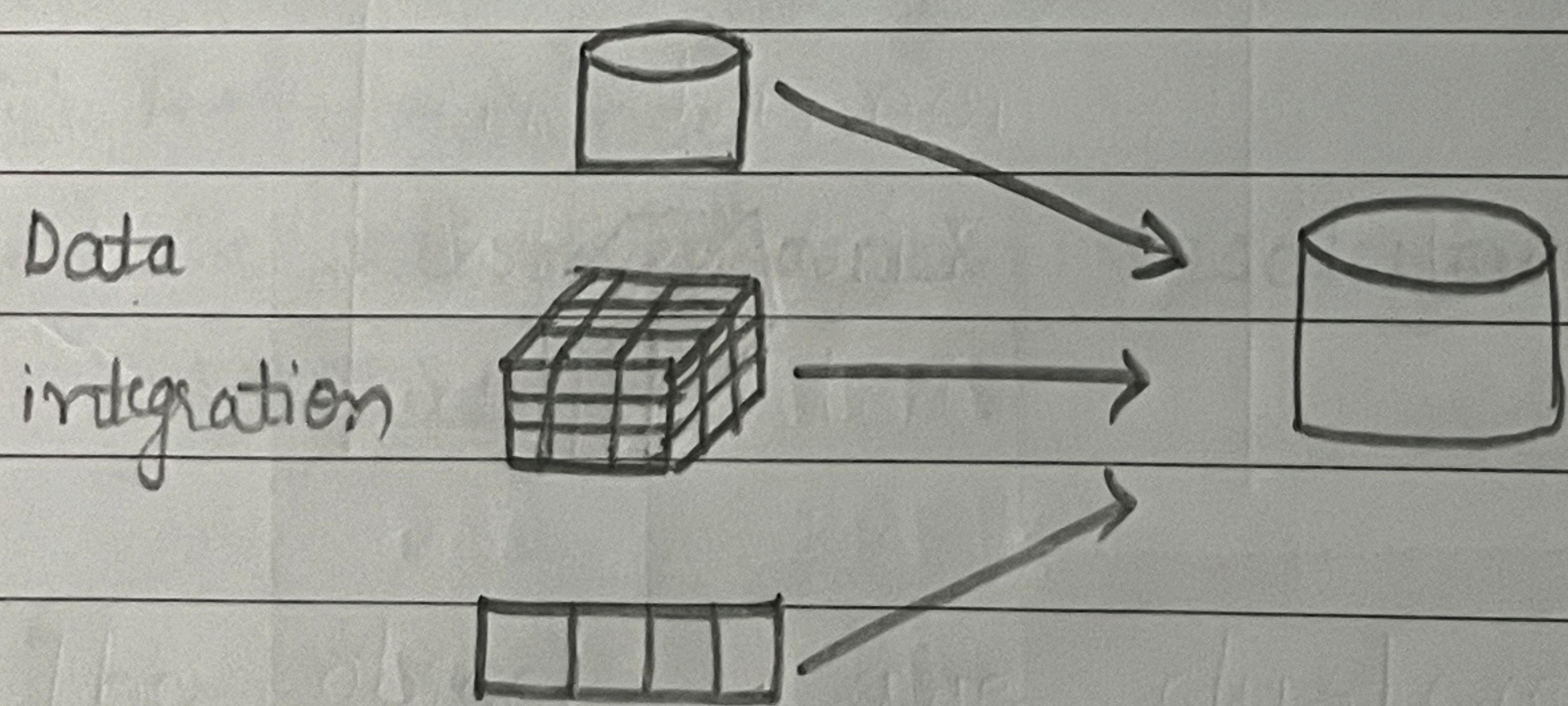
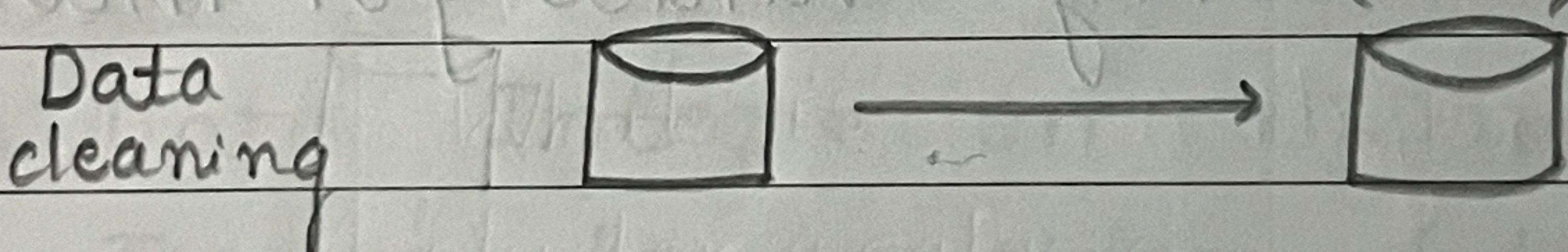


Fig:- Forms of Data preprocessing

Data cleaning:-

* It attempts to fill the missing values, smooth out noise, identifying the outliers and correcting inconsistency in the data.

(i). missing values

- Many tuples have no recorded value for several attributes such as cust_income.

- Methods to fill missing values:

a] Ignore the tuple.

b] Fill the missing values manually.

c] Use a global constant to fill the missing values.

d] Use attribute mean to fill the missing values.

e] Use the most probable value to fill the missing values.

f] Using other customer attributes in a dataset you may construct a decision tree to predict missing values for income.

Data Preprocessing:-

1. Noisy Data:-

- * It is a random error or variance in a measured variable.

2. Clustering:-

- * Outliers may be detected by constraints where similar values are organised into group clusters. values that fall outside of set of clusters considered as outliers.

3. Regression:-

- * Data can be smoothed by fitting data to a function.
- * Linear regression involves finding two lines such that to fit two variables one variable can be used to predict other. Multiple regression where more than two variables are involved.

Data Integration and Transformation:-

* Data Integration Definition:-

Issues to consider data Integration:-

- (i). Schema Integration → Entity identification problem.
- (ii). Redundancy.
- (iii). Detection and resolution of data value conflicts.

Data Trans

- * Data are transformed or consolidated into forms appropriate for mining
It involves,

(i) Smoothing - which works to remove the noise from data.

(ii) Aggregation - where summary or aggregation operations are applied to the data.

- (iii) Generalization where low level data are replaced by higher level concepts through higher level hierarchy.
 street → city or country
 age → young, middleage or senior.
- (iv) Normalization where attribute data are scaled to fall within specified range (0.2 - 1.0)
- (v). Attribute construction:- Where new attributes are constructed and added from given set of attributes.

Data Reduction:-

- * To obtain a reduced representation of the data set.
 - * Maintain integrity of original data.
- (i). Data cube aggregation
 - (ii). Dimension Reduction
 - (iii). Data compression .

Data cube Aggregation:

Year = 1999		Year = 1998		Year = 1997	
Quarter	Sales	Quarter	Sales	Year	Sales
Q1	224			1997	156
Q2	408			1998	235
Q3	350			1999	359
Q4	58				

Data set contains 100's of attributes many of which maybe irrelevant to mining tasks redundant.

Dimensionality reduction reduces the data set size by removing such attributes or dimensions

classmate
Date _____
Page _____

el data are
cepts through

ior
data are
(0.2 - 1.0)
attributess
from

of the

a.

0.0

D

1

(1)

(1)

many
ing tasks

t size

mensions

from it.

* Forward selection

Initial attribute set

[A₁, A₂, A₃, A₄, A₅, A₆]

Initial reduced set

{ }.

→ {A₁}

→ {A₁, A₄}

Reduced → {A₁, A₄, A₆}

* Backward elimination

Initial attribute set

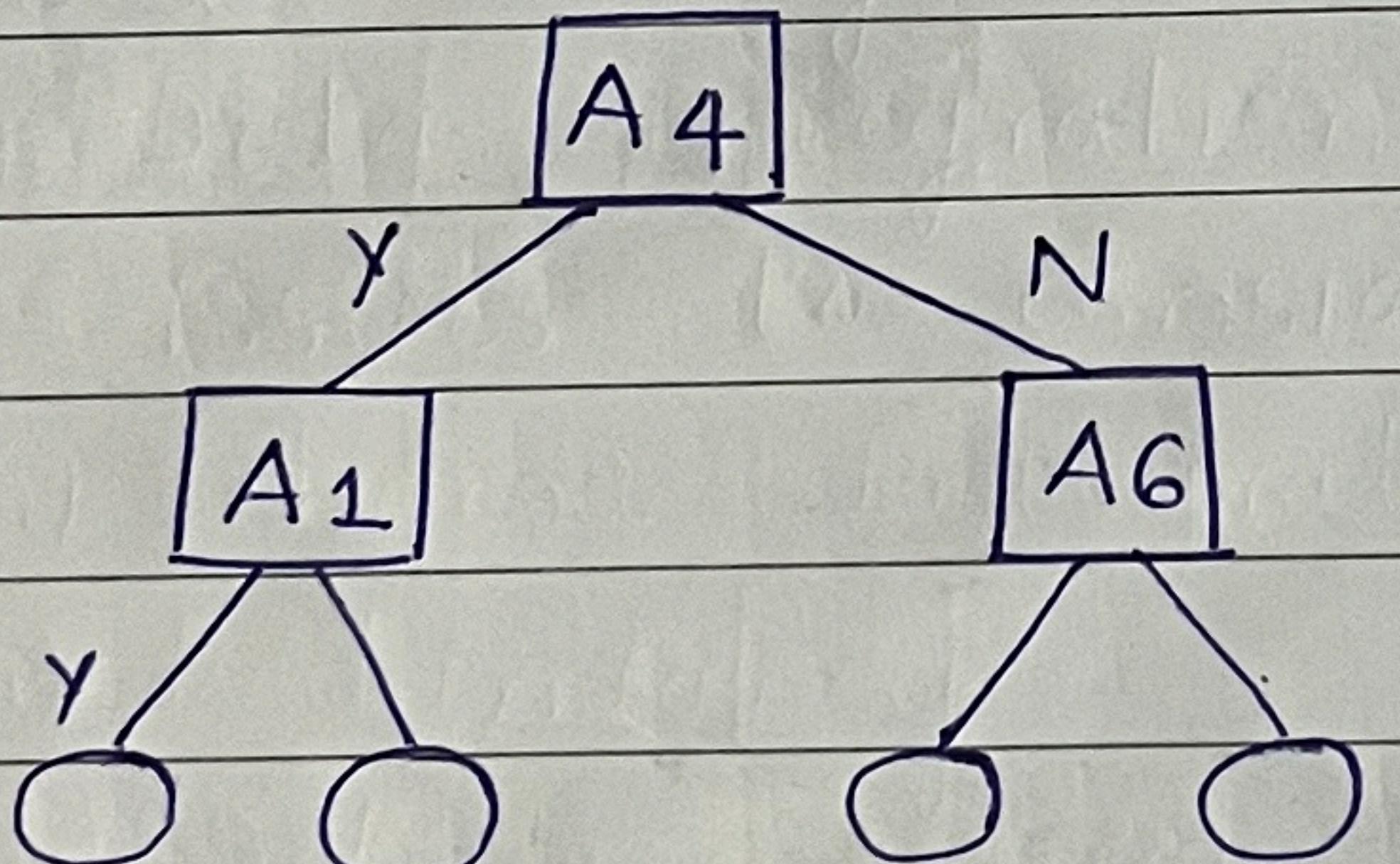
{A₁, A₂, A₃, A₄, A₅, A₆}

→ {A₁, A₃, A₄, A₅, A₆}

→ {A₁, A₄, A₅, A₆}

Reduced → {A₁, A₄, A₆}

* Decision Tree Induction



A₁, A₄, A₆.

Data compression :-

Data encoding or transformations are applied to reduce or obtain compressed representation if the original data can be reconstructed from compressed data without any loss of information is called lossless.

UNIT-II.

DATA WAREHOUSE

Provides tools for business executives to systematically organize, understand and use their data to make strategic decisions

6 marks
**

Key features of Data warehouse:

1] Subject oriented:-

It is organised around major subjects such as customer, supplier, products and sales

2] Data warehouse is constructed by integrating multiple heterogeneous sources.

3] Time variant:-

Data are stored to provide information from historical perspective [past 5-10 years]

4] Non-volatile:-

Data warehouse is a physically separate store of data transformed from application data found in the operational environment.

Difference between operational DB systems and Datawarehouse:-

- OLTP → cover day to day operations.

Eg:- Banking, inventory, purchasing.

- DW → Data analysis and Decision making.

Features	O
Characteristics	T
Orientation	C
User	Ex
Function	D
DB Design	
Access	
No. of records accessed	
No. of users	
DB size	

Features	OLTP	OLAP
Characteristic orientation	operation Processing Transaction.	Informational Processing Analysis
User	Clerk, DBA, DB professional. Ex:- DB Administrator.	Knowledge worker Ex:- Manager, executive, analyst
Function	Day-to-day operations ER Based.	Long term informational requirements
DB Design		Star / Snowflake, subject oriented
Access	Read and write	Mostly read
No. of records accessed	Tens	Millions.
No. of users	Thousands	Hundreds.
DB size	100 mb to gb.	100gb to tb.

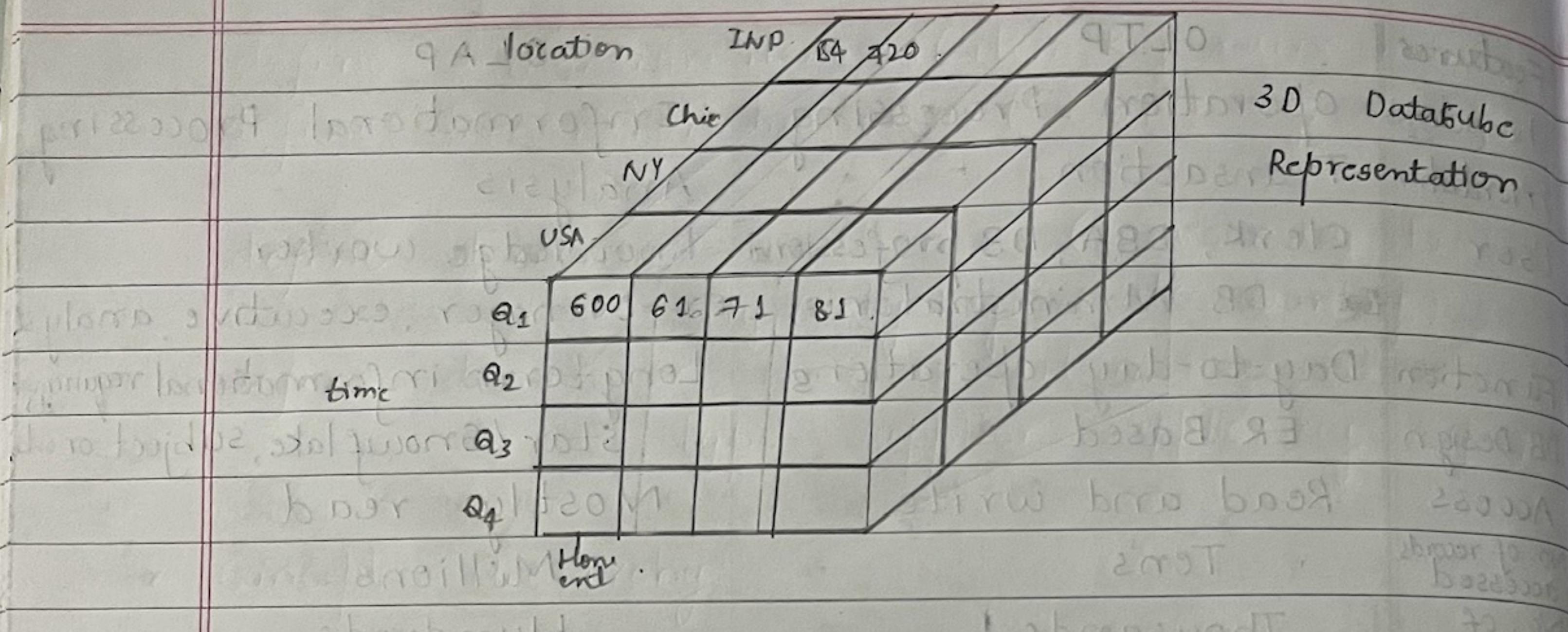
Multidimensional Data model:

- * Views data in form of data cube
- * Data cube - Data to be modelled and viewed in multiple dimensions.
- * Defined by dimensions and facts.
- Ex:- Dimensions → time, item, branch.
- Facts → Dollars sold, units sold.

9-D view of sales data:

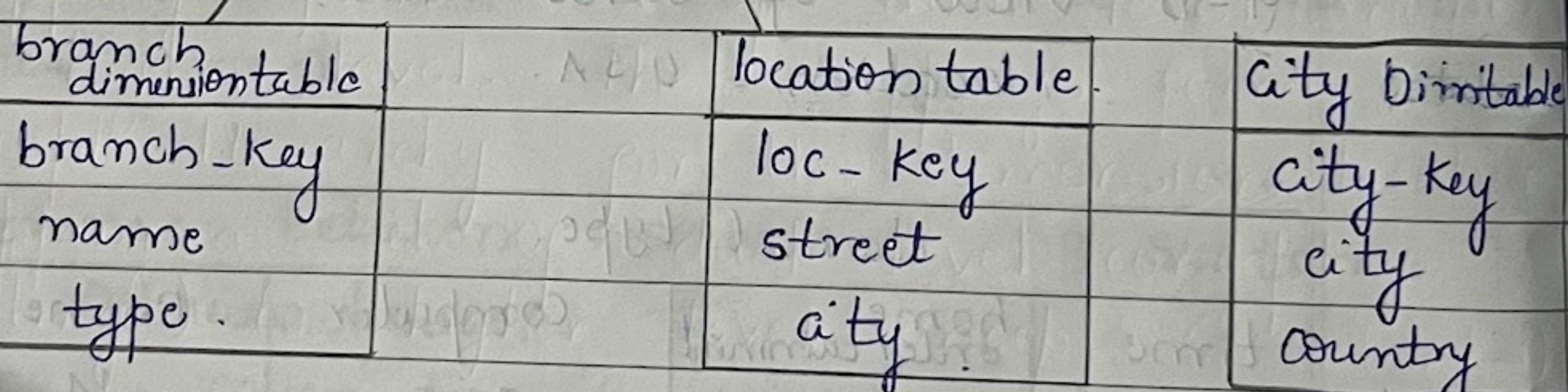
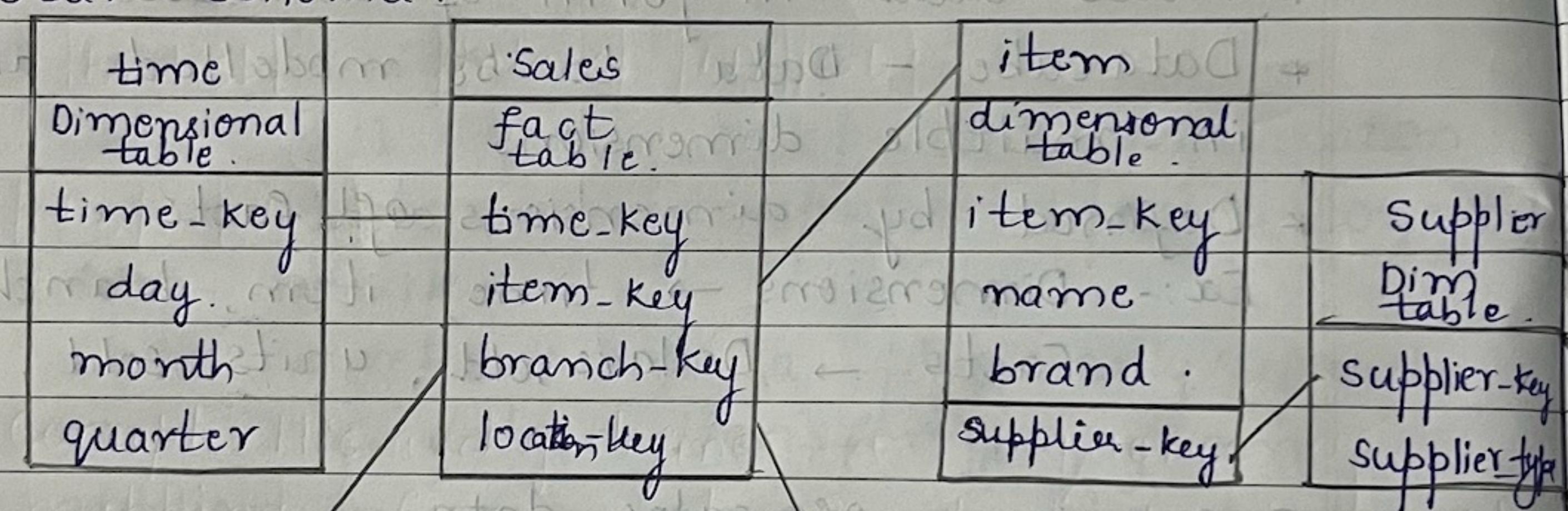
location = USA.

item (type)					
time.	home entertainment	Computer	Phone	Security	
Q1	605	825	14	400	
Q2					
Q3					
Q4					



* Star, Snowflake and fact constellations:-
schemas for multidimensional DB.

1] Star Schema



Definition:- The data warehouse contains .

- * A large central table containing bulk data with no redundancy.
- * A set of smaller attendant tables (dimension table) one for each dimension.

2] Snowflake schema :-
It is the variant of the star schema model where some dimension tables are modified. Thereby further splitting the data into additional tables.

3] Fact constellations schema :-

It requires multiple fact tables to share dimension tables. This kind of schema can be viewed as collection of stars hence it is called as galaxy schema or fact constellation schema.

Sales fact table	Shipping fact table
time-key item-key	item-key time-Key shipper-key

Examples for defining star, snowflake and fact constellation Schema:-

* SQL based DM Query language - DMQL.

* DW defined using 2 language primitives .
- One for cube definition.
- One for dimension definition .

Syntax :-

define cube <cube name><dimension list>: (measure list)
define dimension <dimension name> as (<attribute list>

Eg:- star Schema .

define cube sales_star [time , item, branch, location]
dollars_sold = sum(sales_in_dollars), units_sold

= count(*)

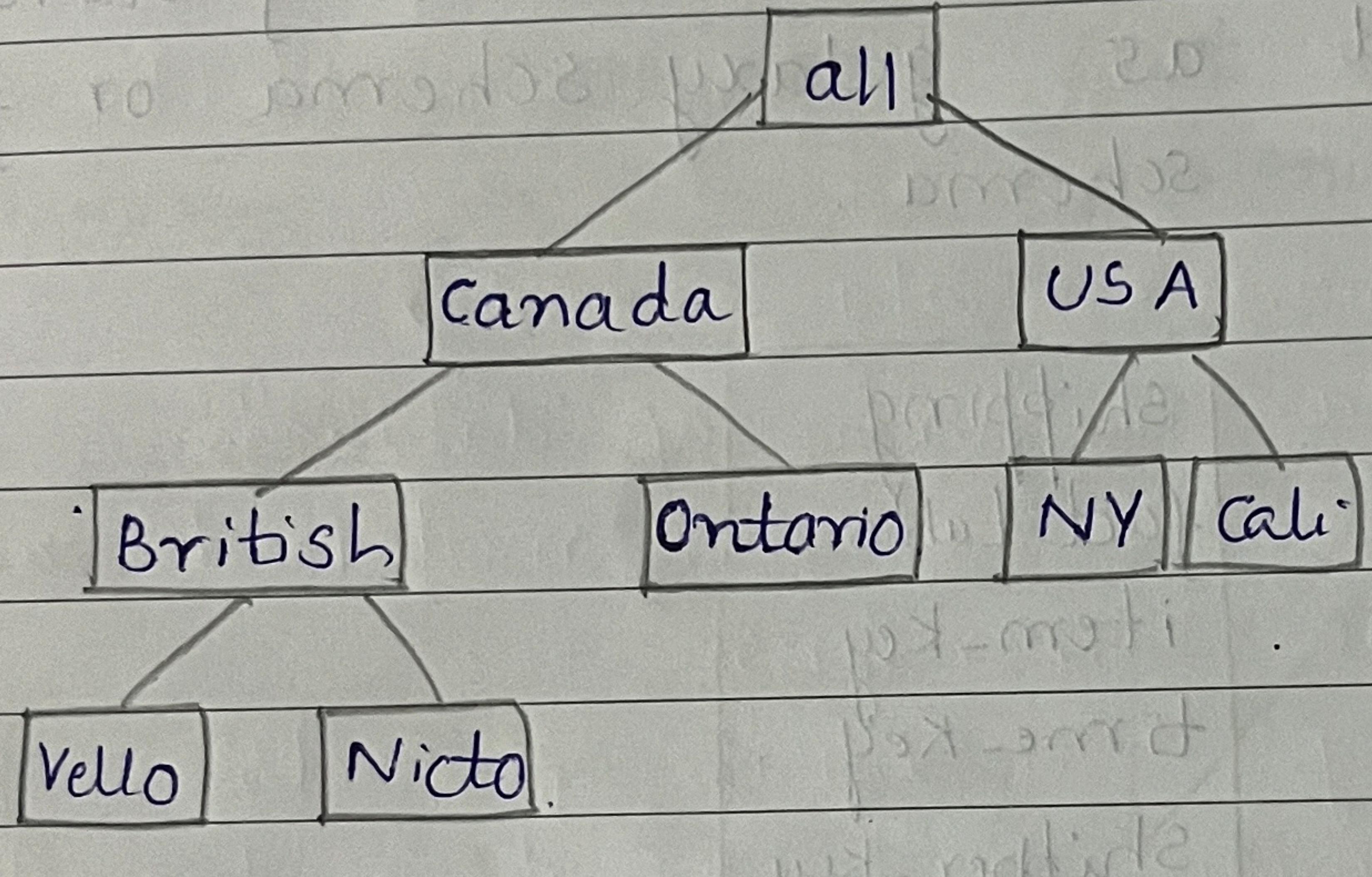
define dimension time as (timekey, day, month, year)
name, brand)

define dimension item as (item-key, name, brand)
 define dimension branch as (branch-key, name)

define dimension location as (location-key, street,
 city (city-key, state, country)).

define dimension item as (item-key, name,
 type, supplier (supplier-key, type))

* Concept Hierarchy:-



Concept Hierarchy defines a sequence of mapping from set of dual level concepts to higher level. The mappings form a concept Hierarchy for the dimension location mapping a set of dual level concepts (settings) to higher level concepts (contents).