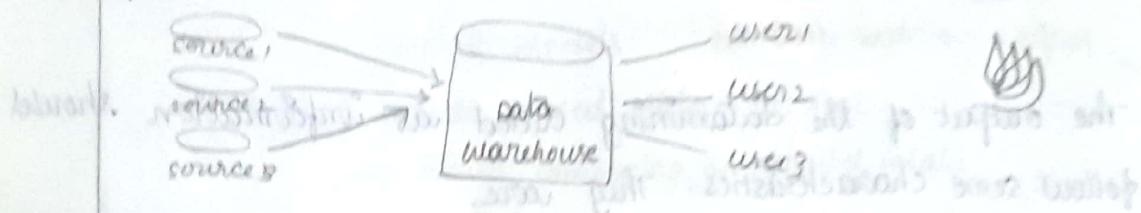


the difference between datawarehouse and database:

In database if number of customers increases, the time taken to get the required information also increases. It is difficult to maintain a data in particular place.

The Datawarehouse is a collection of dissimilar data collected from various sources and stored at particular place.

Data warehouse accepts all formats of data but to store into the Data warehouse it should be transformed into standard formats.



Difference between Datawarehouse and big data.

In Datawarehouse to store the data, it should follow some standard format but in Bigdata it will not follow any standard format.

Data warehouse :

Bill Inmon in 1990's introduced the data warehouse.

According to the Bill Inmon, the collection of data should be called as data warehouse, only if it follows some features i.e The Datawarehouse is

1. Subject Oriented

2. Non-volatile

The data should not be modified or erased. The new data will be added to the collection of data without erasing the old data.

3. Integrated

The data should be collected from different sources such as pdf, doc etc..

4. Time Variant

The data should only be changed in certain period of time but not frequently.

14/12/2018

Data Mining :

Extracting the information from large collection of data which is unknown to the users is called as data mining.

The information that is extracted from data warehouse should possess some characteristics. They are

1. Non-trivial - whatever the information that is getting from the DB that should be relevant to the user (useful)
2. Novel - The o/p should be unique i.e. Any type of algorithm is applied, the o/p should be same
3. useful - The information should be used for future purpose to take decisions.

The output of the datamining called as information should follow some characteristics. They are

1. Non-Trivial

The data should be extracted from warehouse, i.e. the output from the Datawarehouse should

2. Novel

it produces the same output for vary algorithm

3. useful

The information which is extracted from Datawarehouse should be used in future to take decision.

Data mining : (definition).

Extraction of interesting (Non-trivial, Novel, useful) patterns (or) knowledge from huge collection of data.

We can use new technology such as Machine learning, Bigdata, statistics, visualisation, Genetic algorithms

History of data science :

In the early 1960's data collection and database creation
→ file processing.

* early 1960's

data collection and database creation

→ File processing.

* 1970's & early 1980's

→ DBMS

ER models

Indexing & accessing methods

query language - SQL

query processing & optimisation

* mid 1980's to 1990's

→ EER models

→ web based data (ex, XML---)

→ cloud computing & parallel data

* early 1990's to present

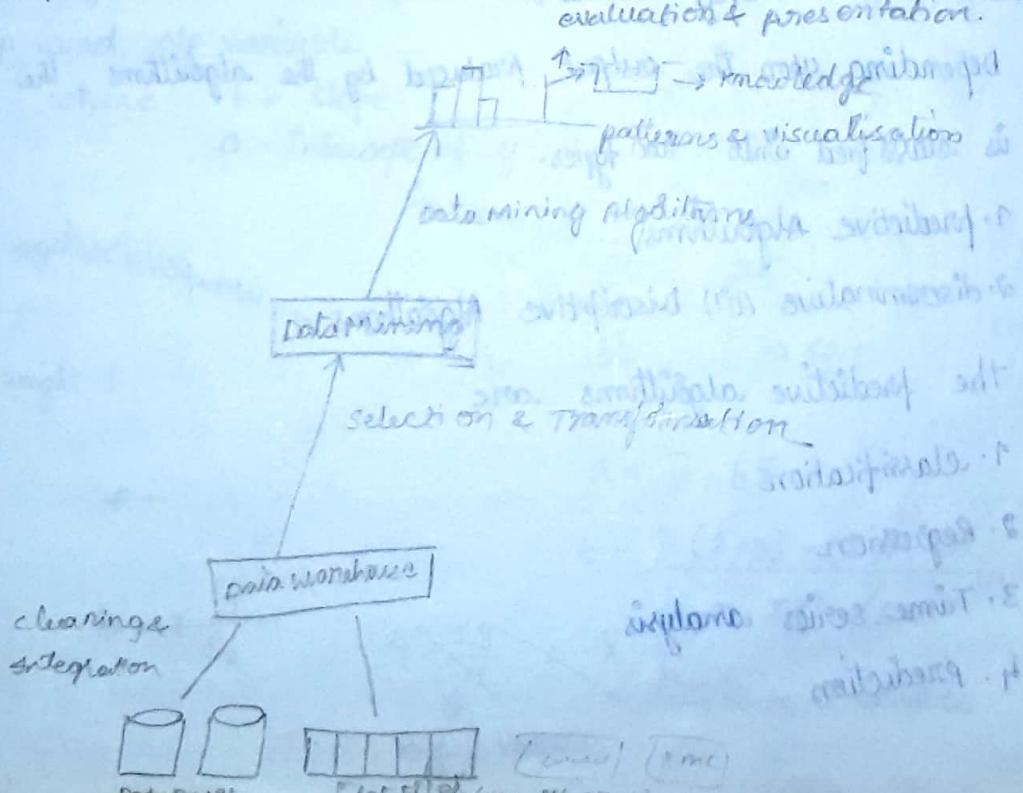
↳ data warehouse

↳ DMT & KDD (datamining techniques & knowledge discovery of data)

↳ DMT in applications.

KDD (knowledge discovery from data).

It indicates the different procedures involved to get the required knowledge from large collection of data.



Data cleaning:

Remove the noise, redundancy data and inconsistent data.

Integration:

collect the information from different data sources which matches to the given query.

Data selection and Transformation:

By using these techniques we select the required information from data warehouse and transform to the format which is understandable by datamining algorithms.

Here we apply normalization procedure and smoothless functions.

patterns and visualisations

identify the interesting information from data mining algorithms and display to the users by applying some visualisation techniques.

Datamining:

Intelligent methods applied to extract the data patterns from data warehouse.

Basic datamining tasks:

Depending upon the output produced by the algorithms, the datamining is classified into two types.

1. predictive Algorithms - supervised

2. discriminative (or) Discriptive Algorithms - unsupervised

The predictive algorithms are

1. classification

2. Regression

3. Time series analysis

4. prediction

The descriptive algorithms are classified as

1. clustering
2. summarisation
3. Association rules
4. sequence discovery.

classification:

predicting the certain output from the given input data lines.
The classification consists of systematic approach and different mathematical techniques which is used to classify the data it is based upon the o/p output which is already known to the users.

Classification algorithms.

1. Decision Trees
2. Neural N/w's
3. Rule based induction.
4. Bayesian N/w
5. Genetic algorithms.

Regression:

it is the technique that is used to predict one or more i/p variable and one o/p variable.

it shows linear relationships b/w i/p & o/p variables.

or $y = a + bx$ is a formula to calculate relationship b/w i/p and o/p variable

where $b = \text{slope}$

$a = \text{intercept of } y$.

$$\begin{array}{ccccccc} \text{age} & 2 & 3 & 10 & 15 & 20 & 25 \\ & & & & & & \\ \text{weight} & 4 & 5 & 30 & 35 & 50 & ? \end{array}$$

$$a = \bar{y} - b \bar{x}$$

$$b = \frac{n(\sum xy) - (\sum x)(\sum y)}{n \cdot \sum x^2 - (\sum x)^2}$$

$$b = r \frac{s_y}{s_x}$$

$$\bar{y} = \frac{\sum y}{n}$$

no. of y variables

$$\bar{x} = \frac{\sum x}{n}$$

no. of x variables.

$$\Sigma x = 50$$

$$b = \frac{5x1848 - 50x124}{5x738 - 2500}$$

$$\Sigma y = 124$$

$$= 2.5$$

$$\Sigma xy = 1848$$

$$a = 24.8 - 2.5 \times 10$$
$$= 24.8 - 25 = -0.2$$

$$\Sigma x^2 = 738$$

$$(\Sigma x)^2 = 2500$$

$$\bar{y} = \frac{\Sigma y}{n} = \frac{124}{5}$$

$$= 24.8$$

$$\bar{x} = \frac{50}{5} = 10$$

$$y = a + bx$$

$$= -0.2 + 2.5x$$

$$\text{if } x = 25 \Rightarrow y = -0.2 + 2.5 \times 25$$

$$\Rightarrow y = 62.3$$

Types of Regression:

1. Linear Regression

2. Multiple Regression

Association Rules:

It is used to analyse and predicts customer behaviour by using if-then statements. It covers the relationships b/w unrelated data in a relational database or other information repository.

The best example is

in a supermarket, the customer who buys the bread chance to buy the butter

bread \rightarrow buys(butter) \rightarrow 90% bread

bread \rightarrow butter [20%, 40%] \rightarrow 40% butter

clustering :

In this the data is divided into different sub classes called as clusters. In clusters the relationship b/w the different object is more.

In cluster the relationship b/w intra class cluster objects is more and interclass cluster relationship is less

the methods in the clusters.

1. partitioning method
2. density based method
3. model based method
4. constraint based method
5. hierarchical method.

examples for the clusters

1. Google search engine

2. marketing loans

3. social N/w analysis

4. credit card and debit card analysis

Time series Analysis

In this technique a series of observations of a variable recorded after successive intervals of time.

The successive intervals of time are usually equal intervals in any unit of time.

Example for Time series Analysis
it is used to study the behaviour of past data and forecast the future data and it ends in business planning and comparative study.

Methods

1. univariate Time Series

Variate

2. Bi-variate Time Series

3. Multi-variate Time Series

Prediction:

Any realworld patterning applications can be seen as predicting future value based on past and current data.

Prediction

Prediction can be

The difference b/w classification and prediction is the classification only works on present data where the off is already known. That prediction works on past and present data & predicts future data.

Examples:

1. Flooding
2. speech recognition
3. Machine learning
4. pattern recognition

18/12/2018

Determining Metrics:

Metrics is used to analyse the algorithm i.e out of two algorithms which one is the better one for the given application.

there are different types of Metrics. They are

1. confusion matrix

2. Cost Matrix

3. ROC curve

4. Accuracy

5. precision

6. Recall

7. F-beta Measures

8. Loss Function

9. Square / Absolute error Methods.

1. Confusion Matrix:

it shows how many relations are correctly classified & how many are not classified correctly.

it gives relationship b/w the correctly classified & not correctly classified relationships classes.

there are 2 classes

1. Actual class

2. predictive class

↳ The o/p of the Alg. is known in advance

↳ This is the o/p of the Algorithm for the same query.

a/b		predictive class		
Actual class		class = yes	class = no	
		class = yes	a	b
		class = no	c	d

True	False	False	True
+ve	-ve	+ve	-ve

among the four regions and are important as they give which are the correctly classified.

Comparison between
Voronoi Regions

Accuracy :

$$\text{Accuracy} = \frac{a+d}{a+b+c+d}$$

Example :-

		+	-
+	+	50	100
	-	50	150

$$\text{ACC} = \frac{a+b}{a+b+c+d} = \frac{60+100}{60+100+50+150} = 0.67$$

$$\frac{d}{d+r} = \frac{50}{50+150} = 0.25$$

$$\frac{a}{a+r} = \frac{60}{60+50} = 0.53$$

$$\frac{c}{c+r} = \frac{100}{100+150} = 0.4$$

Precision :-

$$\text{precision} : \frac{a}{a+c} \quad (\text{True positive Ratio})$$

Precision (π) \Rightarrow accuracy (π)

precision π \Rightarrow accuracy (π)

precision π \Rightarrow accuracy (π)

Cost matrix :

		Predictive class	
		+	-
Actual class	+	$c(+ +)$	$c(+ -)$
	-	$c(- +)$	$c(- -)$
Predictive class	+	$c(+ +)$	$c(+ -)$
	-	$c(- +)$	$c(- -)$

TP
FN
TN
FP

If M1 is used

		+	-
+	+	250	45
	-	5	200

If M2 is used
if ACC is taken \Rightarrow M2 is better

		+	-
+	+	100	40
	-	60	250

$$\text{ACC} = \frac{150+250}{500} = \frac{400}{500} = 0.8$$

$$\text{ACC} = \frac{150+200}{500} = \frac{350}{500} = 0.7$$

$$\text{ACC} = \frac{150+45}{500} = \frac{195}{500} = 0.39$$

		+	-
+	+	250	45
	-	5	200

$$\text{ACC} = \frac{150+45}{500} = \frac{195}{500} = 0.39$$

$$\text{ACC} = \frac{150+100}{500} = \frac{250}{500} = 0.5$$

$$\text{ACC} = \frac{150+40}{500} = \frac{190}{500} = 0.38$$

Now we need to consider acc. to our application
M1 is better

Recall (γ): Recall(γ) = $\frac{\omega}{\omega+b}$ (sensitivity). (TPR).

F-beta measure : F-beta measure = $\frac{2\gamma P}{\gamma + P} = \frac{2\omega}{\omega + b + c}$.

weighted Accuracy :

$$\text{weighted Accuracy} = \frac{w_1 \omega + w_4 d}{w_1 \omega + w_2 b + w_3 c + w_4 d}$$

$$\text{specificity} = \frac{d}{c+d} \Rightarrow (\text{TNR})$$

$$\text{FPR} = \frac{c}{c+d}$$

$$\text{FNR} = \frac{b}{a+b}$$

ROC curve : (receiver operating characteristics).

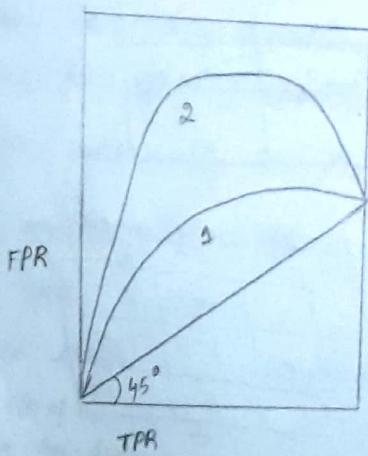
ROC curves are used for visual comparisons of classification models.

These curves come from signal detection theory.

These curves show trade-off b/w true positive rate & false pos. rate.

The area under ROC curve is the measure of accuracy of a model. The area is called as AUC (area under curve). The value is b/w 0.5 to 1.

min distance.



The area of an algorithm for which AUC is more it have more accuracy.

Loss Function:

This function generally used in Regression

$$\text{Absolute error} = |y_i - y_i'| = \sum y_i - \sum y_i'$$

$$\text{Squared error} = (y_i - y_i')^2$$

Test - error :-

$$\text{mean absolute error} = \frac{\sum_{i=1}^d |y_i - y'_i|}{d}$$

$$\text{mean squared error} = \frac{\sum_{i=1}^d (y_i - y'_i)^2}{d}$$

$$\text{relative absolute error} = \frac{\sum_{i=1}^d |y_i - y'_i|}{\sum_{i=1}^d |y_i - \bar{y}|}$$

$$\text{relative squared error} = \frac{\sum_{i=1}^d (y_i - y'_i)^2}{\sum_{i=1}^d (y_i - \bar{y})^2}$$

if we need to measure performance with respect to noise then we use absolute error.

if there are more outlets in our dataset then we use squared error.

19/12/2018
Data Mining in DB perspective:

- 1. Scalability
- 2. Real world data
- 3. Update
- 4. Ease of use
- 5. Security

Scalability :

Algorithms that do not scale up to perform well with massive real world databases are of limited applications. Related to this is a fact that technique should work regardless of amount of available memory space.

Real world data :

The real world data are noisy data and have many missing attribute values.

Algorithm should be able to work even in the presence of these problems.

Update :

Many datamining algorithm work with static datasets.

This is not a realistic approach.

Ease of use :

Although some algorithms may work well, they may not be

well received by the users, if they are difficult to use and understand.

privacy & security:

we are collecting the data from different databases. There may chance to get effected by database worms and viruses.

Data Mining Issues:

1. Human Interaction

2. over fitting

3. outliers

4. Interpretation of result

5. visualisation of results

6. Large datasets

7. High dimensionality

8. Multimedia data

9. missing data

10. irrelevant data

11. noisy data

12. changing data

13. integration

14. Applications of data mining for the soft engineering

• medical diagnosis
• credit evaluation
• fraud detection
• customer segmentation
• market basket analysis
• document classification

• spam filtering
• email filtering
• intrusion detection

• recommendation systems
• information retrieval
• search engines

• decision support systems
• knowledge discovery in databases

• data mining for bioinformatics
• data mining for gene expression analysis

and the best way to solve the problem is also
to understand what are the present issues facing
multiple sectors and projects which
are contributing to the lack of
cooperation between them. In order to
achieve this, there must be a clear
understanding of the different
sectors involved and their
objectives. This will help in identifying
the common goals and interests of all
the stakeholders. Once this is done,
it is important to establish a
collaborative framework that
allows for effective communication
and coordination between the different
sectors. This can be achieved through
regular meetings, joint planning sessions,
and the exchange of information and ideas.
It is also important to recognize
the strengths and weaknesses of each
sector and to work together to
overcome any challenges that may arise.
In addition, it is crucial to build trust
and respect among the different
stakeholders. This can be done by
showing appreciation for the
efforts and contributions of others
and by being open to new ideas and
perspectives. By doing so, it is
possible to create a positive
environment where everyone
feels valued and respected.
Overall, building a successful
collaborative framework requires
a commitment from all
the stakeholders involved to work
together towards a common goal.
It is a challenging task, but with
determination and a willingness to
listen and learn, it is possible to
achieve success.

20/12/18

Data preprocessing:-

Data is nothing but collection of objects and attributes.

Preprocessing means preparing raw data understandable format before applying data mining algorithms.

Importance of Preprocessing:-

Nowadays, we are dealing with real world data and these data in the format of

i) Incomplete means tagging of attribute values of containing aggregate data

Example: $a = b + c$

$b = 20, c = ?$

print(a) # incomplete.

ii) Noisy data - The data contains errors & outliers.

Example: $a = 10$ but $a = -10$ # noisy (unwanted).

iii) Inconsistent data - It containing different codes of names.

Example: if db has

	x_1	x_2		x_1	x_2	
'	0	1	A	0	1	# inconsistent
2	0	1	B	0	1	

Tasks of data preprocessing:-

1) Data Cleaning -
Identify & fill missing data, smoothing the noisy data, remove outliers and remove inconsistency.

2) Data Integration -
We integrate multiple databases, data queues, & flat files.

3) Data Transformation -
We use normalisation & aggregation to transfer the data from one format to another format.

4) Data Reduction -
Reducing the size of data but producing the similar best

analytical results of original data.

5) Data Discretization :-

Some part of data is reduced & replacing by numerical attributes

→ 1. Data Cleaning :-

i) Missing Data : Reasons to missing the data

- * data may not included simply because it was not consider important at that time of entry
- * Relevant data may not be recovered due to misunderstanding or equipment failure.
- * The attributes of interest may not always available.

- Handling Missing Data :

1. Ignoring the tuple :

when a class label is missing we need to delete the tuple but it is an ineffective method unless the tuple contains more number of missing data.

2. Fill in the missing value by manually :

In General this method is the time consuming process & may not be feasible when given a large datasets with many datasets.

3. Use the Global Constant to fill in the missing value:

Replace all missing attributes by some constant, then the mining program may mistakenly think that they form an interesting pattern since, they have a value in common.

4. Use the attribute means to fill in the missing value:

5. Linear (Interpolation) Interpolation : Baban

take the values of up and down of missing value & calculate average & replace it ^{in place of} by missing value.

6. Use the most probable value fitting fill in the missing value :

This may default with inference base tool like decision tree, regression on basis of formula.

ii. Noisy Data:

To deal with the noisy data, we are applying Binning method. The Binning method is again classified into three types. They are

1. Mean

2. Median

3. Boundaries

1. Mean Binning Method:

Step 1: Arrange the data in sorting order.

Step 2: Divide the given data into equal parts called as

bining depth & represented by D and the no. of binning is represented by $\frac{N}{D}$ bins.

Step 3: calculate the mean value of each bin by using the formula

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

Step 4: Replace each value in bin by calculate the mean value of each bin.

Example:

5, 7, 4, 6, 9, 10, 8, 6, 5

Step 1:- Set the data.

4, 5, 5, 6, 6, 7, 8, 9, 10

Step 2:- The depth of each bin is equal to 3

$$D = 3$$

$$\text{no. of bins} = \frac{N}{D} = \frac{9}{3} = 3$$

The first bin contains 4, 5, 5

bin-2 contains 6, 6, 7

bin-3 contains 8, 9, 10

Step 3:- calculate average of each bin.

$$\frac{4+5+5}{3} \Rightarrow \frac{14}{3} = 4.66$$

$$\frac{6+6+7}{3} \Rightarrow \frac{19}{3} = 6.33$$

$$\frac{8+9+10}{3} \Rightarrow \frac{27}{3} \Rightarrow 9$$

step 4

$$\text{the } [4\frac{2}{3}, 4\frac{2}{3}, 4\frac{2}{3}], [6\frac{1}{3}, 6\frac{1}{3}, 6\frac{1}{3}], [9, 9, 9].$$

2. median Binning Method :- (50 percentage)

$$50\% (n+1)$$

$$\frac{50}{100} \times (n+1) \quad \text{where } n \text{ is no. of bins}$$

$$\frac{50}{100} \times (3+1) = \frac{1}{2} \times 4 = 2$$

Replace all the values with "2".

$$[2, 2, 2, 2, 2, 2, 2, 2, 2]$$

3. Boundaries Binning Method :-

calculate the maximum value & min value. In the each bin replace all the values of the bin with the minimum value except the maximum value.

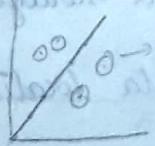
for above example:

$$[4, 5, 5], [6, 6, 7], [8, 8, 12]$$

i) Clustering:
similar values are organised into groups called as clusters. The values that falls outside of clusters called as outliers.

ii) Regression Technique:
data can be smoothen by fitting the data to a function such as with regression.

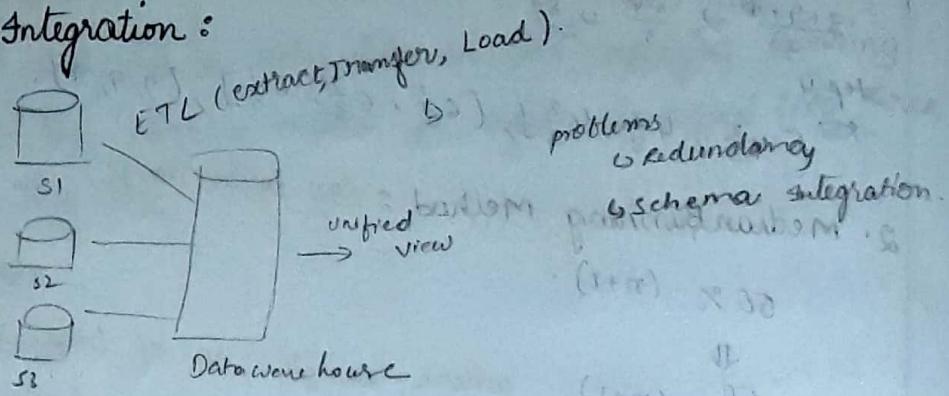
$$y = a + bx$$



iii) Manually we can remove the outliers. but this is not easy.

27/12/18

→ Data Integration :



It involves combining the data from different sources & disparate sources which are stored using various technologies & provide unified view of the data. We use two major approaches to get data from different data sources.

1. tightly coupling :

2. Loose coupling :

1) Tight coupling :

Applies the ETL operation on different data sources & store the results in DW. From that the user will get required information.

2. Loose coupling :

In this we are not applying any ETL operations. The query submitted by the user will directly sent to the diff data sources & the results will transferred to the user.

Advantages of data integration :

1. Faster query processing.

2. Independence (not depending on multiple DBs)

3. Complex query processing

4. High volume data processing.

5. Advance data summarisation and storage.

Disadvantages of data integration :

1. costliest operation because of data localisation, infrastructure and security

2. latency, because data needs to be loaded after ETL operations.

Issues in the data integration :

1. Schema integration :

2. Redundancy

3. detection & resolution of data value conflicts.

1. Schema Integration :

The real world entities from multiple sources be matched is called as the Entity Identification problem.

Example : The Data Analyst or the Computer be sure that customer_id in one DB & cust_id in another DB to represent the same entity.

2. Redundancy :

The Redundancy is occurred by 2-ways
1. object identification, means the same attribute of an object may have different names in different databases.

2. derived data, one attribute may be a derived attribute in another table.

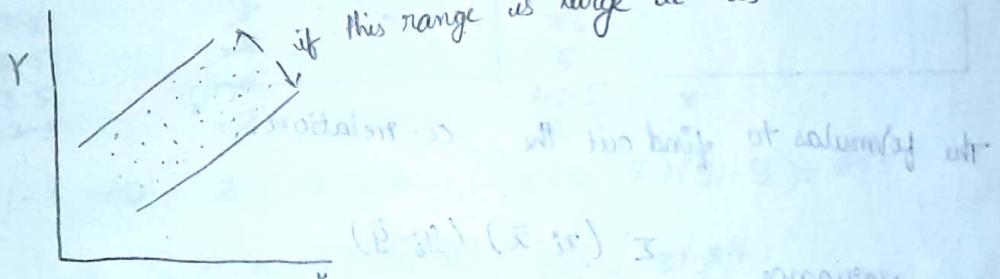
23/12/18

Correlation Analysis

it is used to identify the strength of the association between two variables. In this correlation analysis, there is no causal effect. (numeric data)

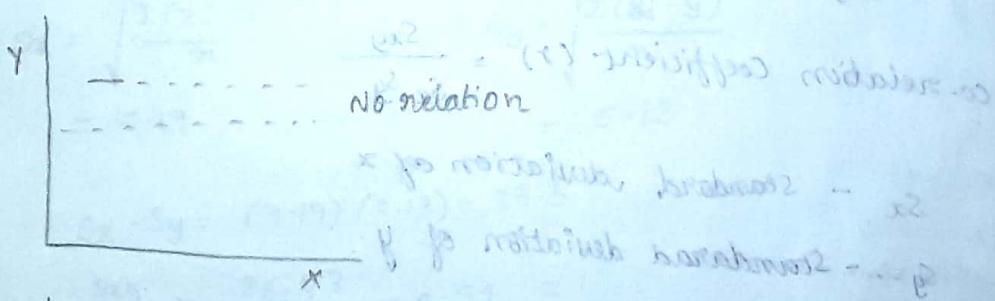
By using scatter plot, we find out the relation b/w the two variables.

if this range is large it is weak relation.



$$\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

completely weak

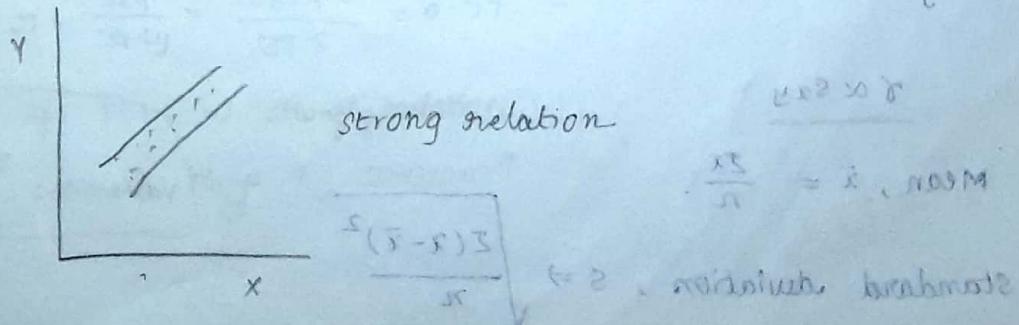


$\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n} = 0$ (no correlation)

X vs running hours1 ... x1

y vs running hours2 ... y2

0.79 ... 0.79



strong relation

$$\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

$$\frac{\sum (x_i - \bar{x})^2}{n} = s_x^2, \text{ no. of M}$$

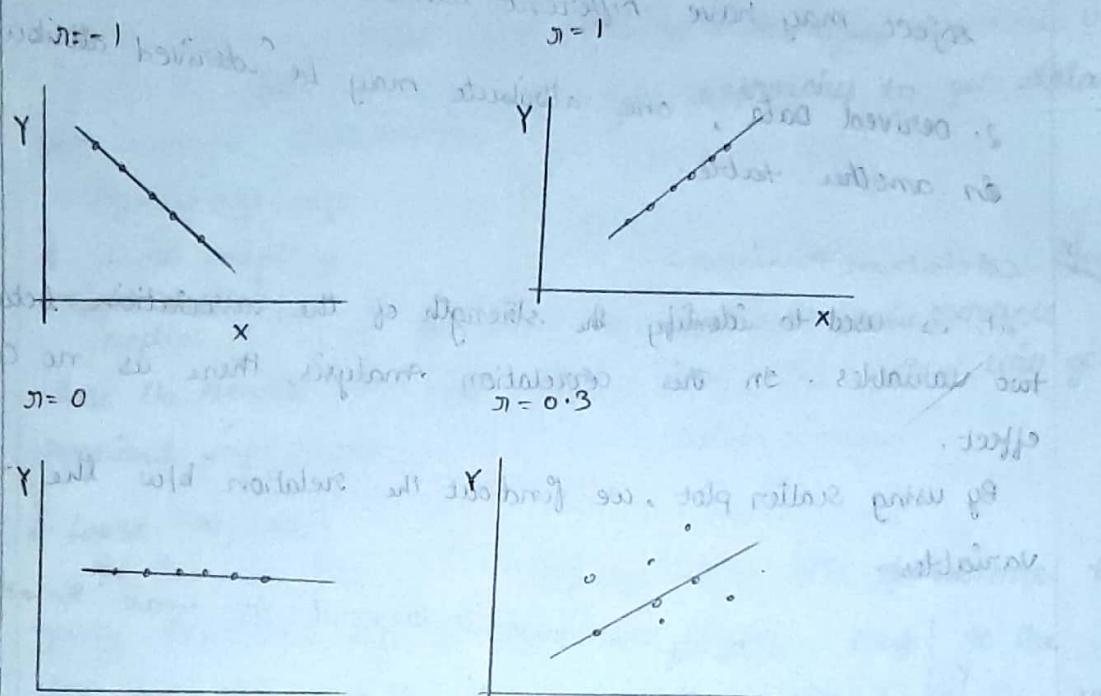
$\frac{\sum (y_i - \bar{y})^2}{n} = s_y^2, \text{ minimum variance}$

In correlation Analysis, the value of 'r' should be b/w -1 & 1.

If $r_{A,B} > 0$, then there is a stronger +ve linear relationship.

If $r_{A,B} = 0$, weaker linear relationship.

If $r_{A,B} < 0$, then there is a stronger -ve linear relationship.



The formulas to find out the co-relation:

$$\text{co-variance} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

$$\text{co-relation coefficient } (r) = \frac{S_{xy}}{S_x S_y}$$

S_x - standard deviation of x

S_y - standard deviation of y

$$r \propto S_{xy}$$

$$\text{Mean, } \bar{x} = \frac{\sum x}{n}$$

$$\text{standard deviation, } s = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

$$\Rightarrow \text{Variance} = \frac{\sum (x - \bar{x})^2}{n}$$

$$\text{correlation coefficient } r = \frac{\text{(181)} n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{[n \Sigma x^2 - (\Sigma x)^2][n \Sigma y^2 - (\Sigma y)^2]}}$$

Example:

temp(x) Icecream(Y)

66	8
72	11
77	15
84	20
83	21
71	11
65	8
70	10

$$\begin{array}{r} 73.5 \\ 66.0 \\ \hline -7.5 \end{array}$$

$$\bar{x} = \frac{\Sigma x}{n} = 73.5 \quad \bar{y} = \frac{\Sigma y}{n} = 13$$

$x_i - \bar{x}$	$y_i - \bar{y}$	$\Sigma (x_i - \bar{x})(y_i - \bar{y})$
-7.5	-5	37.5
-1.5	-2	3
3.5	2	7
10.5	7	73.5
9.5	8	76
-2.5	2	5
-8.5	-8	42.5
-3.5	-3	10.5

$$\Sigma x_i - \bar{x} = 0 \quad \Sigma (y_i - \bar{y}) = 0 \quad \Sigma (x_i - \bar{x})(y_i - \bar{y}) = 255$$

$$s_{xy} = \frac{\Sigma (x_i - \bar{x})(y_i - \bar{y})}{n} = \frac{255}{8} = 31.87$$

$$s_x = \sqrt{\frac{\Sigma (x_i - \bar{x})^2}{n}} \quad s_y = \sqrt{\frac{\Sigma (y_i - \bar{y})^2}{n}}$$

$$= 7.19 \quad = 5.13$$

$$s_x \cdot s_y = (7.19)(5.13) = 37.5$$

$$r = \frac{s_{xy}}{s_x \cdot s_y} = \frac{36.43}{37.5} = 0.99 \approx 1$$

if $r=1 \Rightarrow$ strongly relation.

so, if Temp ↑, Icecream ↑

(3 2) / 12 / 18

Chi-Square : (χ^2 -Square). (18)

χ^2 (used to find the relation b/w two variable. This is used to measure Nominal data. This is the 1900's method by carl pearson).

It is developed in the year 1900 by Carl Pearson. The symbol χ^2 is called as "chi".

This technique only works for the categorical data & the nominal data such as Gender means male & female, Colour means red, yellow, green .., & severity of disease mild, severe. This chi-square technology, it consists of two types of data. i.e. O → original (or) actual (or) observed data. e → expected data.

In order to find out whether 2 objects are strongly correlated or not, we use a degree of freedom to decide.

degree of freedom, $df = (\text{no. of rows}-1)(\text{no. of columns}-1)$.

The formulas to find the expected data :

$$\begin{array}{c|cc|c} & a & b & a+b \\ \hline c & & & \\ & & & \\ \hline & c & d & c+d \\ \hline \text{total} & a+c & b+d & n = a+b+c+d \end{array} \Rightarrow \begin{array}{c|cc|c} & 22F & 3F & 25 \\ \hline 22F & & & \\ & & & \\ \hline 3F & & & \\ \hline 25 & & & \end{array}$$

$$e_{11} = \frac{(a+b)(a+c)}{n}, \quad e_{12} = \frac{(a+b)(b+d)}{n},$$

$$e_{21} = \frac{(c+a)(c+d)}{n}, \quad e_{22} = \frac{(c+d)(b+d)}{n}$$

Example :

		Smoker		Non-Smokers	
		36	14	50	$50 \times \frac{66}{105} = 36$
		$e_{11} = \frac{50 \times 66}{105} = 31.42$	$e_{12} = \frac{60 \times 39}{105} = 28.57$	25	$55 \times \frac{66}{105} = 36$
suffering from lung diseases		36	14	50	$50 \times \frac{66}{105} = 36$
not suffering from lung diseases		30	12	55	$55 \times \frac{66}{105} = 36$
		$e_{21} = \frac{50 \times 55}{105} = 24.57$	$e_{22} = \frac{60 \times 39}{105} = 20.42$	25	$55 \times \frac{66}{105} = 36$
		66	39	85	$85 \times \frac{66}{105} = 54$

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

Calulation - Numerical
Chi-square - Critical

$$\Rightarrow \frac{(36-31.42)^2}{31.42} + \frac{(14-18.57)^2}{18.57} + \frac{(30-34.57)^2}{34.57} + \frac{(25-20.43)^2}{20.43}$$

$$\Rightarrow 0.667 + 1.124 + 0.604 + 1.022$$

$$\Rightarrow 3.417 \Rightarrow \chi^2$$

$$\text{degree of freedom} = 1 \times 1 = 1$$

$p < \chi^2$ (Rejected)
(Strongly correlated)

df	0.05	0.01	0.001
1	3.841	6.635	10.828
2	5.991	9.210	13.816
3	7.815	11.345	16.266
4			

$$P > \chi^2$$

So Accepted
(Independent).
(no relation)

Hence smoking is not the only factor for lung disease.

Tuple Duplication:

If there exists one tuple multiple times in DB it is called

as Tuple Duplication. Then
delete one repeated row to handle tuple duplication.

Data Value Conflicts & Resolution:

If there exists one data in different formats, it leads to conflicts to find out original data (required data) because two values belongs to the same data. This will be resolved by converting one data format into another one to scale the data.

02/01/12 Data Reduction:

It is the technique to reduce the size of the data but provide the same analytical results if you work on original data is called as data reduction. That closely maintains the integrity of data.

Advantages:

To improve the performance of DMT.

To reduce the space.

To reduce the N/W bandwidth to transfer.

Types of Data Reduction Techniques:

1. Dimensionality Reduction

2. Numerosity Reduction

3. Data compression.

Dimensionality Reduction:

It is the process of reducing the no. of random variables

(a) attributes under consideration.

methods in the dimensionality reduction are

1. wavelet Transform { transform original data}
2. principle component Analysis. { sample data (smaller size)}
3. Attribute subset selection.
 - ↳ irrelevant, weakly relevant (or) redundant attribut
 - (or) dimensions are deleted (or) removed.

Numerosity Reduction :

it is a technique to replace the original data size by alternative smaller forms of data representation. There are two types in numerosity reduction. They are

1. parametric numerosity Reduction technique.
2. non-parametric

parametric

In this model, we used to estimate the data, so that the required parameters needs to be stored in memory rather than actual data. The techniques in parametric numerosity reduction technique are

1. Regression

2. Log-Linear models

non-parametric

it is used for storing reduced representation of data.

The techniques in non-parametric reduction are

1. Histograms

2. clustering

3. Sampling

4. Data cube Aggregation

Data Compression :

In this, transformations are applied to obtain reduced & compressed representation of data

it is a reduction in the no. of bits need to be represent the data. There are 2 types of compression techniques

1. Loss-less compression technique

If the original data can be reconstructed from the compressed data without any information loss.

2. Lossy compression technique

Instead of getting original data, we can reconstruct only

an approximation of original data. The techniques in the loss less rare.

1. String compression

2. All the dimensionality reduction and numerosity reduction techniques are examples of data compression.

03/01/2019
7. Wavelet Transforms : it is a digital signal.

The wavelet transformation, is used for data analysis to make more accurate data. The difference b/w Fourier Transformation and wavelet transformation is

* The Fourier Transformation will not produce accurate image (or) data when compared to the wavelet Transformation

* It is not easy to work on continuous signal, But if we work on wavelet Transformation signal, the both areas of the

* In wavelet Transform & -ve will have the common factors. Then we can apply the transformation only one side.

* There are two types of wavelet Transformation

↳ Continuous wavelet Transformation ~~Fourier~~

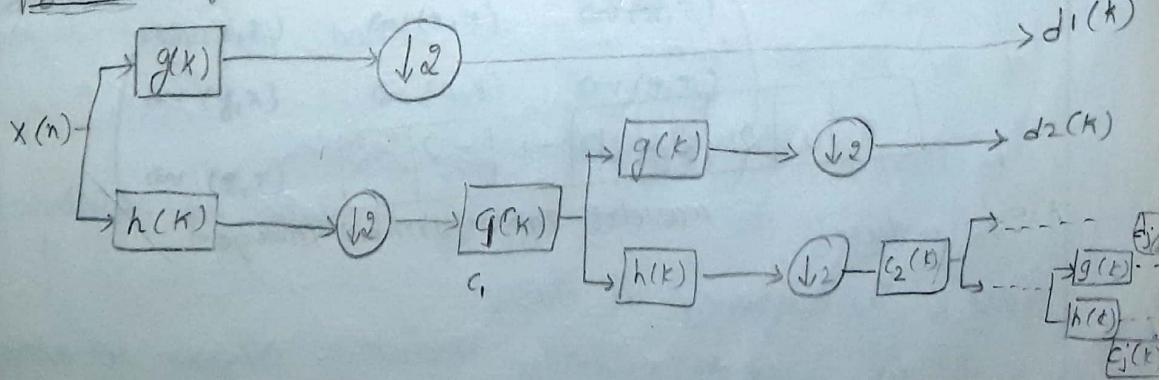
↳ Discrete wavelet Transformation. ~~Fourier~~

Discrete wavelet Transformation : it is a linear signal processing technique

In DWT, we are given an image (a data) of size N & divided into $(x_1, x_2, x_3, \dots, x_n)$ and each data part contains the same length. We call this data parts as wavelet coefficients.

In DWT, there are many algorithms that are used to analyse the data, one among the is pyramid algorithm.

Pyramid Algorithm



The decision tree is used to reduce the dimensionality by removing or not considering the attributes in the original dataset.

parametric data Reduction: There are two techniques.

it is also called regression & log linear methods.
we can also call it as numerical reduction.

These two models are used to reduce the size of data.

Regression:

In linear regression, the data is modeled to fill a straight line

$$y = w_1 x + b$$

where

x is independent variable (predictor variable)
 y is dependent variable (response variable)

w, b are regression coefficients.

If the data which does not fit on the line is called as outlier and this will be deleted from original dataset.

loglinear

This is used to approximate discrete multi-dimensional probability distribution.

Estimate the probability of each point in a multidimensional space for a set of discretized attributes based on a smaller subset of dimensional combination.

$F = a e^{bx}$

where a, b are log linear coefficients.

Non parametric data Reduction:

1. histogram.

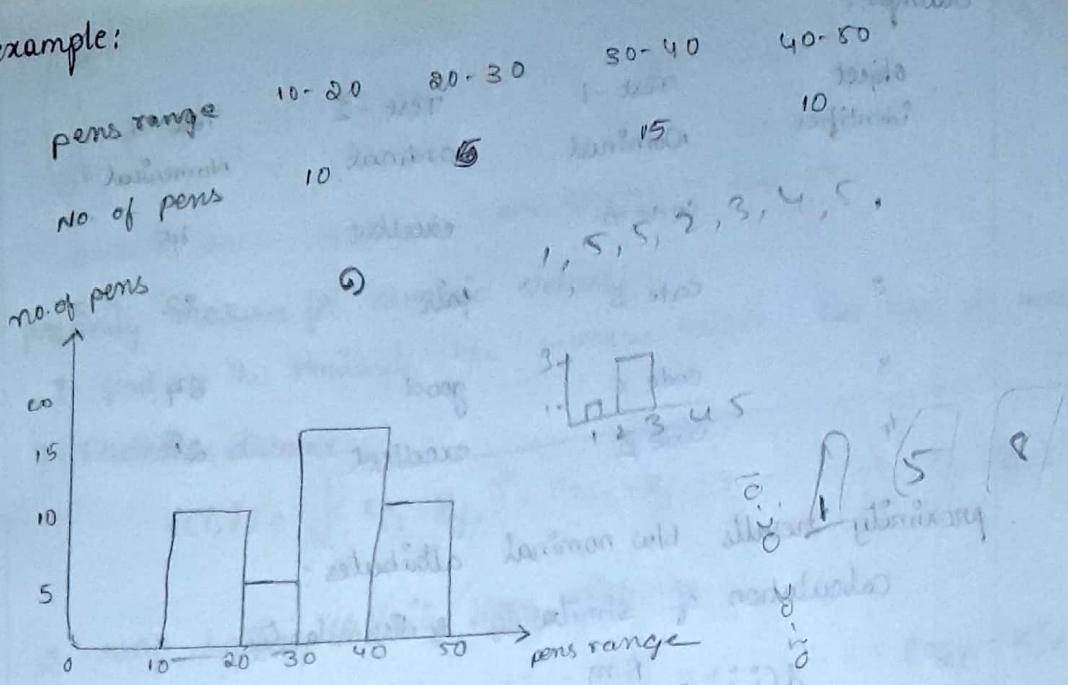
2. clustering

3. Sampling

Histogram:

the graphical presentation of frequency distribution of a continuous series

example:



single ton bucket:

Similarity Measurements: The similarity b/w the two objects will indicate the relationship b/w the two objects. (to make clusters we make similarity & dissimilarity)

$$\text{ex:- } \text{sim}(i, j) = 1 - \text{dis}(i, j)$$

similarity = 0 means there is no relation b/w two objects also indicated by dissimilarity = 1.

similarity = 1 means more relation b/w two objects & also indicated by dissimilarity = 0.

To represent the similarity b/w the objects will use two datacenters

1. data matrix

2. dissimilarity matrix.

(1) clusters

(2) people search

dissimilarity matrix can be represented as
this structure stores a collection of proximities that are available for all pairs of n objects.

$$\begin{bmatrix} 0 \\ d(1,2) & 0 \\ d(2,1) & d(2,3) & 0 \\ d(3,1) & d(3,2) & d(3,4) & 0 \\ d(4,1) & d(4,2) & d(4,3) & d(4,5) & 0 \\ d(5,1) & d(5,2) & d(5,3) & d(5,4) & d(5,6) & 0 \end{bmatrix}$$

examples:

Example :-

object identifier	Test - 1	Test - 2	Test - 3
1	nominal	ordinal	Numerical
2	code A	excellent	45
3	code B	fair	22
4	code C	good	64
	code D	excellent	28

proximity results b/w nominal attributes:

calculation of similarity & dissimilarity

$$d(i,j) = \frac{p-m}{p}$$

where m is the no. of matches:

P is the no. of attributes describe the object (like nominal)

(d) Nominal will describe fine values. (e) Nominal will be used.

دیگری نیستند و این از این دلایل است که این افراد ممکن است در این میان از افرادی باشند که ممکن است از این دلایل خود را بگیرند.

P-17

$$d(z_1, z_2) \Rightarrow \text{Test - 1}$$

P → no. of attributes has one some as nominal

$m = 4$ if 2 → also contains $\omega_0 \wedge \rho_{\text{kin}}$

$$d((x_1), \frac{1-\theta}{1}) = 0$$

$$d(3,1) = \frac{1}{1} = 1 \quad d(4,2) = \frac{1}{1} = 1$$

$$d(3,2) = \frac{1-0}{1} = 1 \quad d(4,3) = \frac{1-0}{1} = 1$$

$$\begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$$

$\Rightarrow d(i,j) = 1$ if $x_i^p = x_j^p$

proximity measures for numeric values:
To find out the similarity b/w numeric values, we use 3 methods

1. Euclidean distance.

$$d(i,j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$$

2. manhattan (city block) distance.

$$d(i,j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

Example:- $x_1 = (1, 2)$
 $x_2 = (3, 5)$.

$$1. \sqrt{(1-3)^2 + (2-5)^2} = \sqrt{2^2 + 3^2} = \sqrt{13}.$$

$$2. |1-3| + |2-5| = 2+3 = 5$$

Weight euclidean distance = weight \times euclidean distance.

3. Minkowski distance?

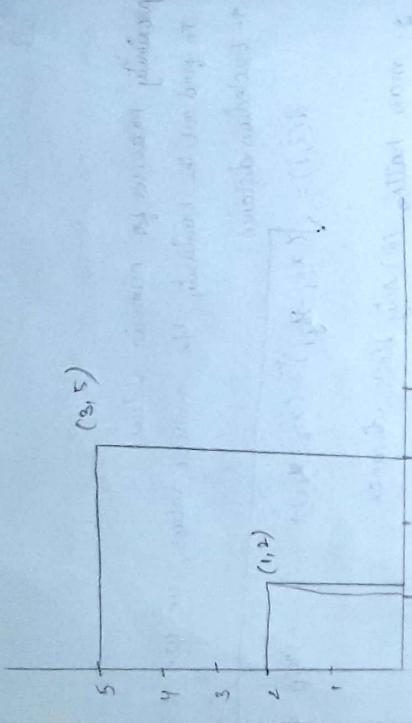
$$d(i,j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

it represents the manhattan when $h=1$
& $h=2$ it represents euclidean distance.

④ Supreme distance: (a) L_{\max} (b) median (c) chebyshev

i.e
 $d(i,j) = \lim_{h \rightarrow \infty} \left(\sum_{p=1}^h |x_{ip} - x_{jp}|^h \right)^{1/h}$

$$h \left(|x_{11} - x_{12}|^h + |x_{21} - x_{22}|^h + \dots + |x_{n1} - x_{n2}|^h \right)^{1/h}$$



$$1. \text{ Euclidean} = \sqrt{4+9} = \sqrt{13}$$

$$2. \text{ manhattan} = |2+3| = 5$$

$$3. \text{ supreme} = |2-3| + |5-2| = 5$$

$$\text{dis}_f^* = \frac{|x_f - x_{if}|}{\max_i x_i - \min_i x_i}$$

$$\text{dis}_m^* = \frac{|x_f - x_{if}|}{\max_i x_i - \min_i x_i}$$

$$\text{Ex:-} \quad \begin{bmatrix} 0 \\ d(3_1) & 0 \\ d(3_1) & d(3_2) & 0 \\ d(4_1) & d(4_2) & d(4_3) & 0 \end{bmatrix} \equiv \begin{bmatrix} 0 & 0.55 & 0 \\ 0.55 & 0 & 0.45 & 0 \\ 0.45 & 0.45 & 0 & 0.65 \\ 0.25 & 0.25 & 0.65 & 0 \end{bmatrix}$$

$$d(2,1) = \frac{|22-45|}{64-22} = \frac{23}{42} = 0.549 \approx 0.55 = (2,1)$$

$$d(3,1) = \frac{|64-45|}{64-22} = \frac{19}{42} = 0.452$$

$$d(3,2) = \frac{|64-22|}{64-22} = 1$$

$$d(4,1) = \frac{|28-45|}{64-22} = \frac{17}{42} = 0.404$$

$$d(4,2) = \frac{|28-22|}{64-22} = \frac{6}{42} = 0.14$$

$$d(4,3) = \frac{|28-64|}{64-22} = \frac{36}{42} = 0.857$$

$d_{1,2}$ are similar

proximity results b/w two ordinal attributes
 ordinal \Rightarrow it contains specific order.

(1) excellent

(2) good

(3) excellent

rank of value

to find the dissimilarity b/w ordinal attributes

$$\text{For ordinal data: } \frac{r_{if} - 1}{m_f - 1}$$

maximum rank

$$① \frac{3-1}{3-1} = 1$$

$$② \frac{1-1}{3-1} = 0$$

$$③ \frac{2-1}{3-1} = \frac{1}{2} = 0.5$$

$$④ \frac{3-1}{3-1} = 1$$

\Rightarrow dissimilarity matrix

$$\begin{bmatrix} 0 & d(2,1) & 0 \\ d(2,1) & 0 & 1 \\ d(3,1) & d(3,2) & 0 \\ d(4,1) & d(4,2) & d(4,3) & 0 \end{bmatrix} = \boxed{\begin{bmatrix} 0 & 0.5 & 0 \\ 0.5 & 0 & 0 \\ 1 & 0.5 & 0 \end{bmatrix}} \quad \text{1 = max. diff.}$$

$$d(2,1) = |\text{difference b/w 2 & 1}| = |0-1| = 1$$

$$d(3,1) = |0.5-1| = 0.5$$

$$d(3,2) = 0.5$$

$$d(4,1) = 0$$

$$d(4,2) = |1-0| = 1$$

$$d(4,3) = |1-0.5| = 0.5$$

$$d(4,4) = 0 \quad i.e. \text{similarity} = 1.0 + 1$$

proximity results for binary attributes.

The binary attribute values may be zero or 1, yes or no.

Object is absent To find the similarity b/w two objects,

we use Jaccard coefficient if

object is present

$$P = \frac{Q + R + S + T}{4} \quad (\text{for no. of attributes})$$

$$q \rightarrow \begin{cases} i=1 \\ j=1 \end{cases} \quad t \Rightarrow \begin{cases} i=0 \\ j=0 \end{cases}$$

$$r \rightarrow \begin{cases} i=1 \\ j=0 \end{cases} \quad s \rightarrow \begin{cases} i=0 \\ j=1 \end{cases}$$

$i = q$

$$\begin{pmatrix} q \\ r \\ s \\ t \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

$$P = q + r + s + t$$

black board 654

Name	gender	fever	cough	test 1	test 2	test 3	test 4
Jack	m	y	n	p.	n	n	n
Jim	m	y	y.	n	n	p	n
Mary	f	y	y	n	p.	n	n

$$P(i, j) = \frac{r+s}{q+r+s+t}$$

$$N=0 \quad d(jack, jim) = \frac{1+1}{1+1+1} = \frac{2}{3}$$

$$= \frac{2}{3}$$

$$d(jim, mary) = \frac{r+s}{q+r+s} = \frac{2}{4} = 0.5$$

$$= 0.66$$

$$d(mary, jack) = \frac{1+0}{2+1+0} = \frac{1}{3} = 0.33$$

$$d(jack, mary) = \frac{0+1}{2+0+1} = \frac{1}{3} = 0.33$$

(distance on ml) $3+2+4+3 = 12$ $P = 9$

Proximity Results for fixed variables

numerical.

$$\begin{array}{c} \text{Nominal values} \\ \left[\begin{array}{ccccc} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{array} \right] \end{array}$$

$$d(i,j) = \frac{\sum_{f=1}^p \delta_{ij}^{(f)}}{\sum_{f=1}^p \delta_{ij}^{(f)}}$$

for all attributes $\delta_{ij}^{(f)} = 1$ if $x_i^{(f)} = x_j^{(f)}$
 $\delta_{ij}^{(f)} = 0$ otherwise

$\delta_{ij}^{(f)} = 1$ means if all the attributes are equal

else $\delta_{ij}^{(f)} = 0$ (del - δ)

ordinal values

$$d(2,1) = \frac{1(1) + 1(1) + 1(0.95)}{1+1+1} = \frac{1+0.5+0.45}{3} = 0.65$$

$$d(3,1) = \frac{1(1) + 1(0.5) + 1(0.45)}{3} = \frac{1+0.5+0.45}{3} = 0.65$$

$$d(3,2) = \frac{1(1) + 1(0.5) + 1(1)}{3} = 0.83$$

$$d(4,1) = \frac{1(1) + 1(0) + 1(0.4)}{3} = 0.13$$

$$d(4,2) = \frac{1(1) + 1(0.5) + 1(0.4)}{3} = 0.71$$

$$d(4,3) = \frac{1(1) + 1(0.5) + 1(0.86)}{3} = 0.785$$



$$(0 + 0.65 + 1 \times 0.5 + 1 \times 0 + 1 \times 0.4 + 0.71 + 0.785 + 0.785 + 0.785) / 12 = 0.65$$

$$2\alpha = 2 + 2 + 2 + 2 =$$

$$2\beta = \sqrt{(1 + 0.65 + 0.71 + 0.785)^2 / 12} = 1.01$$

$$\alpha, \beta = \sqrt{(1 + 0.65 + 0.71 + 0.785)^2 / 12} = 1.01$$

5.12 The cosine similarity :
 if the data is missing in the documents, By using previous methods, we cannot find similarity . i.e the algorithm does not "work".

Sparse matrix \Rightarrow the matrix with more zeroes and less data values. $\begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix}$

It is used to handle the sparse data (the data values are zero). In information retrieval system most of the data objects are absent. for those documents, we will the normal similarity measures will not produce accurate similarity.

The cosine similarity is used to handle sparse data. if resultant cosine value = 1, the angle b/w x, y is 0 then the documents are more similar.

If cosine value is nearer to zero, then documents are dissimilar.

$$\text{Sim}(x, y) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

for example :

document	team	coach	baseball	(soccer)	tennis	(hockey)	win	loss	total
1	5	0	3	2	0	0	2	0	4
2	3	0	1	0	1	0	1	0	2
3	0	7	0	2	0	0	3	0	3
4	0	1	0	0	1	2	2	0	3

term frequency table. (the may time a word is repeated in a document).

$$\begin{aligned} \text{Sim}(x, y) &= 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 0 + 2 \times 1 + 0 \times 0 + 0 \times 1 \\ (1, 2) &= \sqrt{5^2 + 0^2 + 3^2 + 2^2 + 2^2} = \sqrt{25 + 9 + 4 + 4} = 6.48. \\ \|x\| &= \sqrt{9 + 4 + 1 + 1 + 1} = 4.12. \\ \|y\| &= \sqrt{9 + 4 + 1 + 1 + 1} = 4.12. \end{aligned}$$

$$\text{Sim}(x,y) = \frac{25}{6.48 \times 4.12} = 0.936 \approx 0.93$$

$$\text{there the doc1 \& doc2 were more similar}$$

If the document attributes were binary attributes, then the

$$\text{similarity of } x,y = \frac{x \cdot y}{x \cdot x + y \cdot y - x \cdot y}$$

\rightarrow Data Transformation:

The Transf. data are transformed or consolidated into the required format so that, the resultant mining process may be more efficient and the patterns are easier to understand.

Types of Data Transformation:

1. Smoothing

2. Aggregation

3. Generalisation

4. Normalisation $F(0,0) = \frac{25 \cdot 21 - 8}{48 \cdot 2} = 11.75$

5. Attribute Construction

Normalisation: In normalisation, the 2 attribute data are scaled within

a specified range $[(-1 to 1) \text{ or } (0 to 1)]$

In $\Rightarrow 0.1 < 25 \cdot 21 - 8 = 0.8$

Types of normalization:

1. min max normalization

2. Z-score normalization

3. Decimal scaling

Min-Max Normalisation:-

$$N = \frac{Y - \text{min}}{\text{max} - \text{min}} \left((\text{new max} - \text{new min}) + \text{new min} \right)$$

$$\text{Example: marks } 8 \quad 10 \quad 15 \quad 20 \quad \frac{8-0}{20-0} = \frac{8}{20} = 0.4$$

$$\text{new max} = 1$$

$$\text{new min} = 0$$

$$V' = \frac{8-0}{20-0} (1-0) + 0 = \frac{8}{20} = 0.4$$

$$V' = \frac{10-8}{20-8} (1-0) + 0 = \frac{1}{16} = 0.1$$

$$V' = \frac{15-8}{20-8} (1-0) + 0 = \frac{7}{12} \cdot 0.58 + \frac{25}{12} = 0.58$$

$$V' = \frac{20-8}{20-8} (1-0) + 0 = 1 - 0 = 1$$

Z-Score Normalisation:

$$S = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \sqrt{\frac{\sum (x_i - \mu/\sigma)^2}{n-1}}$$

$$Z\text{-score} = \frac{x - \mu/\sigma}{\sigma}$$

Example: Marks: 8, 10, 15, 20
Mean: $\bar{x} = \frac{8+10+15+20}{4} = \frac{53}{4} = 13.25$ and Standard Deviation: $S = \sqrt{\frac{(8-13.25)^2 + (10-13.25)^2 + (15-13.25)^2 + (20-13.25)^2}{3}}$

$$S = \sqrt{\frac{(8-13.25)^2 + (10-13.25)^2 + (15-13.25)^2 + (20-13.25)^2}{3}} = 5.37,$$

$$Z\text{-score for } 8 = \frac{8-13.25}{5.37} = -0.97$$

$$Z\text{-score for } 10 = \frac{10-13.25}{5.37} = -0.60$$

$$Z\text{-Score for } 15 = \frac{15-13.25}{5.37} = 0.32$$

$$Z\text{-Score for } 20 = \frac{20-13.25}{5.37} = 1.25$$

Since our range is not between 0 and 1, we have to scale it.

Exponential Scaling:

$$V' = \frac{V}{10^T}$$

T is no. of digits.
Range: $V_{\min} - V_{\max}$ to $V_{\max} - V_{\min}$

$$V' = \frac{8}{10} = 0.8$$

$$V' = \frac{10}{10} = 1.0 = 0.1$$

$$V' = \frac{15}{10} = 0.15$$

$$V' = \frac{20}{10} = 0.20$$

$$V' = \frac{8}{8} = 1 = 0 + (0-1) \frac{8-0}{8-8} = 1$$

Smoothing: which helps to remove the noise from the data
the techniques that are used are

1. Binning
2. Clustering
3. Regression.

Aggregation: where the summary of Aggregation operations are applied to the given data.

Ex:- Rollup & cube operations in DW.
Generalization: the data where the low level or primitive data are placed by the higher level data by using the concept of Hierarchies.

Ex:- Address.

Attribute Construction: where a new attributes are constructed and added from a given set of attributes to help the mining process.

Ex:- Feature extraction (PCA).