

Data mining refers to extracting or "mining" knowledge from large amounts of data. Unit-II

- ① Data cleaning: (to remove noise and inconsistent data)
- ② Data integration (where multiple data source may be combined)
- ③ Data selection (where data relevant to the analysis task are retrieved from the database)
- ④ Data transformation (where data are transformed or consolidated into forms appropriate for mining by performing summary or aggregation operations)
- ⑤ Data mining (an essential process where intelligent methods are applied in order to extract data patterns)
- ⑥ Pattern evaluation (to identify the truly interestingness patterns representing knowledge based on some interestingness measures.)
- ⑦ Knowledge presentation (where visualization & knowledge representation techniques are used to present the mined knowledge to the user).

Step 1 to 4 are different forms of data preprocessing, where the data are prepared for mining.

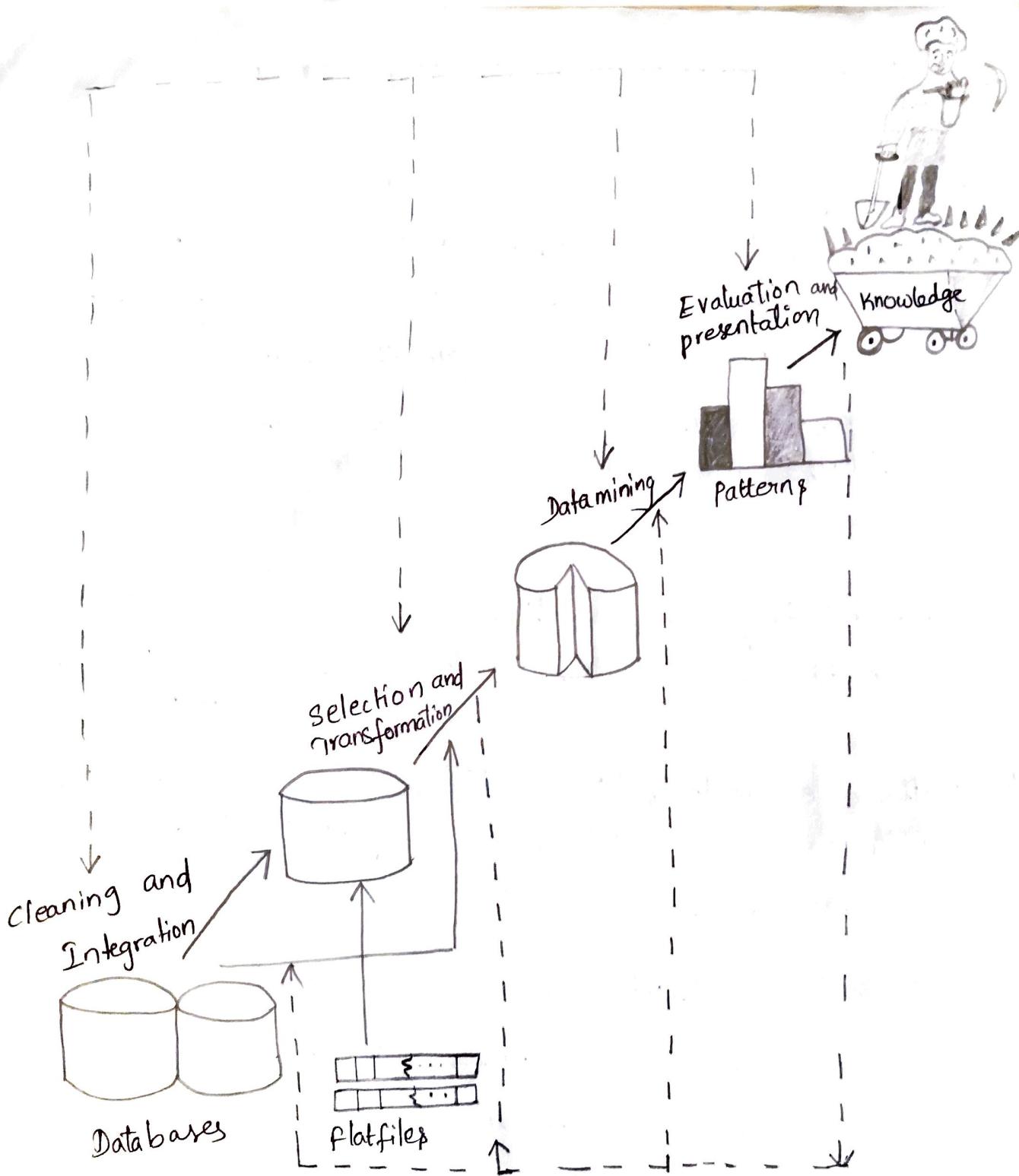


fig: Datamining as a step in the process of Knowledge discovery (KDD)

Database or datawarehouse server ; the database or data warehouse server is responsible for fetching the relevant data, based on the user's data mining request knowledge base ; this is the domain knowledge that is used to guide the search or evaluate the interestingness of resulting patterns. Such knowledge can include concept hierarchies, used to organize attributes or attribute values into different levels of abstraction .

Data mining engine : This is essential to the data mining system and ideally consists of a set of functional modules for tasks such as characterization, association and correlation analysis, classification, prediction, cluster analysis, outlier analysis and evolution analysis

Pattern evaluation module : this component typically employs interestingness measures and interacts with the data mining modules so as to focus the search toward interesting patterns.

User Interface : this module communicates between users and the data mining system, allowing the user to interact with the system by specifying a datamining query or task , providing information to help focus the search and performing exploratory data mining based on the intermediate data mining result.

Based on this view, the architecture of a typical datamining system may have the following major components.

- ① Database, datawarehouse, www, or other information repository: This is one or a set of databases, datawarehouses, spreadsheets or other kinds of information repositories. Data cleaning and data integration techniques may be performed on the data.

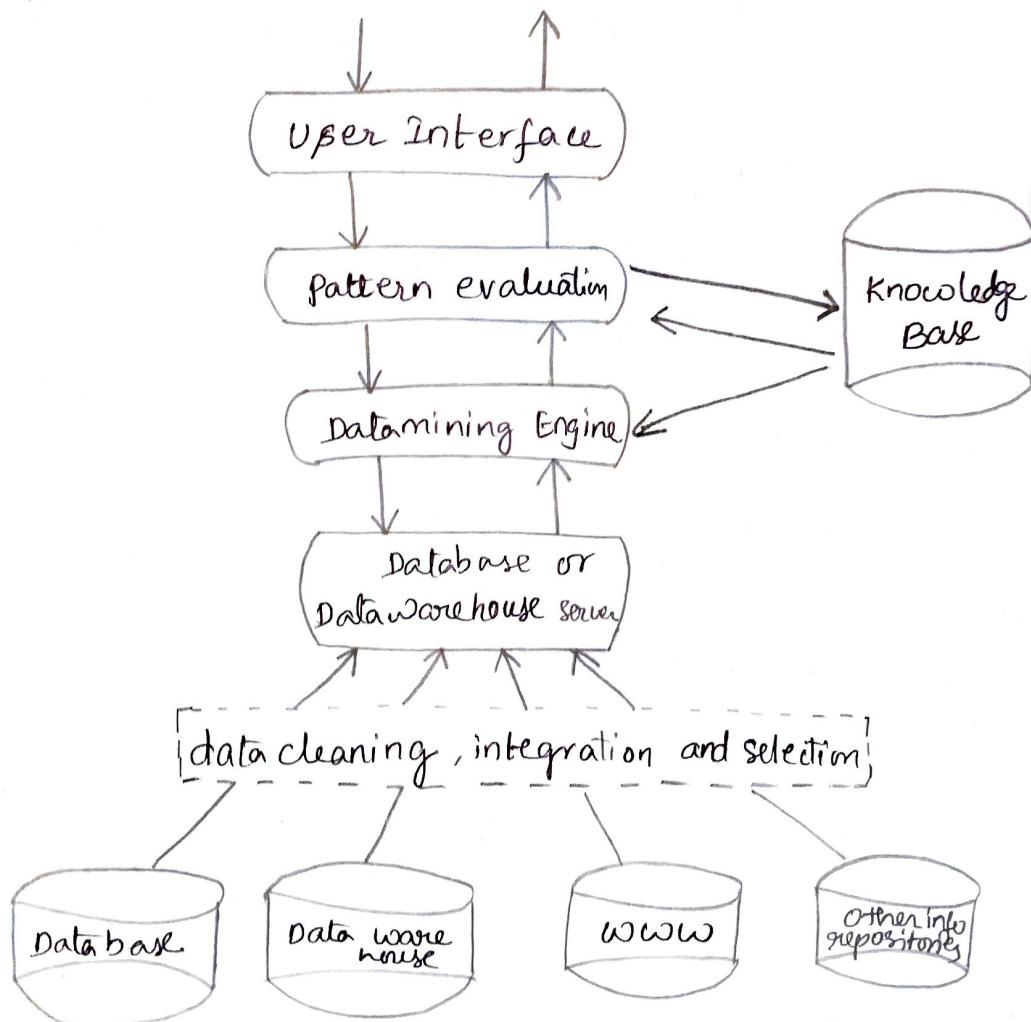


fig : Architecture of a Typical data mining System

## ② Data Mining - On what kind of Data?

- ① Relational Databases
- ② Data warehouses
- ③ Transactional Databases
- ④ Advanced data and information systems and advanced applications

### Object-Relational Databases

Temporal Databases, sequence Databases and Time-series DB's.

Spatial DB's and Spatiotemporal DB's

Text DB's and Multimedia DB's

Heterogeneous DB's and Legacy DB's.

Data streams.

The WWW

### D) Relational DB's

- A Relational DB is a collection of tables, each of which is assigned a unique name.
- Each table consists of a set of attributes (cols or fields) and usually stores a large set of tuples (records or rows).
  - Each tuple in a relational table rep an object identified by a unique key and described by a set of attribute values.
  - An ER data model rep the DB as a set of entities and their relationships.

customer

cust-ID	name	address	age	income	credit-info	category	...
C1	Smith, Sandy	1223 Lake Ave., Chicago, IL	31	\$78000	1	3	...
..	...	...	..	..	..	..	..

## item

item-ID	name	brand	category	type	price	place-made	supplier	...
13	hi-resTV	Toshiba	high resolution TV	TV	\$988.00	Japan	Niko x	\$600.
18	Laptop	Dell	laptop	computer	\$1369.00	USA	Dell	\$983.0
...	...	...	...	...	...	...	...	...

employee (emp-ID, name, category, group, salary, commission)

Branch (branch-ID, name, address)

purchases (transID, cust-ID, emp-ID, date, time, method-paid, amount)

Items\_sold (trans-ID, item-ID, qty)

works-at (emp-ID, branch-ID).

Q

## DW Lt

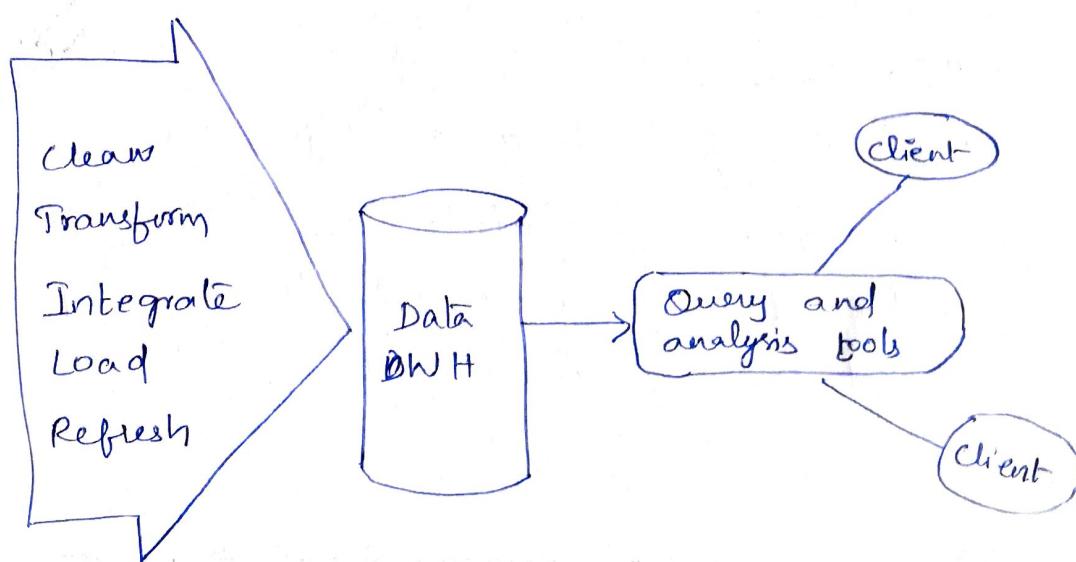


fig: Typical architecture of a dataware house for All Electronics

## Transational Databases

Sales

trans-ID	List of item-IDs
T100	I1, I3, I8, I6
...	....

### ⑤ Advanced DB systems and Advanced DB Applications

→ The new DB applications include handling spatial data (such as maps), engineering design data (such as the design of buildings, system components or integrated circuits), hypertext and multimedia data (including text, image, video and audio data), time related data (such as historical records or stock exchange data) and the www (a huge, widely distributed information repository made available by the Internet).

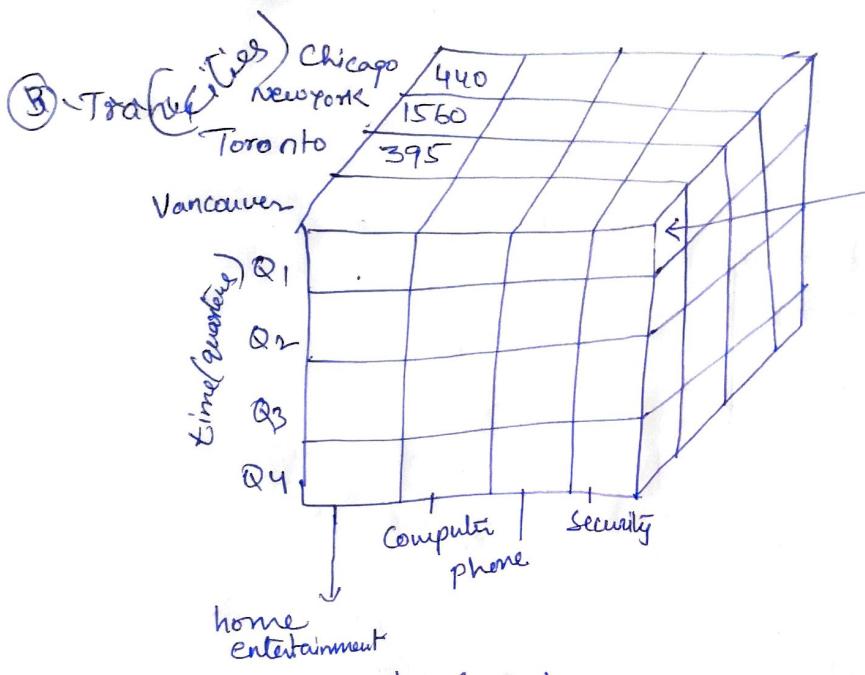
### Object-Oriented DB's

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

→ Each object is associated with a set of variables, messages & methods.

### object-Relational DB's

→ The object-relational model extends the basic relational data model by adding the power to handle complex data model types, class hierarchies, and object inheritance.



< Vancouver, Q1, serial  
 → such as  
 desk  
 Temp

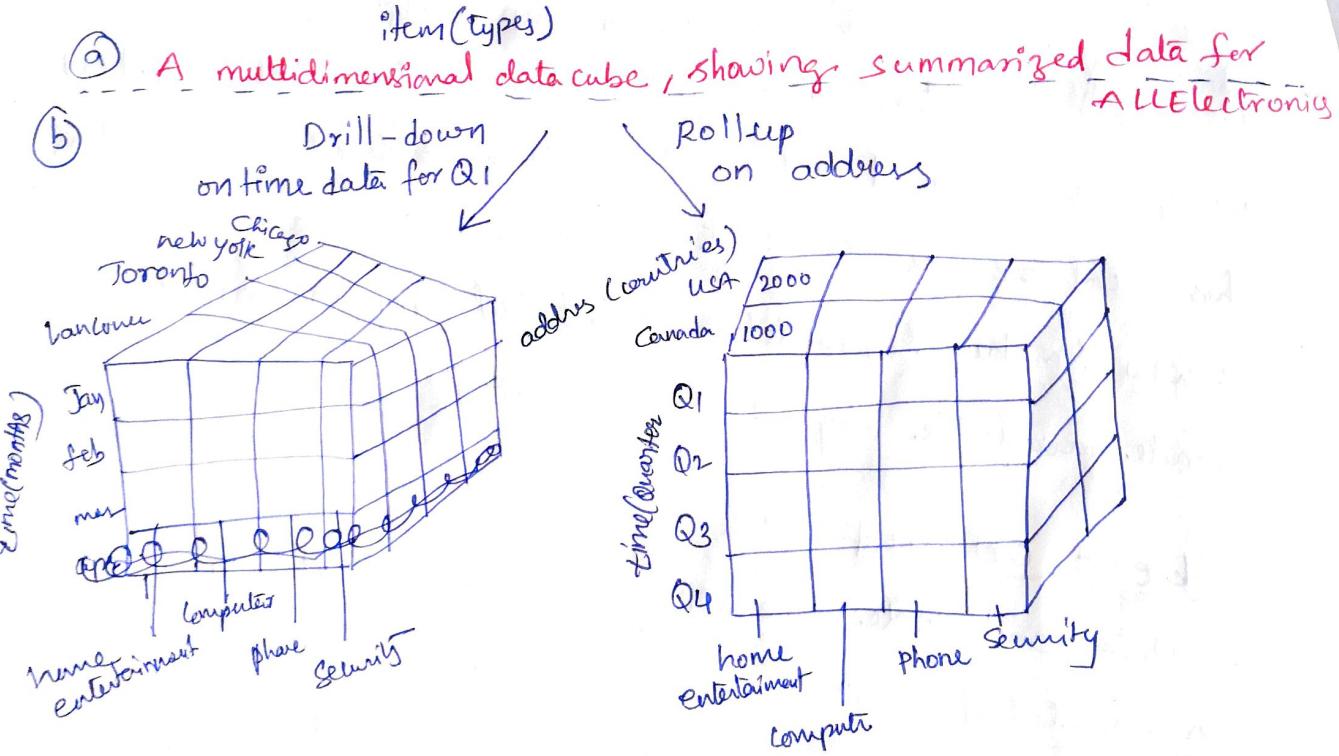


fig: showing summarized data resulting from  
 drill-down and roll-up operations on the  
 cube in .

Spatial DB's: contain spatial-related information.  
→ such DB's includes geographic (map) DB's, VLSI chip design DB's and medical and satellite image DB's.

### Temporal DB's and Time-series DB's

→ Both store time-related data.

→ A time-series DB stores sequences of values that change with time, such as data collected regarding the stock exchange.

→ A temporal DB usually stores relational data that include time-related attributes (timestamps).

### Text DB's and Multimedia DB's

→ Text DB's are DB's that contain word descriptions for objects. (Ex: keywords or paragraphs).

→ Multimedia DB's store image, audio and video data.

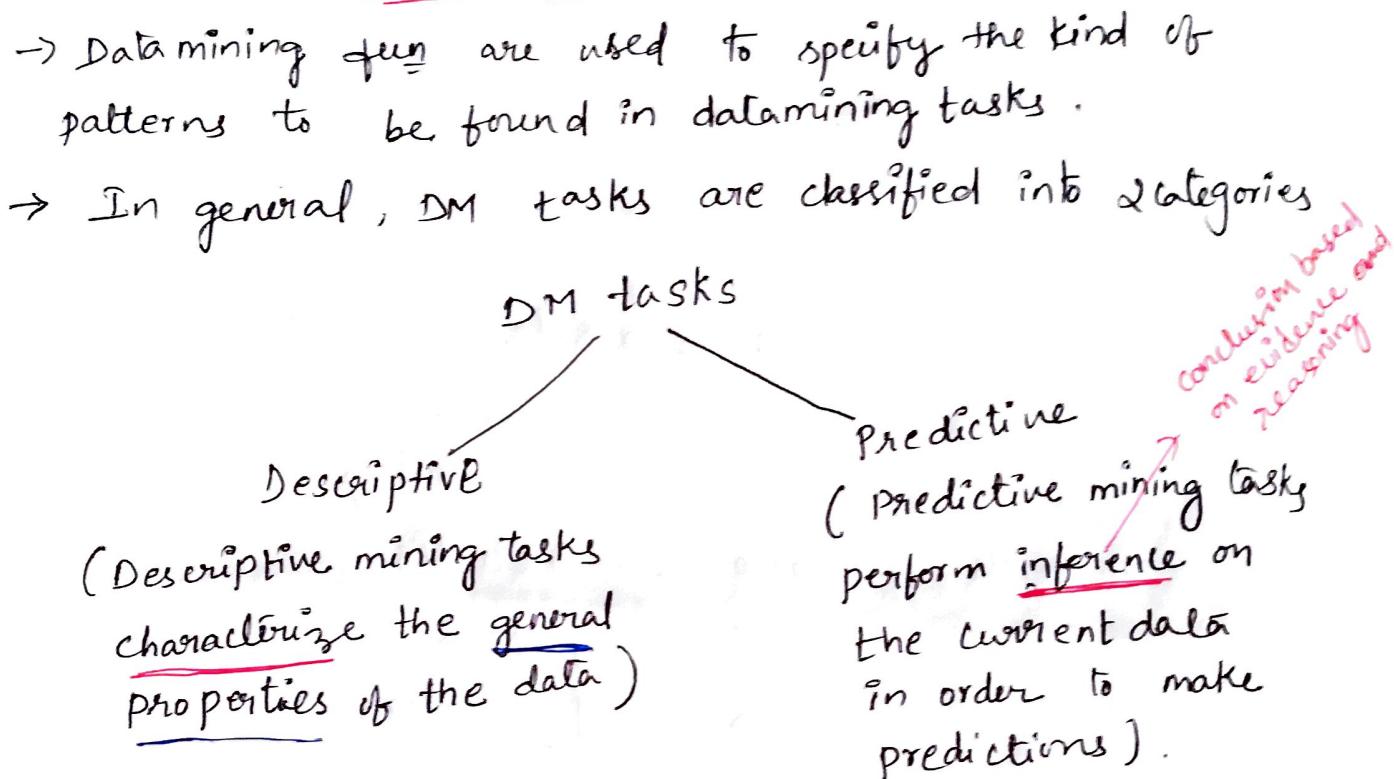
Ex: picture content-based retrieval, voicemail systems, video on-demand sys, the www systems, speech-based user interfaces that recognize spoken commands.

### Heterogeneous DB's & Legacy DB.

A legacy db is a group of heterogeneous db's that combines diff kinds of data systems;

The www:

## Data Mining Functionalities - What kind of patterns can be mined



→ DM systems should be able to discover patterns at various granularities (i.e. diff levels of abstraction).

DM fuc and kinds of patterns they can discover are:

- ① concept / class description: Characterization & Discrimination
- ② Mining frequent patterns, correlations and Associations.
- ③ classification and prediction
- ④ cluster Analysis
- ⑤ outlier "
- ⑥ Evolution "

# ① concept / class description?

characterization and Discrimination

→ Data can be associated with classes & concepts

Ex: All Electronics store

classes of items for sale include

- computers
- printers

Concepts of customers includes

- Big spenders
- budget spenders.

→ Data characterization is a summarization of the general characteristics or features of a target class of data.

Ex: To study the characteristics of S/W products whose sales increased by 10%.

In the last year, D/P rep in pie charts, bar charts, curves, multidimensional data cubes and tables

→ Data discrimination is a comparison of the general features of target class data objects with the objects from one or a set of contrasting classes

Data characterization

Generalization (common class  
to all customer)  
i.e A & B belongs to the same

Data discrimination

comparative class  
i.e B is better than A

in  $y$ ".

E2: Given the AllElectronics relational database, a datamining system may find association rules like

$\text{age}(x, "20\text{--}29") \wedge \text{income}(x, "20K\text{--}29K") \Rightarrow \text{buys}(x, "cdplayer") [S=2\%, C=60\%]$

where  $x$  is a variable , rep a customer. The rule indicates that of the AllElectronics customer under study , 2%. (support) are 20 to 29 years of age with an income of 20K to 29K and have purchased a cd player at AllElectronics .

→ there is a 60% probability (confidence, or certainty) that a customer in this age and income group will purchase a CD player.

NOTE :- That this is an association b/w more than one attribute, or predicate (i.e. age, income and buys).

→ Adopting the terminology used in multidimensional databases, where each attribute is referred to as a dimension, the above rule can be referred to as a multidimensional association rule.

→ Suppose, as a marketing manager of AllElectronics, you would like to determine which items are frequently purchased together with in the same transactions.

→ An example of such a rule is

contains (T, "Computer")  $\Rightarrow$  contains (T, "Software")

[support = 1%, confidence = 50%]

meaning that if a transaction, T, containing "Computer", there is a 50% chance that it contains "Software" as well and 1% of all of the transactions contain both.

→ This association rule involves a single attribute or predicate (i.e. contains) that repeats.

→ Association rules that contain a single predicate are referred to as single-dimensional association rules.

→ Dropping the predicate notation, the above rule can be written simply as

\* computer  $\Rightarrow$  software [ 1% - 50% ]

Classification and prediction, is the process of finding set of models (or fun's) that describes and distinguish data classes or concepts for the purpose of being able to use the model to predict the class of objects whose class label is unknown.

→ the derived model is based on the analysis of a set of training data (i.e data objects whose class label is known)

→ the derived model may be rep in classification ( IF - THEN ) rules various form

decision Trees

mathematical formulae  
neural networks

→ A decision tree is a flowchart like tree structure where each node denotes a test on an attribute value, each branch represents an outcome of the test, and tree leaves rep classes or class distributions.

→ Decision trees can be easily converted to classification rules.

→ A neural netw is when used for classification, is typically a collection of neuron-like processing units with weighted connections b/w the units.

- Classification can be used for predicting the class label of data objects.
- Users may wish to predict some missing or unavailable data values rather than class label.
- This is usually the case when the predicted values are numerical data and is often specific referred to as prediction.

Ex: age ( $x$ , 29..35) and income ( $x$ , high)  $\rightarrow$  class (?)

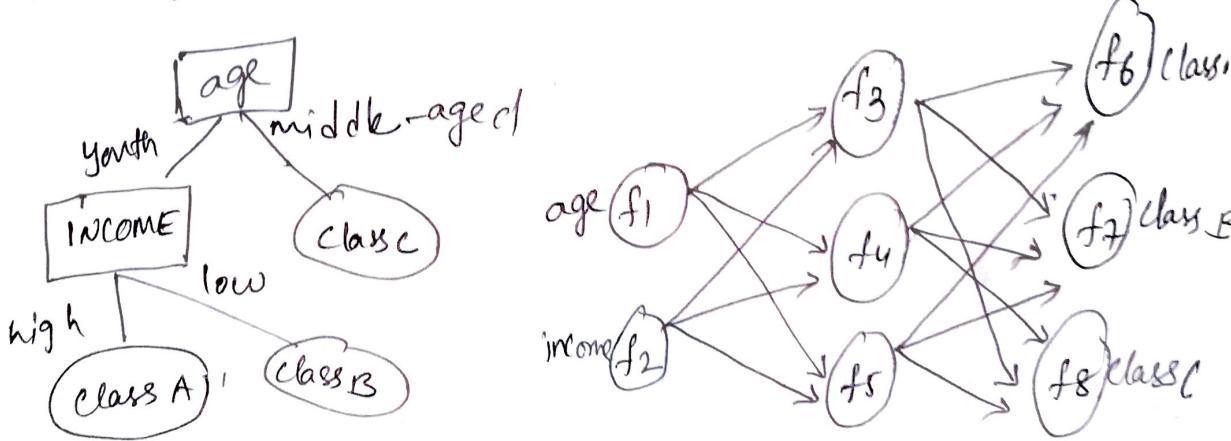


fig: A classification model can be represented in various forms such as (a) IF-THEN rules, (b) a decision tree, or a (c) neural network.

(4) Cluster Analysis: Unlike classification and prediction which analyze class-labeled data objects, clustering analyzes data objects without consulting a known class label.

→ In general, the class labels are not present in the training data simply because they are not known to begin with.

- clustering can be used to generate such labels.
- The objects are clustered or grouped based on the principle of maximizing the intra-class similarity and minimizing the inter-class similarity.
- clustering can also facilitate taxonomy formation that is, the organization of observations into a hierarchy of classes that group similar events together.

Ex: cluster analysis can be performed on AllElectronics customer data in order to identify homogeneous subpopulations of customers.

→ These clusters may represent individual target groups for marketing.

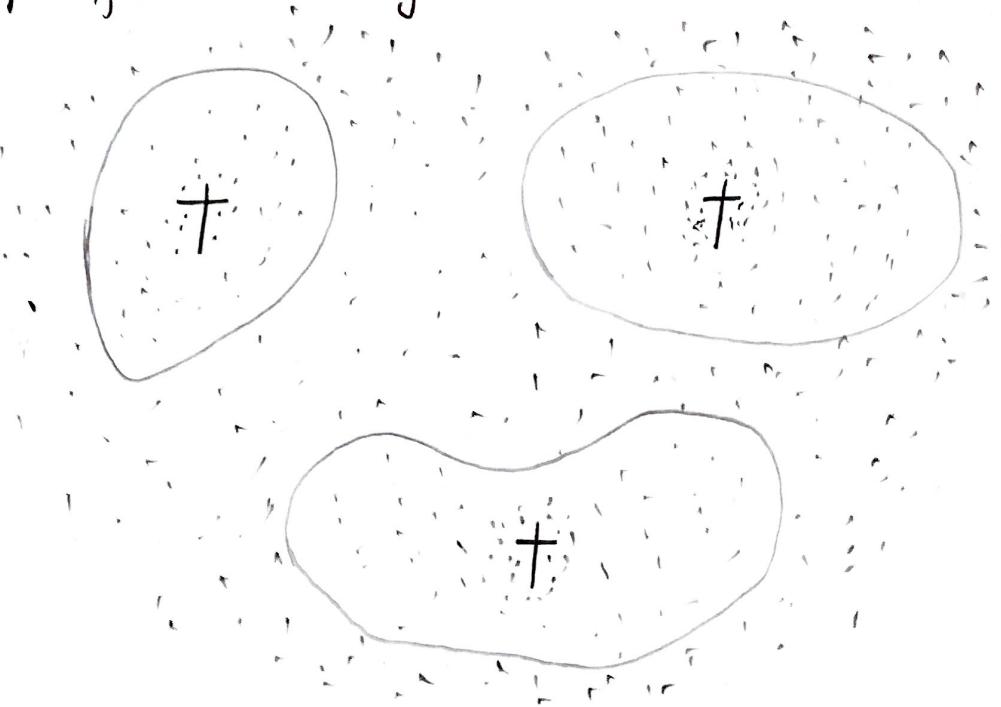


fig: A 2-D plot of customer data with respect to customer locations in a city, showing 3 data clusters. Each cluster "center" is marked with a "+".

## ⑤ outlier Analysis (Fraud Detection)

class

A database may contain data objects that do not comply with the general behavior or model of the data. These data objects are outliers.

- Most datamining methods discard outliers as noise or exceptions.
- The analysis of outliers data is referred to as outlier mining.
- Outliers may be detected using statistical tests that assume a distribution or probability model for the data.

## ⑥ Evolution Analysis

Data evolution analysis describes and models regularities or trends for objects whose behavior changes over time. Although this may include characterization, discrimination, association and correlation analysis, classification, prediction or clustering of time related data, distinct features of such an analysis include time-series data analysis, sequence or periodicity pattern matching, and similarity-based data analysis.

## Classification of Datamining System

→ Because of the diversity of disciplines contributing to DM, Data mining research is expected to generate a large variety of datamining systems.

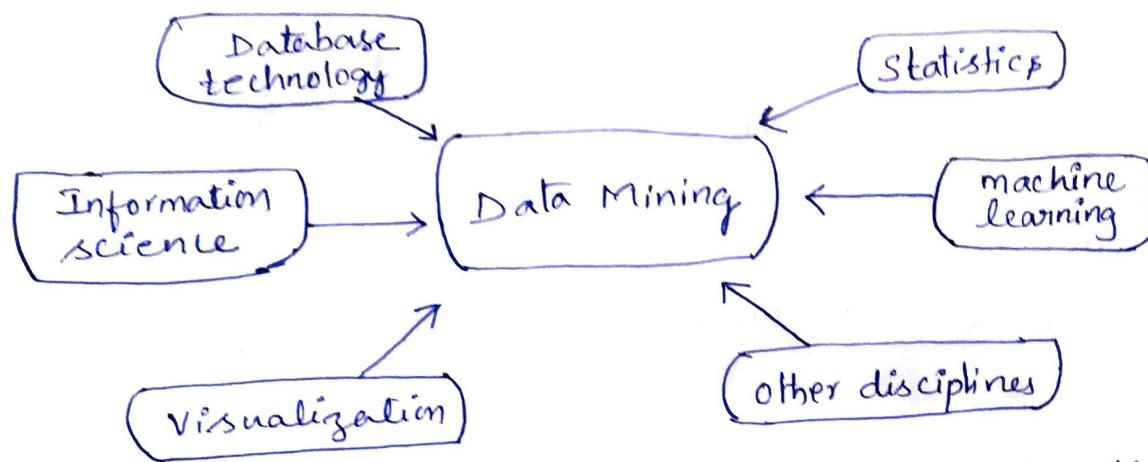


fig: Data mining as a confluence of multiple disciplines.

① classification according to the kinds of DB's mined.

~~if classifying according to Datamodels~~ relational, transactional, object-oriented, object-relational or datawarehouse mining system.

~~if classifying according to the special types~~ spatial, time-series, text, or multimedia data mining system or a www mining system.

② classification acc to the kinds of knowledge mined

→ Based on DM functionalities, such as characterization, discrimination, association, classification, clustering, outlier analysis and evolution analysis.

③ classification according to the kinds of techniques utilized

→ Database-oriented or datarehouse-oriented techniques, machine learning, statistics, visualization, pattern recognition, neural networks and so on.

④ classification according to the applications adapted

→ Finance; telecommunications, DNA, stock markets, e-mail and so on.

## Major Issues in Data Mining

### ① Mining Methodology and user interaction issues

- Mining diff kinds of knowledge in DB's.
- Interactive mining of knowledge at multiple levels of abstraction.
- Incorporation of background knowledge (info regarding the domain under study)
- DM query languages and adhoc data mining → SQL
- Presentation and visualization of DM results. → DMQL
- Handling noisy or incomplete data → easily understood by humans
- Pattern evaluation - the interestingness problem.

### ② Performance issues

- Efficiency and scalability of DM algorithms.
- Parallel, distributed and incremental mining algorithms → huge size of many DB's, wide distribution of data.

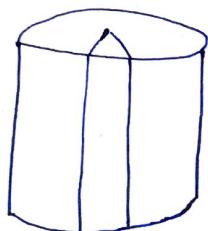
### ③ Issues relating to the diversity of database types:

- Handling of relational and complex types of data.
- Mining information from heterogeneous DB's and global info systems.

## Data Mining task primitives

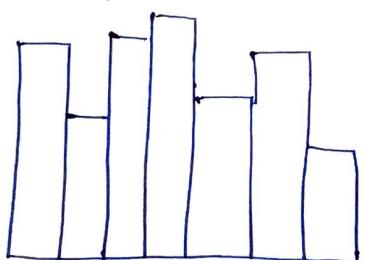
- ① Task-relevant data: what is the data set that I want to mine?
- ② What kind of knowledge do I want to mine?
- ③ What background knowledge could be useful here?
- ④ What measures can be used to estimate pattern interestingness?
- ⑤ How do I want the discovered patterns to be presented?

- ① The set of task-relevant data to be mined



### Task-relevant data

Database or Data warehouse name  
Database tables or data warehouse cubes  
Conditions for data selection  
Relevant attributes or dimensions  
Data grouping criteria.



### Knowledge type to be mined

Characterization

Discrimination

Association

Classification / Prediction

Clustering

A user studying the buying habits of all Elec customers may choose to mine association rules of the form

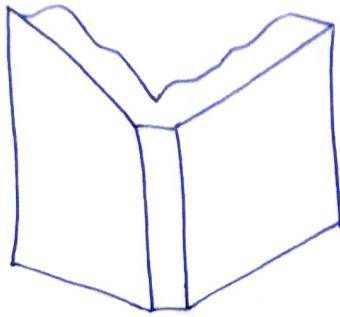
$$P(x: \text{customer}, w) \wedge Q(x, y) \Rightarrow \text{buys}(x, z)$$

P, Q  $\rightarrow$  predicate vars

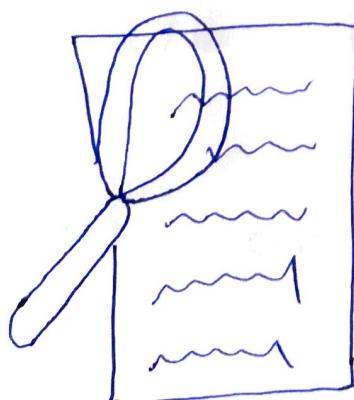
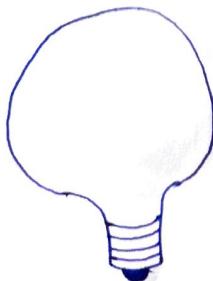
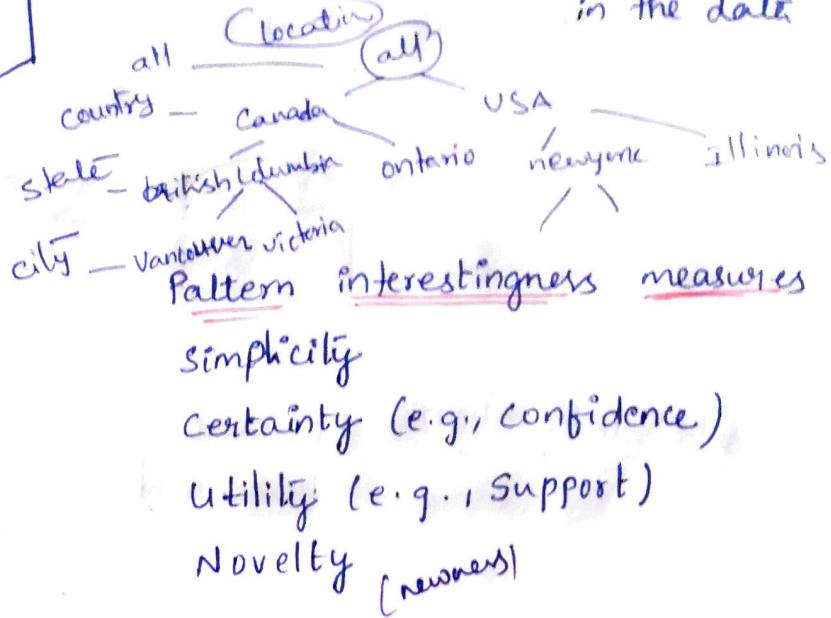
w, y, z are obj vars

age(x, "30..39")  $\wedge$  income(x, "40k..49k")  $\Rightarrow$  buys(x, "VCR")

and occupation(x, "stud")  $\wedge$  age(x, "20..29k")  $\Rightarrow$  buys(x, "comp") [2.2%, 60%]  
[1.4%, 70%]



Background Knowledge  
concept hierarchies → see from a set of low-level concepts to higher-level user beliefs about relationships in the data  
Location (11) mapping



## Visualization of discovered patterns

Rules, tables, reports, charts,  
graphs, decision trees and  
algorithms

Drill-down and roll-up.

## fig: Primitives for specifying a data mining task.

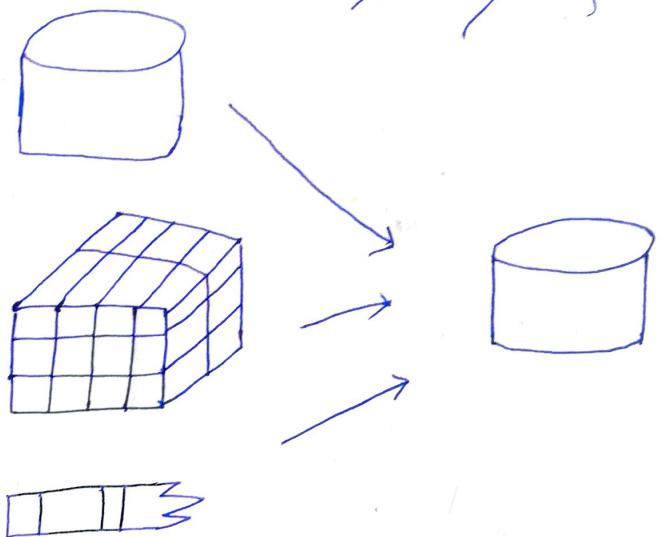
# Data Preprocessing

→ Data Preprocessing is an imp issue DW & DM as real-world data tend to be incomplete, noisy and inconsistent.

## Data cleaning



## Data Integration



## Data transformation

-2, 32, 100, 59, 48 → -0.02, 0.32, 1.00,

## Data reduction

		attributes																
		A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	...	A <sub>126</sub>												
transactions	T <sub>1</sub>																	
	T <sub>2</sub>																	
	T <sub>3</sub>																	
	T <sub>4</sub>																	
	...																	
	T <sub>2000</sub>																	

		attributes																
		A <sub>1</sub>	A <sub>3</sub>	...	A <sub>115</sub>													
transactions	T <sub>1</sub>																	
	T <sub>4</sub>																	
	...																	
	T <sub>1456</sub>																	

fig: Forms of data preprocessing

## Descriptive Data summarization

### ① Measuring the central tendency.

The ~~bar~~ most common and most effective numerical measure of the "center" of a set of data is the (arithmetic) mean.

The mean of this set of values is

$$\bar{x} = \frac{\sum_{i=1}^N x_i}{N} = \frac{x_1 + x_2 + \dots + x_N}{N}$$

$$\text{Median} = L_1 + \left( \frac{\frac{N/2 - (\sum \text{freq})l}{\text{freqmedian}}}{\text{width}} \right)$$

variance and standard deviation

the variance of  $N$  observations,  $x_1, x_2, \dots, x_n$  is

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \frac{1}{N} \left[ \sum x_i^2 - \frac{1}{N} (\sum x_i)^2 \right].$$

$L_1$  is the lower boundary of the median interval,

no. of values in the entire data set,

$(\sum \text{freq})l$  is the sum of the freq's of all of the intervals that are lower than the median interval,

$(\text{freqmedian})$  is the freq of the median interval,

width is the width of the median interval.

## Data cleaning

### ① Missing Values

- Ignore the tuple: → when class label is missing
- fill in the missing value manually! - Time consuming
- Use a global constant to fill in the missing value [unknown or os]
- Use the attribute mean to fill in the missing value. [Avg value fitting]
- Use the attribute mean for all samples belonging to the same class as the given tuple
- Use the most probable value to fill in the missing value
  - ↳ determined with regression, inference-based tools using a Bayesian formalism or decision tree induction

### ② Noisy Data

Noise is a random error or variance in a measured variable.

- ① Binning
  - ② clustering → outliers may be detected by clustering.
  - ③ combined computer and human inspection; - helps to identify outlier patterns.
  - ④ Regression → linear regression involves finding Ex: mislabeled characters
    - ↳ the "best" line to fit two variables.  
So that one variable can be used to predict the other.
- Binning: Binning methods smooth a sorted data value by consulting its "neighborhood", that is, the values around it.
- The sorted values are distributed into a no. of "buckets" or bins.
- Bcuz binning methods consult the neighborhood of values, they perform local smoothing.

Ex: = sorted data for price (in dollars):

4, 8, 15, 21, 21, 24, 25, 28, 34.

Partition into (equidepth) bins:

Bin 1: 4, 8, 15

Bin 2: 21, 21, 24

Bin 3: 25, 28, 34

Smoothing by bin means:

Bin 1: 9, 9, 9

Bin 2: 22, 22, 22

Bin 3: 29, 29, 29

Smoothing by bin boundaries:

Bin 1: 4, 4, 15

Bin 2: 21, 21, 24

Bin 3: 25, 25, 34

fig: Binning methods for data smoothing

③ Inconsistent data: inconsistent in data recording for some transactions.  
→ Some data inconsistencies may be corrected manually using external references.

Ex: Errors made at data entry may be corrected by performing a paper trace.

## Data Integration and Transformation

- The merging of data from multiple data stores.
- Redundancy is an important issue.
- Some redundancies can be detected by correlation Analysis.

$$\chi^2 = \sum_{i=1}^C \sum_{j=1}^R \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$E_{ij} = \frac{\text{count}(A=a_i) \times \text{count}(B=b_j)}{N}$$

Ex: A  $2 \times 2$  contingency table for the data.

are Gender and Preferred - Reading correlated?

	Male	Female	Total
Fiction	250 (90)	200 (360)	450
non-fiction	50 (210)	1000 (840)	1050
Total	300	1200	1500

$$E_{11} = \frac{300 \times 450}{1500} = 90.$$

$$\begin{aligned} \chi^2 &= \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} \\ &= 284.44 + 121.90 + 71.11 + 30.48 = 507.93. \end{aligned}$$

Issues Schema integration & object matching can be tricky

- (1) Schema integration can be detected by correlation analysis.
- (2) Redundancy, can be detected by correlation analysis and data integration is the detection and resolution of data value conflicts.
- (3) Data integration is the detection and resolution of data value conflicts.

## Data Transformation

→ The data are transformed or consolidated into forms appropriate for mining.

- (1) smoothing - remove noise from the data. such techs incl. Binning, regression and clustering.
- (2) Aggregation - summary or aggregation Ex: daily sales may be aggregated so as to compute monthly & annual total amounts.
- (3) Generalization - Techs: Data cube for analysis of data at multiple granularities. (P) (raw) data can be replaced by higher level concepts. Techs: Concept hierarchy.
- (4) Normalization
- (5) Attribute construction.

→ (feature construction), where new attributes are constructed and added from the given set of attributes to help the mining process.

where the attribute data are scaled so as to fall within a small specified range, such as -1.0 to 1.0  
or  
0.0 to 1.0

Min-max normalization performs a linear

transformation on the original data.

→ Min-Max normalization maps a value,  $v$  of  $A$  to  $v'$  in the range

$[new_{-}min_A, new_{-}max_A]$  by computing

$$v' = \frac{v - min_A}{max_A - min_A} (new_{-}max_A - new_{-}min_A) + new_{-}min_A$$

Suppose that the minimum & maximum values for the attribute income are \$12,000 & \$98,000, resp. we would like to map income to the range [0.0, 1.0]. By min-max normalization, a value of \$73,600 for income is transformed to.

$$\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$$

Z-score normalization or zero-mean normalization, the values for an attribute, A are normalized based on the mean and standard deviation of A. A value, v of A is normalized to  $v'$  by computing.

$$v' = \frac{v - \bar{A}}{\sigma_A}$$

where  $\bar{A}$  &  $\sigma_A$  are the mean & SD.

Z-Score normalization suppose that the mean and standard deviation of the values for the attribute income are \$54,000 & \$16,000, respectively, with z-score normalization, a value of \$73,600 for income is transformed to

$$\frac{73,600 - 54,000}{16,000} = 1.225$$

## Data Reduction

→ Applied to obtain a reduced representation of the data set, that is much smaller in volume, yet closely maintains the integrity of the original data.

- ① Data cube aggregation, where aggregation operations are applied to the data in the construction of a data cube.
- ② Attribute subset selection, where irrelevant, weakly relevant or redundant attributes or dimensions may be deleted and removed.
- ③ Dimensionality reduction, where encoding mechanisms are used to reduce the data set size.  
compressed or encrypted without losing the identity.
- ④ Numerosity reduction, where the data are replaced or estimated by alternative, smaller data representations clustering, sampling & the use of histograms.
- ⑤ Discretization and concept hierarchy generation is a form of numerosity reduction that is very useful for the automatic generation of concept hierarchies.

## ① Data cube aggregation

→ Data cubes store multidimensional aggregated information.

The diagram illustrates the process of data cube aggregation. On the left, three separate 2D tables are shown, each representing sales data for a specific year (2002, 2003, and 2004). Each table has 'Quarter' and 'Sales' as columns. An arrow points from these three tables to a single 1D table on the right, which contains the total sales for each year, with the years listed in chronological order.

Year	Sales
2002	\$ 1,568,000
2003	\$ 2,356,000
2004	\$ 3,594,000

fig: AllElectronics for the years 2002 to 2004

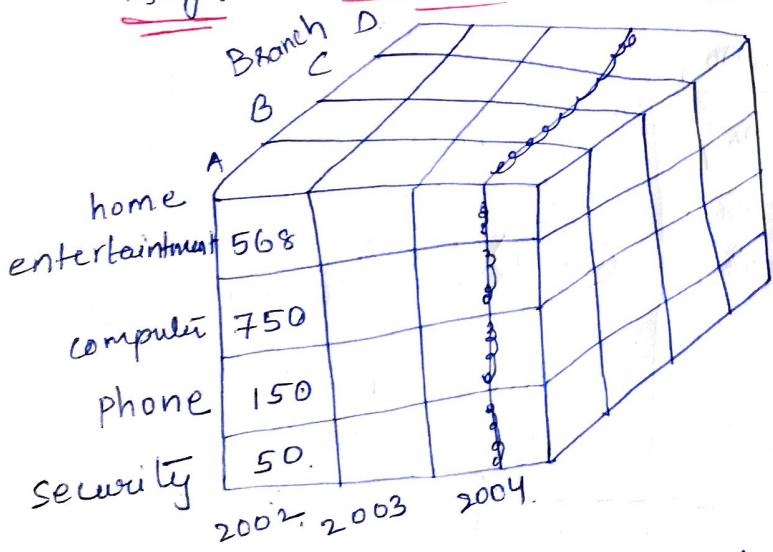


fig: A data cube for sales at AllElectronics

## Attribute subset selection

- Attribute subset selection reduces the data size by removing irrelevant or redundant attributes (dimensions).
- The goal of attribute subset selection is to find a minimum set of attributes such that the resulting probability distribution of the data class is as close as possible to the original distribution obtained using all attributes.

Forward selection	Backward elimination	Decision tree induction
Initial attribute set: $\{A_1, A_2, A_3, A_4, A_5, A_6\}$ Initial reduced set: $\{ ? \}$ $\Rightarrow \{A_1\}$ $= \{A_1, A_4\}$ $\Rightarrow$ Reduced attribute set: $\{A_1, A_4, A_6\}$	Initial attribute set: $\{A_1, A_2, A_3, A_4, A_5, A_6\}$ Initial reduced set: $\Rightarrow \{A_1, A_3, A_4, A_5, A_6\}$ $\Rightarrow \{A_1, A_4, A_5, A_6\}$ $\Rightarrow$ Reduced attribute set: $\{A_1, A_4, A_6\}$	Initial attribute set: $\{A_1, A_2, A_3, A_4, A_5, A_6\}$ <pre> graph TD     A4[A4 ?] -- Y --&gt; A19[A19]     A4 -- N --&gt; Ab2[Ab2]     A19 -- Y --&gt; Class1((Class 1))     A19 -- N --&gt; Class2((Class 2))     Ab2 -- Y --&gt; Class3((Class 1))     Ab2 -- N --&gt; Class4((Class 2))   </pre> $\Rightarrow$ Reduced attribute set: $\{A_1, A_4, A_6\}$

- ① Stepwise forward selection: The procedure starts with an empty set of attributes as the reduced set.  
→ 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 elimination: The procedure starts with the full set of attributes. At each step, it removes the worst attribute remaining in the set.
- ③ Combination of forward selection and backward elimination: the stepwise forward selection and backward elimination methods can be combined so that, at each step, the procedure selects the best attribute and removes the worst from among the remaining attributes.
- ④ Decision tree induction: Decision tree alg's, such as ID<sub>3</sub>, C4.5 and CART, were originally intended for classification. Decision tree induction constructs a flowchart-like structure where each internal (nonleaf) node denotes a test on an attribute, each branch corresponds to an outcome of the test, and each external (leaf) node denotes a class prediction.

## Dimensionality Reduction

- In Dimensionality Reduction, data encoding or data transformations are applied so as to obtain a reduced or "compressed" rep of the original data.
- If the original data can be reconstructed from the compressed data without any loss of information, the data reduction is called lossless.
- If, instead, we can reconstruct only an approximation of the original data, then the data reduction is called lossy.
- Two popular and effective methods of lossy dimensionality reduction:
  - ① wavelet transforms
  - ② Principal component analysis.
- ① wavelet transforms  
The discrete wavelet transform (DWT) is a linear signal processing tech that, when applied to a data vector  $X$ , transforms it to a numerically diff vector  $x'$ , of wavelet coefficients.  
→ The two vectors are of the same length.

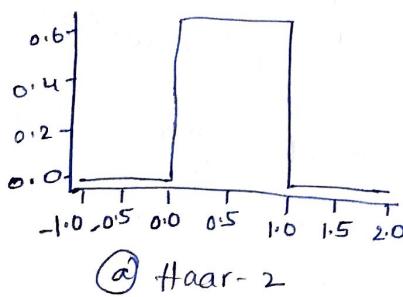
→ When applying this tech to data ~~for redu~~, we consider each tuple as an  $n$ -dimensional vector, that is  $X = (x_1, x_2, \dots, x_n)$ , depicting  $n$  measurements made on the tuple from  $n$  attributes.

→ The usefulness lies in the fact that the wavelet transformed data can be truncated.

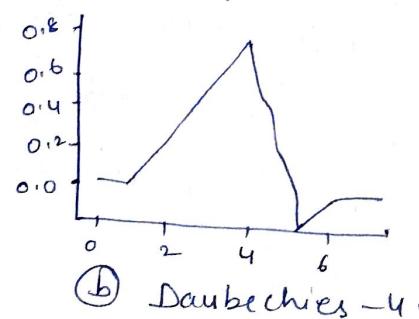
→ A compressed approximation of the data can be retained by storing only a small fraction of the strongest of the wavelet coefficients.

Ex:

→ All wavelet coefficients larger than some user-specified threshold can be retained.  
 → All other coefficients are set to 0.  
 → The resulting data representation is therefore very sparse, so that operations that can take adv of data sparsity are computationally very fast if performed in wavelet space.



② Haar-2



③ Daubechies-4.

## Principal components Analysis

- PCA as a method of dimensionality reduction.
- suppose that the data to be reduced consists of tuples or data vectors described by 'n' attributes or dimensions.
- PCA also called Karhunen -Loeve, or K-L method), searches for 'K' n-dimensional orthogonal vectors that can best be used to represent the data, where  $K \leq n$ .
- the original data are thus projected onto a much smaller space, resulting in dimensionality reduction.

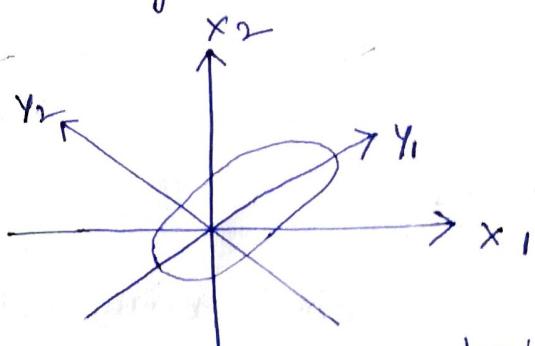


fig: principal components analysis  $y_1, y_2$  are the first two principal components for the given data.

## Basic Procedure

- ① The i/p data are normalized, so that each attribute falls within the same range. This ensures step helps ensure that attributes with large domains will not dominate attributes with small domains.
- ② PCA computes  $K$  orthonormal vectors that provide a basis for the normalized i/p data. These are unit vectors that each point in a direction perpendicular to the others.

These vectors are referred to as the components. The i/p data are a linear combination of the principal components.

③ The principal components are sorted in order of decreasing "significance" or strength. The principal components essentially serve as a new set of axes for the data, providing info about variance. That is, the sorted axes are such that the first axis shows the most variance among the data, the second axis shows the next highest variance and so on.

For ex: In the above fig show the first two principal components  $y_1$  and  $y_2$  for the given set of data originally mapped to the axes  $x_1$  and  $x_2$ . This info helps identify groups or patterns within the data.

④ Be<sup>cuz</sup> the components are sorted according to decreasing order of "significance", the size of the data can be reduced by eliminating the weaker components, that is those with low variance.

## Numerosity Reduction

Techniques of numerosity reduction are used to reduce the data volume by choosing alternative,

"smaller" forms of data representation.

These techniques may be parametric or nonparametric.

① parametric methods is used to estimate the data, so that typically only the data meters need to be stored; outliers may also be stored.

Ex! Log-linear models which estimate discrete multidimensional

probability distributions, are an example.

② Non-parametric methods for storing reduced representations of the data include histograms, clustering and sampling.

Regression and log-linear models

→ Used to approximate the given data.

→ In linear regression, the data are modeled to fit a straight line.

Ex: A random variable,  $y$  (called a response variable),

can be modeled as a linear fun of another

variable,  $x$  (called a predictor variable),

with the eqn =

$$y = wx + b$$

where the variance of  $y$  is assumed to be constant. In the context of datamining,

$x$  and  $y$  are numerical database attributes.

The coefficients,  $w$  and  $b$  (called regression coefficients).

specify the slope of the line and the Y-intercept, respectively.

→ These coefficients can be solved for by the method of least squares, which minimizes the error b/w the actual line separating the data and the estimate of the line.

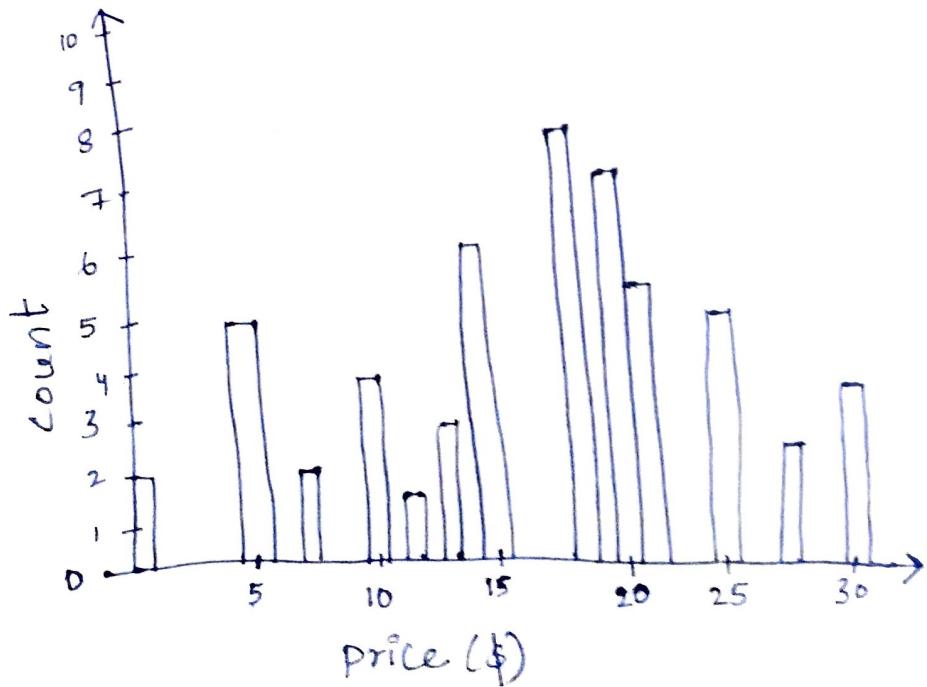
→ Multiple linear regression is an extension of (simple) linear regression, which allows a response variable  $y$ , to be modeled as a linear function of two or more predictor variables.

→ Log-linear models approximate discrete multidimensional probability distributions.

→ Given a set of tuples in  $n$  dimensions (e.g., described by  $n$  attributes), we can consider each tuple as a point in an  $n$ -dimensional space.

## Histograms

- Histograms use binning to approximate data distributions and are a popular form of data reduction.
  - A histogram for an attribute,  $A$ , partitions the data distribution of  $A$  into disjoint subsets, or buckets.
  - If each bucket represents only a single attribute-value / frequency pair, the buckets are called singleton buckets.
- Ex:- The following data are a list of prices of commonly sold items at AllElectronics (rounded to the nearest dollar). The numbers have been sorted : 1, 1, 5, 5, 5, 5, 5, 8, 8, 10, 10, 10, 10, 10, 12, 14, 14, 14, 15, 15, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 18, 18, 18, 20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 25, 25, 25, 25, 28, 28, 30, 30, 30.



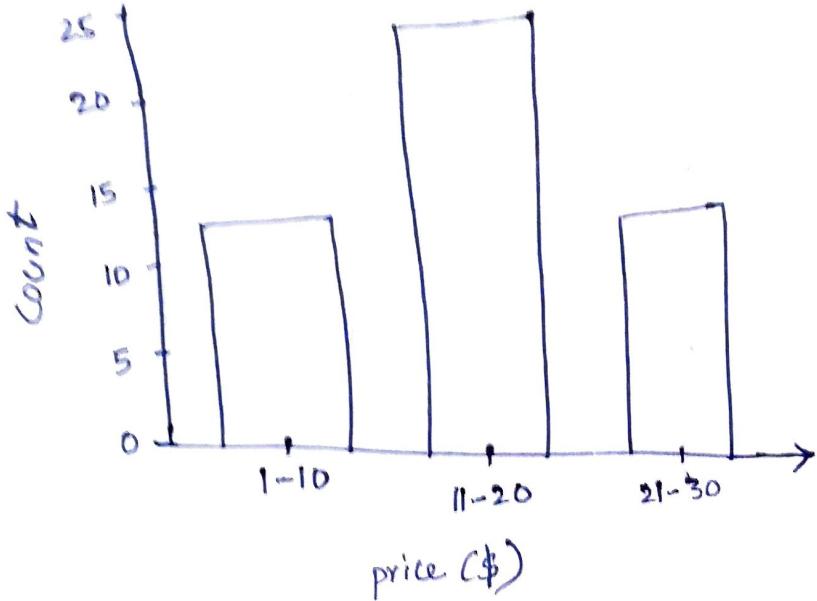


fig: An equal-width histogram for price, where values are aggregated so that each bucket has a uniform width of \$10.

"How are the buckets determined and the attribute values partitioned?"

### Partitioning rules:-

- ① Equal-width:- The width of each bucket range is uniform.
- ② Equal-frequency:- The buckets are created so that, roughly, the freq of each bucket is constant.
- ③ V-Optimal:- If we consider all of the possible histograms for a given no. of buckets, the V-optimal histogram is the one with the least variance.
- ④ MaxDiff:- We consider the diff b/w each pair of adjacent values. A bucket boundary is established b/w each pair for pairs having the  $B-1$  largest differences, where  $B$  is the user-specified no. of buckets.

Clustering :- Clustering techniques consider data tuples as objects. They partition the objects into groups or clusters, so that objects within a cluster are "similar" to one another and "dissimilar" objects in other clusters.

→ the "quality" of a cluster may be represented by its diameter, the maximum distance b/w any two objects in the cluster.

→ centroid distance is an alternative measure of cluster quality and is defined as the avg distance of each cluster object from the cluster centroid.



fig: A 2-D plot of customer data

→ A concept hierarchy for a given numerical attribute defines a discretization of the attribute.

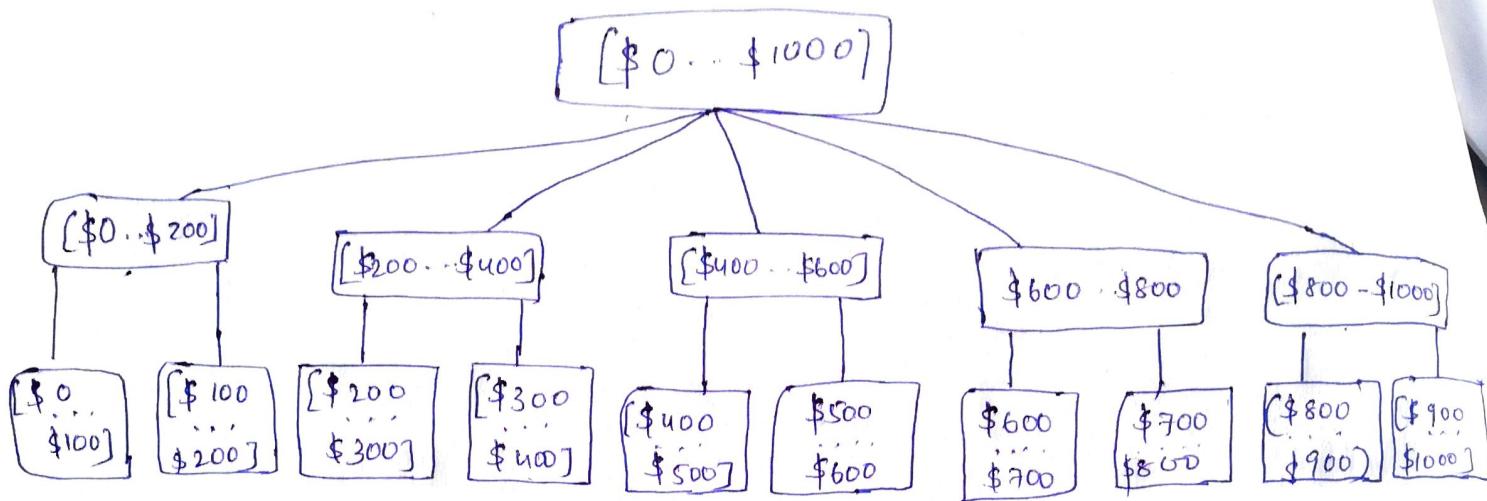


fig: A concept hierarchy for the attribute price, where an interval  $(\$x \dots \$y)$  denotes the range from  $\$x$  (exclusive) to  $\$y$  (inclusive).

→ concept hierarchies can be used to reduce the data by collecting and replacing low-level concepts (such as numerical values for the attribute age) with higher-level concepts (such as youth, middle-aged or senior).

## Smoothing by Binning

(3)

using the data for age

13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 26, 26,

30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 50, 70

So

① sort the data  
 $N = 27$ ,  $D_{\text{depth}} = 3$

$$\textcircled{2} \quad \frac{N}{D} = \frac{27}{3} = 9$$

$\textcircled{3}$  mean or median  $\rightarrow$  (50th percentile)  
 or boundaries [min & max val]

$$\text{mean} = \frac{1}{N} \cdot \sum_{i=1}^N x_i \quad \text{Step 3} \\ \text{Bin 1: } \frac{1}{3} \cdot (44) = 14 \cdot \frac{2}{3} = 14 \frac{2}{3}$$

$$= \frac{1}{27} \quad B_2: \frac{1}{3} (55) = 18 \cdot 33 = 18 \frac{1}{3}$$

$$B_3: \frac{1}{3} (63) = 21$$

smoothing by bin means  
step 1: sorted, by incrementing

$$B_4: \frac{1}{3} (72) = 24$$

$$B_5: \frac{1}{3} (81) =$$

$$B_6: \frac{1}{3} (90) = 30 \cdot 66 \frac{2}{3} = 33 \frac{2}{3}$$

$$B_7: \frac{1}{3} (105) = 35$$

$$B_8: \frac{1}{3} (121) = 40 \cdot 33 = 40 \frac{1}{3}$$

$$B_9: \frac{1}{3} (138) = 46$$

step 4: Bin 1 =  $14 \frac{2}{3}, 14 \frac{2}{3}$ .

$$\therefore \frac{N}{D} = \frac{27}{3} = 9 \text{ bins}$$

Bin 1: 13, 15, 16, ~~16~~, 19, 20

Bin 2: ~~16~~, 19, 20

Bin 3: 20, 21, 22

|| 4: 22, 25, 25

|| 5: 26, 25, 30

|| 6: 33, 33, 35

|| 7: 35, 35, 35

|| 8: 36, 40, 45

|| 9: 46, 50, 70

Ex:

= 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34.

Sol partition into three equal-depth bins

Step 1:  $\text{Bins} = \frac{N}{3} = \frac{12}{3} = 4$

$$\text{Bin 1} = 4, 8, 9, 15 = \frac{36}{4} = 9$$

$$\text{Bin 2} = 21, 21, 24, 25 = \frac{91}{4} = 22.75 = 23$$

$$\text{Bin 3} = 26, 28, 29, 34 = \frac{117}{4} = 29.25 = 29$$

Step 2: smoothing by bin means

Bin 1: 9, 9, 9, 9

Bin 2: 23, 23, 23, 23

Bin 3: 29, 29, 29, 29

smoothing by bin boundaries

Step 3: Bin 1: 4, 4, 4, 4, 15

Bin 2: 21, 21, 25, 25

Bin 3: 26, 26, 26, 34

bin 1 = 4, 8, 15

bin 1 = 4, 4, 15

4, 8  
15

Ex: Suppose that the values for a given set of data are grouped into intervals. The intervals and corresponding frequencies are as follows.

age	freq
1-5	200
5-15	450
15-20	300
20-50	1500
50-80	700
80-110	44

$$\text{lower median value} = \frac{\sum \text{freq}}{2} = 950$$

$$\text{width} = \frac{\text{width of median interval}}{20-50} = 30.$$

compute an approximate median value for the data.

$$\text{Sof} \quad \text{median} = L_1 + \left( \frac{\frac{N}{2} - (\sum \text{freq})_t}{\text{freq}_{\text{median}}} \right) \text{width}$$

we have  $L_1 = 20$  lower boundary

$$\text{freq}_{\text{median}} = 1500$$

$$N = 3194$$

$$\text{width} = 30$$

$$(\sum \text{freq})_t = 950$$

$$\text{median} =$$

100  
20-50

$$\text{median} = 20 + \left( \frac{\frac{3194}{2} - 950}{1500} \right) 30$$

$$= 20 + (0.431) 30 = 20 + 12.94$$

$$= 32.94 \text{ years.}$$

Ex: Min-max normalization.

suppose that the min and max values for the attribute income are \$12000 and \$98000, resp. we would like to map income to the range [0.0, 1.0]. By min-max normalization, a value of \$73,600 for income is transformed to

$$v' = \frac{73,600 - 12000}{98000 - 12,000} (1.0 - 0) + 0 \\ = 0.716$$

$$v = 73,600$$

$$\min_A = 12000$$

$$\max_A = 98000$$

$$\boxed{v' = \frac{v - \min_A}{\max_A - \min_A} (\text{newmax}_A - \text{newmin}_A) + \text{newmin}_A}$$

Ex: Z-score normalization suppose that the mean and standard deviation of the values for the attribute income are \$54000 and \$16,000 resp. with Z-score normalization, a value of \$73,600 for income is transformed to

$$v' = \frac{v - \bar{A}}{\sigma_A} = \frac{73600 - 54000}{16000} = 1.225.$$

Ex: suppose that the data for analysis includes the attribute age. The age value for the data tuples are (in increasing order)

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

- to smooth the data, means to smooth the above data, using a bin depth of 3.
- (a) use smoothing by bin means to smooth the above data, using a bin depth of 3. Illustrate your steps comment on the effect of this technique for the given data.
- (b) How might you determine outliers in the data?
- (c) what other methods are there for data smoothing?

Sol Step 1: sort the data

Step 2: partition the data

into equal - freq bins

Bin 1: 13, 15, 16

Bin 5: 25, 25, 30

Bin 2: 16, 19, 20

Bin 6: 33, 33, 35

Bin 3: 20, 21, 22

Bin 7: 35, 35, 35

Bin 4: 22, 25, 25

Bin 8: 36, 40, 45

Bin 9: 46, 52, 70

Step 3: calculate the  $\bar{m}_{th}$

# Discretization and concept hierarchy generation for numerical data

- Methods :-
- (1) Binning
  - (2) Histogram analysis
  - (3) Entropy-based discretization
  - (4)  $\chi^2$ -merging
  - (5) Cluster analysis
  - (6) discretization by intuitive partitioning.

Binning : Is a top-down splitting tech based on a specified no. of bins used

→ Binning methods are used for data smoothing.

Ex : attribute values can be discretized by applying equal-width or equal-freq binning, and then replacing each bin value by the bin mean or median, as in smoothing by bin means or smoothing by medians, resp.

## Histogram Analysis

→ Like binning, histogram analysis is an unsupervised discretization tech becoz it does not use class info.

→ Histograms partition the values for an attribute A, into disjoint ranges called buckets.

## Entropy-Based Discretization

→ Entropy-based discretization is a supervised, top-down splitting tech.

→ To discretize a numerical attribute, A, the method selects the value of A that has the minimum entropy as a split-point, and recursively partitions the resulting intervals to arrive at a hierarchical discretization.

### Interval merging by $\chi^2$ Analysis

→ Which employs a bottom-up approach by finding the best neighboring intervals and then merging these to form larger intervals, recursively. This merge proceeds as follows. Initially, each distinct value of a numerical attribute A is considered to be one interval.

→  $\chi^2$  tests are performed for every pair of adjacent intervals.

→ Adjacent intervals with the least  $\chi^2$  values are merged together, becoz low  $\chi^2$  values for a pair indicate similar class distributions.

→ This merging process proceeds recursively until a predefined stopping criterion is met.

### Cluster Analysis

→ A clustering alg can be applied to discretize a numerical attribute, A, by partitioning the values of A into clusters or groups.

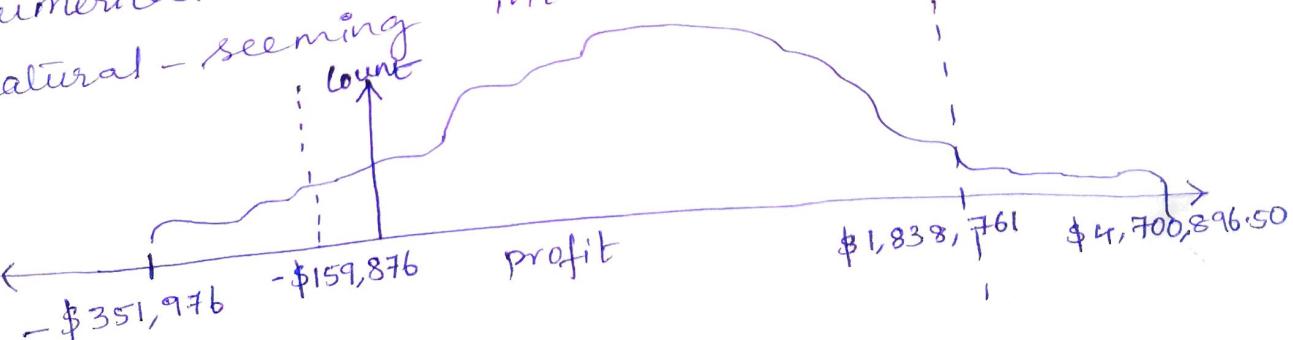
Discretization by Intuitive partitioning

→ Above discretization methods are useful in the generation of numerical hierarchies, many users would like to see numerical ranges partitioned into relatively uniform, easy-to-read intervals.

that appear intuitive or "natural".

Ex: Annual salaries broken into ranges like (\$50,000, \$60,000) are often more desirable than ranges like (\$51,263.98, \$60,872.34), obtained by, say, some sophisticated clustering analysis.

→ The 3-4-5 rule can be used to segment numerical data into relatively uniform, natural-seeming intervals.



Q: Suppose a hospital tested the age and body fat data for 18 randomly selected adults with the following result.

age	23	23	27	27	39	41	47	49	50
%fat	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2
age	52	54	54	56	57	58	58	60	61
%fat	34.6	42.5	28.8	33.4	30.2	34.1	32.9	41.2	35.7

- (a) calculate the mean, median and standard deviation of age and % fat.
- (b) Draw the boxplots for age and % fat.
- (c) Normalize the two variables based on z score normalization.

Sol calculate the mean, median & standard deviation

- (a) calculate the mean, median & standard deviation of age and % fat.
- for the variable age the mean is 46.44  
 the median is 51  
 $SD = 12.85$