

B, B2

IT

27\05\21

AD : ESE

Devansh Shah
1914078
~~Deshab.~~ ①

Q) 1) a

2) d

3) c

4) d

5) a

6) c

7) a

8) a

9) d

10) d

Q1) (B)

(A) OLTP

OLAP

- It stands for online transaction processing
- Used in day to day operations like purchasing, inventory, banking, etc.
- The database design is application oriented
- Usage is repetitive
- access is read/write index/hash on primary key
- data is current, up to date detailed, flat relational isolated
- metric is transaction throughput
- It stands for online analytical processing
- Used in data analysis and decision making in data warehouse system
- The database design is subject oriented
- Usage is adhoc
- access is to lots of scans
- data is historical, summarized multidimensional, integrated, consolidated
- metric is query throughput

(d) The advantages of hybrid OLAP are:-

- 1) HOLAP has better accessibility in comparison to both ROLAP and MOLAP models.
- 2) Higher processing ability when related to ROLAP and MOLAP system.
- 3) Cubes are smaller than MOLAP since only precise data is fetched for processing.
- 4) It has the volume to occupy huge data as it is a relational database in HOLAP.

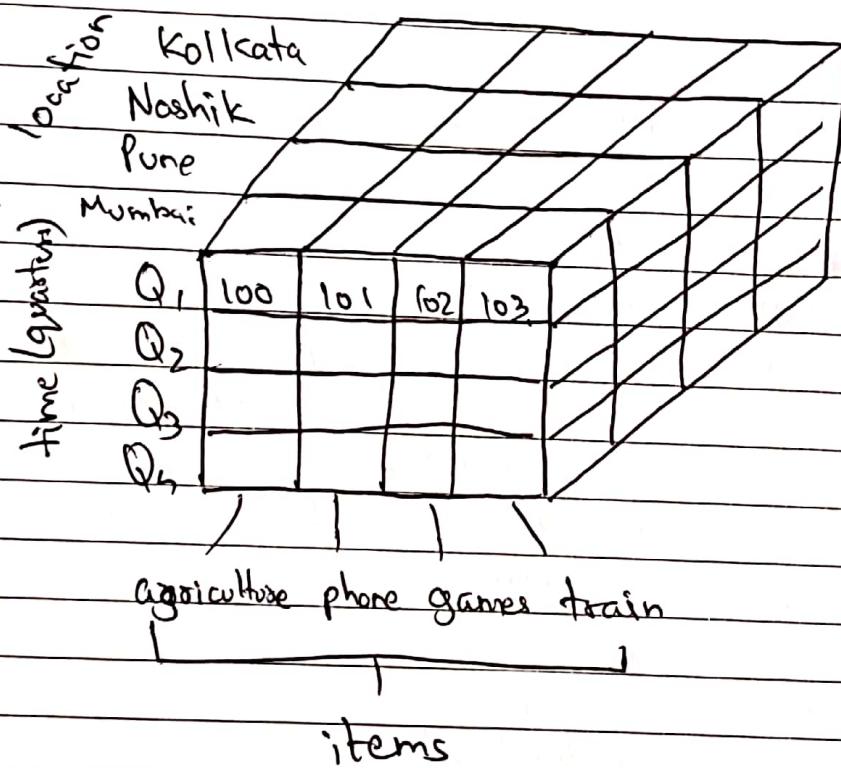
(e) The limitations of ROLAP are:-

- 1) There are many performance problems associated to the processing of complex queries that require multiple passes through relational data.
- 2) It is difficult to maintain aggregate tables in data warehouse.
- 3) Development of middleware to facilitate development of multidimensional applications is required.
- 4) Does not have complex functions that are provided by OLAP tools.

Devansh
1914078
Deshub

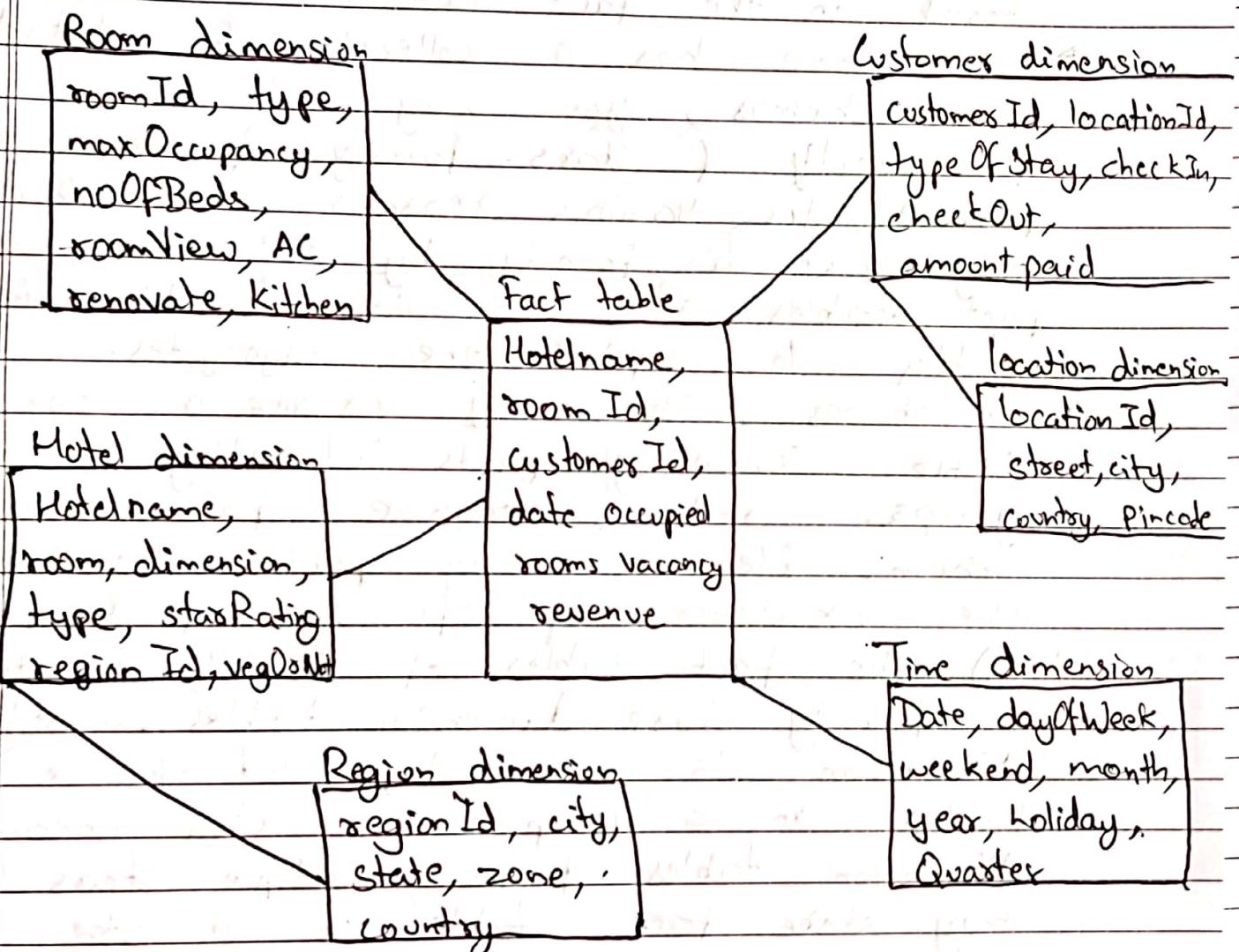
- (c) Multidimensional databases are specialized data stores that organise facts by dimensions such as geographical region, product line, sales person, time.

Below is a neat sketch of multidimensional OLAP which represents quantities of a product sold by specific retail locations during certain time periods.



- (f) Real world applications of OLAP are:-
- Data mining, marketing
 - management reporting, agriculture
 - budgeting, forecasting
 - Business reporting for sales.

- (Q2) i) Snowflake Schema for hotel management system
- The schema has one or more normalized dimensions
 - Fact table has no. of occupied rooms, no. of vacant rooms and revenue.
 - Dimensions are room, time, customer, hotel, region, location



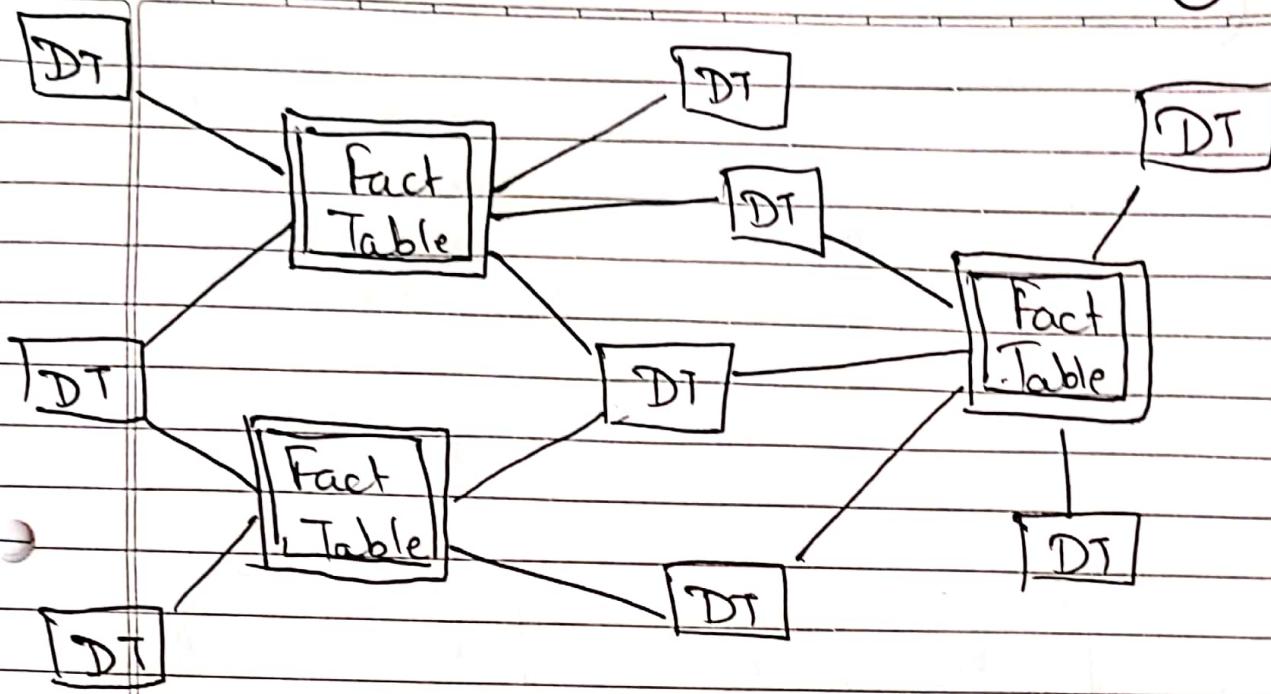
Here, region and location dimension have been normalized.

Q2) iii) When you look at single star schema with its fact table and the surrounding dimension tables, you know that it is not the extent of a data warehouse. Almost all the time, all the data warehouse contain multiple star schemas. Each star schema serves a particular purpose to track the measure stored in the fact table.

- When one has a collection of related star schemas, you may call it a collection of family of stars. Family of stars are formed for various reasons such that one may form a family by just adding aggregate fact tables and the derived dimension tables to support the aggregates.
- Sometimes, one may create a core fact table containing facts interesting to most users group and customize fact tables for them. Hence many factors lead to the existence of family of stars.
- eg) The fact tables of family of stars in a family share dimension tables, mostly time dimension is shared by most of the fact tables in a group. This also means dimension tables form multiple stars that may share fact table of a star.

Devansh
1914078
Dishub

(7)



DT : dimension table.

03) i)

Noise removal :-

- Noise is a random error or variance in measured variable. Noisy data may be due to faulty data collection instruments, data entry problems and technology limitations.
- Multiple ways to handle data are
 - i) Binning - sort data by consulting the values around it.

eg) Price = 4, 8, 15, 21, 21, 24, 25, 28, 35
 Bin a : 4, 8, 15
 Bin b : 21, 21, 24
 Bin c : 25, 28, 35

(equally partitioned on equal frequency)

- i) Clustering
- ii) Smoothing
- iii) regression

Data discretization :-

- Data discretization converts a large number of data values into smaller ones so that the data evaluation and data management becomes very easy

eg) Age (before discretization) : 10, 11, 17, 19, 38, 40, 42,
 after 70, 73, 80

After discretization.

10, 11, 17, 19 - young, 38, 40, 42 - mature, 70, 73, 78 - old

(9)

(iii) Data Smoothing :-

Data smoothing refers to a statistical approach of eliminating outliers from data sets to make the patterns more noticeable. It's achieved by manipulating data to remove or reduce any volatility or other type of noise.

e.g.) an economist can smooth out data to make seasonal adjustments for certain indicators like retail sales by reducing the variations that may occur each month like holidays or gas prices.

(iv) Min-max normalization :-

In this technique of data normalization, linear transformation is performed on the original data. Minimum and maximum value from data is fetched and each value is replaced according to the formula:

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

v' : new value

v : old value

(10)

Eg) min = 8
max = 20

marks : 8, 10, 15, 20

new Max = 16

new Min = 0

for marks = 18 $\therefore V_1 = \frac{18 - 8}{20 - 8} \times (16 - 0) = 10$

$$\text{Min Max} = \frac{V_1 - \text{min marks}}{\text{Max marks} - \text{min marks}} (\text{new max, new min})$$

$$= \frac{18 - 8}{20 - 8} \times (16 - 0)$$

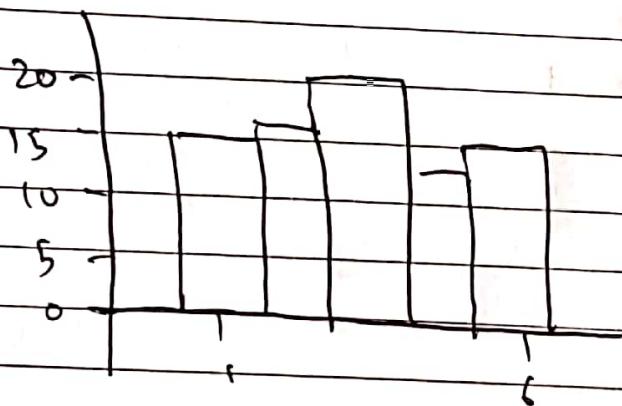
$$= \frac{0}{12} \times 16$$

$$= 0.$$

v) Z-scores

- Also known as standardized scores, they are scores or data values that have been given a common standard.
 - This standard is a mean of zero and a standard deviation of 1.
 - Z-scores are not necessarily normally distributed.
- eg) Scores of 100 people on 2 tests where max of test 1 is 6 and min = 1

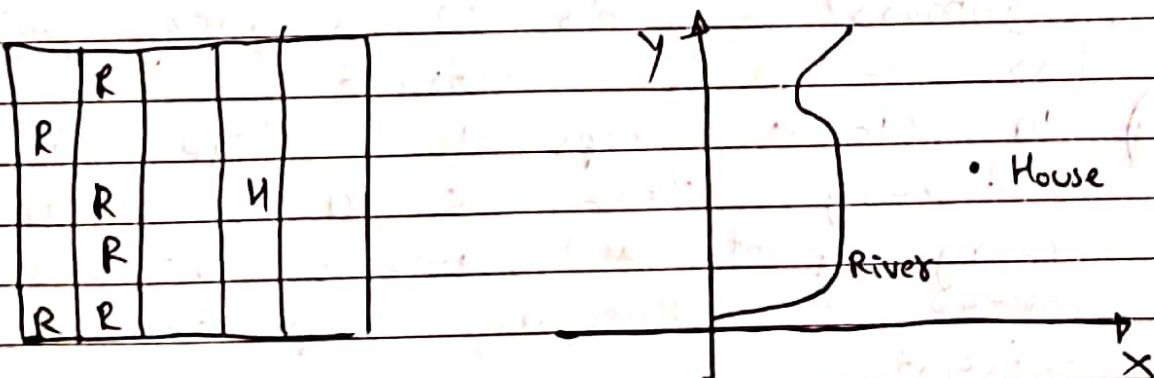
It's histogram goes like



∴ not normally distributed.

- Q5) A spatial database is a spatial type of database. It is used for storing and querying of data which are defined in a geometric space such as linear, points or other polygons.
- One can use spatial database to get the view of the real world as we see it, in the form of a database. We can either capture the entire world or a portion of it.
 - As we are visualizing the real world so all the measurements input into the database must be as complete and accurate as possible. It should be able to represent real world.
 - There are peculiar things about spatial databases as they change on the basis of what attributes of the world we capture, when we are capturing them and how much we captured. They are termed as attributes, time period and study sites.

Spatial databases are represented in 2 ways:-
 Raster model Vector model



(13)

→ Types of queries :-

i) Range queries :-

- It is used to find all desired selection in a particular range. Its base is a region. Hence also called as regional queries.
- Range queries usually need 3 things to be satisfied.

i) Object of interest

ii) Query point

iii) region

- eg) find all the rivers in a 10 km radius to Mumbai city. So, here all the rivers will be objects of interest, Mumbai will be query point, and 10km radius circle with origin at mumbai as region.

2) Nearest Neighbour :-

- The nearest neighbour query, also known as K nearest neighbours (KNN) gives K results as output. It works on the basis of, it finds out k of the nearest (neighbours) elements to a query point. It has no region constraint.

- eg) 10 closest shops to college which sell stationary, then $K=10$, college is query point and nearest distance from college to stationary shops are objects of interest.