

→ 1.7 Need for preprocessing the Data:

Data in the real world is incomplete, noisy & inconsistent

Incomplete — Lacking attribute values of certain attributes of interest or containing only aggregate data

Noisy — Containing errors or outlier values

Inconsistent — Containing discrepancies

All these things occur in large databases & Data-

Warehouses.

Reasons for Incomplete data:-

- > Attributes of interest may not be available
- > Data is considered as unimportant at time of entry
- > Relevant data is not recorded due to misunderstanding
- > Data may be deleted.
- > Modifications of data may be overlooked.

Reasons for Noisy Data:-

- > Data collection instruments may be faulty
- > Error in data transmission
- > Technology limitations — limited buffer size

\* Incorrect data may result from inconsistencies in naming conventions of data code uses or inconsistent format for all fields

\* Data Cleaning — Cleans data by filling in missing values, smoothing noisy data, removing outliers, resolving inconsistencies. If data is dirty mining results in unreliable output. Some mining methods deal with incomplete, noisy data but they are not always robust so preprocessing step is to clean data through some cleaning routines.

\* Data Integration — Integrating multiple databases, data out of file. Attributes have diff names in diff databases, causing inconsistencies and redundancies

cust-so in one data etc. 122 and Customers in movies show  
 so naming inconsistencies arise.

> Large no of redundant data slow down of complex discovery process  
 so Data Integration need to be performed to detect and  
 remove redundancies.

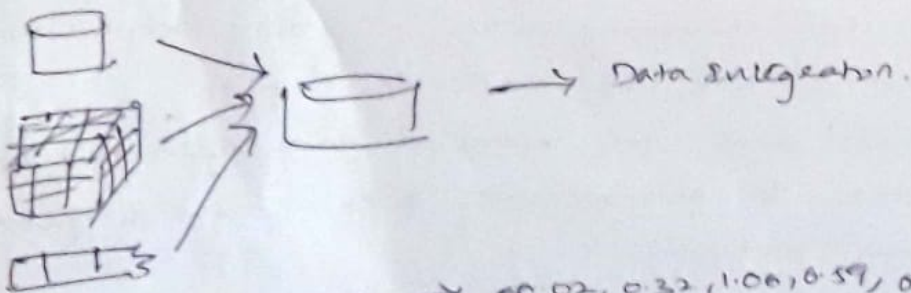
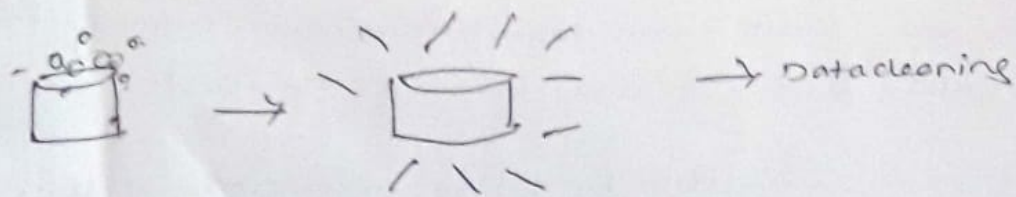
Data Transformation

Neural NW, Nearest neighbour classifiers,  
 clustering provide better results if the data is normalized  
 i.e scaling to  $[0, 1]$ . Data Transformation operations are  
 normalization, aggregation.  $\hookrightarrow$  provides ~~to~~ moves towards success  
 of mining process

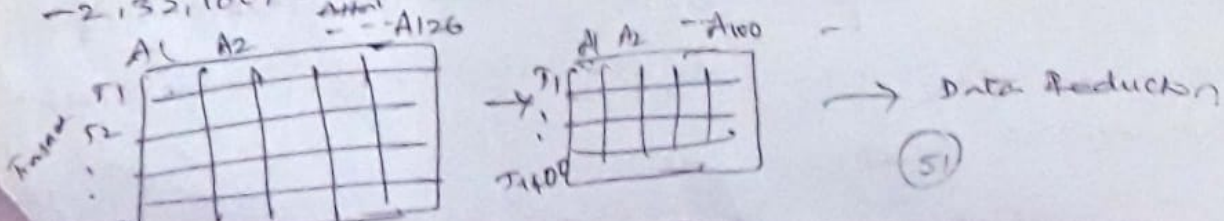
Data Reduction

Obtains a reduced representation of data set i.e  
 much smaller in volume yet produces same analytical  
 results. Includes Data aggregation, — building data cube,  
 Attribute subset selection — removing irrelevant attributes  
 dimensionality reduction — uses encoding schemes  
 Numerosity reduction — replace data by clusters of personnel

Data can also be reduced by generalization with concept hierarchies  
 i.e low level concepts are replaced with high level concepts  
 Data discretization is also form of reduction.



$-2, 132, 1001, 59, 48 \rightarrow 0.02, 0.32, 1.06, 0.59, 0.48 \rightarrow$  Data Transformation



### 1.8 DATA CLEANING:-

This process fills missing values, smooths data, detect outliers, correct inconsistency in data.

#### 1.8.1 Missing values:-

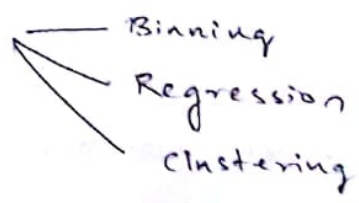
There are some methods for filling missing values

- 1) Ignore tuple → done when class label is missing
  - > not effective unless several attributes are missing
  - > poor when percentage of missing values per attribute varies considerably.
- 2) Fill in Missing values manually →
  - > Time consuming
  - > Not feasible in large data set with many missing values
- 3) Use global constant to fill missing values →
  - > Replace missing attributes with same constant like unknown or -∞
  - > But use may mislead by decision process
  - > Simple, not foolproof
- 4) Use attribute mean to fill missing values →
  - > Calculate mean of attribute and replace missing value of that attribute with that mean
- 5) Use attribute mean for all samples belonging to same class as given tuple →
  - > If we classify customers according to credit risk we can replace mean of income attribute value to credit risk category of given tuple
- 6) Use most probable value to fill missing values →
  - > determined with regression, Bayes formalism or Decision tree Inductors
  - > Most probable one

#### 1.8.2 Noisy Data:-

Noise is a random error var. in a measured variable

Some of Smoothing techniques



Binning :- 4, 8, 15, 21, 21, 24, 25, 28, 34  
 partition into equal frequency bins

> Binning method smooths data  
 sorted data by considering neighbours.

- Bin 1 : 4, 8, 15
- Bin 2 : 21, 21, 24
- Bin 3 : 25, 28, 34

> sorted values are distributed into bins  
 > performs local smoothing

- Smoothly by bin means
- Bin 1 : 9, 9, 9
  - Bin 2 : 22, 22, 22
  - Bin 3 : 29, 29, 29

Each bin is replaced by mean value of bin

- Smoothing by bin boundaries
- Bin 1 : 4, 4, 15
  - Bin 2 : 21, 21, 24
  - Bin 3 : 25, 25, 34

Min, max values in bin are identified as bin boundaries. Each value is replaced by closest boundary value

\* Smoothing by bin median — each bin value is replaced by bin median

Regression :-

- \* Data is smoothed by fitting data to some fn (Regression)
- \* Linear Regression — finding best line to fit 2 attributes so that one attribute is used to predict other.
- \* Multiple linear regression — more than 2 attributes are involved and data are fit to multidimensional surface

Clustering :-

- \* Outliers are detected with clustering where similar values are grouped into cluster.
- \* values that fall outside clusters are outliers.

data cleaning as a process :-

- \* First step is Discrepancy Detection
- \* Discrepancies are caused due to several factors
  - > poorly designed data entry forms that has more fields
  - > Human Error in entry
  - > deliberate errors
  - > data decay (outdated addresses)
  - > Inconsistent data representation
  - > Inconsistent use of code
  - > Error in instrument devices
  - > System Errors
  - > in data integration is many connections

\* Dealing discrepancy detection — we use metadata

- What are domain, datatype of attribute
- What are acceptable values of attribute
- What is range of length of values
- Do all values fall within range

\* We should look at inconsistent code and representation, Field overloading is a source of error: — new attributes are squeezed into unused portion of already defined attribute

- \* Data should be examined regarding unique rule, consecutive rule, Null rule
- unique rule — each value in attribute must be diff from others
- consecutive rule — there should be no missing values b/w lowest & highest values of attribute
- Null rule — handling of blank, question marks, special characters.

Data Inspection

\* There are number of tools for discrepancy detection.

Data Scrubbing tools — use simple domain knowledge to detect and correct errors in data.

Tools with grammar, context — rely on parsing & fuzzy matching for detection of syntax errors & fuzzy matching techniques when cleaning

Data Auditing tools — find discrepancies by analyzing data to discover rules, relationships & detect data that violate such conditions

2nd step

- \* Some data inconsistencies are corrected manually but require data Transformation. — After finding discrepancies we need to define & apply series of transformations to correct them
- Data Migration tools — used for transformations.
- ETL tools — specify it through GUI

\* All above process is error prone & time consuming

- \* There are many new approaches to data cleaning
- Potter's Wheel tool — public tool
- Integrally discrepancy detect & TX
- Discrepancy checking is done automatically
- Users can develop & refine TX
- \* Development of declarative languages for TX operators is a new approach

## 1.9 DATA INTEGRATION AND TRANSFORMATION:

DM requires integration — merging of data from multiple data stores. Data need to be transformed into forms appropriate for mining.

### Data Integration:

- Combines data from multiple sources such as multiple databases, data cubes or flat files
- \* There are no of issues to consider during integration
  - > Schema Integration, object matching are tricky
  - > Entity Identification problem — how entities from one multiple sources are matched?
- ie how analyst is sure of Customer-id in one database and Cust-number in another database are same
- Meta data is used to solve problem. It for each attribute it include name, meaning, data type, range of values permitted for attribute, Null values for handling zero, or null values.

\* Metadata is used to help transform data

Redundancy — Attribute is redundant if it can be derived from another attribute & set of attributes.

Cause of redundancy — Inconsistency in attribute or dimension naming inconsistent

\* Redundancies are detected by correlation analysis ie they measure how strongly one attribute implies other.

For Numerical attr → we use Correlation Coefficient  
(or)  
Pearson product moment Coefficient

$$r_{AB} = \frac{\sum_{i=1}^N (a_i - \bar{A})(b_i - \bar{B})}{N \sigma_A \sigma_B} = \frac{\sum_{i=1}^N (a_i b_i) - N \bar{A} \bar{B}}{N \sigma_A \sigma_B}$$

- N — no of tuples
- $a_i, b_i$  — values of A, B
- $\bar{A}, \bar{B}$  — mean of A, B
- $\sigma_A, \sigma_B$  — SD of A, B
- $\sum (a_i b_i)$  — sum of AB Crossproduct

\*  $-1 \leq r_{AB} \leq +1$

2)  $r_{AB} > 0$  A, B are +vely correlated  
Higher value stronger correlation.  
Higher value indicates A or B may be removed as redundant

$r_{AB} = 0$  A, B are Independent  
ie no correlation

$r_{AB} < 0$  A, B are -vely correlated Data Transfer

\* Scatter plots are Used to view Correlation

\* Not always A causes B or B causes A.  
Ex in demographic database if A is no of hospitals & B is no of thefts in area. There A will not cause B But both are linked to 3<sup>rd</sup> attribute population.

Correlation Categorical Attributes

Correlation is discovered by  $\chi^2$  (Chi square) test

if A has c distinct values  $a_1, a_2, \dots, a_c$   
B has r " " "  $b_1, b_2, \dots, b_r$

Contingency table — has c values of A as columns  
r values of B as rows

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

Pearson  $\chi^2$  statistic

$O_{ij}$  — observed frequency  
 $E_{ij}$  — Expected "

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

## 2x2 Contingency table

$(r-1) \times (c-1)$   
 $\downarrow$   
 degree of freedom

	male	female	Total
Fiction	250 (90)	200 (360)	450
Nonfiction	50 (210)	1000 (840)	1050
Total	300	1200	1500

$\bar{e}_{ij}$  in parameters — expected frequencies  
 we can verify expected frequency

$$e_{ij} = \frac{\text{count}(\text{male}) \times \text{count}(\text{fiction})}{N} = \frac{300 \times 450}{1500} = 90$$

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.93$$

## Third Issue: Detection and resolution of data value conflicts

Attribute value from diff sources may differ  
 i.e. weight is stored in metric units in one system and British Imperial units in another.

## Data Transformation:-

data are transformed into appropriate form for mining.

It involves

- > Smoothing — remove noise
  - ← Missing
  - ← Regression
  - ← Clustering
- > Aggregation — summary of aggregation operations are applied.  
 Daily sales is aggregated to compute monthly and annual total amount.
- > Generalization — low level data are replaced by higher level by their concept hierarchies.  
 Street is generalized to City, Country  
 age " " " young, Middle, Senor
- > Normalization — data are scaled so as to fall within specific range i.e. -1.0 to 1.0 or 0.0 to 1.0
- > Attribute Construction — New attributes are constructed and added  
 (57) — to help mining process.

Smoothing: user specify Transformations to correct data inconsistencies.

Normalization — Attribute is normalized by scaling its value so they fall within specified range such as 0.0 to 1.0.  
useful for classification alg involving nn, nearest neighbors classifiers and clustering.

In NN normalizing step for each attribute speed up learning phase.  
In Distance based methods — larger range values are scaled to smaller range (ie binary attributes)

Min-Max Normalization

performs linear x on original data  
 $min_A, max_A$  — minimum and maximum value of attribute A  
A value of A is mapped to  $v'$  in range  $[new-min_A, new-max_A]$

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

Ex min, max of income are \$12,000 and \$98,000  
Range is [0.0, 1.0]

\$73,600 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

values are normalized based on mean & standard deviation of A  
v is mapped to  $v'$  by computing

$$v' = \frac{v - \bar{A}}{\sigma_A} \quad \bar{A} - \text{mean} \quad \sigma_A - \text{SD}$$

this method is useful when min, max are unknown

Ex — Mean of A \$54,000  
SD of A \$16,000

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

Normalization by decimalScaling:

normalizes by moving decimal points of value of attribute A  
No of decimal points moved depend on maximum absolute value of A. v is transformed to  $v'$

$$v' = \frac{v}{10^j} \quad j = \text{smallest integer such that } \text{Max}(|v'|) < 1$$

Ex - Values of A range from -986 to 917.

Min absolute value of A is 986

$\therefore d = 3$  (1000)

-986 is normalized to -0.986  
 917 " " " 0.917

Attribute Construction - new attributes are constructed from given attributes and are added to improve accuracy & understandability of structure in high dimensional data

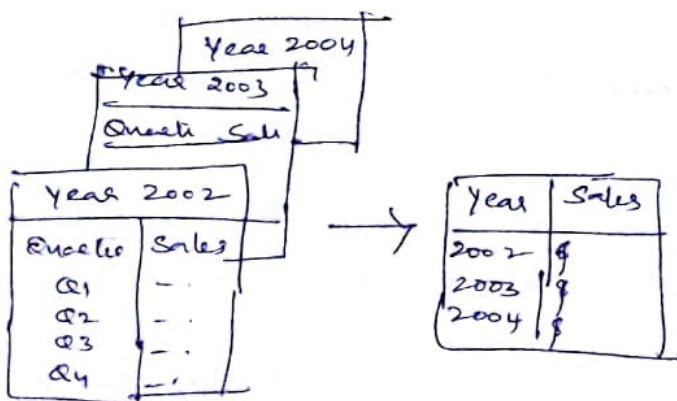
## 1-10 DATA REDUCTION :-

Complex data analysis and mining on huge amounts of data can take a long time, making such analysis impracticable or infeasible. Data reduction is applied to obtain a reduced representation of data set is much smaller in volume yet maintain integrity of original dataset. Thus mining on reduced data set should be more efficient yet produce same analytical results

Strategies of Data reduction

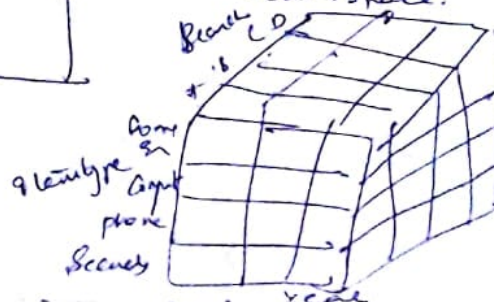
Data Cube Aggregation :- Aggregation operations are applied to the data in construction of data cube

Sales data from 2002 to 2004 is present which is in sales per quarter but we are interested in annual sales. So data can be aggregated. Resulting data is smaller in volume w/o loss of information necessary for analysis.



Data Cube store multidimensional aggregated information

\* Each cell holds aggregated data value corresponds to data point in multidimensional data space.



\* Concept hierarchies allows data to be analyzed at multiple levels of abstraction. (19)

Data Cube provides fast access to precomputed, summarized data which benefits OLAP & Data Mining

Cube at low level — Base Cuboid  
↳ Correspond to Individual Entry

at High level — Apex Cuboid

Ex-It may give total sales for all 3 years for all item types & for all branches

Data Cube — called as lattice of Cuboids

### Attribute Subset Selection:

Data set may contain hundreds of attributes, many of which may be irrelevant to task or redundant. The domain expert can pick useful attributes but which is difficult & time consuming.

Attribute subset selection reduces data size by removing irrelevant or redundant attributes

Goal — To find minimum set of attributes such that resulting probability distribution of data is as close as possible to original distribution.

How to find good subset of attributes —

for  $n$  attributes there are  $2^n$  possible subsets

\* Heuristic methods are used for Attribute Subset Selection

↳ Greedy methods

↳ Make locally optimal choice which lead to globally optimal solution.

\* Best (or worst) attributes are determined using statistical tests which assume attributes are independent or one another

Some of heuristic methods are

Stepwise forward Selection — Starts with Empty set of attributes  
Best attributes are determined & added to reduced set

(20) At each step best of remaining attributes is added to set

Ex Initial attribute set  $\{A_1, A_2, A_3, A_4, A_5, A_6\}$

Initial reduced set:  $\{ \}$

$\{A_3\}$

$\{A_1, A_4\}$

reduced attribute set:

$\{A_1, A_4, A_6\}$

Stepwise Backward Elimination — starts with full set of attributes  
— at each step removes worst attribute in set

Ex — Initial attribute set  $\{A_1, A_2, \dots, A_6\}$

$\Rightarrow \{A_1, A_3, A_4, A_5, A_6\}$

$\Rightarrow \{A_1, A_4, A_5, A_6\}$

$\Rightarrow \{A_1, A_4, A_6\}$

Combination of forward selection & Backward Elimination — The procedure selects the best attribute & removes worst attribute from any remaining attributes

Decision Tree Induction — alg such as ID3, C4.5, CART are used for classification.  $\downarrow$  constructs flowchart like structure where each internal node denotes test on attribute. Each branch corresponds to the outcome of test & each leaf node denotes class prediction. At each node alg chooses "best attribute" to partition data into individual classes.

Attribute in tree form reduced subset of attribute.

### Dimensionality Reduction:-

Encoding &  $\neq$  are applied to get reduced representation of original data

Lossless — if original data is reconstructed from compressed data w/o any loss of info then it is lossless

Lossy — if we reconstruct only an approximation of original data then it is lossy

Two popular and effective methods of lossy dimensionality reduction are wavelet transforms

& principal component analysis.

Curse of dimensionality  
when dimens  $\uparrow$  data become sparse

Density, distance b/w points are critical to clustering, outliers etc.

(5)

Dimensionality reduction

- avoid curse of dimension
- reduce time & space in
- eliminate irrelevant features
- Reduce noise

## Wavelet Transforms :- (Used for Image Compression) 22

\* Discrete wavelet Transform (DWT) is linear signal processing technique when it is applied to data vector  $X$  it is transformed to  $X'$

$X' \rightarrow$  with wavelet Co-efficients

$X$  and  $X' \rightarrow$  same length

$X = (x_1, x_2, \dots, x_n)$  is  $n$  dimensional vector ( $n$  attributes)

\* Both vectors are of same length but still this method is useful for data reduction

> Here wavelet transformed data can be truncated

> Compressed data is retained only with strongest wavelet coefficients for  $n$  coefficients larger than user specified threshold are retained, all other coefficients are set to zero (0).

> Now data is sparse and operations are performed very fast in wavelet space

> wavelet  $px$  is used to remove noise and smoothing.

\*  $\Rightarrow$  DWT is closely related to DFT - discrete Fourier Transform which involves Sines & Cosines.

DWT gives better lossy compression i.e. DWT provides more accurate approximation of original data than DFT, & also requires less space than DFT.

\* There are several families of DWT such as Haar-2 and Daubechies-4

\* General procedure uses hierarchical pyramid algorithm and is given below

① Length  $L$  of 2D vector must be integer power of 2  
If necessary pad vector with zeros as necessary

② Each TX involves two fn  $\left\{ \begin{array}{l} \text{data smoothing} \\ \text{weighted difference} \\ \text{(bring detailed feature of data)} \end{array} \right.$

(62)

- ③  $2$   $f^n$  are applied to all pairs of data which results in  $2$  sets of data of length  $1/2$
- ④  $2$   $f^n$  are recursively applied to data obtained in previous loop until datasets are of length  $2$ .
- ⑤ Selected values obtained from above loops are designated as wavelet Co-efficients of  $x^k$  data.

\* Matrix multiplication is applied to  $2$ lp data to get Wavelet Co-efficients

Matrix — orthogonal

↳ cols are unit vectors

↳ orthogonal re inverse is its transpose

\* Atq has complexity of  $O(n)$   
 $n$  —  $2$ lp vector

\* wavelet  $f^k$  are applied to datacube. First apply  $f^k$  to  $1^{st}$  dimension, then to second & so on

### Principal Component Analysis (PCA)

- \* Data consist of vectors described by  $n$  attributes  $\Rightarrow$  dimension
- PCA also called as Karhunen - loève or K-L method
- Search for  $k$   $n$ -dimensional orthogonal vectors ( $k \leq n$ )
- \* Original data is projected into smaller space

#### Basic procedure

①  $2$ lp data is normalized so attribute falls within <sup>same</sup> range  
 so that attributes with large domain will't dominate attribute with smaller domains.

② PCA computes  $k$  orthogonal vectors  
 ↳ unit vectors  $\perp$  to each other.  
 known as principal components

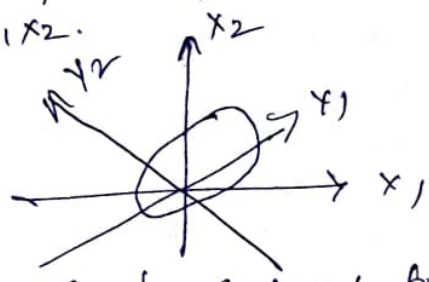
$2$ lp data is linear combination of principal components

3) Principal Components are sorted in decreasing significance of strength. These serve as new set of axes providing info about variance.

1st axis — More variance

2nd axis — next highest variance

$y_1, y_2$  — principal components of data originally mapped to  $x_1, x_2$ .



4) Size of data can be reduced by eliminating weaker components. i.e. with low variance

\* PCA is computationally inexpensive

\* Can be applied to ordered & unordered data sparse and skewed data

Numerosity Reduction:-

parametric

Nonparametric

Do not assume model for storing reduced representation of data include histogram, clustering & sampling.

a model is used to estimate data. only parameters are stored instead of actual data

Regression and log linear Models:-

used to approximate given data

linear Regression — data is modeled to fit in st. line

$y$  (Random variable) is modeled as linear fn of  $x$  (random variable Predictor / Response)

$y = wx + b$

$w$  &  $b$  regression co-efficients

$x$  &  $y$  numerical database attributes

Co-efficients are solved by method of least squares

Multiple linear regression — linear fn of 2 or more predictor variables

Loglinear models — approximate discrete multidimensional probability distribution

$n$  dimensions are present

Each tuple is considered as point in  $n$ -dimensional space

\* These models are used to estimate probability of each point in multidimensional space

\* used to construct higher dimensional dataspace from lower dimensional space

\* Also used for dimensionality reduction, data smoothing  
↓  
lower dimensional point occupy less space than original

\* used on sparse data

Histograms:-

\* Use binning to approximate data distributions

\* attribute  $A$  is partitioned into disjoint subsets of buckets  
If each bucket has single value then it is called singleton buckets

Partitioning rules :-  
 $w = \frac{B-A}{N}$   
↓  
low & high attribute

① Equal-width — width of each bucket range is uniform

② Equal frequency or equal depth — frequency of each bucket is constant (same no. of samples)

③  $v$ -optimal —  $v$ -optimal histogram is one with least Variance. Histogram variance is weighted sum of original values of each bucket.

④ Max-Diff — consider all bin adjacent values. Bucket boundary is established s.t. each pair for pairs have  $p-1$  largest difference

$p$  — user specified no. of buckets

(65)

\* V-optimal, MaxDiff — more accurate

Clustering:-

- \* Considers data-tuples as objects. partition objects into clusters - objects within clusters are similar to each other and dissimilar to other objects in other clusters.
- \* Similarity is based on distance fn  
 Quality of cluster — its diameter — mean distance b/w object in clusters  
 — Cluster centroid (average point in space for cluster)
- \* clustering is more effective for data that can be organized into distinct clusters than for smeared data
- \* Multidimensional index trees are used for providing fast data access which is used for hierarchical data reduction and also provide approximate answers to queries
- \* Index tree recursively partitions multidimensional space with root node representing entire space. Tree are balanced. parent node has keys and pointers to child nodes. leaf node has pointers to data-tuples they represent
- \* Index trees store aggregate as well as detailed data at diff levels of abstraction, which provides hierarchy of clusters.  
 Each child node is bucket then tree is hierarchical histogram
- \* There are many measures for defining clusters and cluster quality.

Sampling:-

↳ one form of Data Reduction  
 • larger dataset is represented by smaller random sample (subset) of data

Data set → D (56)  
 Total tuples → N

We draw 's' sample of N tuples from D

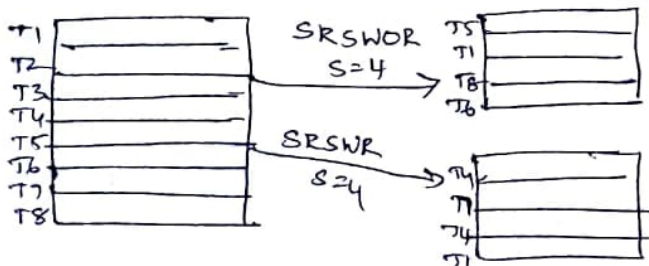
$$s < N$$

probability of drawing any tuple is  $\frac{1}{N}$

ie all tuples are equally likely to be sampled

Simple Random Sample with Replacement (SRSWR) of size s :-

Similar to SRSWOR but each time a tuple is drawn from D it is recorded and replaced. ie it placed back in D and so we can draw it again.



Cluster Sample :-

D is grouped to M disjoint clusters then SRS of 's' clusters are obtained where  $s < M$

Ex - tuples in database are retrieved one page at a time so each page is considered as clusters.

Stratified Sample :-

- > D is divided into mutually disjoint parts called "strata"
- > Apply SRS on each stratum and generate stratified sample
- > useful when data is skewed

Ex - a stratified sample is obtained from customer data where each stratum is created for each customer age group.

Advantage

- > Cost of obtaining sample is proportional to size of sample s instead of N
- > Sampling complexity is "sublinear" to size of data  
↳ increases as <sup>(67)</sup> no of dimensions increases.
- > Sampling is refinement of reduced data set. Such set ... sample size.

## 1.11 DATA DISCRETIZATION AND CONCEPT HIERARCHY

### GENERATION :-

- \* Data discretization techniques — reduce no of values for given continuous attribute by dividing range of attribute into intervals. Labels are used to replace actual data values. So that it simplifies original data.
- \* Supervised discretization      \* unsupervised discretization
  - ↓ uses class info
  - ↓ doesn't use class info
- \* Top down discretization (or) splitting
  - ↳ starts with one or few points (split point) to split range and repeat this recursively on resulting intervals
- \* Bottom up discretization (or) merging — contrast to splitting
- \* Discretization can be performed recursively on attribute to provide hierarchical partitioning of value known as Concept hierarchy which is useful for mining at multiple levels of abstraction
  - ↳ one form of data reduction where low level concepts are replaced by higher level concepts
- \* Manual definition of Concept hierarchies is tedious task and time consuming and so several discretization methods are used to automatically generate Concept hierarchies for Numerical Attributes and for Categorical it can be generated automatically.

### Discretization and Concept Hierarchy Generation

#### for Numerical Data.

Concept hierarchies for numerical attributes can be constructed automatically based on data discretization.

## Binning :-

- \* Top-down splitting technique based on no of bins
- \* Attribute values are discretized by applying equal-width (or) equal-frequency and then applying replace each bin value by bin-mean or median
- \* Apply these techniques recursively to resulting partitions to generate Concept hierarchies
- \* Doesn't use class information and so called as unsupervised discretization technique
- \* Sensitive to  $\left\{ \begin{array}{l} \text{presence of outliers} \\ \text{user specified no of bins} \end{array} \right.$

## Histogram Analysis :-

- \* unsupervised discretization technique
- \* Histogram partition value of A into buckets with equal widths (or) equal frequency.  
range  $\downarrow$  same no of tuples
- \* Apply histogram analysis algorithm recursively to generate <sup>multiple</sup> Concept hierarchy. Recursive procedure stops when it reaches pre-specified no of Concept levels (or) when it reaches minimum interval size per level.  
 $\downarrow$  no of values for each partition.

## Entropy Based Discretization

- \* Entropy — one discretization measure by claud Shannon  
Supervised, Top down Splitting
- \* To discretize a numerical attribute 'A' method selects value of 'A' that has minimum Entropy as a splitpoint & recursively partition the resulting intervals to get hierarchical discretization.

(69)

is further refined by procedure samples

\*  $D$  <sup>has</sup> set of data tuples <sup>(24)</sup> & class label attribute <sup>✓</sup>  
 provides class info per tuple

\* Basic method for Entropy based discretization,

① Considers splitpoint to partition range of  $A$   
 $A$  can partition tuples in  $D$  into 2 subsets satisfying  
 Condition  $A \leq \text{splitpoint}$  and  $A > \text{splitpoint}$ , which  
 creates a binary discretization.

② If we want to classify tuples in  $D$  by partitioning on a  
 attribute  $A$  & some splitpoint

2 classes  $\left\{ \begin{array}{l} C_1 \\ C_2 \end{array} \right.$

1st partition

Contains tuples of  $C_1$  but  
 also some of  $C_2$

Still how much info is needed for perfect classification.

This info is called Expected information for classifying tuple  
 in  $D$  based on  $A$ .

$$\text{Info}_A(D) = \frac{|D_1|}{|D|} \text{Entropy}(D_1) + \frac{|D_2|}{|D|} \text{Entropy}(D_2)$$

$D_1, D_2$  correspond to  $A \leq \text{splitpoint}$  &  
 $A > \text{splitpoint}$

$|D|$  — no of tuples in  $D$

on classes  $C_1, C_2, \dots, C_m$

$$\text{Entropy}(D) = - \sum_{p=1}^m p_p \log_2(p_p)$$

$p_c$  — probability of class  $C$  in  $D$

\* process of determining a splitpoint is recursively applied  
 to each partition until stopping criteria is met.



minimum info requirements on all  
 candidate splitpoint  $<$  threshold  $\epsilon$

Or no of intervals  $>$  threshold max. interval

## Interval Merging by $\chi^2$ Analysis:

- \*  $\chi^2$  based method, Bottom-up approach - find best neighbouring intervals & merge to form larger intervals recursively
- \* If two adjacent intervals have similar distribution of classes then intervals can be merged. Else separate procedure
- \* each value of A is considered as one interval
- \*  $\chi^2$  test is performed on all adjacent intervals. We merge intervals with least  $\chi^2$  value. (Similar class distribution)  
It is repeated until some stopping criteria is met.
- \*  $\chi^2$  used in data integration and we can construct a contingency table for our data
  - └ 2 columns (2 adjacent intervals)
  - └ m rows (no. of distinct classes)

We use  $O_{ij}$ ,  $E_{ij}$

### \* Stopping criteria

- ①  $\chi^2$  values for all pairs exceeds some threshold  
high value of significance — over discretization  
low " " — under discretization  
Significance level ↓ 0.10 to 0.01
- ② No. of levels intervals cannot be over prespecified max. (10 to 15)
- ③ class frequencies must be considered within interval

## Cluster Analysis:

- \* Discretize numerical attribute A by partitioning values of A into clusters
  - └ considers closeness of datapoints so that it produces high quality discretization methods,

- \* follows either Top-down or Bottom up strategy  
Each cluster forms a node
- \* Initial clusters may be further decomposed into several subclusters forming lower level of hierarchy.  
con clusters are formed by repeatedly grouping neighbouring clusters to form high level concepts

### Discretization by Intuitive partitioning:-

- \* we would like to discretize numerical ranges partitioned into relatively uniform, easy-to-read intervals  
Annual salaries broken into ranges like ( $\$50,000, \$60,000$ ) are desirable than ranges like ( $\$51,263.98, \$60,872.34$ )
- \* 3-4-5 rule — to segment numerical data into relatively uniform, natural intervals  
partitions given range of data into 3, 4 or 5 equal width intervals recursively and level by level based on MSB

Rule is as follows

- > If interval covers 3, 6, 7 or 9 distinct values at MSB then partition range into 3 intervals  
for 3, 6, 9  $\rightarrow$  3 equal width intervals

- > If it covers 2, 4, or 8 distinct values at MSB then partition to 4 equal width intervals

- > If it covers 1, 5 or 10 distinct values at MSB then partition range into 5 equal width intervals.

Rule is recursively applied to each interval creating concept hierarchy

$\rightarrow \$351,976.00$

$\$4,700,896.50$

(72)

\* user need automatic generation of hierarchy of profit

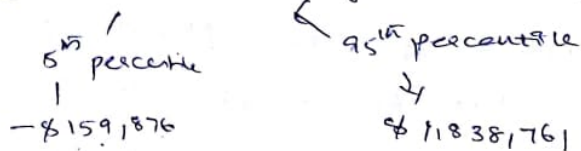
$[-\$1,000,000 \dots \$0]$

data b/n  $5^{th}$  and  $95^{th}$  percentile are  $\$159,876$  &  $\$1,838,761$

Results are .

- ① Min =  $-\$351,976.00$   
 MAX =  $\$4,700,876.50$

5<sup>th</sup> low and high values for 1<sup>st</sup> level of discretization are



- ② msd is at million digit position  
 msd = 1,000,000

rounding low =  $-\$1,000,000$   
 high =  $+\$2,000,000$

- ③ We have to do 3 partitions

$$\frac{\text{HIGH} - \text{LOW}}{\text{msd}} = \frac{2,000,000 - (-1,000,000)}{1,000,000} = 3$$

3 equal width segments

$(-\$1,000,000 \dots \$0)$   $(\$0 \dots \$1,000,000)$  &  $(\$1,000,000 \dots \$2,000,000)$

### CONCEPT HIERARCHY GENERATION FOR CATEGORICAL DATA:-

Categorical data — histogram — discrete — data

- ↓
- finite no of distinct values
- No ordering among values

## Specification of a partial ordering of attributes Explicitly at the Schema Level by Users or Experts :-

> Involve group of attributes

> We can define concept hierarchy by specifying partial (or) total ordering of attributes at schema level

location dimension of Datawarehouse contain following attributes — Street, city, province or state and Country.

Total ordering is  $\text{Street} < \text{city} < \text{province or state} < \text{Country}$ .

## Specification of position of hierarchy by explicit data grouping

In large databases it is difficult & unrealistic to define entire concept hierarchy. So we specify explicit groupings of ~~for~~ a small portion of intermediate level data.

Ex After specifying province and country form hierarchy at schema level we can use some intermediate levels such as { Alberta, Saskatchewan, Manitoba } & { prairies-Canada and { British Columbia, prairies-Canada } & Western Canada.

## Specification of set of attributes but not of partial ordering:

User specify set of attributes forming concept hierarchy but omit their partial ordering. System automatically generate attribute ordering to construct hierarchy.

ie Based on no. of distinct values per attribute in given set it automatically generates hierarchy

Country <sup>with</sup> small no distinct values than Street.  
Attribute with most distinct values — at Low level of hierarchy  
" " " — top level

Specification of only partial set of attributes:

Sometimes users have vague idea of what should be included in hierarchy, so he can include only a small subset of relevant attributes in the hierarchy specification.

For ex he may include only street, city. To handle such partially specified hierarchies we should embed data semantics in database, so so attribute with tight semantic connections can be pinned together.

7th 1st home B  
100 (0, 61, 62, 71)  
221 77, 81, 82, 83,  
88, 89, 92, 93,  
94, 95, 96, 97, 98, 99

85, 86

18th 2nd home  
100  
59, 58, 60, 65, 69, 76,  
77, 87, 88, 92, 93,  
94, 95, 97, 98, 99, 100

(75)

## ARCHITECTURE OF DATA MINING SYSTEMS:-

A good system architecture will facilitate the system to make best use of software environment which accomplish DM tasks in an efficient and timely manner interoperate and exchange the information with other information systems be adaptable to user requirements and evolve with time.

### Desired Architecture for DM systems:-

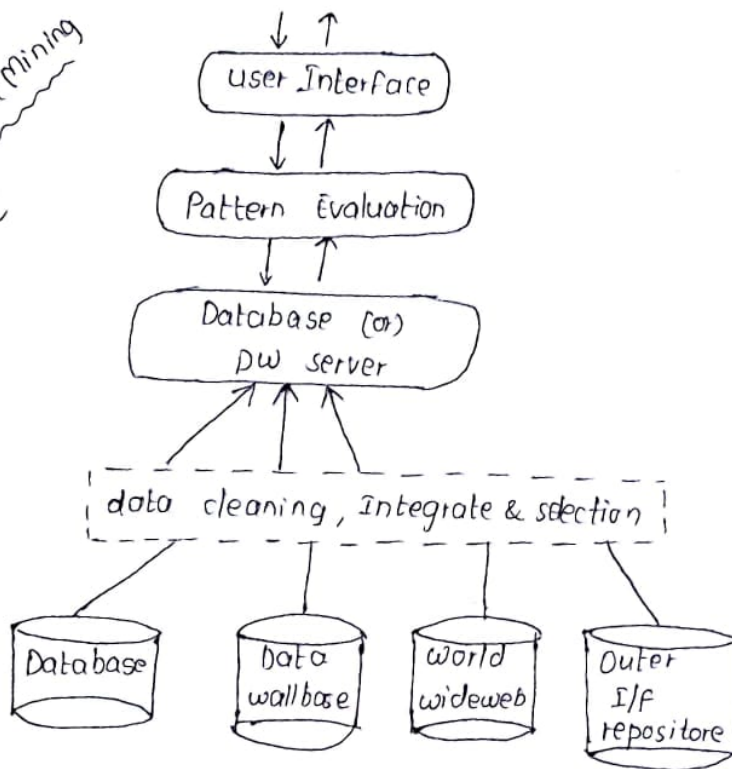
Based on different architecture designs a DM system can be integrated with a DB/DW system using following coupling systems.

- a) No Coupling:- A DM system will not utilize any function of a DB (or) DW system. It fetches data from particular source process data and by using some mining algorithm and then store the mining results in the another file.
- b) Loose Coupling:- Loose coupling means that a DM system will use some facilities of a DB (or) DW system, fetching data from a data repository managed by these system performing data mining and then storing the mining results either in a file (or) in a designated place in database (or) DW.
- c) Semitight Coupling:- Semitight coupling means that besides linking a DM system to DB/DW system efficient implementation of a few essential data mining primitives can be provided in DB/DW system. These primitives can include sorting, and precomputation of some essential statistical measures such as sum, count, max and min so on.
- d) Tight coupling:- Tight coupling means that a DM systems is smoothly integrated into DB/DW system. The DM subsystem is treated as functional component of an

information system. DM queries and functions are optimized based on mining query analysis, data, structure, indexing schema & query process methods of DB/DW system.

DM, DB and DW system will evolve and integrate together as one I/f system with multiple functionalities which provides a uniform I/f processing system.

Architecture of  
Typical data mining  
System



Knowledge Base:- It is the knowledge which is used to guide the search (or) evaluate the interestingness of resulting patterns. There are diff kinds of knowledge such as concept hierarchies  
 - To organise attribute values into different abstractions  
 - To access pattern's interestingness.

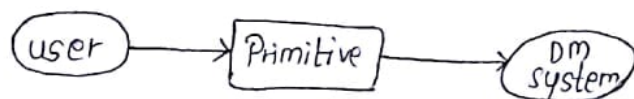
Data Mining Engine - It consists of set of functional modules for tasks such as characterization, association & correlation analysis, classification, prediction, clusters analysis, outlier analysis and so on...

Pattern-Evaluation Module - It employs interestingness measures and interacts with DM module so as to focus the search towards interesting pattern.

User Interface: - This module communicates b/w user and data mining system allowing user to interact with system by specifying on Query (or) task, providing i/p to help focus search and performing exploratory data mining based on the intermediate DM results.

### DATA MINING TASK PRIMITIVES:

- \* Each user will have a determining task in mind i.e some form of data analysis that he (or) she would like to have performed.
- \* DM task is specified in form of DM Query which is i/p to DM system.
- \* DM Query is defined in terms of DM task primitives.
- \* These primitives allow user to interactively communicate with the DM system during discovery in order to direct the mining process.



Incorporating these primitives in an DM Query Language

- More flexible user interaction
- foundation for design of GUI
- Standardization of DM industry and practise

- \* DM primitives specify following
  - Task-relevant data
  - Type of knowledge to be mined
  - Background knowledge
  - Pattern interestingness (or) measurements
  - Visualization of discovered pattern

### Task-relevant data to be mined

\* It specifies position of database (or) set of data in which is interested.

- It includes
- ✓ Database (or) Datawarehouse name
  - ✓ Database tables (or) datawarehouse cubes
  - ✓ Condition for data selection
  - ✓ Relevant attributes (or) dimensions
  - ✓ Data grouping criteria

\* Task relevant data in general known as Mined view. The portion of database to be mined is called mined view.

### Example

If determining task is to study association b/w items frequently purchased of All electronics by customers in Canada, the task relevant data can be specified as

- name of database (or) DW to be used (eg. All electronics db)
- name of tables (or) data cubes containing relevant data (eg. item, customer, purchases, item sold)
- conditions for selecting relevant data (eg. relative data pertaining to purchases made in Canada for current year.

→ schema hierarchy (order among attributes in DB schema)

eg. street < city < province < country  
or state

→ set-grouping hierarchy (given attribute/dimension groups)

eg {20...39} = young

{40...59} = middleaged

→ operation derived hierarchy

email address: dmbook@cs.sfu.ca

login name < dept < university < country

operations specified by user, operations can include the decoding of int encodes string, int extraction complex data objects and data cluster

→ Rule based hierarchy

low-profit-margin(x) < = price(x, P1) and  
Cost(x, P2) and (P1 - P2) < \$50

Information about domain to be mined usefully in discovery process.

⇒ Measurements of Pattern Interestingness

They are used to guide mining ~~process~~ process.  
Different kinds of knowledge may have different interesting measures

Ex Interesting-measures for association rule include Support and Confidence

Rules whose Support & Confidence value are below user-specified threshold are considered uninteresting.

information system ...

- relevant attributes (or) dimensions  
(eg name & price from itemtable  
& income age from costumertable)

Type of knowledge to be mined:-

It determines data mining function to be performed such as

characterization

Discrimination

Association

Classification/predictor

Outlier analysis

other DM tasks

**Example**

A user studying buying habits of All Electronics, customers may choose to mine association rule of the form

$P(X: \text{customer}; w) \cap Q(X; Y) \text{ buys } (X; Z)$

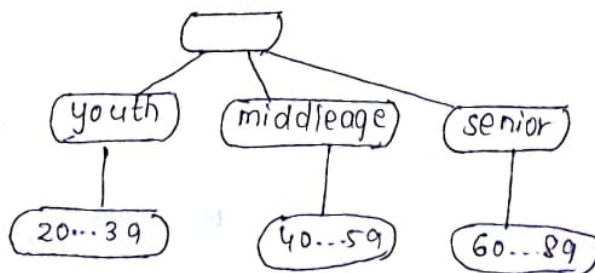
Eg:

$\text{age}(X; [30\ 39]) \cap \text{income}(X; [40\ 50k]) \text{ buys } (X; [VCR])$

[22%, 60%]

Background knowledge

This kind of knowledge is useful for guiding the knowledge discovery process and for evaluating patterns found. concept hierarchies are popular form of Background knowledge which allow data to be mined at multiple levels of abstraction.



## Simplicity

Eg: (Association) rule length, decision tree size,  
Can be viewed as function of pattern structure,  
defined pattern size in bits, (or) no of attributes  
or operations appeared in pattern

✓ Certainty - which assess validity or trustworthiness  
of pattern.

Eg:- Confidence  $p(B|A)$  - for Association Rule  
classification reliability or Accuracy

Certainty factor,

Rule Strength

Rule Quality

Discriminatory weight etc.

✓ Utility — Eg Support =  $p(A \cup B) = \frac{\text{tuple count of } A \cup B}{\text{total tuple}}$

Potential Usefulness

Support (Association)

Noise Threshold (Description)

✓ Novelty — that contributes new info or increased  
performance to give not previously known  
Surprising (used to remove redundant rule).

# Visualization of Discovered patterns

This refers to form in which discovered patterns are to be displayed i.e. d/f Backgrounds/ usages many require d/f forms of presentation.

Eg:- Rules, Tables, Crosstabs, pie/bar chart etc

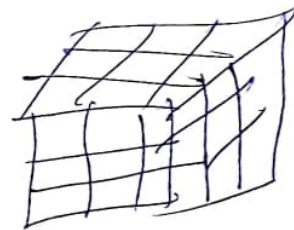
## Rules

age(x, "young") and income(x, "high")  $\Rightarrow$  class(x, "A")  
 age(x, "young") and income(x, "low")  $\Rightarrow$  class(x, "B")  
 age(x, "old")  $\Rightarrow$  class(x, "C")

## Tables

age	income	class
young	high	A
young	low	B
old	high	C
old	low	C

## DataCube



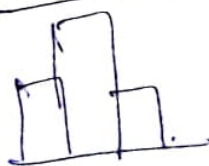
## Crosstab

age	income		class		
	high	low	A	B	C
young	83	-	1		
old	-	-			
Count					

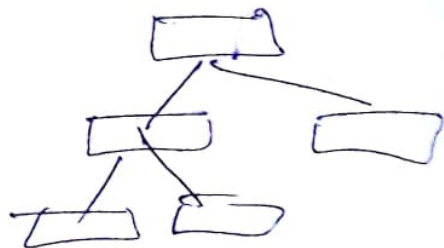
## Pie-chart



## Barchart



## Decision tree



Visualization of discovered patterns in various forms can help users with different backgrounds to identify patterns of interest and to interact can guide the system in future discovery ~~form~~ of visual forms as rules, tables, cross tabs, pie chart, bar chart, decision tree.

⑤

• relevant attributes (or) dimensions

(eg name & price from item table & income, age from customer table)

Type of knowledge to be mined :-

It determines data mining function to be performed

- Such as
- Characterization
- Discrimination
- Association
- Classification/prediction
- Clustering
- Outlier analysis
- Other DM tasks



Example

A user studying buying habits of All electronics customers may choose to mine association rule of the form

$$p \subseteq X : \text{customers} ; N \cup \{ (X ; Y) \} \text{ buys } (X ; Z)$$

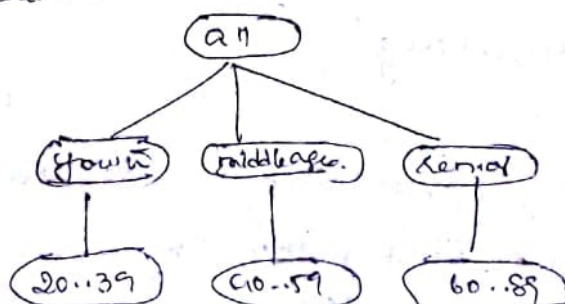
eg:

$$\text{age } (X ; \{30-39\}) \wedge \text{income } (X ; \{40-50k\}) \Rightarrow \text{buys } (X ; \{vcr\})$$

[2.2% | 60%]

Background knowledge

This kind of knowledge is useful for guiding the knowledge discovery process and for evaluating patterns found. concept hierarchies are popular form of Background knowledge which allow data to be mined at multiple levels of abstraction.



✓ schema hierarchy (order among attributes in DB schema)

Eg Street < city < province < country  
or etc

✓ set grouping hierarchy (given attribute / dimension groups)

Eg \$20-39y = young  
\$40-59y = ~~Senior~~ middle aged

✓ operation derived hierarchy

email address : dmbook @ cs.sfu.ca

loginname < dept < university < country

Operations specified by users, operations can include the decoding of int encodes string, int extraction etc. Complex data objects and data clusters

✓ Rule based hierarchy

low-profit-margin (X) < = price (X, P1) and cost (X, P2) and (P1 - P2) < \$50

Df about domain to be mined useful in discovery. Proc

Measurements of pattern Interestingness :-

They are used to guide the mining process.

Df. kinds of knowledge may have dift interesting measures.

Ex Interesting measures for association rule include support & confidence

Rules whose support & confidence values are below user-specified threshold are considered uninteresting.

✓ Simplicity

Eg (association) rule length, (decision) tree size  
simplicity = CNF - no of conjuncts

Can be viewed as m<sup>th</sup> of pattern structure, defined pattern size in both no of attributes (or) operations appear in rule

\* D1

(4)

✓ Certainty - which assess validity (or trustworthiness) of pattern

Eg Confidence  $P(A|B)$  - for Association rules  
classification reliability (or accuracy)  
Certainty factor  
rule strength  
rule quality  
Discriminating weight etc

✓ utility - Eg Support =  $P(A \cup B) = \frac{\text{tuple count of } A \cup B}{\text{total tuple}}$

min lot of errors  
min supp etc

potential usefulness  
support (association), noise threshold (descriptor)

✓ Novelty - that contribute new info (or increased performance to give pattern set)

not previously known  
Surprising (used to remove redundant rules)

### Visualization of Discovered patterns

This refers to forms in which discovered patterns are to be displayed i.e. diff backgrounds (images) may require diff forms of presentation

Eg :- Rules, tables, Correlabs, pie bar chart etc

#### Rules

age (x, "young") and income (x, "high")  $\Rightarrow$  class (x, "A")  
age (x, "young") and income (x, "low")  $\Rightarrow$  class (x, "B")  
age (x, "old")  $\Rightarrow$  class (x, "C")

#### Table

age	income	class	count
young	high	A	1242
young	low	B	1342
old	high	C	856
old	low	C	752

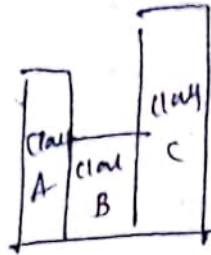
#### Correlab

age	income		class		
	high	low	A	B	C
young	83				
old	85				
Count					

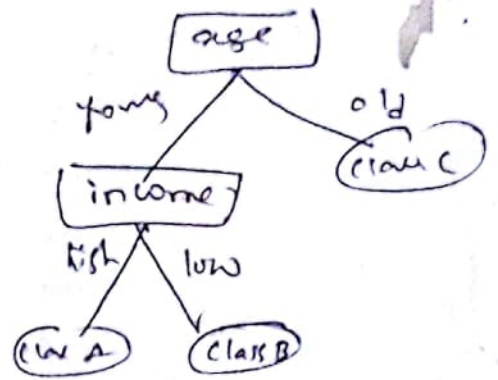
### Pie chart



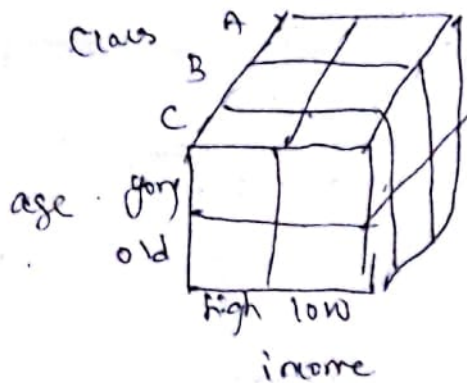
### Bar chart



### Decision tree



### Data cube



Visualization of discovered patterns in various forms can help users with diff backgrounds to identify patterns of interest and to interact with the system in further discovery. Some of the visual forms are rules, tables, cross tables, pie chart, bar chart, decision tree etc.

## DATA MINING QUERY LANGUAGE (DMQL)

- \* DMQL can provide the ability to support ad-hoc and interactive data mining in order to facilitate flexible & effective DM duty.
- \* By providing standardized language like SQL DMQL has hope to achieve similar effect like SQL on relational database.
- ✓ foundation for system development & evolution.
- ✓ facilitate its exchange, technology transfer, commercialization & wide acceptance etc.

\* DMQL is designed with the primitives described earlier

## Syntax for DMQL

Syntax for specification of

- task relevant data
- kind of knowledge to be mined
- interestingness measures
- pattern presentation & evaluation

\* If we put all these primitives together it is called a DMQL Query

\* In an inductive database (IDB) ordinary queries can be used to access and manipulate data while inductive queries can be used to generate, manipulate and apply patterns. IDB becomes an Extended Querying process

## Features of DMQL

- Data mining query language satisfy the closure property
- A classical approach in database theory
- Relation algebra
- Has Optimization Strategies
- Supports Association rule mining process.

Syntax Syntax is Specification of Task relevant data

- Use database database-name (or) Use datawarehouse datawarehouse-name
- from relation(s)/cube(s) [where condition]
- in reference to attr - con dim-list
- Order by order-list
- group by grouping-list
- having condition

Example: This Example shows how to use DML to specify last relevant data for mining of association on items frequently purchased at AllElectronics by Canadian Customers w.r. to customer income & age. In addition the user specifies that she would like data to be grouped by date. The data are retrieved from relational database

Use database AllElectronics-db  
 in relevance to I.name, I.price, c.income, C.age  
 from Customer C, item I, purchase P, itemsold S  
 where I.item-ID = S.item-ID and  
 S.trans-ID = P.trans-ID and  
 P.Cust-ID = C.Cust-ID &  
 C.address = 'Canada'

group by P.date

Syntax: kind of knowledge to be mined:

■ Characterization — specifies characteristic descriptions to be mined

Mine-knowledge-Specification ::=

mine characteristics [as pattern-name]

analyze measure(s)

■ Discrimination

Mine-knowledge-Specification ::=

mine comparison [as pattern-name]

for target-class where target-condition

& versus contrast-class-1 where contrast-condition-1 &

analyze measure(s)

⑥

Eg

Mine comparison as purchase groups  
for big spenders where  $\text{avg}(\text{I-price}) \geq \$100$   
versus budget spenders where  $\text{avg}(\text{I-price}) < \$100$   
analyze count.

### Association

Mine-knowledge-specification ::=

Mine associations [as patternname]  
[matching < metapattern >]

Eg:

Mine associations as buying habits  
matching  $P(X:\text{custom}, W) \wedge Q(X, Y) \Rightarrow \text{buys}(X, Z)$

### Classification

Mine-knowledge-specification ::=

Mine classification [as patternname]  
analyze classifying-attribute-con-dimension

### Other patterns

Clustering, outliers analysis, predictor

### Syntax: Concept Hierarchy Specification

\* To specify what concept hierarchies to use

use hierarchy < hierarchy > for < attribute-or-dimension >

\* we use the syntax to define the type of hierarchies.

### Schema hierarchies

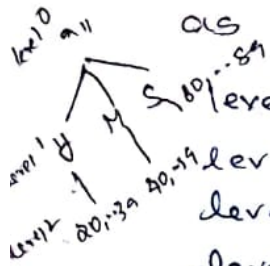
define hierarchy time-hierarchy on date as [date, month, quarter, year]



predefined concept hierarchy for schema  
date (day, month, quarter, year) as structure

## Set-grouping hierarchies

define hierarchy age-hierarchy for age on Customer as



level 1 : { young, middleaged, Senior } < level 0 : all

level 2 : { 20, ..., 39 } < level 1 : young

level 2 : { 40, ..., 59 } < level 1 : middle-aged

level 2 : { 60, ..., 89 } < level 1 : Senior

## Operation-derived hierarchies

define hierarchy age-hierarchy for age on Customer as

{ age-category (1), ..., age-category (5) } :=

clusters (default, age, 5) < all (age)

These clustering alg is to be performed on all of ages values in Customer table to form 5 clusters. age-categories! } (cluster is) runs

## Rule based hierarchies

define hierarchy profit-margin hierarchy on item as

level-1 : low-profit-margin < level-0 : all

if (price - cost) < \$ 50

level -1 : medium-profit margin < level-0 : all

if (price - cost) < \$ 50 and ((price - cost) <= \$ 250)

level -1 : high-profit margin < level-0 : all

if (price - cost) > \$ 250

## Specification of Interestingness Measures

Interestingness measures and thresholds can be specified by user with Strml

with <interest-measure-name> threshold = threshold-value

eg with Support threshold = 0.05

with confidence threshold > 0.7

## Specification of pattern presentation

- \* specify the display of discovered pattern  
display as  $\langle \text{result-form} \rangle$
- \* To facilitate interactive viewing at diff concept levels  
we use rollup, drill down options  
the following syntax is used.

Multilevel\_Manipulation := rollup on attribute (or) dimension  
| drill down on attribute (or) dimension  
| add attribute (or) dimension  
| drop attribute (or) dimension

### Putting it all together : A DMQL Query

Use database AllElectronics-db

Use hierarchy location-hierarchy for B. address  
mine characteristics as customer purchasing  
analyze count%

in relevance to c.age, I.type, I.place-made  
from customer C, item I, purchase P, Items-Sold S,  
Works-at W, branch.

Where I.item\_ID = S.item\_ID and S.trans\_ID =  
P.trans\_ID and P.cust\_ID = C.cust\_ID and  
P.method\_paid = "AmEx"  
and P.Empl\_ID = W.Empl\_ID and W.branch\_ID =  
B.branch\_ID and B.address = "Canada"  
and I.price >= 100

with noise threshold = 0.05

display as table

Based on all primitives we design a Query language re  
DMQL

## Other Data Mining Languages & Standardization Efforts:

### Association Rule Language Specifications

✓ MSRL (Imielinski & Viswanath 1997)

✓ MineRule (Meopoulos and Ceri '96)

✓ Query floccle based on Datalog Syntax (Gsur et al '98)

### MSRL - Mini Structured Query Language:

#### Syntax

GetRules (C) [INTO <rulebase name>]

[WHERE <rule constraints>]

[SQL - groupby clause]

[USING encoding-clause]

Where C - source table

rulebase name - name of the object

rule constraints - conditions on desired rules

### MINE Rule

It extracts association rule b/n values of attributes in relational table

#### Syntax:

MINE RULE <tablename> AS

SELECT DISTINCT [<cardinality>] <Attributes>

AS BODY [<cardinality>] <Attributes>

AS HEAD [, SUPPORT] [, CONFIDENCE]

FROM <table> [WHERE <where clause>]

GROUP BY <Attributes> [HAVING <having clause>]

[CLUSTER] & V <Attributes>

[HAVING <having clause>]]

EXTRACTING RULES WITH

SUPPORT : <real>, Confidence <real>

# DMQL Syntax

(8)

Use database < databasename >

Use hierarchy < hierarchyname >

<DMQL> ::= <DMQL-statement>; & <DMQL-stmt>}

<DMQL-statement> ::= <Data-Mining-statement>  
<concept-hierarchy-definition-statement>  
<visualization-and-presentation>

<Data-Mining-statement> ::=

Use database <databasename> use datawarehouse  
<datawarehouse-name>

Use hierarchy <hierarchy name> for <attribute-or-dimension>

Use mine-knowledge-specification

in relevance to <attribute-or-dimension list>

from <relation(s) | cube>

[where <condition>]

[order by <order list>]

[group by <grouping-list>]

[having <condition>]

Use with [ <interest-measure-name> ] threshold =

<threshold-value> [for <attribute(s)>]

<Mine-knowledge-specification> ::= <Mine-char> | <Mine-discr> |  
<Mine-assoc> | <Mine-class> | <Mine-pred>

all above  
which are given before

The design of effective DMQL requires a deep understanding of power, limitation and underlying mechanisms of various kinds of DM tasks

• OLEDB for DM (Microsoft' 2000)

- ✓ Based on OLE, OLEDB, OLEDB for OLAP
- ✓ Integrating DBMS, datawarehouse & DM

• CRISP - DM (Cross-Industry standard process for DM)

- ✓ providing a platform and process structure for effective data mining
- data mining
- ✓ Emphasize on deploying DM technology to solve the business problems

\* OLE DB for DM is an extension of OLEDB Aps to access database systems.

Syntax:-

```
CREATE MINING MODEL <DMM name>
  ( <columns definition> )
  USING <algorithm>
  [ ( <algorithm parameters> ) ]
```

\* DMQL adopts SQL-like syntax so that it can be easily integrated with relational query language

Syntax of DMQL is defined in an Extended BNF grammar where [ ] represents zero (0) or one occurrence

{ } - 0 or more occurrences

- DM Query is parsed to form an SQL query which retrieves data and then perform the operation.

# UNIT-4

\* Data Mining

- descriptive Data Mining
- predictive Data Mining

①

Descriptive data mining — describes dataset in concise & summarize manner and present interesting & general properties of data

predictive Data Mining — Analyzes data in order to construct one or more set of models and attempts to predict the behaviour of data sets

\* Databases usually store large amount of data and users would like to view the data in concise & summarized (succinct), descriptive terms. These descriptions provide overall picture of data, and can

\* Users like the ease and flexibility of having datasets described at different levels of granularity and from different angles. Such descriptive data mining is called Concept description

## CONCEPT DESCRIPTION :- (CD)

\* It is simplest kind of descriptive data mining.

Concept — collection of data such as frequent buyers, graduate students & so on

↳ general descriptions for characterization & comparison of data

\* CD sometimes called as class description when concept to be described refers to class of objects.

Characterization — provides concise & summarization of given collection of data.

Discrimination — provide description comparing 2 or more collection of data.

CD involve both characterization and discrimination

- \* CD is closely related with generalization. i.e. when is large amount of data we are able to describe data in concise and generalized rather than at low levels of abstraction.
- \* If we allow data to be generalized at multiple levels of abstraction we can examine general behaviour of data

Example

Instead of viewing or examining individual customer's sales manager may prefer to view data generalized to higher levels such as customer groups according to geographic regions frequency of purchase per group and so on  
 multilevel data generalization = multidimensional data analysis in DW

- \* The differences b/w CD in large database & OLAP involve the following
  - Complex datatypes and Aggregation
  - User Control Versus Automation

Complex datatypes & Aggregation :-

\* DW & OLAP tools are based on multidimensional data model that view data in form of data cube which consist of dimension (or attributes), & measures. Datatypes are measure are restricted to some extent. Many current OLAP systems confine dimensions to nonnumeric data. Measures such as counts, sums, averages apply only to numeric data

\* In CD database attributes are of various type such as numeric, nonnumeric, spatial, text or images

\* Aggregation of attributes in database include complicated datatypes.

\* CD handles complex datatypes of their attributes & aggregation if necessary.

User Control Versus Automation

\* OLAP in DW is usercontrolled process i.e. all operations are directed & controlled by users. Although it is userfriendly it requires good understanding of role of each dimension.

\* CD strives for more automated process which helps to determine which dimension should be included in analysis, degree to which

\* If a data set should be generalized in order to produce summarization of data (2)

## DATA GENERALIZATION AND SUMMARIZATION - Based CHARACTERIZATION

The data & objects in database contains detailed int at primitive levels. i.e. for example item selection in a sales database contain items no, brand, category, place, price & so on. So we have to summarize large set of data and present at higher conceptual levels.

Data generalization is a process that abstracts large set of task-relevant data in database from low conceptual level to higher conceptual levels.

Data generalization approach include

- a) Data cube - based Data aggregation
- b) Attribute Oriented Induction.

(a) Data cube based on aggregation approach was discussed in 2nd unit.

### (b) Attribute - Oriented Induction for Data characterization

It is to first collect the task-relevant data using a database query and then perform generalization based on no of distinct values of each attribute in relevant set of data.

Generalization is perform by either Attribute removal  
(or)  
Attribute generalization

\* Attribute oriented induction approach is Query oriented, generalization based, on-line data analysis technique.

\* The essential operation of Attribute oriented Induction is data generalization.

### Attribute Removal :-

If there is large set of distinct values for an attribute in initial working relation, either

- (a) There is no generalization operator on attribute (i.e. there is no concept hierarchy defined for attribute)
- (b) Its higher level concepts are expressed in terms of other attributes, then attribute should be removed from working relation.

An attribute value pair represents a conjunct in a generalized tuple (or) rule. The removal of conjunct eliminates a constraint and thus generalizes the rule.

In case (a) — If there is no generalization operator on attribute should be removed  $\because$  it can't be generalized and has large set of disjuncts.

### Attribute Generalization :-

If there is large set of distinct values for an attribute in initial working relation and there exist set of generalization operators on attribute, then generalization operator should be selected and applied it to attribute.

Use of generalized tuple operators to generalize an attribute within tuple will cover more of original data tuples. This corresponds to rule known as climbing generation trees in learning from examples (or) concept tree ascension.

The control of how high an attribute should be generalized is typically quite subjective. Control of this process known as Attribute Generalization Control.

If attr is generalized too high it is overgeneralized and rules may not be very informative.

\* If attribute is not generalized to sufficiently high level then it is undergeneralized.

So some balance should be attained in Attribute Oriented Induction

These are mainly 2 ways to control a generalization process. They are

- attribute generalization threshold control
- generalized relation threshold control

Attribute generalization threshold control :-

- \* We can set one generalization threshold for all attributes (or set one threshold for each attribute)
- \* If no of distinct values is greater than of attribute is greater than attribute threshold further attribute removal (or attribute generalization) should be performed.
- \* DM systems have typically default threshold value b/w 2 to 8. We can modify these threshold values

If we want to further generalize relation we can reduce the threshold (roll up along attribute)

If we feel generalization reaches too high we can increase the threshold (drill down along attribute)

Generalized relation threshold control :-

- \* This technique sets threshold for generalized relation. If
- \* If no of tuples in generalized relation is greater than threshold then further generalization should be performed.
- \* Default threshold value is b/w 10 to 30. Also it may be adjustable.

These two techniques are applied in sequence, i.e. 1st apply attribute threshold control technique to generalize each attribute and then apply relation threshold control to further reduce size of generalized relation.

Many users are interested in statistical info at diff levels. i.e. we need to accumulate (total) of counts & other aggregated values. For each tuple we have count value and an initial working relation it is set to '1'.

Now we perform Attribute removal (or) generalization and then Merge all identical tuples. This group count value is equal to sum of all tuples merged.

### Example of Attribute oriented Induction:-

Let us take initial working relation is a collection of task relevant data

name	Gender	major	birth-place	birthdate	residence	phone #	gpa
X	M	CS	vanowca	8-12-76	AAAA-	6871598	3.67
Y	M	CS	Newyork	28-2-70	BBBBB	853-9106	3.76
Z	F	physic	leade	21-3-75	CCCCC	420-5232	3.83
:	:	:	:	:	:	:	:

for each attribute generalization proceeds as follows

1) name → It has large no of distinct values and there is no generalization operator defined for it so this attribute is removed.

2) Gender → It has only 2 distinct values and so it is retained and no generalization is performed.

3) Major → let us think major is generalized to values of arts & science, engineering, business? and let threshold is set to 5 and there are more than 20 distinct values for major in initial working relation. major is generalized by climbing up concept hierarchy.

Birth-place → It has large no of distinct values so we need to generalize it. let us think concept hierarchy exist city & province (or) state & country. If no of distinct values for country is  $>$  attribute generalization threshold then attribute should be removed.

If no. of distinct values for country is  $\leq 4$  less than threshold then birth-place can be generalized to country.

5) Birth-date  $\rightarrow$  Suppose hierarchy exists that can generalize birth-date to age and age to age-range. If no. of age range's interval is too small w.r. to threshold then generalization should take place.

6) Residence  $\rightarrow$  Residence is defined by no., street, city, residence state, residence country. No. of distinct value for street, no. are high since they lie at very low level.  $\therefore$  attribute no., street should be removed.  $\therefore$  residence is generalized to residence-city which has fewer distinct values.

7) Phone #  $\rightarrow$  It has too many distinct values and therefore removed in generalization.

8) GPA  $\rightarrow$  Suppose concept hierarchy exist for gpa that groups values for grade point average into numerical value intervals like  $\{3.75-4.0, 3.5-3.75\}$  etc. These are grouped into descriptive values such as {excellent, very good}. The attribute can be generalized.

ii Generalized relation is.

<u>gender</u>	<u>major</u>	<u>birth-country</u>	<u>age-range</u>	<u>residence-city</u>	<u>gpa</u>	<u>count</u>
M	Science	Canada	20-25	Richmond	very-good	16
F	Science	Foreign	25-30	Burnaby	excellent	22
:	:	:	:	:	:	:

Identical tuples are grouped into one and count is accumulated.

## Efficient Implementation of Attribute Oriented Induction :-

Algorithm: Attribute Oriented Induction - Mining generalized characteristics in a relational database given user's DM request.

Input:

- ✓ DB - relational database
- ✓ DMQuery - a data mining query
- ✓ a-list - a list of attributes
- ✓ GenList - set of concept hierarchies with generalization operators on attribute  $a_i$ .
- ✓ a-gen-threshold( $a_i$ ) - attribute generalization threshold for each  $a_i$ .

Output:

P - a prime generalized relation.

Method:

Step 1:  $W \leftarrow \text{get-task-relevant data (DMQuery, DB)}$

Let W be working relation that hold task relevant data.

Step 2: prepare for generalization(W)

- (a) Scan W and collect distinct values for each attribute  $a_i$ . If W is large take sample of W
- (b) For each attribute  $a_i$  determine whether  $a_i$  should be removed and if not compute min desired level  $L_i$  based on given threshold and determine mapping pair  $(v, v')$ 
  - $v$  - distinct value of  $a_i$  in W
  - $v'$  - corresponding generalized value at level  $L_i$

Step 3:  $P \leftarrow \text{generalization}(W)$

P is derived by replacing each value  $v$  in W by  $v'$  in mapping rule while accumulating count & calculating of other aggregated values

- (a) For each generalized tuple, insert the tuple into sorted prime relation P by binary search - If tuple is already in P simply increase count and other aggregated values & insert into P

(b) Mostly no. of distinct values at prime relation level is small  
So it can be coded as m-dimensional array where m is no. of attributes in P. and each dimension has generalized values.

Efficiency of algorithm is as follows:-

- Step 1 — It is a query to collect data relevant data into W.  
Its processing efficiency depends on query processing methods used
- Step 2 — collect statistics on W. It requires scanning W at most once.  
Cost for entire step is dependent on no. of distinct values for each attribute
- Step 3 — There are N tuples in W and P tuples in P. A B  
for each tuple t in W we substitute value based on mapping.  
This results in generalized tuple t'  
✓ If variation (a) is adopted each t' takes  $O(\log P)$  to find location for count increment. ∴ total time complexity is  $O(N \times \log P)$  for all generalized tuples.  
✓ If variation (b) is adopted each t' takes  $O(1)$  to find tuple for count increment. ∴ overall time complexity is  $O(N)$  for all generalized tuples.

Presentation of Derived Generalization:

\* Generalized descriptions are most commonly displayed in form of generalized relation (or table).

\* Descriptions can be visualized in form of cross tabulation or cross-tabulation.

In cross tab

each row — value from another attribute

column — value from another attribute

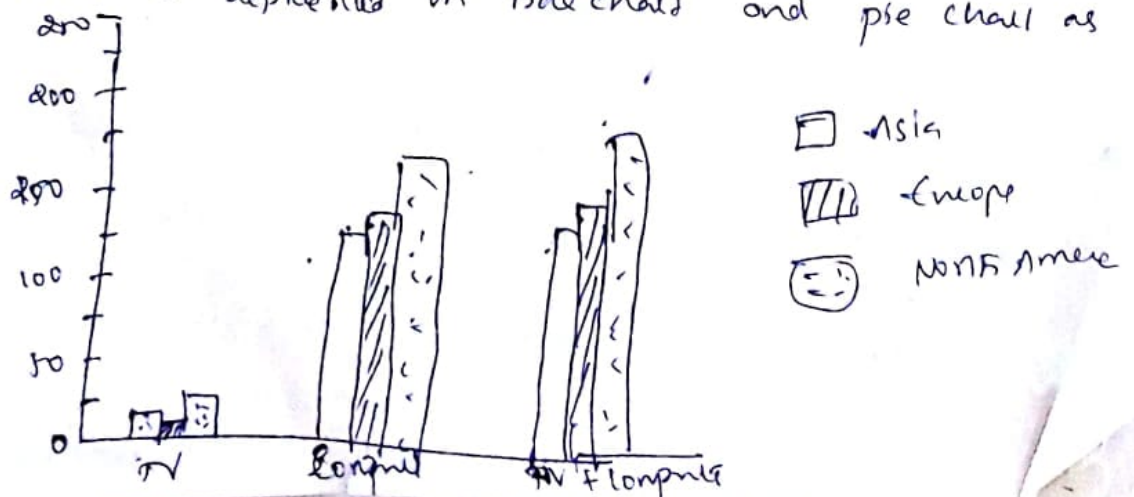
Location	Items	Sales (in million dollars)	Cost (in thousand)
Asia	TV	15	800
Europe	TV	12	250
North America	TV	28	450
Asia	Computer	120	1000
Europe	Computer	100	1200
North America	Computer	200	1800

A Crosstab for sales in 2004

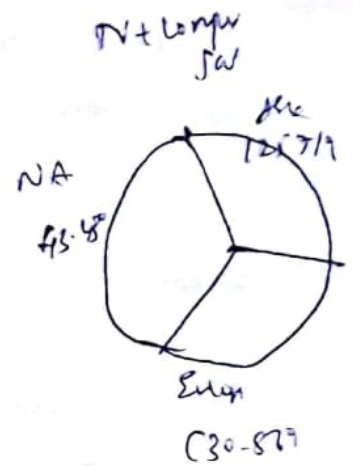
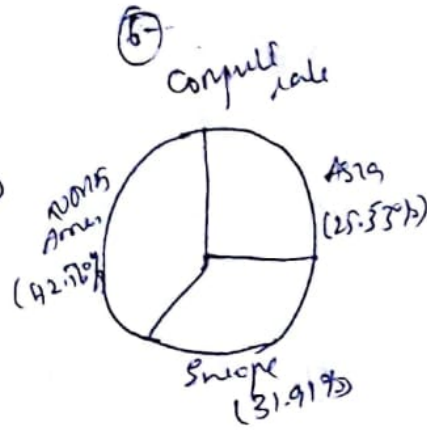
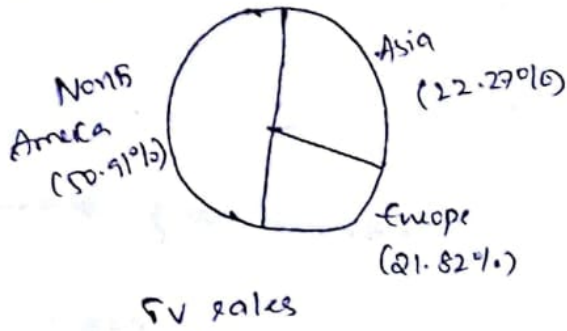
Location	Item					
	TV		Computer		Total	
	Sale	Cost	Sale	Cost	Sale	Cost
Asia	15	800	120	1000	135	1800
Europe	12	250	100	1200	112	1450
North America	28	450	200	1800	228	2250
All regions	45	1000	420	4000	525	5000

\* Generalized data can be presented graphically using bar charts, pie charts and curves. Visualization with graphs is popular in data analysis. Such graphs can represent 2D or 3D data.

The above table is represented in bar chart and pie chart as follows



## Pie-Chart



## Cube View

The Cube is shown for the dimensions item, location and cost.

The size of cell represents the count of corresponding cell. The brightness of cell can be used to represent another measure of the cell such as sum(sales).

\* pivoting, drilling, slicing operations can be performed on the cube browser by mouse clicking.

## Quantitative Rule:-

\* Data in a large database usually span a diverse range of distributions, a single generalization tuple is unlikely to cover (or represent 100% of) initial working relation tuples on cases.

\* Quantitative information such as percentage of data tuples that satisfy left and right hand side of rule should be associated with each rule.

\* A logic rule that is associated with Quantitative information is called a Quantitative rule.

$$t\text{-weight} = \text{count}(a_n) / \sum_{i=1}^n \text{count}(a_i)$$

where

$t$ -weight as each tuple in corresponding generalized relation

$a_n$  be the generalized tuple describing target class

$n$  — no of tuples for target class in the generalized relation.

The range for  $t$ -weight is  $[0.0, 1.0]$  or  $[0\%, 100\%]$

### Quantitative characteristic rule

(1) In logic form by associating corresponding  $t$ -weight value with each disjunct covering target class.

(2) In relational table can convert form by changing the count values in these tables for tuples of the target class to corresponding  $t$ -weight values.

\* Each disjunct of rule represents a condition. The disjunction of these conditions form a necessary condition of target class.

\* All tuples of target class must satisfy the condition.

Rule should be expressed in the form

$$\forall x, \text{target\_class}(x) \Rightarrow \text{Condition}_1(x) [t:w_1] \vee \dots \vee \text{Condition}_m(x)$$

$x$  — target class

$[t:w_m]$  as

$w_i$  — ~~the~~  $t$ -weight value for condition

$i \in \{1, \dots, m\}$

new

(7)

\* The Crosstab shown in the table can be transformed into logic rule form.

$\forall x, \text{item}(x) = \text{"Computer"} \Rightarrow$

$(\text{location}(x) = \text{"Asia"}) [t: 25.00\%] \vee (\text{location}(x) = \text{"Europe"})$

$[t: 30.00\%] \vee (\text{location}(x) = \text{"North America"}) [t: 45.00\%]$

First t-weight value of 25.00% is obtained by 1000 - the value ~~count~~ of count slot for ("Asia, computer") divided by 4000 (total no of items sold)

## ANALYTICAL CHARACTERIZATION:-

### Analysis of Attribute Relevance:-

\* We may not be sure enough that which attribute should be included for class characterization and comparison. We may specify too many attributes and slow down the system.

Measures of Attribute Relevance can be used to identify weak attributes and can be excluded from the concept description process.

\* The incorporation of preprocessing step into class characterization and comparison is referred to as the Analytical Characterization

### Why perform Attribute Relevance Analysis:-

The first limitation of class characterization or multi-dimensional data analysis in DW and OLAP is handling of complex objects.

The second limitation is lack of automated generalization process → the user must tell system which dimensions should be included in class characterization and how high each dimension should be generalized. But it is not difficult as he can set attribute threshold too specify which level a given dimension should reach using command "generalize dimension location to the country level"

x By default there is threshold value b/w 2 to 8. If necessary we can <sup>do</sup> roll up or drill down operations. Selection of attributes should not be too high or too low. If it is too high - low i.e. if we include few attributes it will be incomplete.

# Analytical Characterization: An Example

If the mined concept description involve many attributes Analytical characterization should be performed. This procedure first removes irrelevant less weakly relevant attributes prior to performing generalization let us Examine Example

Example 1:- Mining general characteristics describing graduate students using Analytical characterization.

Given attributes Name (birth-place  
Gender, birth-date  
major | phone# gpa

Step 1 :- Target class data are collected which consist of set of graduate students  
Contrasting class data is also required to perform Relevance analysis. ie let us take as set of undergraduate students

Step 2 :- Relevance Analysis is performed via by Attribute oriented Induction

Name } are removed :- the no of distinct values  
phone# } Exceeds their threshold value.

Concept hierarchies are used to generalize birth-place to birth-country and birth-date to age-range.

Major gpa are generalized to still higher levels

∴ Attributes remaining are  
gender  
major  
birth country  
age range  
gpa.

Candidate relation obtained for Analytical characterization: - target class (graduate students):

gender	major	birth-country	age-range	gpa	count
M	Science	Canada	21...25	Very good	11
F	Science	Foreign	26...30	Excellent	22
M	Engineering	Foreign	26...30	Excellent	18
F	science	Foreign	26...30	Excellent	25
M	science	Canada	21...25	Excellent	21
F	Engineering	Canada	21...26	Excellent	18

Candidate relation obtained for Analytical characterization: - Contrasting class (undergraduate students):

gender	major	birth-country	age-range	gpa	count
M	Science	Foreign	< 20	Very good	18
F	Business	Canada	< 20	Fair	20
M	Business	Canada	< 20	Fair	22
F	Science	Canada	21...25	Fair	24
M	Engineering	Foreign	21...25	Very good	22
F	Engineering	Canada	< 20	Excellent	26

Step 3: Attributes are evaluated using relevance analysis measure such as Information gain

C1 be a class corresponding to graduate - 120 samples

C2 correspond to class undergraduate - 130 samples

$$I(S_1, S_2) = I(120, 130) = -\frac{120}{250} \log_2 \frac{120}{250} - \frac{130}{250} \log_2 \frac{130}{250} = 0.9988$$

Next we need to compute entropy of each attribute let us do for major. ie we need to look at the distribution of graduate and undergraduate students for each value of major.

let us calculate Expected information for each distribution

For major = "Science"

$$S_{11} = 84 \quad I(S_{11}, S_{21}) = 0.9183$$

$$S_{21} = 42$$

For major = "Engineering"

$$S_{12} = 36 \quad I(S_{12}, S_{22}) = 0.7892$$

$$S_{22} = 46$$

For major = "Business"

$$S_{13} = 0 \quad I(S_{13}, S_{23}) = 0$$

$$S_{23} = 42$$

Expected information needed to classify a given sample

$$E(\text{major}) = \frac{126}{250} I(S_{11}, S_{21}) + \frac{82}{250} I(S_{12}, S_{22}) + \frac{42}{250} I(S_{13}, S_{23})$$

$$= 0.7873$$

∴ Gain is

$$\text{Gain}(\text{major}) = I(S_1, S_2) - E(\text{major})$$

$$= 0.2115$$

By we have to calculate gain for all attributes. Then the values are sorted in increasing order.

- gender - 0.0003
- district - 0.0407
- major - 0.2115
- gpa - 0.1115
- age - 0.5971

\* Let relevance threshold = 0.1 to identify weakly relevant attributes:

Gain of gender, district are < threshold. So they are removed.

The contrasting class is also removed. result is the initial target class working relation

candidate - . . .

Step 4:- Attribute oriented induction is applied to Introl target class working relation.

## Mining Class Comparison: Discriminating between Different classes

In many applications users are not interested in a single class described can characterized they would prefer to some description that compares and distinguish one class from other comparable class.

\* Target class and contrasting class must be comparable i.e. Both class should have similar dimensions and attributes.

Ex 1 classes person, address, item are not comparable.

Ex 2 sales in last 3 years are comparable.

\* The techniques developed can be extended to handle class comparison across several comparable classes.

i.e. generalization is performed synchronously among the classes compared so that attributes in all classes to be generalized at same level of abstraction.

### Class Comparison Methods and Implementations:-

General procedure is as follows.

- 1) Data collection:- The set of relevant data in database is collected by query processing and is partitioned into target class and one or set of contrasting classes.
- 2) Dimension Relevance Analysis:- If there are many dimensions and analytical comparison is desired then dimension relevance analysis should be performed and only highly relevant dimensions are included in further analysis.

(11)

Presentation of class comparison descriptions,

Generalized relations, Crosstabs, bar chart, pie chart, curves, etc.  
Rules.

Quantitative Discriminant Rule - The discriminative features of the target and contrasting classes of a comparison description can be described quantitatively by a quantitative discriminant rule which associates a statistical interesting measure d-weight with each generalized type in description.

$$d\text{-weight} = \text{Count}(a_i \in C_f) / \sum_{i=1}^m \text{Count}(a_i \in C_f)$$

where  $a_i$  be generalized type  
 $C_f$  be target class

The rule of given comparison description is written in the form

$$\forall x, \text{target-class}(x) \leftarrow \text{Condition}(x) [d: d\text{-weight}]$$

We can present both characterization & comparison in the same table and in same rule.

Quantitative description Rule

A Quantitative characteristic rule and Quantitative discriminant rule for the same class can be combined to form a Quantitative description rule for the class which displays t-weights and d-weights associated with corresponding characteristic and discriminant rule.

A Quantitative characteristic rule provides a necessary condition for the given target class since it presents a probability measurements for each property that can occur in the class. Such a rule is of form

$$x, \text{target-class}(x) \Rightarrow \text{Condition}_1(x) [t:w_1] \vee \dots \vee \text{Condition}_m(x)$$

where each condition represents a property of target class (the rule is of form). The rule indicates that

if  $X$  is in target class, the probability that  $X$  satisfies condition  $i$  is value of  $t$ -weight,  $w_i$  where  $i \in 1 \dots m$

A Quantitative discriminant rule provides a sufficient condition for target class since it presents a Quantitative measurement of properties that occur in target class versus that occur in contrasting class. Such rule is of form

$$\forall X, \text{target\_class}(X) \Leftrightarrow \text{Condition}_1(X) [d:w_1] \wedge \dots \wedge \text{Condition}_m(X) [d:w_m]$$

The rule indicates if  $X$  satisfies condition  $i$ , there is a probability of  $w_i$  ( $d$ -weight) that  $X$  is in target class where  $i \in 1 \dots m$

A Quantitative characteristic rule and Quantitative discriminant rule for a given class can be combined as follows to form Quantitative description rule

For each condition we have to show both  $t$ -weight and  $d$ -weight

A bidirectional arrow should be used between the given class and conditions.

Such a rule is of form

$$\forall X, \text{target\_class}(X) \Leftrightarrow \text{Condition}_1(X) [t:w_1, d:w_1'] \theta \dots \theta \text{Condition}_m(X) [t:w_m, d:w_m']$$

where

$\theta$  represents logical disjunction / conjunction

The rule indicates that for  $i$  from 1 to  $m$ , if  $X$  is in target class, there is a probability of  $w_i$  that  $X$  satisfies condition  $i$  and if  $X$  satisfies condition  $i$ , there is probability of  $w_i'$  that  $X$  is in target class.

Ctry	TX			computer			base - lists		
	Count	t-weight	d-weight	Count	t-weight	d-weight	Count	t-weight	d-weight
-Europe	80	25%	40%	80	75%	30%	320		
NA	120			50					
both region	200			80					

t-weight for (Europe, TX) = 9s 25% =  $\frac{80}{(80+240)} = \frac{80}{320} = 25\%$

d-weight for (Europe, TX) = 40% =  $\frac{80}{200} = 40\%$

## Mining descriptive Statistical Measures in Large

### Databases:-

Users are interested to learn some data characteristics regarding both central tendency and data dispersion.

Measures of central tendency include mean, median, mode, and range

Measures of data dispersion include - Quartiles, Outliers, and variance

### Measuring the central tendency

\* The most common and most effective numerical measure of center of center of set of data is mean,

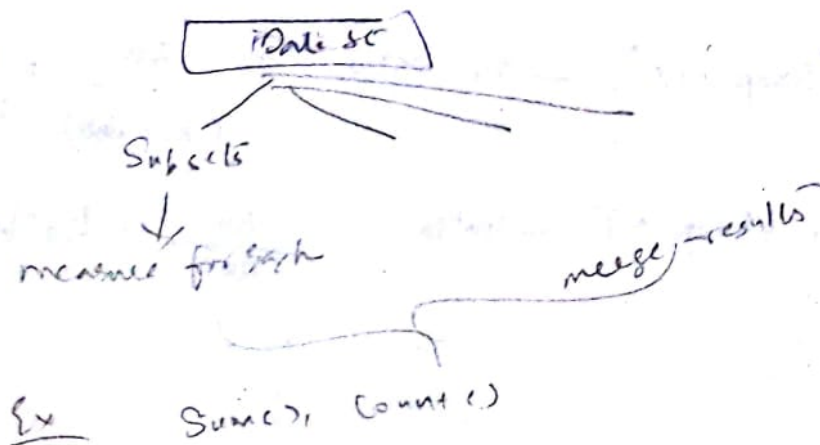
let  $x_1, x_2, \dots, x_n$  be set of  $N$  values for observations.

Mean of set of values is

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

This corresponds to avg() in SQL

A distributive measure is a measure that can be computed for given data by partitioning data into smaller subsets, computing measure for each subset and then merging results in order to arrive at measure value.



### Algebraic Measure

by applying algebraic fn to distributive measure  
 ex - avg(x) — Sum(x) / count(x)

$$\bar{x} = \frac{\sum_{i=1}^N h_i x_i}{\sum_{i=1}^N h_i}$$

→ weighted arithmetic mean

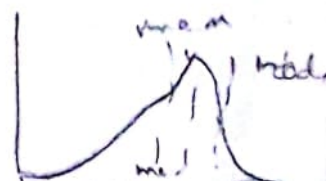
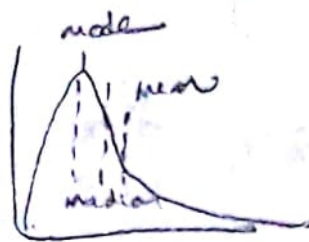
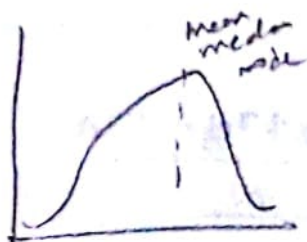
### Statistic

Computed on entire data set as whole

ex - median =  $h_l + \left( \frac{N/2 - \sum \text{freq}}{\text{freq}_{\text{median}}} \right) \text{width}$

lower boundary of median interval

Sum of frequencies of all intervals lower than median interval



Mode

- \* The mode for set of data is value that occurs most frequently in the set.
- \* Sometimes there occurs more than one mode  $\therefore$  greatest frequency correspond to several values, which results in more than one mode.

$$\text{mean} - \text{mode} = 3 \times (\text{mean} - \text{median})$$

Measuring Dispersion of Data :-

The degree to which numerical data tend to spread is called

- the dispersion, or variance of data.

Most common measures are range, five no. summary (based on quartiles)

Interquartile Range (IQR), Standard deviation

- \* Let  $x_1, x_2, \dots, x_n$  be set of observations for some attribute

The range is difference b/w max & min.

- \* Let us assume data are sorted in the increasing numerical order

kth percentile

The kth percentile of set of data in numerical order is value  $x_k$

- having property  $k\%$  percent of data entries lies at or below  $x_k$

Median is 50th percentile

Other than median most commonly we use (percentiles) Quartiles.

Quartile  $Q_1$  is first quartile i.e. 25th percentile

$Q_3$  is 75th percentile.

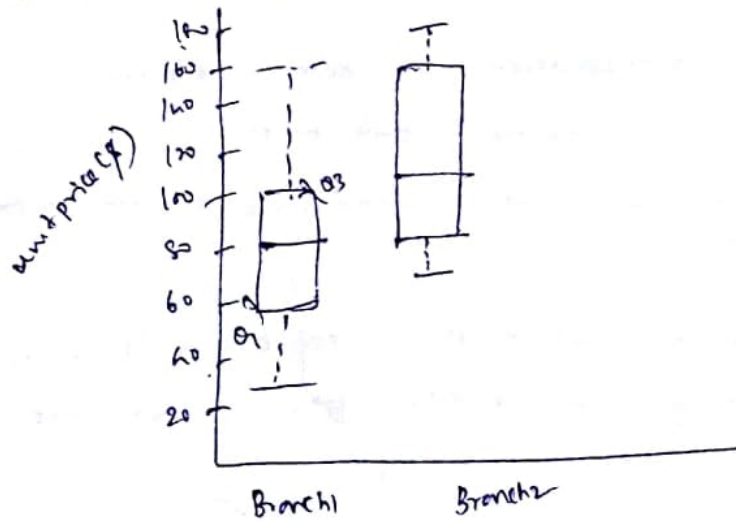
- \* Quartiles, median give some shape to distribution

- \* Distance b/w 1st & 2nd quartile is a simple measure that gives range covered by middle half of data. This distance

is called Interquartile Range (IQR)

$$IQR = Q_3 - Q_1$$

- \* A common rule of thumb for identifying suspected outliers is to highlight values falling at least  $1.5 \times IQR$  above the 2nd quartile or below first quartile.
- \* Using  $Q_1, Q_3$ , median we don't have complete information if we don't have full shape of distribution. It is obtained by giving min & max values. This is known as Five-number Summary.  
The five-number summary consist of median,  $Q_1, Q_3$ , smallest and largest individual observations. i.e. Minimum,  $Q_1$ , <sup>Median</sup>  $Q_3$ , Maximum
- \* Boxplots are popular form of visualizing distribution. It consist of
  - \* ends of box are quartiles i.e. box length is equal to IQR
  - \* Median is marked by line within box
  - \* Two lines (whiskers) outside box extend to min & max observations



The whiskers are extended to extreme low & high observations with  $H$  values are less than  $1.5 \times IQR$ .

## Variance and Standard deviation

Variance of  $N$  observations  $x_1, x_2, \dots, x_N$  is

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

$\bar{x} = \text{mean}$

Standard deviation &  $\sigma$  is same root of variance.  
 Variance, Standard deviation are algebraic measures. They are computed from distributive measures.

↓  
 computed data by partition data set.

## Graphic displays of Basic Descriptive Data Summaries

Other than bar chart, pie chart and line graph there are popular types of graph for display of data summaries & distributions.

These include Histograms

Quantile plots

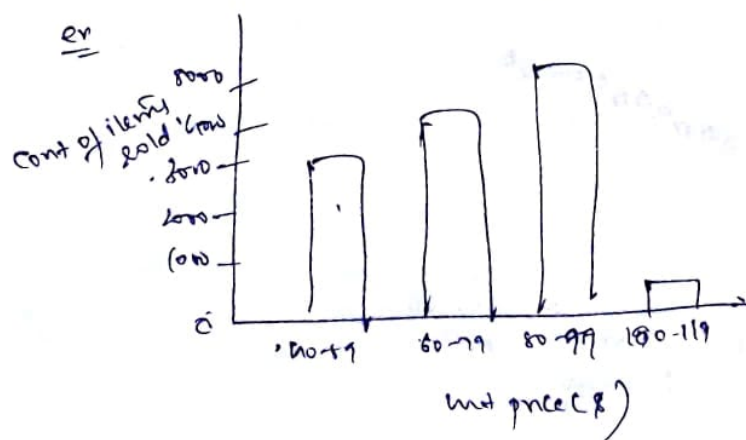
r-r plots

Scatter plots

Box curves.

### Histogram

- \* It is a graphical method for summarizing distribution of given data attribute (A)
- \* It partitions the distribution of A into buckets. width of each bucket is typically uniform for variables
- \* Each bucket is represented by rectangle whose height is equal to relative frequency of each bucket.



Height = Relative frequency of bucket

# Quantile plot

- \* It is simple & effective way for univariate data distribution
- \* 1<sup>st</sup> it displays all of data for given attribute.
- \* 2<sup>nd</sup> it plots Quantile information

ie let  $x_i$  for  $i = 1$  to  $N$  be data in increasing order.

$x_1$  - smallest observation

$x_N$  - largest observation.

Each  $x_i$  is paired with  $f_p$  which indicates approximately 100 $f_p$ % of data are below or equal to value  $x_i$ .

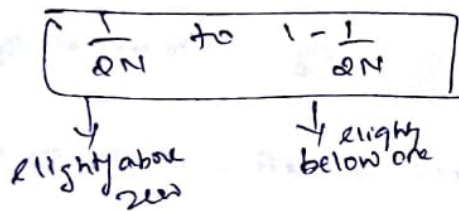
0.25 quantile -  $Q_1$

0.50 ——— median

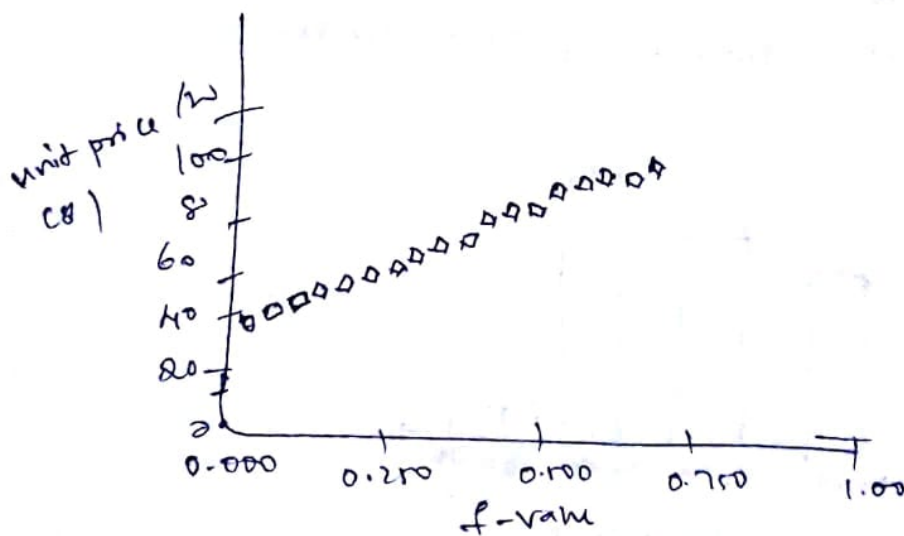
0.75 ———  $Q_3$

$$f_p = \frac{i - 0.5}{N}$$

Number increase in equal steps of  $\frac{1}{N}$  ranging from the



On a Quantile plot  $x_i$  is graphed against  $f_i$



## Quantile-Quantile plot or Q-Q plot

It graphs Quantiles of one univariate distribution against the corresponding Quantiles of another.

Let we have 2 sets of distributions for variable unit price, data from two different branch locations.

Let  $x_1, x_2, \dots, x_n$  be data from first branch

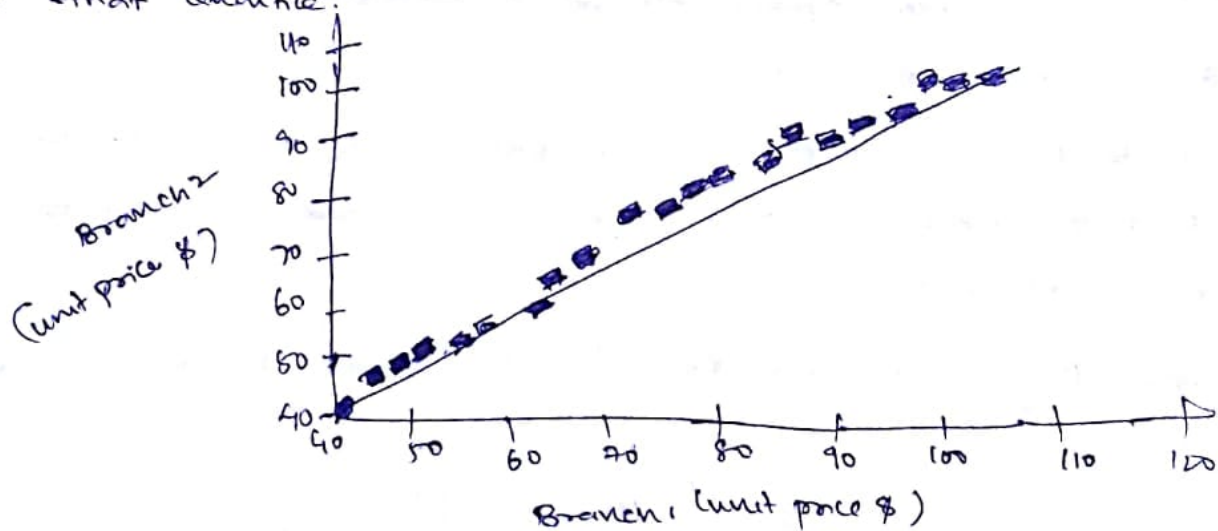
$y_1, y_2, \dots, y_m$  be data from second

$M=N$  & where each data set is sorted in increasing order.

↳ we simply plot  $y_i$  against  $x_i$  where  $x_i, y_i$  are both  $\frac{i-0.5}{N}$  Quantiles of their respective data

- If  $M \neq N$  there can be only  $\min$  points on Q-Q plot, if  $M > N$   $\frac{i-0.5}{M}$  Quantile of  $y$ -data which is plotted against  $\frac{i-0.5}{N}$  Quantile of  $x$ -data

The fig shows unit price of items sold at branch 1 versus branch 2 for that Quantile.



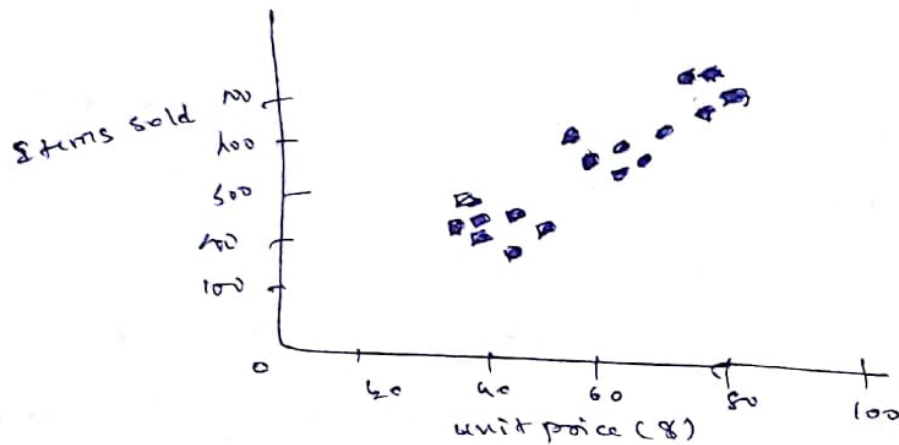
- \* Here lowest point in left corner corresponds to 0.03 quantile. At this quantile unit price of each of items sold at branch 1 was slightly less than at branch 2. In other words 3% of items sold at branch 1 were less than (or) equal to \$40 while 3% of items at branch 2 were less than (or) equal to \$45.
- \* At the highest Quantile we see that unit price of item at branch 1

candidate analysis

Where less than 60% tend to branch 1.

Here there is a shift in distribution of branch 1 w.r to branch 2 in unit price of items sold at branch 1 tend to lower than that of branch 2

### Scatter plot



A scatter plot is one of the most effective graphical methods for determining if there appear to be a relationship or trend b/w 2 numerical attributes.

To construct a scatter plot, each pair of values is located as a pair of coordinates in algebraic senses and are plotted as points in plane.

From this figure we can see clusters of points and outliers.



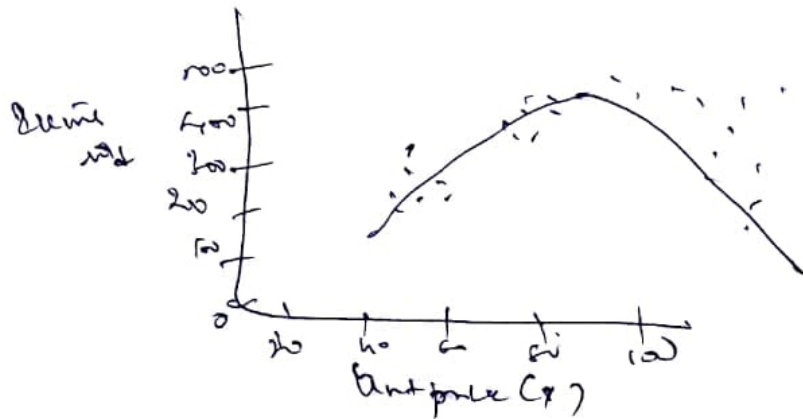
Positive Correlation



Negative Correlation

- Scatter plots are used to find out correlation b/w attributes
- \* Given n attributes a scatter plot contains  $n \times n$  grid of scatter plots - that provides visualization of each attribute
- \* It becomes less effective as no. of attributes grow

Loess Curve



This adds a smooth curve to a scatter plot in order to provide better perception of pattern of dependence. The word "loess" is short for "local regression".

To fit loess curve values need to be set for 2 parameters -

- $\lambda$  - smoothing parameter
- $k$  - degree of polynomial that are fitted by regression.

$\lambda$  - any +ve no. typically b/w 1/4 and 1

$k$  can be 1 (or) 2

Choosing  $\lambda$  is imp as it produces smooth curve. The curve becomes smoother as  $\lambda$  increases.

∴  $\lambda$  is very small underlying pattern is tracked.