

Types of databases

Section - 2

We can apply data mining techniques to the following types of databases.

- (1) Relational databases.
- (2) Data warehouse
- (3) Transactional databases.
- (4) Spatial databases.
- (5) Time-series databases.
- (6) Text databases.
- (7) multimedia database.
- (8) world wide web.

Advanced
database
system.

Relational databases

- A relational database is a collection of tables, each having a unique name.
- Each table consists of set of attributes (columns or fields) and stores a large set of tuples (records or row).
- Each tuple represents an object identified by a unique key.
- An E-R model is used to show the relationships.
- With the help of SQL, we can write database queries to access the relational data.

Customer (Table)

Customer ID	Name	Address	Age	Income	Credit Info
-------------	------	---------	-----	--------	-------------

Item (Table)

Item ID	Name	Brand	Category	Type	Price	Player Model	Supplier	Cost
---------	------	-------	----------	------	-------	--------------	----------	------

Employee (Table)

Employee ID	Name	Category	Group	Salary	Commission
-------------	------	----------	-------	--------	------------

Branch

Branch ID	Name	Address
-----------	------	---------

Purchasing

trans-ID	cust-ID	empl-ID	date	time	method paid	Amount
----------	---------	---------	------	------	-------------	--------

items_sold

trans-ID	item-ID	qty.

works_at

empl-ID	branch-ID

Given

A query as transformed into a set of relational operations such as join, selection and projection. One is then optimized for processing.

→ A query allows extraction of subsets of the data.

→ For example, Show ~~all~~ a list of all items that were sold in the last quarter can be done

→ Relational languages also include aggregate functions such as sum, average, count, max, and min.

→ This help us to show the total sales of the last month grouped by branch or how many trans. occurred in Dec. etc.

→ When data mining is applied to relational databases, we can search for trends or data patterns.

Ex- Data mining systems can analyze customer data to predict the credit risk of new customers based on their income, age and previous credit information.

→ Data mining also detect deviations actual comparison of sales amount of this and expected sales

such as items whose sales are less from those expected in comparison to the previous year. Such deviations can further be investigated.

Data Warehouse

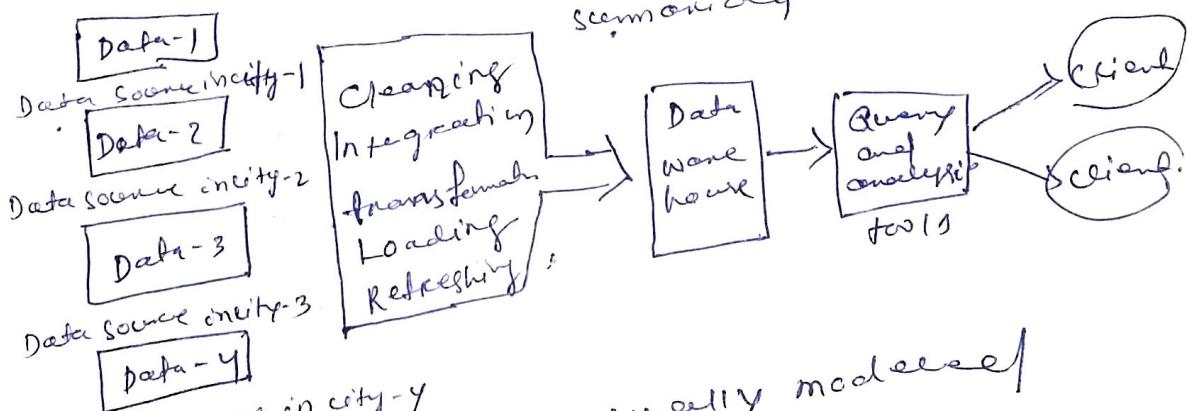
Let us say an organization is having branches around the world. Each branch has its own set of databases.

The president of the organization needs to provide an analysis of the company's sales / per item / per branch for the three quarters.

This is an difficult job, since the relevant data are spread out over several databases, physically located at different sites.

That is the reason why we are using data warehouse.

A datawarehouse → is a repository of information collected from multiple sources, stored under a unified schema, and resides at a single site
 → data stored from the past 5-10 years and are typically summarized.



- A data warehouse is usually modeled by a multidimensional database structure.
- Each dimension corresponds to an attribute or a set of attributes in the schema and each cell stores the value of some aggregate measure e.g. count or sales amount.

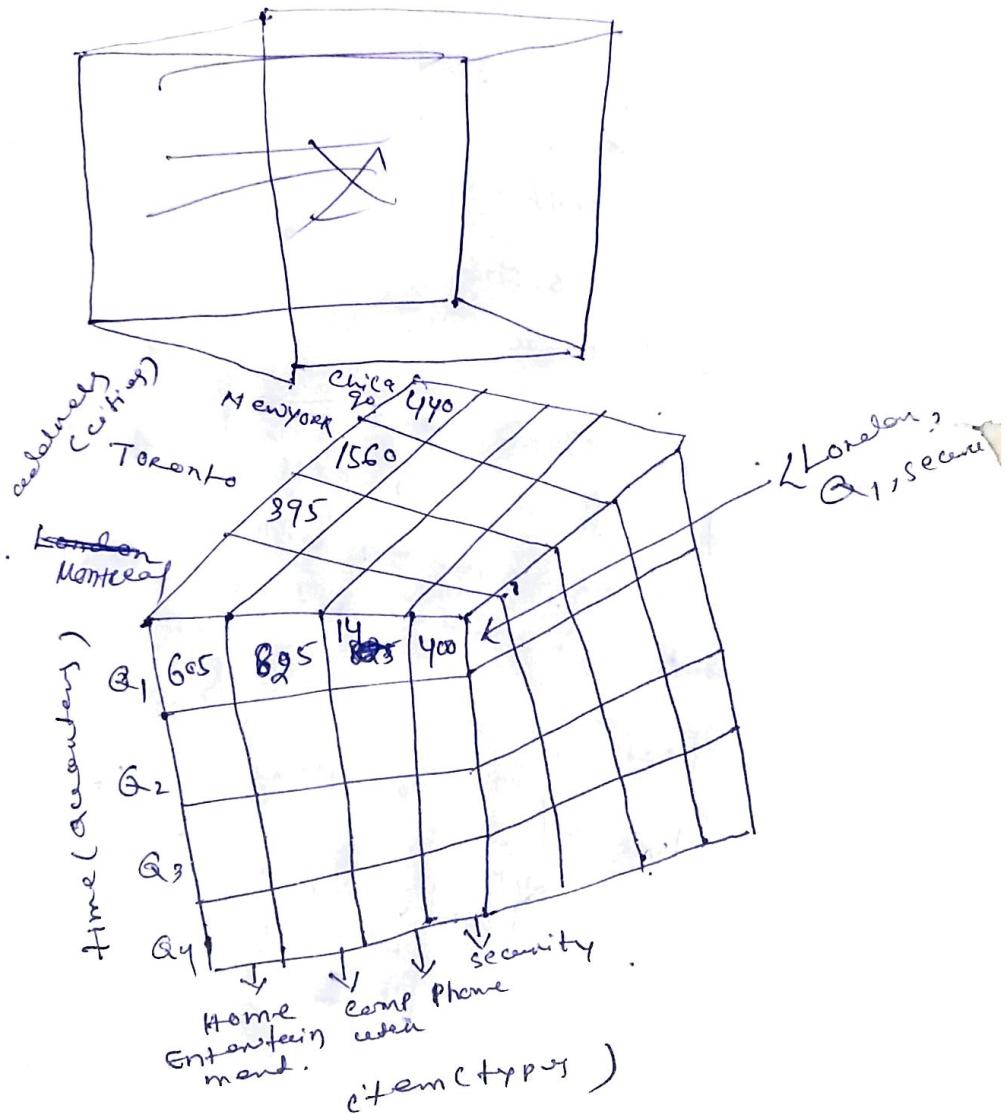
→ The actual physical structure of a data warehouse is a multidimensional data cube.

→ A data cube gives a multidimensional view of data and allows the pre-computation and fast accessing of summarized data.

Data warehouse →
collects information about subjects that cover an entire organization.

Data mart →
is a dept. subset of a datawarehouse.
it focuses on selected subjects and its scope is dept.-wide.

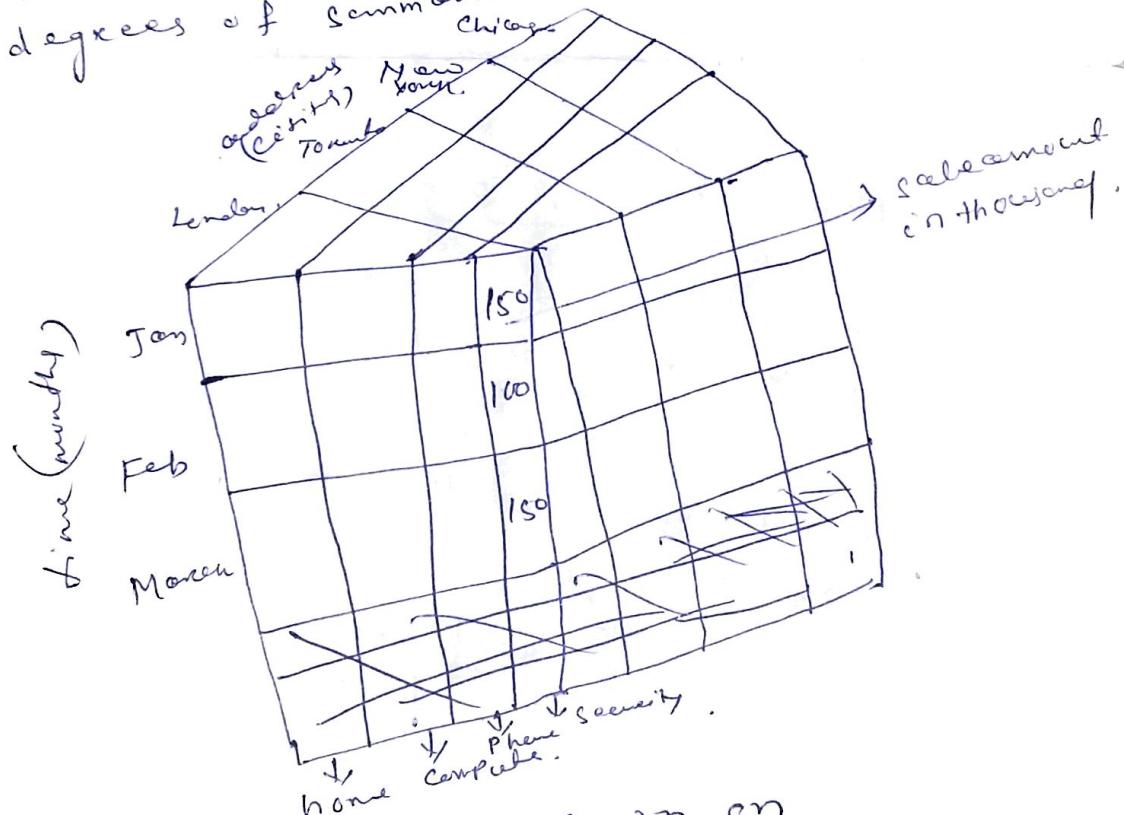
~~605
604
603
602~~



A multidimensional data cube showing summarized data for an organization.

most year - such deviations can be investigated.

- The cube has three dimensions: address, time and item type. (3)
- The aggregate value stored in each cell of the cube is sales-amount
- Ex- The total sales for the 1st quarter Q1, for items relating to security in London is \$400,00.
- by providing multidimensional data views and the pre-computation of summarized data, data ware house is well suited for OLAP - (online analysis present)
- OLAP operations include drill-down and roll-up, which allow the user to view the data at different degrees of summarization.



(Drill-down on
time data for Q1)

Drill-down on sales data summarized by quarter to see the data summarized by month.

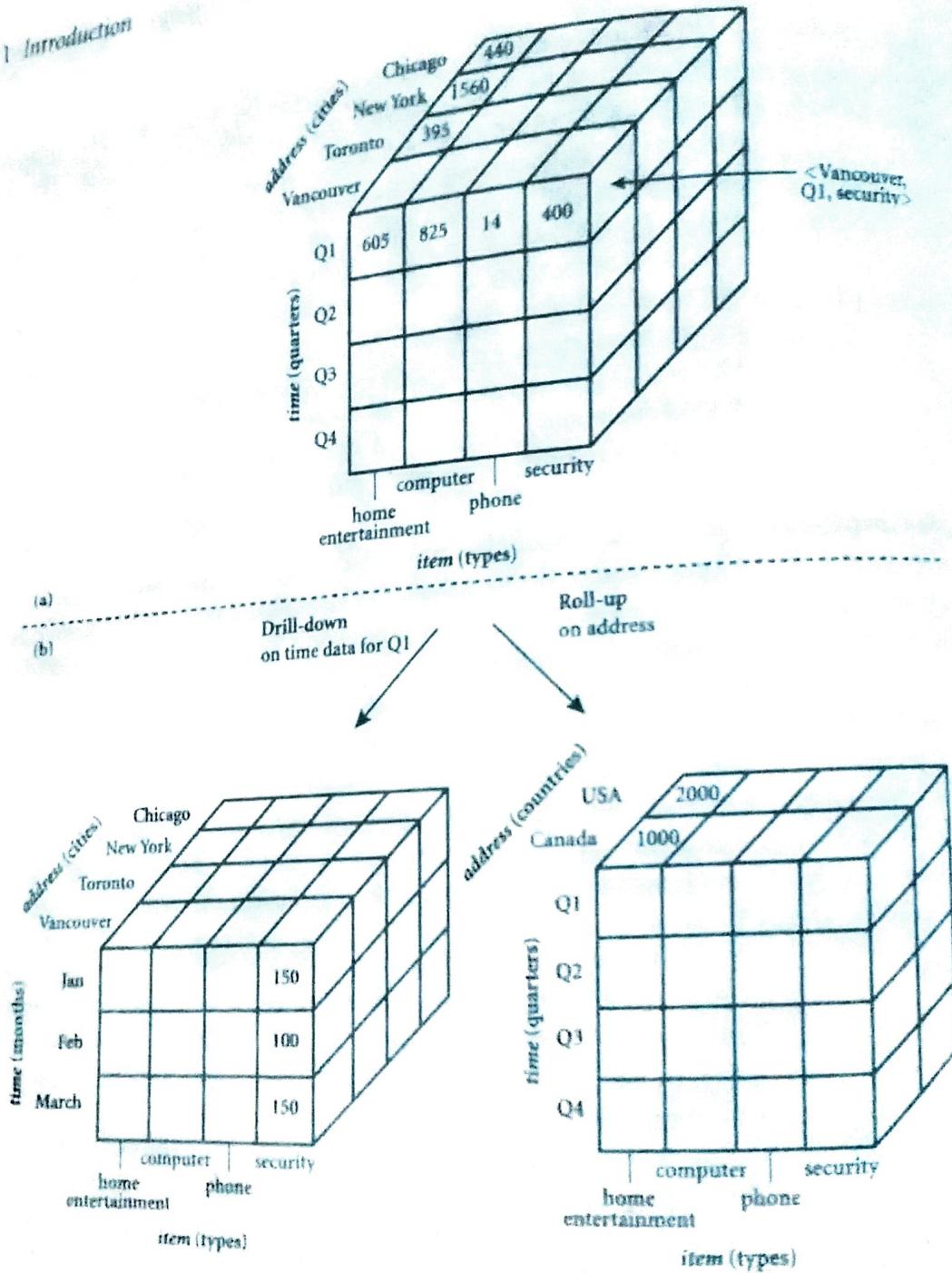
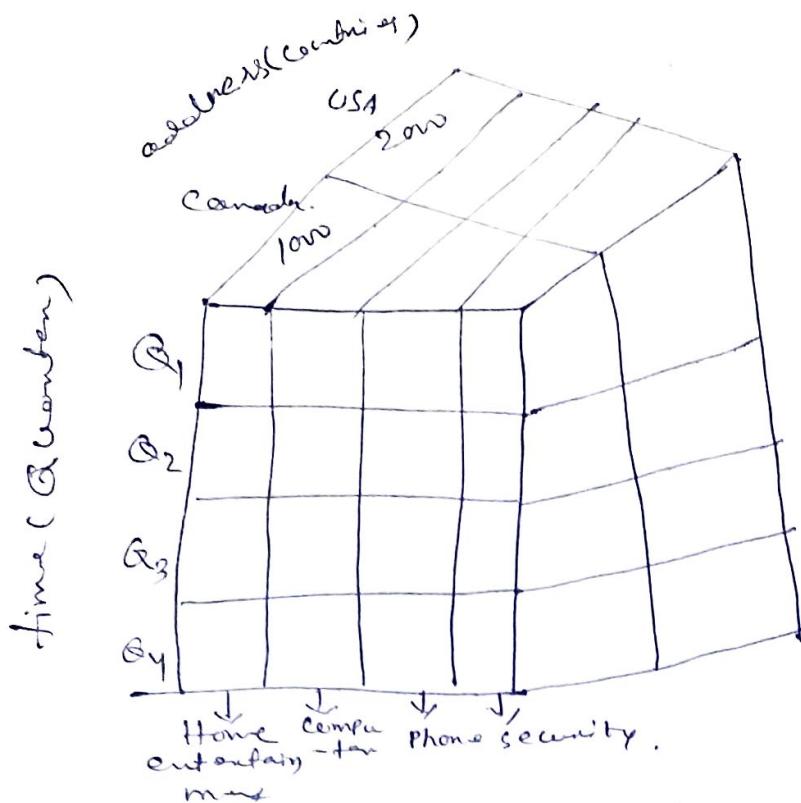


Figure 1.8 A multidimensional data cube, commonly used for data warehousing. (a) showing summarized data for AllElectronics and (b) showing summarized data resulting from drill-down and roll-up operations on the cube in (a). For improved readability, only some of the cube cell values are shown.

1.3.3 Transactional Databases

In general, a transactional database consists of a file or files that store data in a structured form.



(roll-up on address)

roll-up on sales data summarized by category to view the data summarized by country -

Transactional Databases

- It consists of a file where each record represents a transaction.
- A transaction including a unique transaction identity number and a list of items making the transaction.
- sales (table) →

trans-ID	list of items_IDs
T100	11, 13, 18, 116
T200	12, 18

Fragment of a transaction database for sales
- It may have additional tables associated with it, which contain other information regarding the sale, for example date of transaction, customer ID no., the ID no. of sales person and so on.

The sales table in the above figure is a nested relation because the attribute list of item_ids contains a set of items.

- Relational database systems do not support nested relational structures.
- Hence, the transactional database is usually either stored in a flat file or a format similar to that of the ^{sales} table.
- Or unfolded into a standard relation in a format similar to that of the items_sold table (Relational table).
- If we want to know how many transactions include item no. 13? We can answer it by ^{joining} ~~scraping~~ ^{active} database ^{and} sales table.
- Suppose we ask, "Which items sold well together?"
- Market basket data analysis would enable us to bundle groups of items together as a strategy for maximizing sales.
- Ex: The printers are commonly purchased together with computers. We could offer an expensive printer at a discount rate to customers buying computers, in the hope to increase the selling of the expensive printers.
- Data mining systems can identify frequent itemsets, i.e. sets of items that are frequently sold together.

Object Relational (OR) bases

→ This type of databases are considered based on an object-relational data model.

→ ~~The object-relational data model inherits the essential concepts of object-oriented databases, where each entity or instance is considered as an object.~~

Example → employees, customers or items are objects.

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

→ Each object has associated with it the following:-

- (i) A set of variables that describe the objects.
- (ii) A set of messages that the object can use to communicate with other objects or rest of the database.

~~(iii)~~ A set of methods, where each method holds the code to implement a message. After receiving a message, the method returns a value in response.

→ Objects that share a common set of properties can be grouped into an object

class.

→ Each object is an instance of its class.

→ Object classes can be organized into subclasses so that each class objects represents properties that are common to objects in that class.

Ex An employee class can contain variables like name, address and DOB. Suppose the class salesperson is a subclass of the class employee.

A sales-person object would inherit all of the variables pertaining to its superclass employee. In addition, it has all

(5)

of the variables that pertain to salesperson. Determining techniques needed for handling complex object structures, complex data types, class and subclass hierarchy, property inheritance etc.

Booking

Meeting ID	Room	Date	Time
DB Group 1331	PC	6/11/2016	10am

Temporal databases → stores relational

data that include time-related attributes. Example: Booking (meeting, room, time)

Ex ~~customer shopping sequence~~
A tuple (m, r, t) denotes the fact that meeting m is in room r at time t . sequences of

Sequence Databases → stores sequences of ordered events, with or without a concrete notion of time.

Ex customer shopping sequences; biological sequences → formation of DNA

Time series database → stores sequences of values or events obtained over repeated measurements of time (i.e. hourly, daily, weekly).

hourly, daily, weekly).

Stock exchange data.

Ex Inventory control
Change in Temperature can be used to find data mining techniques can be used to find the trend of changes for objects in the database.

such information can be useful in decision making and strategy planning.

Spatial database

It contains spatial related information.

It contains spatial (^(map)) related databases, VLDB, Examples → geographic databases, medical and satellite images

The data represents objects defined in geometric shape.

Temporal Database

EMP ID	Name	Dept	Salary	Valid from Start	Valid time end
10	John	Research	11000	1985	1990
10	John	Sales	11000.	1990	1993
10	John	Sales	12000	1993	JULY 2000
11	paul	Resear	10000.	1988	1995
12	George	R&D	10500.	1991	JULY
13	Niraj.	Soft	15500.	1998	JULY

time period attached to the data.

expressing when it was valid

Conventional database do not keep track of past or future DB

Benefits:

*(As few
few
a rectangular
grid of pixels
image & etc.)*

Spatial data, ^{one} represented in vector format consisting of n-dimensional bit maps are pixel maps.

Ex - 2-D satellite image \rightarrow where each pixel represents the rainfall in a given area.

Maps are represented in vector format where roads, bridges, buildings and lakes are represented as unions of basic geometric shapes.

Geographic databases have numerous applications :-

- (1) Forestry and ecology planning.
- (2). providing public service information regarding the location of telephone and electric cables, pipes and sewage system.
- (3) Vehicle navigation and dispatching system.

Ex \rightarrow taxi can store a city map with all information about the city, as well as the current location of each driver.

Data mining may discover patterns which describe the climate of mountains areas

Text Database

These contain word descriptions for objects. These words are not simple keywords but rather long sentences or paragraphs.

Ex - Product specification,
error reports,
warning messages.

Text databases → highly unstructured ex - Web page
→ semi structured → ex - e-mail
→ structured → ex : library catalogue (6)

Multimedia Database

These store image, audio and video data.
Applications → voice-mail, video on demand, www, speech based user interface

Heterogeneous Database

- It consists of a set of interconnected, autonomous component databases.
- The components communicate in order to exchange information and answer queries.
- Objects in one component may differ from objects in other component.

Legacy database

It is a group of heterogeneous database that combines different kinds of data

systems

Ex - relational or object oriented database.
Network database.
Spread sheets.
Multimedia databases.

- Information exchange occurs such databases is difficult, as they are in different form.

Data Streams

Here data flow in and out of a window dynamically or continuously.

Ex : - time series data
power supply.

stock exchange
telecommunications.

weather or environment monitoring

www

(1)

- users + exercise from one page to another page.
- understanding user access patterns will help us to improve marketing decision.
we can put some advertisement.
- Capturing user access patterns in such distributed environment environment is called weblog mining.

Q.1 What is data mining?

2. Describe the steps involved in data mining when viewed as a process of knowledge discovery (KDD)
3. suppose your job as a S/W engineer at AGO's to design a data mining system to examine the university course database, which contains the following information: the name, address and status of each student, the courses taken and the cumulative grade point average (CGPA). Describe the architecture you would choose. What is the purpose of each component of this architecture?
- ④ How is data warehouse different from database?
How they are similar?
- ⑤ Define each functionality of data mining.
- ⑥ Noisy data \rightarrow salary = -10
Incomplete data \Rightarrow salary = " "
inconsistent data. \Rightarrow age = 29, DOB = 24/2/1986



Second Edition



Data Mining

Concepts and Techniques

Jiawei Han and Micheline Kamber

