

INTRODUCTION TO DATA MINING

Data mining is the extraction of data for processing data or finding data, ~~is~~ that can be used.

It is knowledge discovery from data.

Also called

- KDD (knowledge discovery in databases)
- data archeology

* KDD. (5 M)

Collecting data

Data cleaning

Finding if data is missing or any other data is required.

Data Integration

Data Warehouse

↓ Data Transformation

All the information is in a particular format.

choosing the data that is required.

← Selection

↓ Data Mining

Pattern Evaluation & presentation

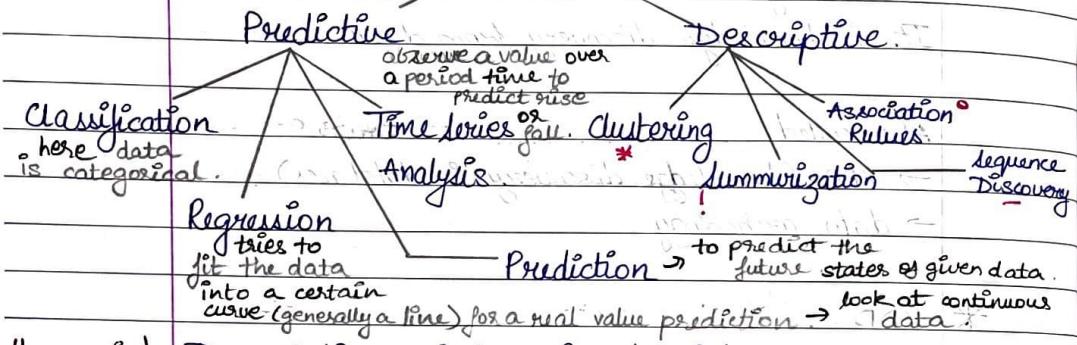
↓

Knowledge

Data mining is required to find information that can be related to each other, while query processing can be done to only find the information from already available data.

* DATA MINING MODELS AND TASKS

Data Mining



Unsupervised Descriptive Mining takes into account only the data that is currently available and gives information. It does not require prior information.

Supervised Predictive Mining takes into account prior information. It allows to predict future information or to get data that must be missing.

* Partitioning of data on certain condition set by us.

! describing the data or summarizing or generalisation of data.

• relationship between different data

eg. Market - basket analysis used in supermarkets same time to increase the sale of items.

- to determine a pattern between different data.

to make this is done over a period of time (eg. a month).

* TYPES OF DATA

1. DATASET

Consists of datapoints or data object. Can be mapped in any dimension. Usually mapped in any two dimension. Each datapoint have attributes or dimensions.

a₁ a₂ a₃ ... a_n ← attributes

d ₁				
d ₂				
d _n				

↑ data objects

a_m is generally
label
class

d_i = <a_{i1}, a_{i2}, a_{i3}, .. a_{in}>
↑ FEATURE/ATTRIBUTE VECTOR

Observed values for a given attribute are known as observations.

A set of attributes used to describe a given object is called attribute vector / feature vector.

* DATA OBJECTS & ATTRIBUTE TYPE

1. Nominal / Categorical attributes

Related to names ; Values are symbols or names. Cannot find mean, median but mode can be found easily.

eg. Hair styles, colors

Nominal values may be numeric.

They are not quantitative.

2. Binary Attributes

Can only have 2 states (0 / 1)

0 - value absent = 1 - value present

Two types

→ symmetric ⇒ both have same weight, are equally probable.

→ A ~~symmetric~~ binary \Rightarrow both states are not equally valuable. Here one state is more favourable than the other! eg. In medical test, more favourable ones is ~~more than~~ ~~than~~ less than.

3. Ordinal attributes

In the form of ranks

eg. Grades \rightarrow A+, A, A-, B, ...

They have a definite order.

4. Numeric attributes.

Continuous data; Measurable Quantity.

Two types

\Rightarrow Interval scaled.

Can be compared, find the difference.

Here '0' means absence of ~~the~~ value for attribute.

\Rightarrow Ratio scaled

Cannot be compared, values are ordered.

Mean, median & mode can be found.

Here '0' has a meaning on the scale.

eg. Temperature

DISCRETE \rightarrow finite number of values.
eg. medical test, points.

CONTINUOUS \rightarrow ~~countable~~ numeric values.

eg. age,

COUNTABLE INFINITY \rightarrow values can grow infinite but still countable.

eg. Zipcode.

* MEAN, MEDIAN, MODE

Mean, $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Weighted mean $\bar{x} = \frac{\sum w_i x_i}{\sum w_i}$

\Rightarrow Median = $L_1 + \left[\frac{N/2 - (\sum f_{\leq m})}{f_{m \text{ median}}} \right] \text{ width}$

$L_1 \rightarrow$ lower boundary of the median interval

$N \rightarrow$ no. of values.

$f_{\leq m} \rightarrow$ sum of frequencies of intervals lower than the median interval.

$f_{m \text{ median}} \rightarrow$ frequency of median interval.

width \rightarrow width of median interval.

Mode \rightarrow values that occurs most frequently.

Types.

\rightarrow Unimodal \rightarrow Bimodal

\rightarrow Trimodal \rightarrow Multimodal

Midrange \rightarrow Average of largest and smallest values in the set.

Grouped mean = $\frac{\sum (f_i \times x_{im})}{n}$

$n \rightarrow$ total number of observations.

$f_i \rightarrow$ frequency of i^{th} observation

$x_{im} \rightarrow$ midpoint of i^{th} x_m .

Grouped median :

We first find where $\frac{N}{2}$ observation lies in which range. & use the above median formula.

$$\text{Grouped mode} = L + \left[\frac{f_m - f_{m-1}}{(f_m - f_{m+1}) + (f_m - f_{m-1})} \right] \times \text{width}$$

$L \rightarrow$ lower limit of the group with the mode
(group with highest frequency)

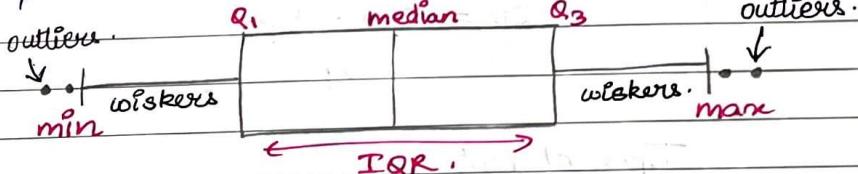
$f_m \rightarrow$ frequency of group with mode

$f_{m-1} \rightarrow$ frequency of group before f_m .

$f_{m+1} \rightarrow$ frequency of group after f_m .

width \rightarrow width of the groups.

Box plot

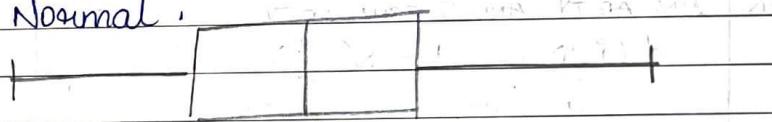


Outliers \rightarrow beyond min & max.

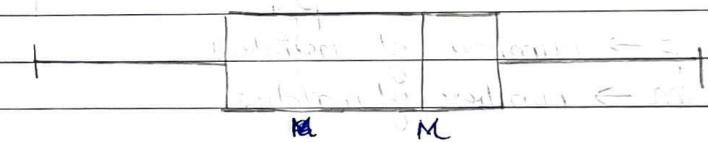
Minimum $\rightarrow Q_1 - 1.5 * IQR$

Maximum $\rightarrow Q_3 + 1.5 * IQR$

Normal



Negatively skewed



NEGATIVELY
SKEWED

NORMALLY

POSITIVELY
SKEWED

MEASURING THE DISPERSION DATA

Range \rightarrow difference between largest max & smallest min values.

Quartiles \rightarrow points taken at regular intervals in a distribution.

4-points : Q_1 , Q_2 , Q_3

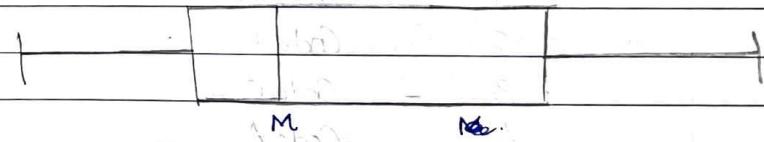
25% 50% 75%

Interquartile range $\rightarrow IQR = Q_3 - Q_1$.

Used to find if there are any outliers in the distribution.

Five no. summary \rightarrow min, Q_1 , M, Q_3 , max.

Positively skewed



Variance

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2$$

Standard Deviation

$$\delta = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Z-score

$$Z_{ij} = \frac{x_{ij} - m_{ij}}{s_j}$$

Types of plot:

1. Quantile Plot.
2. Q-Q Plot.
3. Histograms.
4. Scatter plot.
5. Loess curve.

* SIMILARITY AND DISSIMILARITY.

$$d(i, j) = 1 - s(i, j)$$

DISTANCE MATRIX DISSIMILARITY	
0	
$d(2,1)$	0
$d(3,1)$	$d(3,2)$
\vdots	\vdots
$d(n,1)$	$d(n,2)$
	... 0

PROXIMITY MEASURE FOR NOMINAL / CATEGORICAL DATA

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

$p \rightarrow$ number of variables.

$m \rightarrow$ number of matches.

e.g. Obj. Id. Test 1

- 1 Code A
- 2 Code B
- 3 Code C
- 4 Code A.

$$p = 1$$

$m = 0 \rightarrow$ if the mat do not match.

1, if they match.

$$d(1,2) = \frac{1-0}{1+1} = 1$$

PROXIMITY MEASURE FOR BINARY VARIABLES.

Contingency table

Object j

	1	1 - 0	sum
Object i	a	b	a+b
0	c	d	c+d
	a+c	b+d	-p

$a \rightarrow$ no. of attributes. 1 & 1

$b \rightarrow$ no. of attributes. 1 & 0.

$c \rightarrow$ no. of attributes. 0 & 1.

$d \rightarrow$ no. of attributes. 0 & 0.

Distance measure for symmetric binary variables

$$(1) d(i, j) = \frac{b+c}{a+b+c+d}$$

Distance measure for asymmetric binary variables

$$d(i, j) = \frac{b+c}{a+b+c+d}$$

Jaccard coefficient (similarity measure for asymmetric binary variables)

$$\text{Sim}_{\text{Jaccard}}(i, j) = \frac{a}{a+b+c}$$

DISSIMILARITY OF NUMERIC DATA.

EUCUDIAN DISTANCE

$$d(i, j) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_p - y_p)^2}$$

MINKOWSKI DISTANCE

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

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

PROXIMITY MEASURE FOR ORDINAL DATA

→ Has to be normalized using z-score.

$$z = \frac{x_{ij} - \bar{x}_j}{\text{rank } \bar{x}_j}$$

$$\bar{x}_j = \frac{1}{n_j}$$

$$\text{Total no. of ranks} = n_j$$

$$b_1 < b_2 < b_3 < b_4$$

For ratio-scaled, we find its log values and then apply Manhattan / Euclidian formula

For numeric data, we can find the minimum & maximum and then use the difference between them as denominator.

$$\text{Formulas} \leftarrow d(i, j) = \frac{\sum_{f=1}^p s_{ij}(f) d_{if}(f)}{\sum_{f=1}^p s_{ij}(f)}$$

* DATA PREPROCESSING

Data collected should not be

→ incomplete → noisy → inconsistent

Data collected should be preprocessed for the following reasons.

→ accuracy.

→ completeness.

→ consistency.

→ timeliness.

→ believability.

→ interpretability.

→ accessibility.

→ value added. → e.g., giving meta data with the data for source information.

Major tasks in Data Preprocessing.

1. DATA CLEANING

→ Identifying incomplete data.

→ resolve inconsistencies.

→ smooth noisy data.

2. DATA INTEGRATION

Integrate into multiple databases.

3. DATA TRANSFORMATION

Normalization & reduction.

4. DATA REDUCTION

5. DATA DISCRETION

How to handle Missing Data?

→ Ignore the tuple : usually done when class label is missing. This may result in loss of imp. data

→ Fill it automatically :

For this we can find the most probable value - regression, inference based such as Bayesian formula or decision tree.

→ Filling manually can be tedious & infeasible

→ Using a constant value for replacing missing data may result in the formation of a new class

How to handle Noisy data?

Bining →

Data preprocessing method used to minimize the

effects of small observation errors.

smooth the data by looking at neighbouring data.

Regression →

smooth by fitting the data into regression functions.

Clustering →

detect and remove outliers. (clusters use centroids)

Combined computer & human inspection.

* BINNING METHOD FOR SMOOTHING DATA.

Sorted data → 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29

34

Partition into equal-frequency (equi-depth) bins.

Bin 1 : 4, 8, 15

Bin 2 : 21, 21, 24

Bin 3 : 25, 28, 34

Smoothing by bin means.

Bin 1 : 9, 9, 9

Bin 2 : 22, 22, 22

Bin 3 : 29, 29, 29

Smoothing by bin boundaries.

Bin 1 : 4, 4, 15

Bin 2 : 21, 21, 24

Bin 3 : 25, 25, 34

* DATA CLEANING.

→ Data discrepancy detection.

- make use of meta data (eg. domain, range, etc.)
- check field overloading.
- check unique rule, consecutive rule & null rule.
- use commercial tools.

→ Data scrubbing : use simple domain knowledge to detect errors and make corrections.

→ Data auditing : by analysing data to discover rules & relationship to detect violations.

→ Data migration & integration.

* DATA INTEGRATION.

→ Combine data from multiple sources into a coherent store.

→ First step. Entity Identification problem.

- To identify the different attributes and their relationship.
- Mapping of data
- To associate similar data.

→ Second step. Redundancy & Correlation Analysis.

- Redundancy can be found using correlation analysis.

CORRELATION ANALYSIS.

χ^2 - TEST. (CATEGORICAL DATA).

$$\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$$

Larger the value of χ^2 , more likely the variables are related.

PEARSON'S PRODUCT OF MOVEMENT COEFFICIENT (NUMERICAL DATA)

$$\rho_{A,B} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(n-1)\sigma_A \sigma_B}$$

$$= \frac{\sum (AB) - n(\bar{A}\bar{B})}{(n-1)\sigma_A \sigma_B}$$

$\rho_{A,B} > 0 \rightarrow$ positive correlation

$\rho_{A,B} < 0 \rightarrow$ negative correlation

$\rho_{A,B} = 0 \rightarrow$ no relation

$$\rho_{A,B} = \frac{\text{COV}(A, B)}{\sigma_A \sigma_B}$$

$$\text{COV}(A, B) = \frac{1}{n} \sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})$$

* DATA REDUCTION

HISTOGRAMS

→ equal-frequencies

→ equal-width

histogram frequencies

clustering → combined into bins

binning → overlapping bins

sampling → random selection

cluster or stratified sampling

DATA FATTENING

→ adding noise

but good

add noise to training set and test set

CLASSIFICATION AND PREDICTION

What is classification?

→ Data is generally discrete.

→ Data ~~is also~~ can also be categorical

→ It is supervised learning.

→ Data is classified based on previous historical data.

* SUPERVISED LEARNING

Classification

The training set → known labels

Test set → unknown labels

* CLASSIFICATION - A TWO-STEP PROCESS

1. MODEL CONSTRUCTION

→ describing a set of predetermined classes

→ model is created on a training set.

2. MODEL USAGE

→ for classifying future or unknown object

→ model is chosen based on its accuracy

* ISSUES

Data preparation.

→ Data cleaning.

→ Relevance analysis.

→ Data Transformation

Evaluating Classification Methods.

→ Accuracy

Classifier accuracy - predicting class label

Predictor accuracy - giving value of predicted att.

→ Speed.

→ Robustness.

→ Simplicity.

→ Interpretability.

* BAYESIAN CLASSIFICATION.

Based on Baye's Theorem.

→ Prior knowledge is required for probability as with new entries entries probability may increment or decrement.

$X \rightarrow$ data sample : class label is unknown.

$H \rightarrow$ hypothesis that X belongs to class c .

$P(H|X) \rightarrow$ probability that the hypothesis holds given the observed data sample X .

$P(H) \rightarrow$ prior probability, initial probability.

$P(X) \rightarrow$ probability the sample data is observed

$P(X|H) \rightarrow$ posteriori probability.

$$P(H|X) = \frac{P(X|H) \cdot P(H)}{P(X)}$$

Informally, it can be written as.

posteriori = likelihood \times prior evidence

Classification is done on the basis of maximum posteriori.

EXAMPLE ① - For categorical/discrete data.

$P(C_i) \rightarrow$ Prior probability of class each class.

$$P(\text{buys_comp} = \text{"yes"}) = \frac{9}{14} = 0.643$$

$$P(\text{buys_comp} = \text{"no"}) = \frac{5}{14} = 0.357$$

Compute likelihood of $P(X|C_i)$ of each attribute.

$$P(\text{age} = \text{"<=30"} | \text{buy_comp} = \text{"yes"}) = \frac{2}{9}$$

$$P(\text{age} = \text{"< 30"} | \text{buy_comp} = \text{"no"}) = \frac{3}{5}$$

$$\begin{aligned}
 & P(X) = P(X|C_1)P(C_1) + P(X|C_2)P(C_2) + P(X|C_3)P(C_3) \\
 & (H)9 \cdot (H|X)9 = (X|H)9 \\
 & (X)9
 \end{aligned}$$

$$\begin{aligned}
 P(X|C_i) &= P(X|buys_comp = "yes") \\
 &= 0.044.
 \end{aligned}$$

$$P(C_i|X) = P(X|C_i) * P(C_i)$$

$$\begin{aligned}
 P(buys_comp = "yes" | X) &= 0.028. \\
 &= 0.028.
 \end{aligned}$$

$$\begin{aligned}
 P(buys_comp = "no" | X) &= 0.007. \\
 &= 0.007
 \end{aligned}$$

$P(\text{buys_comp} = \text{yes} | X) > P(\text{buys_comp} = \text{no} | X)$
 $\therefore X$ belongs to class ($\text{buys_comp} = \text{"yes"}$).

If any instance of a class is missing, the probability will be zero, does it is directly classified into other class.

EXAMPLE ② - For continuous attributes.

$$P(\text{class} = \text{"short"}) = \frac{4}{15}$$

$$P(\text{class} = \text{"medium"}) = \frac{8}{15}$$

$$P(\text{class} = \text{"tall"}) = \frac{3}{15}$$

$$X = (\text{gender} = M, \text{height} = 1.95)$$

$$P(\text{gender} = M | \text{class} = \text{"short"}) =$$

$$P(\text{gender} = M | \text{class} = \text{"medium"}) =$$

$$P(\text{gender} = M | \text{class} = \text{"tall"}) =$$

$$\begin{aligned}
 P(x_k | C_i) &= g(x_k, \mu_{ci}, \sigma_{ci}) \\
 &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}
 \end{aligned}$$

continuous.

$$P(\text{height} = 1.95 | \text{class} = \text{"short"}) = \text{Range}(1.6m - 1.7m)$$

$$\begin{aligned}
 \mu &= 1.65 & \sigma &= 0.05 \\
 g(x, \mu, \sigma) &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}
 \end{aligned}$$

$$= \frac{1}{\sqrt{2\pi} \times 0.05} e^{-\frac{(1.95-1.65)^2}{2 \times 0.05^2}}$$

$$= 1.2154 \times 10^{-7}$$

$$P(\text{height} = 1.95 \mid \text{class} = \text{"medium"}) = \text{Range}(1.75-1.95)$$

$$\mu = 1.85, \sigma = 0.1$$

$$g(x, \mu, \sigma) =$$

$$P(\text{height} = 1.95 \mid \text{class} = \text{"tall"}) = \text{Range}(2m-2.2m)$$

$$\mu = 2.1, \sigma = 0.1$$

$$g(x, \mu, \sigma) =$$

$$\mu = 1.29.$$

$$\sigma = 0.1$$

Posteriori Probability of C conditioned on X.

Example ③ If $x = \text{brown, tall, average, no}$

$$X = \langle \text{brown, tall, average, no} \rangle$$

Prior probability

$$P(\text{class} = \text{"Sunburnt"}) = \frac{6}{10}$$

$$P(\text{class} = \text{"None"}) = \frac{4}{10}$$

$$P(\text{brown} \mid \text{sb}) = \frac{2}{6} = \frac{1}{3}$$

$$P(\text{brown} \mid \text{none}) = \frac{3}{6} = \frac{1}{2}$$

$$P(\text{tall} \mid \text{sb}) = \frac{1}{6}$$

$$P(\text{tall} \mid \text{none}) = \frac{1}{6}$$

$$P(\text{avg} \mid \text{sb}) = \frac{2}{6}$$

$$P(\text{avg} \mid \text{none}) = \frac{2}{4}$$

$$P(\text{no} \mid \text{sb}) = \frac{5}{6}$$

$$P(\text{no} \mid \text{none}) = \frac{3}{4}$$

Posteriori Probability of class on X.

$$P(X \mid \text{sb}) = \frac{6}{10} \times \frac{1}{6} \times \frac{2}{6} \times \frac{5}{6} \times \frac{6}{10}$$

$$= 0.0154 \times 9.34 \times 10^{-3}$$

$$P(X \mid \text{none}) = \frac{3}{10} \times \frac{1}{4} \times \frac{2}{4} \times \frac{3}{4} \times \frac{4}{10}$$

$$= 0.0703 \times 0.028$$

So X belongs to class = "none"

Between 2 data, no conditional dependency is present. If 2 data are conditionally dependent then Naive Bayes cannot be used.

LAPLACE CORRECTION

Avoiding the 0-Probability Problem.
We add a tuple ~~to~~ for each type of attribute.

For example .

$$\text{Income (low)} = 0$$

$$\text{income (med)} = 990$$

$$\text{income (high)} = 10.$$

We add 1 to each case.

$$\text{Prob (Income (low))} = 1/1003.$$

$$\text{Prob (Income (med))} = 991/1003$$

$$\text{Prob (Income (high))} = 11/1003.$$

The probability will be very small , but posterior probability won't decrease .

* DECISION TREE INDUCTION

ID 3 - INTERACTIVE DICHOTOMIZER 3.

Internal Nodes - Attributes

Leaf Node - Classes

Branches - Values of Attribute

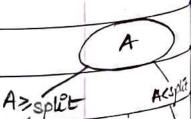
For categorical,

→ Branches are based on the value of attributes.



For continuous,

→ A split point is found



Algorithm

Attribute Selection : Information Gain ID3.
Select attribute with the highest information gain.

$p_i^o \rightarrow$ probability that an arbitrary tuple in D belongs to class C.

$$p_i^o | C_i, p_i^o$$

|D|.

Expected Information needed to classify a tuple in D.

$$\text{Info}(D) = - \sum_{i=1}^m p_i^o \log_2(p_i^o)$$

Information needed (after using A to split D into v partitions) to classify D.

$$\text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Info}(D_j^o).$$

Information gained by branching on attribute A.

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D).$$

Example ①

Class → buys - comp = "yes"

buys - comp = "no".

$$\text{Info}(D) = - \frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

$$= 0.940 \text{ bits.}$$

age	p_i	n_i	$I(p_i, n_i)$
≤ 30	2	3	0.971
31 - 40	4	0	0
> 40	3	2	0.971

$$\text{Info age (D)} = \frac{5}{14} I(2, 3) + \frac{4}{14} I(4, 0) + \frac{5}{14} I(3, 2)$$

5 $I(2, 3) \rightarrow$ means "age ≤ 30 " has 5 out of 14 sample, with 2 yes's & 3 no's.

$$I(2, 3) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.971$$

$$\text{Info age (D)} = \frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971$$

= 0.694 bits.

$$\text{Gain (age)} = \text{Info (D)} - \text{Info age (D)}$$

$$= 0.940 - 0.694$$

$$= 0.246$$

$$I(\text{Income (D)}) = \frac{4}{14} I(2, 2) + \frac{6}{14} I(4, 2)$$

$$= \frac{4}{14} I(2, 2) + \frac{4}{14} I(4, 2)$$

M	T	W	T	F	S	S
Page No.:						
Date:						YOUVA

* the metadata in the data warehouse should be same for every user viewing it
 (VIVA Q's → Evolution)

M	T	W	T	F	S	S
Page No.:						
Date:						YOUVA

INTRODUCTION TO DATA WAREHOUSE.

* DATA WAREHOUSE

→ single*, complete & consistent store of data obtained from a variety of different sources, made available to end users in a way they can understand & use in a business context.

* EVOLUTION

1960's → Batch reports.

was hard to find & analyse information.

1970's → Terminal based Decision systems (DSS) and Executive information systems (EIS).

1980's → Desktop data access

1990's → Data warehousing.

Terabytes - 10^{12} eg. Walmart data.

Petabytes - 10^{15} Geographic Information sys.

Exabytes - 10^{18} .

Zettabytes - 10^{21}

Zottabytes - 10^{24}

Data Warehouse is the collection of historic data
 Database is the collection of data that is used for day to day transactions

Data warehouses are not replacing databases

OLTP → Online Transaction Process
 OLAP → Online Analytic Processing.

M	T	W	T	F	S	S
Page No.:						YOUVA
Date:						

What - is (Banks)

Type of account present

Number of transactions

Withdrawals & Deposits done.

Data granularity

→ data should be
in summarized
format.

DEFINITION *

A data warehouse is at.

→ subject-oriented eg. sales, product, customer.

→ integrated. → data inconsistency is removed.

→ time-varying historical data, snapshots of data,

allows analysis of past, relates info to present.

→ non-volatile → data cannot be changed.

collection of data that is used primarily in organizational decision making.

* INFORMATION PACKAGE

What were the sales of Godrej refrigerators of white colour in the year 2010 January in Tulu area? → Query.

A. Involves ~~historical data~~ ^{DIMENSIONS}

Product	Time	Region
Colour(white)	Year (2010)	Juhu
	Month (Jan.)	South Africa

Since data collected is historical, 'time' dimension must always be included.

ER

IP

Entity → Dimensions

Relationships → Facts.

Pg 120 → DW by Tonniah - TP examples.

(Databases)

OPERATIONAL (OLTP)

→ Current values

→ Optimized for transaction.

→ High access frequency.

→ Read, update, delete

→ Predictable, repetitive usage.

→ Response time in sub-seconds

→ Large number of users

M T W T F S S

Page No.:						
Date:						YOUVA

M T W T F S S

Page No.:						
Date:						YOUVA

(Data Warehouses).

INFORMATIONAL (OLAP)

Archived, derived, summarized

Optimized for complex queries

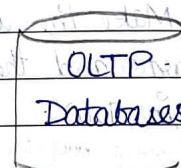
Medium to low

Read.

Adhoc, random heuristic usage.

Several seconds to minutes.

Relatively small.



Add, delete, insert with update.

Decision making

3 data granularity levels in the bank.

Daily

Monthly

Quarterly

For a data cube with dimensions more than 3, are called hypercube.

Identify the subjects for Airline

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:	YOUVA					

Airline

Airline Name

Airline sales

No. of passengers

Cost.

Departure Time

Arrival Time

Type of seat

Domestic / Int.

Departure

Destination

No. of stops in between

No. of seats left

Flight time

No. of kg of luggage

Services

Time seat was booked

Top down Approach

→ Make the DW and
then break it down.
(usually should be
done)

Bottom Up Approach.

Make the small
parts and then integrate
in DW.
(data marts → DW),

DATA CUBE

A visualization technique

IP for Airlines

In paper
why we
have
chosen the
given
dimension,
subject of
analysis
for flight

TIME	FLIGHT	CUSTOMERS	BOOKING	ROUTE
Year	No. of seats booked	Name	Company	Departure
Month	Type of flight	Gender	Booking Date	Arrival
Day of the Month	Type of class	Age		No. of stops
Day of the Week	Cost	No. of luggage		Time taken
Hour	Services			
Minute	Total no. of seats			
Season				
Quarterly				

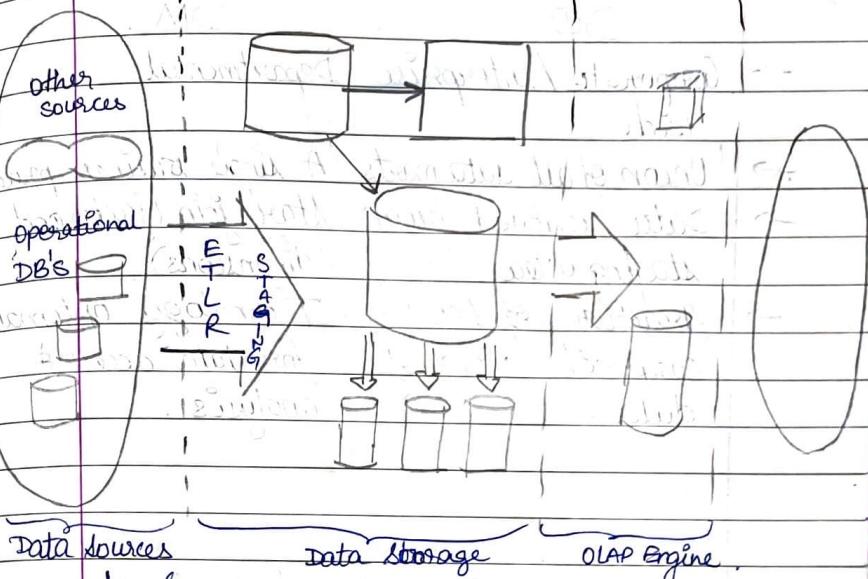
FACTS : No. of seats vacant, most used route,

Time flight is booked, Time when most seats are booked, Class in which seats are booked more, if catering was chosen

Metadata - data above data.

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:	YOUVA					

* DATA WAREHOUSE ARCHITECTURE



1. SOURCES

2. STAGING

→ ETL Process

Extract Transform Load Refresh.

Metadata - varies from business to business.

It stores the rules to handle data present correctly.

Data Mart - Every individual business processes have a mini data warehouse in large organizations. These are called data marts. These data refer to one large data warehouse.

3. OLAP

4. Analysis

All the data present are for AI, DS are applied in Stage 3 & 4

~~TOP~~

* DIFFERENCE B/w DATA WAREHOUSE & DATA MART

DW

DM

- Corporate / enterprise wide
- Union of all data marts
- Data received from staging area.
- Structure for corporate view of data.
- Departmental
- A single business process
- Star join (facts and dimensions)
- Technology optimal for data access & analysis.

* DESIGN DECISIONS

- Choosing the process.
 - select subjects from IP.
- Choosing the grain.
 - level of detail is important
- Identifying and conforming the dimensions.
- Choosing the facts.
 - selecting metrics, units of measurement.
- Choosing the duration of database.
 - how much historical data is required.

* CRITERIA FOR COMBINING TABLES INTO DIMENSION TABLE

- Model should provide best data access.
- Whole model should be query centric.
- Optimized for query and analysis.
- Every dimension can interact equally with fact table.
- Should allow table drill up and drill down, along dimensional hierarchy.

* STAR SCHEMA

A star schema represents the relationship between the different dimensions ~~and it~~. The dimensions directly interact with the fact table.

* POSSIBLE ETL FAILURES

- A missing source file.
- A system failure.
- Inadequate metadata.
- Poor mapping information.
- Inadequate storage planning
- A source structural change.
- No contingency plan.
- Inadequate data validation.