

07/04/21  
Tuesday

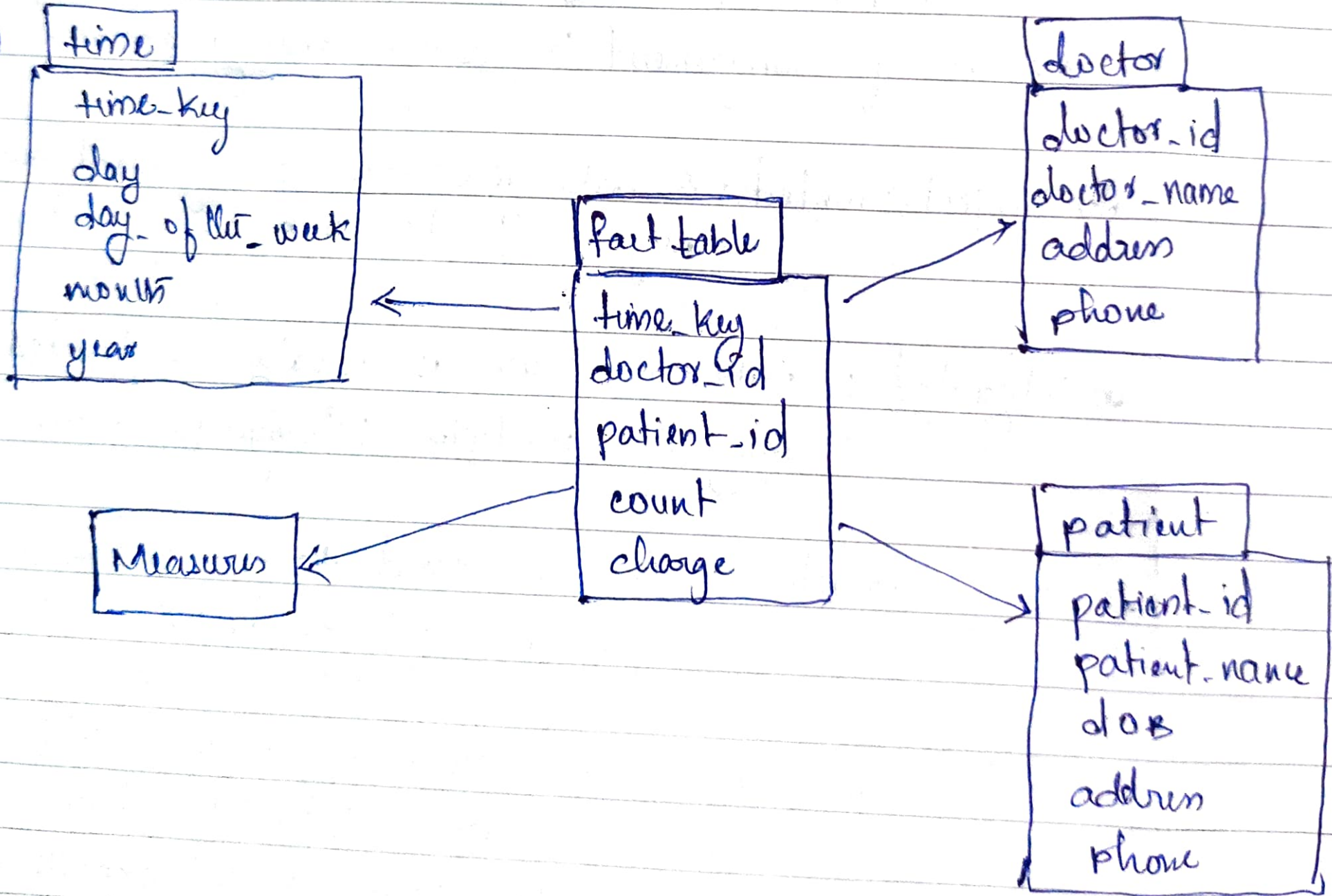
DATA MINING AND DATAWAREHOUSING  
[CS 402]

Christy Varghese  
S8 CSE A  
Roll No. 34.

- ① @ Data warehouse provides architectures and tools for business executives to systematically organize, understand & use their data to make strategic decisions. Data warehouse is a subject-oriented, time variant, non-volatile & integrated collection of data in support of management's decision making process.
- ② Subject-oriented : A data warehouse is organized around major subjects, such as : customer.
- ③ Integrated : Data warehouse is usually constructed by integrating multiple heterogeneous data.
- ④ Time-variant : Data are stored to provide information from a historical perspective.
- ⑤ Non-volatile : A data warehouse is always a physically separate store of data transformed from the application data found in the operational env.

①

②



②

### OLTP

- ① Operational data, OLTPs are the original source of data
- ② short & fast inserts & updates by end users.
- ③ A snapshot of on-going business processes.

### OLAP

- ① Consolidation data, OLAP data comes from various OLTP databases.
- ② Periodic long running batch jobs refresh the data.
- ③ Multidimensional views of various kinds of business activities



④ To control and run fundamental business tasks

④ Relatively standardized and simple queries returning relatively few records.

④ To help with planning, problem solving decision support.

④ often complex queries involving aggregations.

③ ① Stepwise forward selection : The procedure starts with an empty set of attributes as the reduced sets. The best of the original attributes is determined and added to the reduced set.

At each subsequent iteration or step, the best of the remaining original attributes is added to the set.

② Stepwise backward selection : The procedure starts with the full set of attributes.

At each step, it removes <sup>the</sup> worst attribute remaining in the set.

③ (iii) A combination of forward selection & backward elimination:

The stepwise forward selection & backward elimination methods can be combined so that, at each step, the procedure selects the best attribute & removes the worst from among the remaining attributes.

④ ⑤ \* The entropy of  $s$  is defined as:

$$\text{Entropy}(s) = \sum_{i=1}^c -p_i \log_2(p_i)$$

where  $s$  is a segment of dataset having a number of class labels.

$p_i$  is the proportion of examples in  $s$  having the  $i$ th class label.

④ ⑥ 5, 9, 11, 13, 15, 9, 8, 12, 11, 13, 18, 19, 18, 19, 13.

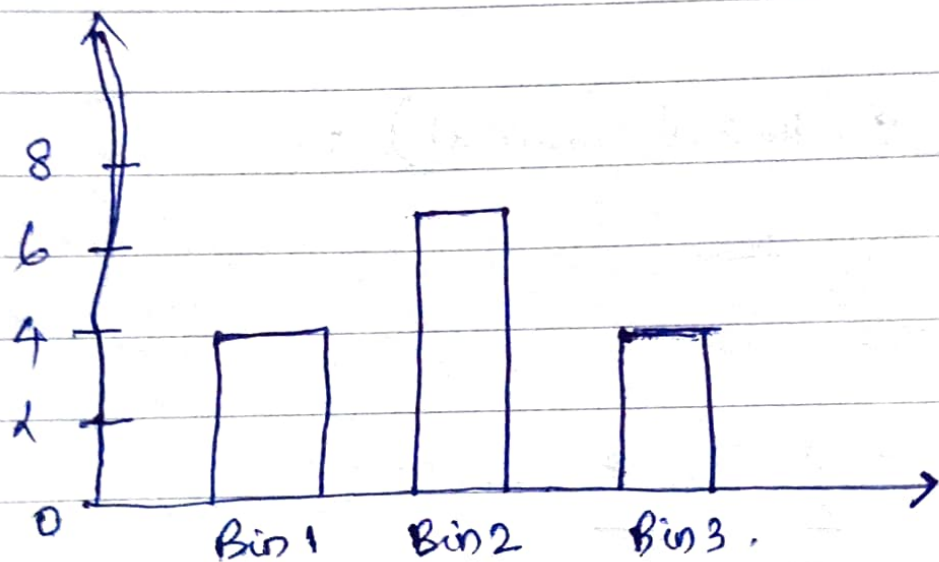
$$\text{width} = \frac{\text{max} - \text{min}}{n}.$$

$$\text{width} = \frac{19 - 5}{3} = \frac{14}{3} = 4.666 = \underline{\underline{5}}.$$

Bin 1: <sup>Range</sup> (5 to 10) : 5, 8, 9, 9

Bin 2: (10-15) : 11, 11, 12, 13, 13, 13, 15

Bin 3: (15-20) : 18, 18, 19, 19.





⑤

## OLAP (Online Analytical Processing).

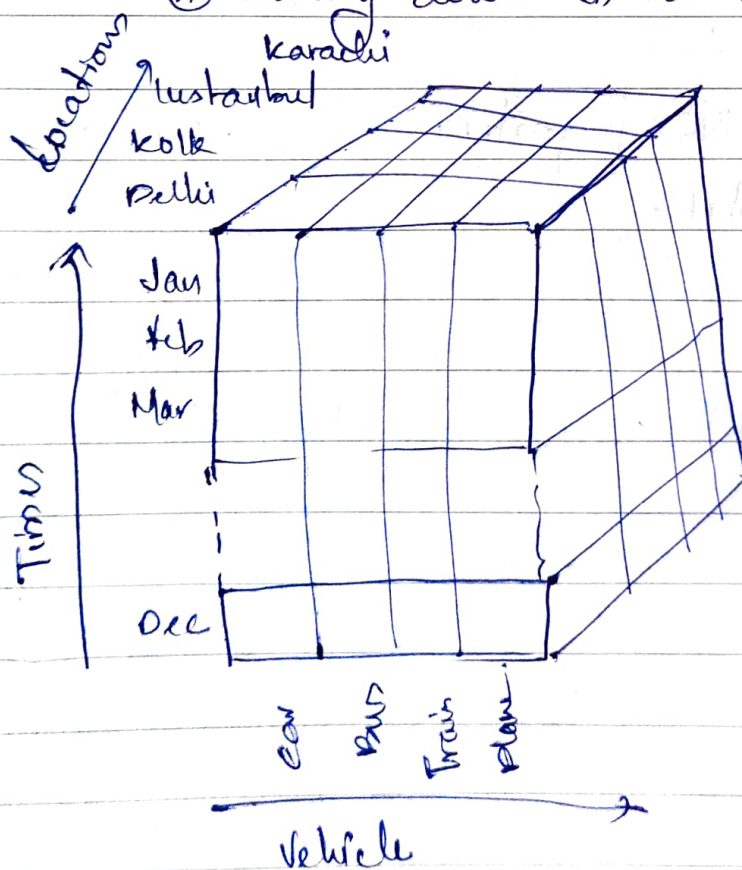
It is a slow technology that allows users to analyze information from multiple database systems at the same time.

### OLAP operations

#### (i) Drill Down

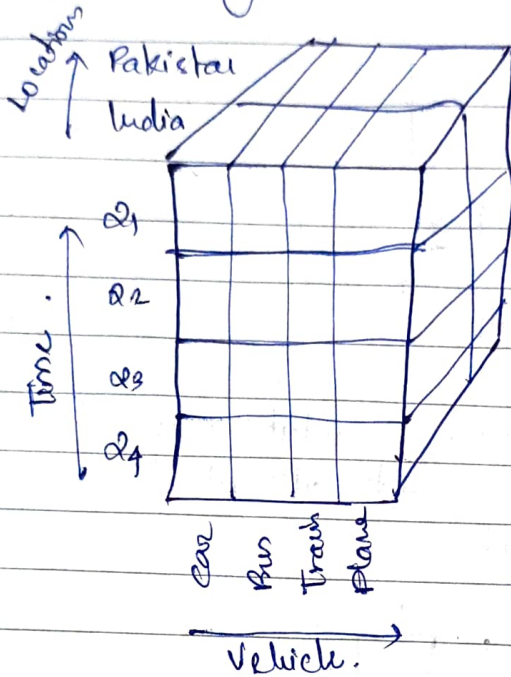
In this, the less detailed data is converted into highly detailed data.

⊕ Moving down in the concept hierarchy.



(ii) Rollup : It is just opposite of the drill down op.  
It performs aggregation on the OLAP cube.

- ⊗ Climbing up the concept hierarchy.
- ⊗ Reducing the dimensions.

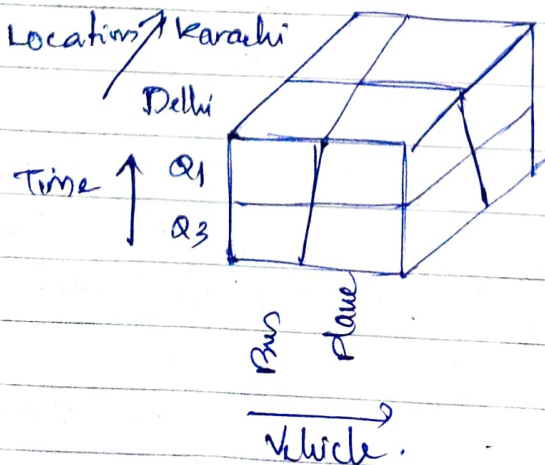


(iii) Slice : It selects a sub-cube from the OLAP cube by selecting 2 or more dimensions

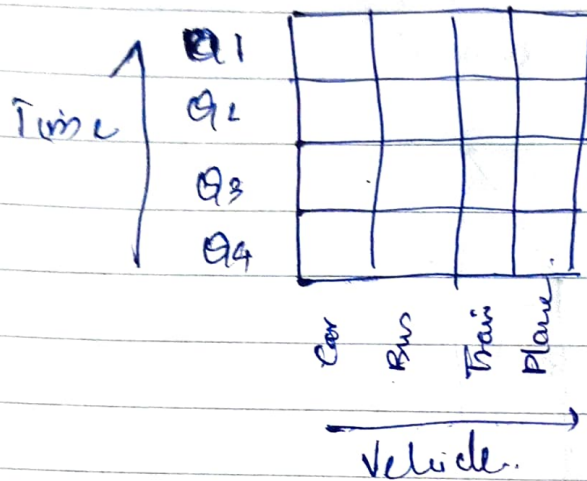
⊙ Location : Karachi, Delhi

⊙ Time : Q1, Q3

⊙ Vehicle : Bus, Plane

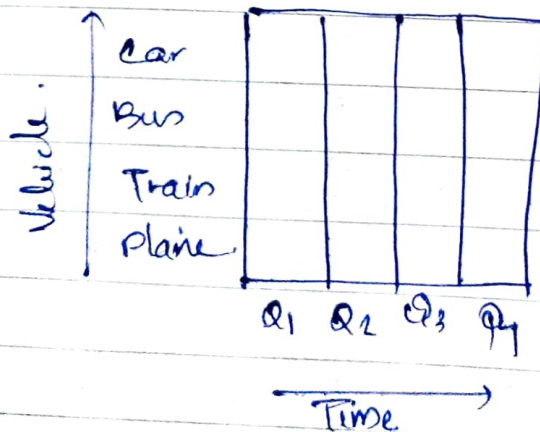


(iv) Slice: It selects a single dimension from OLAP cube which selects new sub-cube creation.



Location = Delhi.

(v) Pivot: It is also known as rotation operation as it rotates the current view to get a new view of the representation.



⑥ ① min\_max normalization.

$$\text{newMax} = 1 \quad \text{newMin} = 0$$

$$\text{data} = 200, 400, 500, 700, 1000$$

$$V' = \frac{V - \text{min}A}{\text{Max}A - \text{min}A} (\text{newMax}_A - \text{newMin}_A) + \text{newMax}_A.$$

$$V' = \frac{200 - 200}{1000 - 200} (1 - 0) + 1 = 0$$

9x



(2)

$$V' = \frac{400 - 200}{1000 - 200} (1 - 0) + 1 = \frac{200}{800} (1) = \underline{\underline{0.25}}$$

$$V' = \frac{500 - 200}{1000 - 200} (1 - 0) + 1 = \frac{300}{800} (1) = \frac{300}{800} = \underline{\underline{0.375}}$$

$$V' = \frac{700 - 200}{1000 - 200} (1) = \frac{500}{800} (1) = \underline{\underline{0.625}}$$

$$V' = \frac{1000 - 200}{1000 - 200} = \underline{\underline{1}}$$

⑥ ② z-score normalization

$$\Rightarrow V' = \frac{V_i - \bar{A}}{\sigma_A}$$

$$A' = \frac{200 + 400 + 500 + 700 + 1000}{5} = \underline{\underline{560}}$$

$$\sigma^2 = \frac{1}{5} \left[ (200^2 - 560^2) + (400^2 - 560^2) + (500^2 - 560^2) + (700^2 - 560^2) + (1000^2 - 560^2) \right]$$

$$\Rightarrow \frac{1}{5} \left[ -273600 + -153600 + -63600 + 176400 + 686400 \right]$$

$$= \frac{1}{5} [ 372000 ] = \underline{\underline{74400}} = \underline{\underline{272.764}}$$

x-score

$$200 \Rightarrow \frac{200 - 560}{272.764} = \underline{\underline{-1.319}}$$

$$400 \Rightarrow \frac{400 - 560}{272.764} = \underline{\underline{-0.795}}$$

$$500 \Rightarrow \frac{500 - 560}{272.764} = \underline{\underline{-0.298}}$$

$$700 \rightarrow \frac{700 - 560}{272.764} = \underline{\underline{.5132}}$$

$$1000 \rightarrow \frac{1000 - 560}{272.764} = \underline{\underline{2.186}}$$

$$\text{Ans} \Rightarrow -1.317, -0.795, -0.298, 0.5132, 2.186.$$

⑥ (iii) Normalization of decimal scoring.

$$\Rightarrow V' = \frac{V}{10^J} \quad J=4$$

$$\Rightarrow V' = \frac{200}{10000} = 0.02$$

$$\Rightarrow V' = \frac{400}{10000} = 0.04$$

$$\Rightarrow V' = \frac{500}{10^4} = 0.05$$

$$V' = \frac{700}{10^4} = 0.07$$

$$V' = \frac{1000}{10000} = \underline{\underline{0.1}}$$



⑦ ① I.G (S, aerial, animal)

$$A \Rightarrow 3 \quad M = 2 \quad R = 1$$

$$B = 2 \quad f = 2$$

$$\text{Entropy (S)} = -\frac{3}{10} \log_2 \frac{3}{10} - \frac{2}{10} \log_2 \frac{2}{10} - \frac{2}{10} \log_2 \frac{2}{10} - \frac{2}{10} \log_2 \frac{2}{10} - \frac{1}{10} \log_2 \frac{1}{10}$$

$$= 2.246 \Rightarrow \underline{\underline{2.25}}$$

$$\text{Entropy (Syes)} = -\frac{2}{2} \log_2 \frac{2}{2} = 0$$

$$\text{Entropy (Sno)} \Rightarrow -\frac{2}{8} \log_2 \frac{2}{8} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{2}{8} \log_2 \frac{2}{8} - \frac{3}{8} \log_2 \frac{3}{8}$$

$$\Rightarrow \underline{\underline{1.905}}$$

$$\text{I.G (S, aerial animal)} \Rightarrow 2.246 - \left(\frac{2}{10} \times 0\right) - \left(\frac{8}{10} \times 1.905\right)$$

$$\Rightarrow \underline{\underline{0.722}}$$

②

⑦ ⑪ 14 C & has legs

$$\begin{aligned} \text{Entropy (legs)} &= -\frac{2}{7} \log_2 \frac{2}{7} - \frac{3}{7} \log_2 \frac{3}{7} - \frac{2}{7} \log_2 \frac{2}{7} \\ &= 1.0327 + 0.5238 = \underline{\underline{1.5567}} \end{aligned}$$

$$\begin{aligned} \text{Entropy (SNO)} &= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \\ &= 0.8899 + 0.52832 \\ &= \underline{\underline{0.9183}} \end{aligned}$$

$$\begin{aligned} \text{res (C & has legs)} &\Rightarrow 2.246 - \left( \frac{7}{10} \times 1.5567 \right) - \left( \frac{3}{10} \times 0.9183 \right) \\ &\Rightarrow \underline{\underline{0.88082}} \end{aligned}$$