

1) → Discrimination and Classification :-

- Difference :- Discrimination refers to the comparison of the general features of the target class object with the general features of objects from one or a set of differentiating classes.

But, classification refers to the way of finding a set of functions that depict and recognize data classes for the purpose of being able to use the model to predict the class of objects whose class label is unknown.

- Similarity: They both deal with the analysis of class data objects.

→ Characterization and Clustering :-

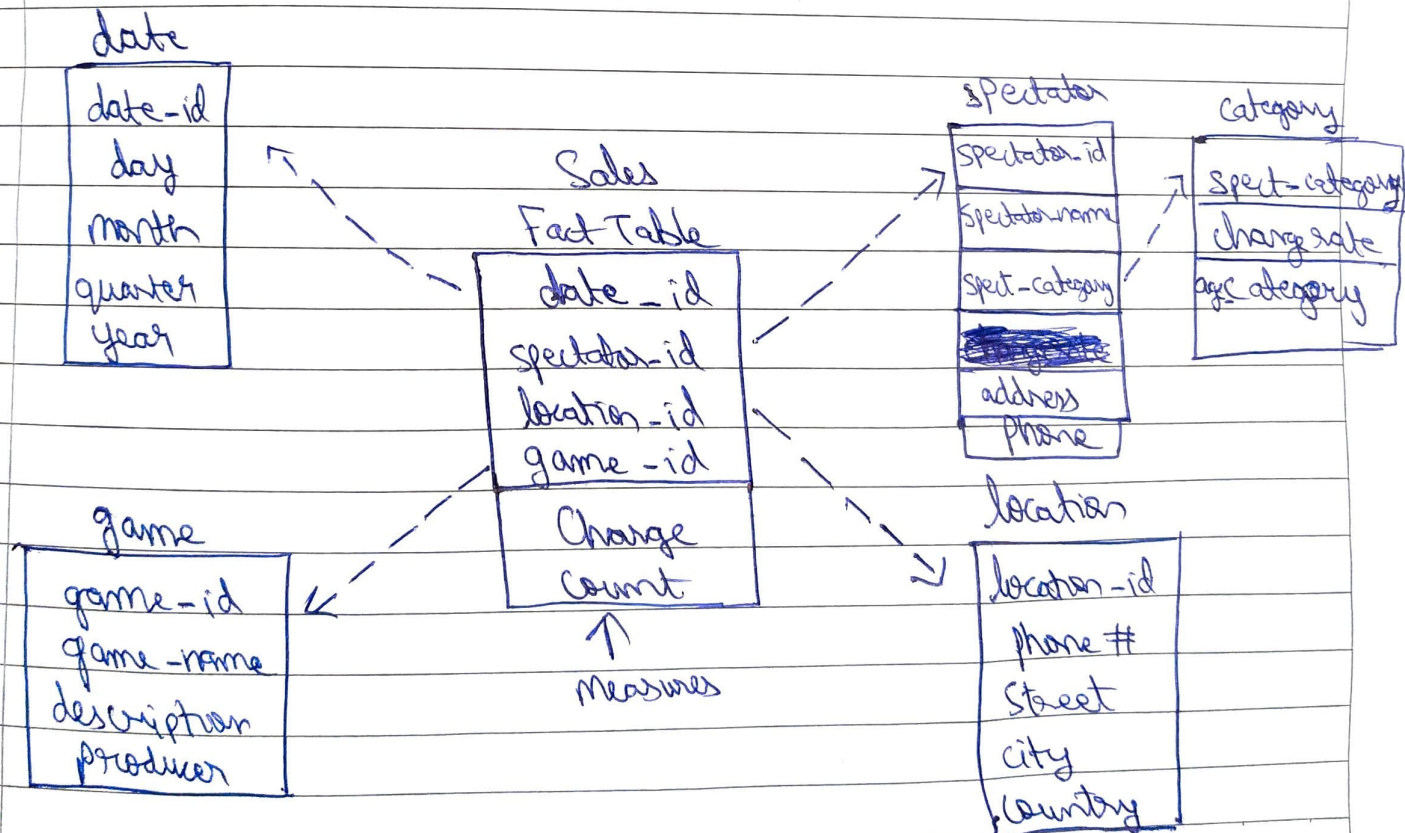
- Difference :- Characterization refers to a summarization of the general characteristics or features of a target class of data while clustering deals with the analysis of data objects without consulting a known class label.

- Similarity :- They both deal with grouping together objects of data that are related or have high similarity in comparison to one another.

→ Classification and regression:-

- Difference :- Classification predicts categorical labels which are discrete and unordered. Regression on the other hand predicts missing or unavailable and often numerical data values.
- Similarity :- This pair of tasks is similar in that they both are tools for prediction.

2.) The ~~set~~ snowflake schema diagram for the data warehouse is as follows:-



b) The DMQL for the corresponding schema would be :-

define cube sales_snowflakes [date, game, spectator, location] :

charge = ~~sum~~ charge_rate,
count = count (*)

define dimension date as (date-id, day, month, quarter, year)

define dimension game as (game-id, game-name, description, producer)

define dimension spectator as (spectator-id, spectator-name, category (spectator-category, charge-rate, age-category), address, phone)

define dimension location as (location-id, phone number, street, city, country);



3.7 Age: 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70

Number of data points = 27

a) i) Smoothing by bin means:-

	Mean
bin 1: 13, 15, 16	14.7
bin 2: 16, 19, 20	18.3
bin 3: 20, 21, 22	21
bin 4: 22, 25, 25	24
bin 5: 25, 25, 30	26.7
bin 6: 33, 33, 35	33.7
bin 7: 35, 35, 35	35
bin 8: 36, 40, 45	40.3
bin 9: 46, 52, 70	56

∴ by smoothing, we get

bin 1: 14.7, 14.7, 14.7
 bin 2: 18.3, 18.3, 18.3
 bin 3: 21, 21, 21
 bin 4: 24, 24, 24
 bin 5: 26.7, 26.7, 26.7
 bin 6: 33.7, 33.7, 33.7
 bin 7: 35, 35, 35
 bin 8: 40.3, 40.3, 40.3
 bin 9: 56, 56, 56

Answer

ii) Smoothing by bin boundaries :-
(take the value which is down)

Bin 1: 13, 16, 16

Bin 2: 16, 20, 20

Bin 3: 20, 20, 22

Bin 4: 22, 25, 25

Bin 5: 25, 25, 30

Bin 6: 33, 33, 35

Bin 7: 35, 35, 35

Bin 8: 36, 36, 45

Bin 9: 46, 46, 70

b) Min-max:- Minimum = 13
Maximum = ~~20~~ 70

$$\text{Formula:- } V' = \frac{V - \text{min}_A}{\text{max}_A - \text{min}_A} \times (\text{newmax}_A - \text{newmin}_A)$$

So, Old data: 13, 15, 16, 16, 19, 20

So, normalization:-

$$13 \rightarrow \frac{13 - 13}{57} = 0$$

$$15 \rightarrow \frac{15 - 13}{57} = \frac{2}{57} = 0.0350877$$

$$16 \rightarrow \frac{16 - 13}{57} = \frac{3}{57} = 0.0526316$$

$$16 \rightarrow \frac{16 - 13}{57} = \frac{3}{57} = 0.0526316$$

$$19 \rightarrow \frac{19-13}{57} = \cancel{0.857143} \quad 0.105263$$

$$20 \rightarrow \frac{20-13}{57} = \cancel{0.122807}$$

∴ Normalized data:-

$$0, 0.285714, 0.428571, 0.428571, 0.857143$$

~~Answer~~

Normalized data:-

$$0, 0.0350877, 0.0526316, 0.0526316, 0.105263$$

0.122807

Answer

4.) Since we have 2-dimensional data, formula is

$$D = \sqrt{(iA_1 - jA_1)^2 + (iA_2 - jA_2)^2}$$

Given new data point $x = (1.4, 1.6)$

i) Euclidean distances are:-

$$x_1 \rightarrow 0.1414 = \sqrt{(1.5-1.4)^2 + (1.7-1.6)^2}$$

$$x_2 \rightarrow 0.6708 = \sqrt{(2-1.4)^2 + (1.9-1.6)^2}$$

$$x_3 \rightarrow 0.2828 = \sqrt{(1.6-1.4)^2 + (1.8-1.6)^2}$$

$$x_4 \rightarrow 0.2236 = \sqrt{(1.2-1.4)^2 + (1.5-1.6)^2}$$

$$x_5 \rightarrow 0.6083 = \sqrt{(1.5-1.4)^2 + (1.0-1.6)^2}$$

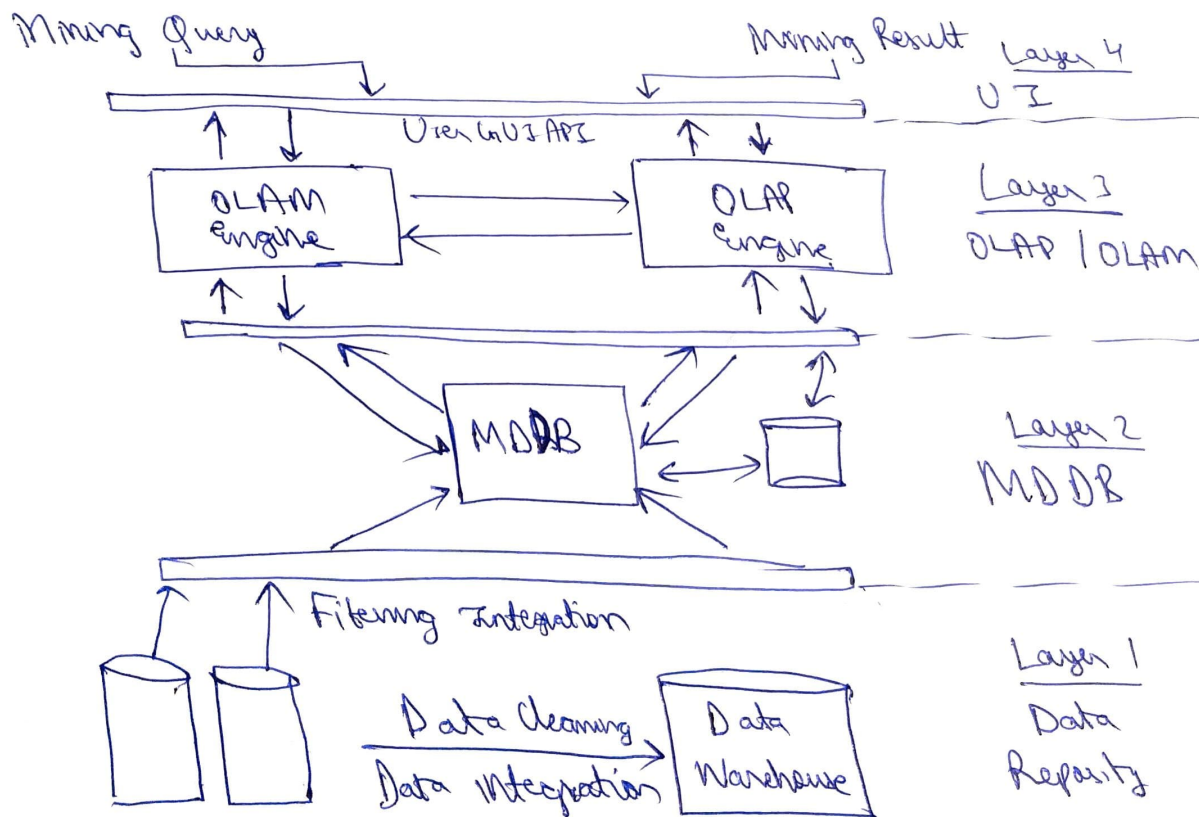
ii) formula for cosine similarity is

$$S(x, y) = \frac{x^T \cdot y}{\|x\| \|y\|}$$

So,

$X_1 \rightarrow$	0.99999	} <u>Answer</u>
$X_2 \rightarrow$	0.99575	
$X_3 \rightarrow$	0.99997	
$X_4 \rightarrow$	0.99903	
$X_5 \rightarrow$	0.96536	

S:->



- For banking and financial institutions, there is a requirement of huge data mining and data processing capabilities.
- It would use a transactional database where each record represents a transaction.
- OLAP & OLAM mechanisms accept custom online requests through a graphical user interface (GUI) application programming interface.
- We work with the data cube using cube API. data cube is created by integrating multiple databases.
- OLAP is the technology behind many business intelligence ~~(BI)~~ (BI) applications. OLAP is a powerful technology for data discovery, including capabilities for limitless report viewing, complex analytical calculations and predictive scenarios (budget, forecast).
- These technologies enable fast aggregations and calculations of underlying data sets, one can understand its usefulness in helping business leaders make better informed decisions.

- OLTP is customer oriented and used for query & transaction processing. It also enables management of current data for decision making.
- Enables short transactions with atomic consistency control.
- Hence, useful for banking and financial institutions.