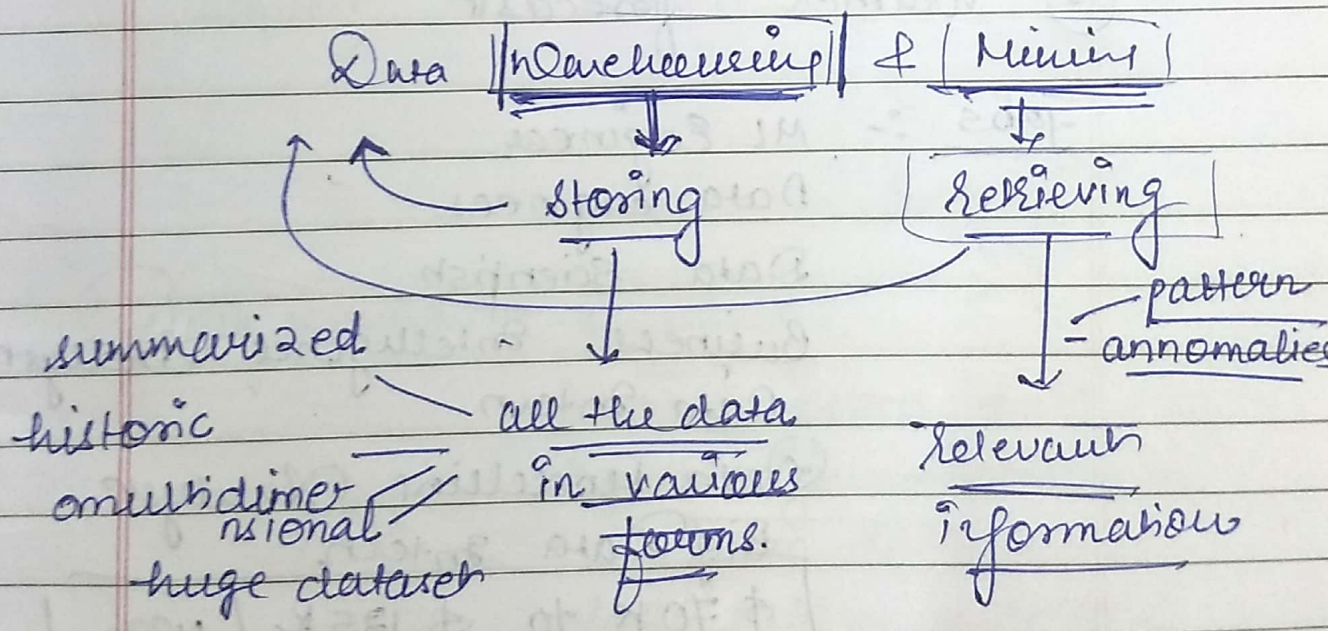


Data Mining Concepts & Techniques
Lecture 1

contents :-

- Introduction to the Subject
- Importance of the Subject
- Applicability of the Subject
- Unit 1
 - Importance of DM & DM Architecture
 - ~~Data Mining Functionalities~~
 - ~~Classification of DM systems.~~

Introduction to subject



IoT + ML + AI → Increasing data, over the cloud.

huge data collection so as to analyze and find insights. Data collected on a data warehouse, mining is cleaning - data filling - visualizing the findings

From the huge sets of data, looking for patterns, anomalies, associations, with the goal of extracting value.

- ① market analysis
- ② online transaction theft detection
- ③ improving user experience based on the filters based on user choice.
- ④ strategic decisions in companies
- ⑤ energy saving of appliances based on usage
- ⑥ weather forecast

jobs :- ML Engineer
 Data engineer
 Data Scientist
 Business Intelligence Analyst
 Data Intern

① Data Modelling Analyst
 Big Data Intern

\$ 70K to \$ 125K / year

tools : { Spark } | Java - Language
 { R } | Python
 { Hadoop }

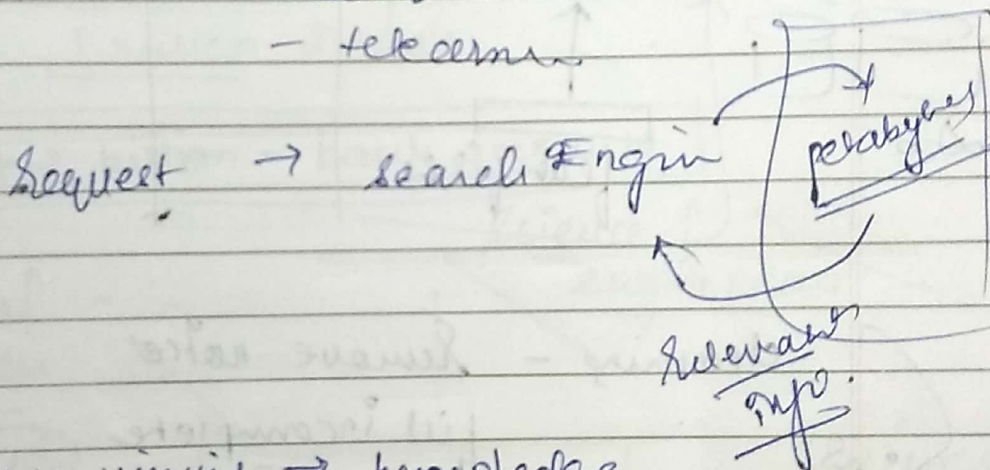
Micromasters Big Data Certificate Program

size of the digital universe will double every 8 years

50x growth rate of data.

Exponential growth rate towards 2020.

- business
- medicine
- science & engineering
- social media
- telecom



Data mining → knowledge

mining from data.

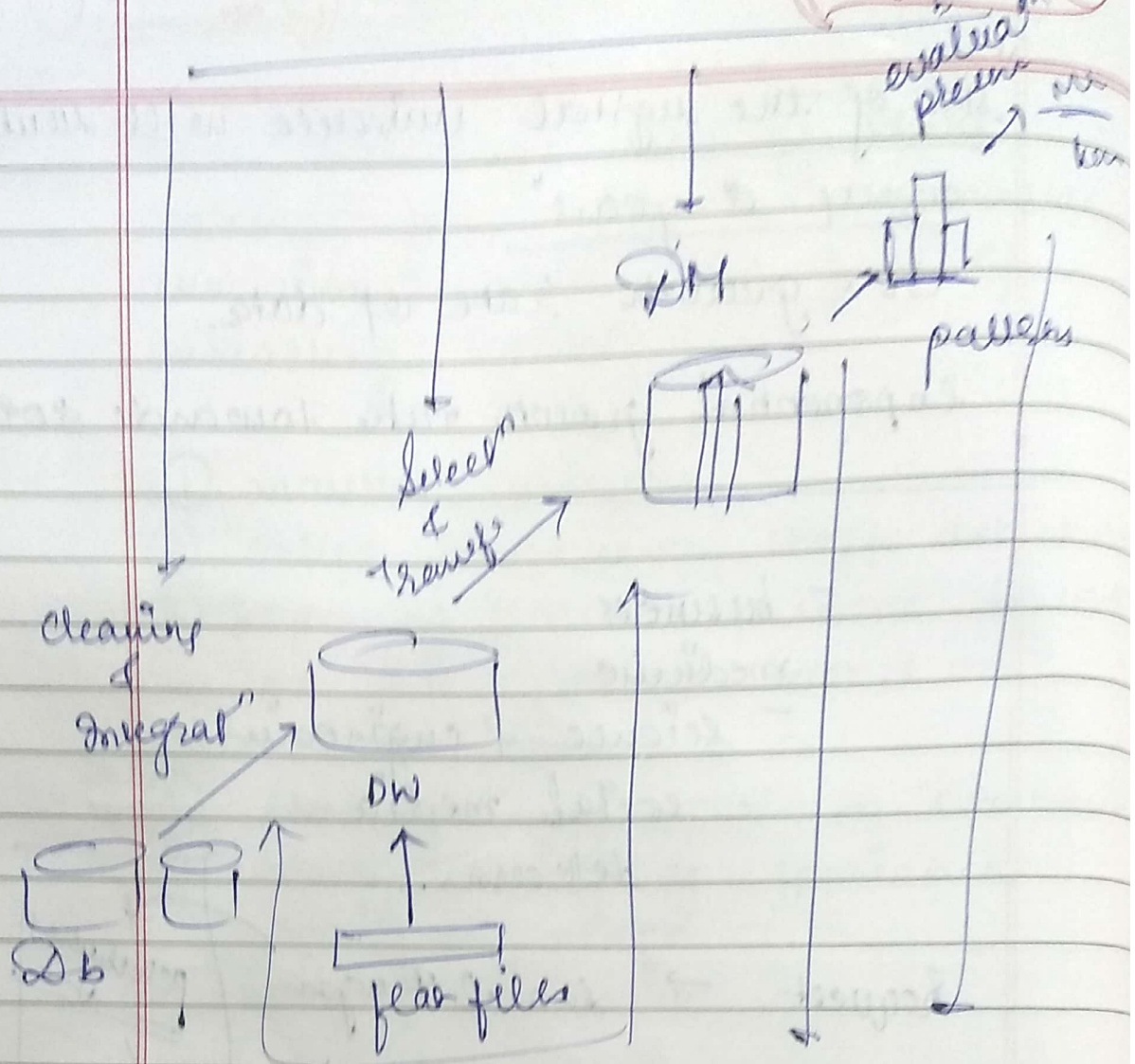
knowledge discovery from data

KDD

data pattern analysis

data dragging

data archaeology



data processing

cleaning - remove noise
 fill incomplete
 remove inconsistency

integrat - multiple sources
 combined

select - relevant to analysis
 data retrieved

transform - data transformed &
 consolidated into forms
 appropriate for mining.

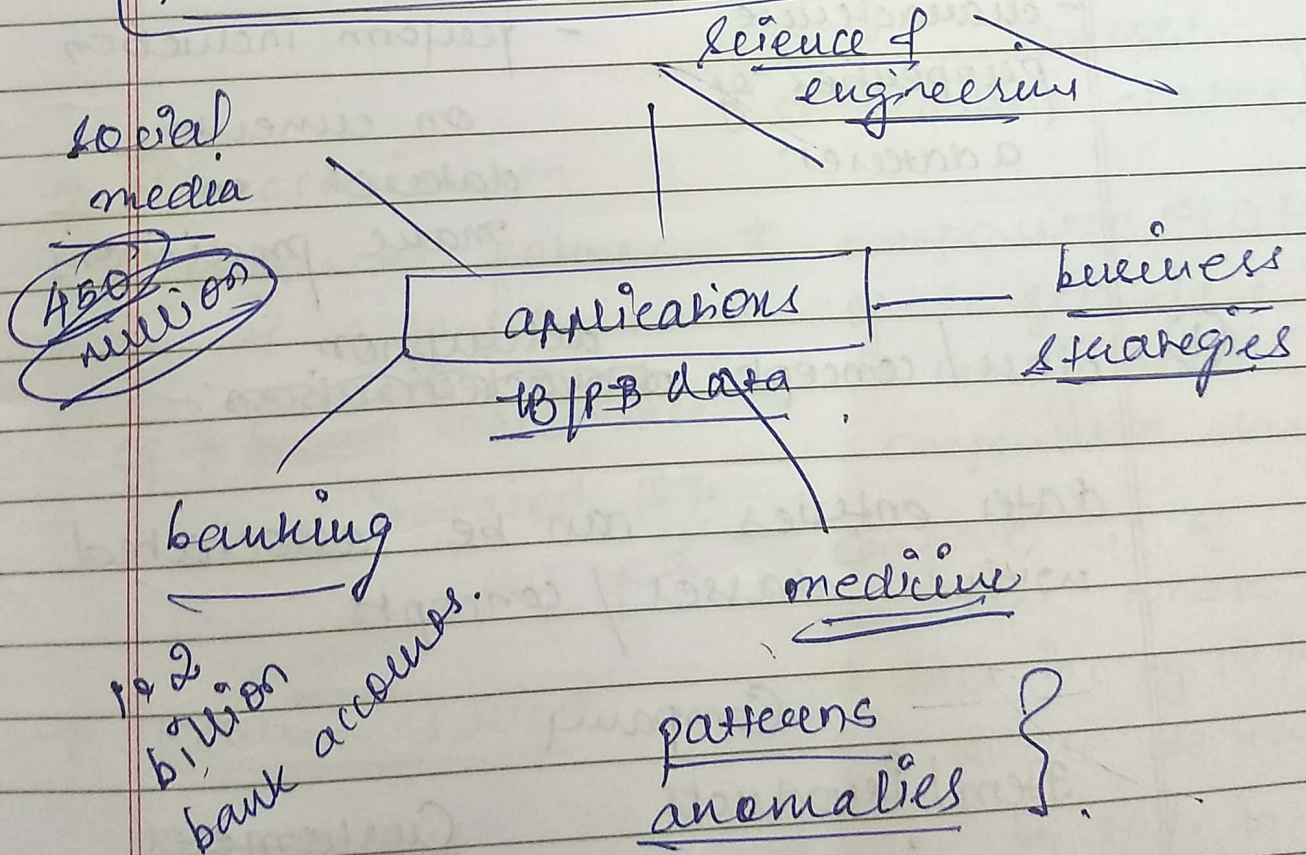
pattern evaluation - to identify interesting patterns rep. cancelled, ~~based on~~

knowledge representation - visualization

Whatsapp

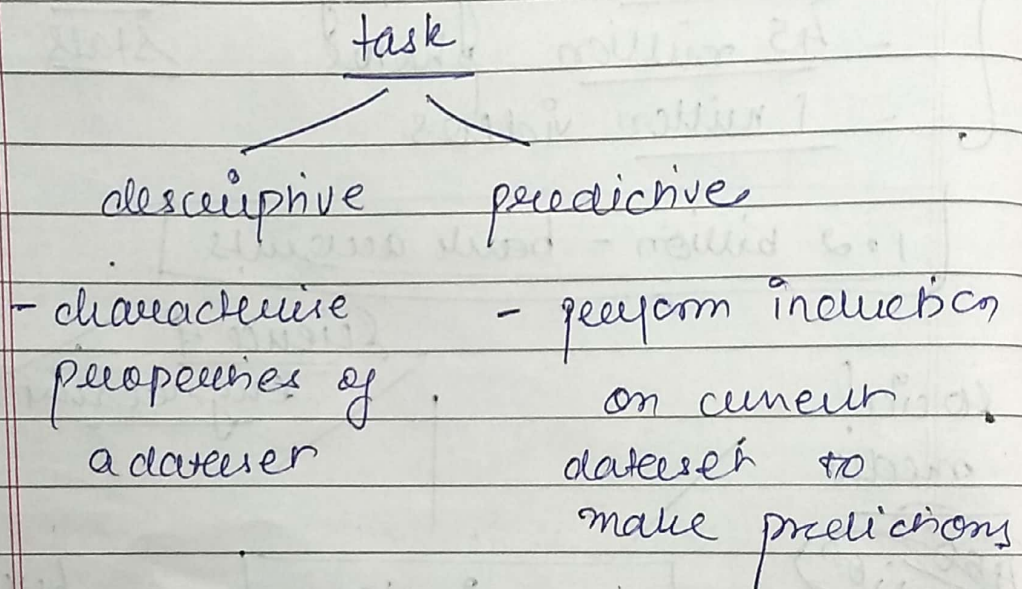
- 450 million users
 - 55 billion msg/day
 - 45 million photos
 - 1 million videos.
- 2018
stats

1.2 billion - bank accounts



Data Mining Functionalities

- ↳ characterization } class / concept description
- ↳ discrimination }
- ↳ mining frequent patterns
- ↳ associations & correlations
- ↳ classification
- ↳ regression
- ↳ cluster analysis
- ↳ outlier analysis



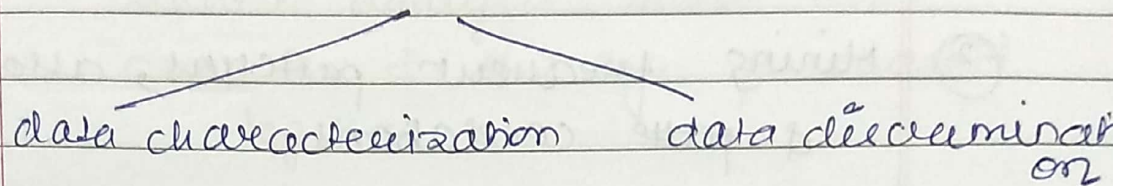
① class / concept ^{description} ~~characterization~~ :-

data entries can be associated with classes / concepts

Company	
Item / Products	<u>Customers</u>
printers	big spenders
computers	Budget spenders.

This helps defining a group in the form of a concept in a summarized, precise form.

↳ this description is called class / concept description.



summarizing the data following a class in general terms (target class)

by comparison of our target class with a set of comparative classes.

(contrasting classes)

→ characterization

of general features of a target class

eg → ~~product~~ product whose sale increased 10% this year

→ comparison of CF of target class to CF of other comparative class.

compare CF of products whose sale inc by 10%

& whose decreased by 30% in last year

→ techniques used

① attribute oriented induction

② roll up operation in OLAP

discrimination rules

Output - pie chart,
bar chart, MD cubes,
MD tables,
generalised
relations.

② Mining frequent patterns, associations
& ~~paths~~ correlations

→ Frequent pattern - appears frequently
in your dataset.

frequent pattern: ~~set of~~ i

frequent itemset - set of items that
appear together in a transaction

(sequential) frequent subsequence - an item purchase
followed by another item.

frequent substructure - in the form
of a graph/tree is substructure

↳ associations & correlations

association analysis -

buys (X, "computer") → buys (X, "sw")

[support = 1% confidence = 30%]

confidence 50% means if customer buys a laptop then are 50% chances that he will buy sw as well.

support 1% means 1% transactions under analysis show that computer & sw are purchased together.

by single dimensional Rules

dependent on 1 attribute

multidimensional Rules

dependent on multiple attributes

age (x, "29") \wedge income (x, 80K) \rightarrow

buys (x, laptop)

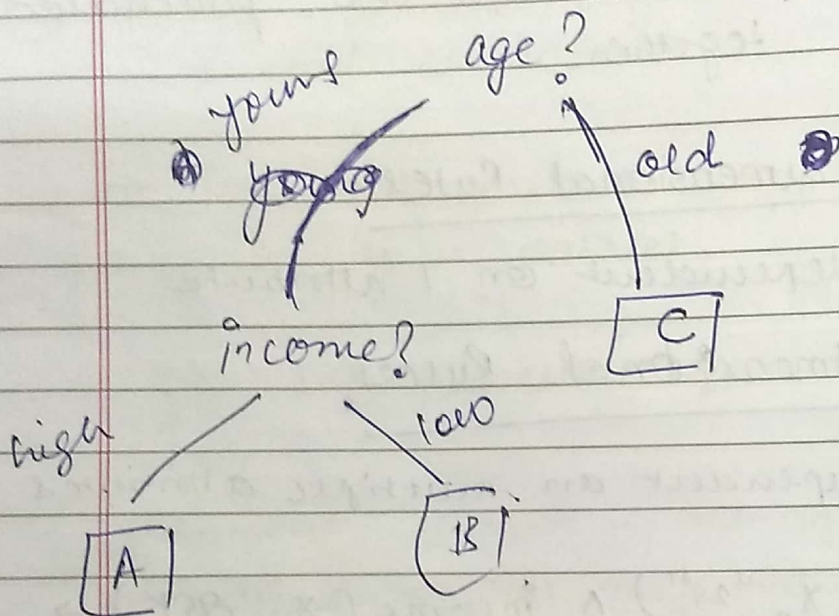
minimum support & confidence threshold must be satisfied.

② Classification & Regression.

classification - finding a model that describes & distinguishes data classes / concepts.

Model derived from analysis of the ^{training} dataset.

This model is used to then predict class labels of objects whose labels are unknown - classification rules, decision trees, neural n/w, mathematical formulae



$\text{age}(x, \text{"young"}) \rightarrow \wedge \text{income}(x, \text{"high"}) \rightarrow \text{class}(x, \text{"A"})$
 \downarrow
 $\text{income}(x, \text{"low"}) \rightarrow \text{class}(x, \text{"B"})$

$\text{age}(x, \text{"middle"}) \rightarrow C$

$\text{age}(x, \text{"senior"}) \rightarrow C$

classification

categorical (unlabelled data, discrete)

Regression - models continuous values to predict missing / unavailable numerical data values

Regression analysis is used for numeric prediction.

It encompasses ^{id of} distribution of trends over data

Both preceded by Relevance Analysis

identifies attributes relevant to classification / regression process.

4) Cluster Analysis

1) - classification & regression analyze labelled data whereas cluster analysis is constructing a class label.

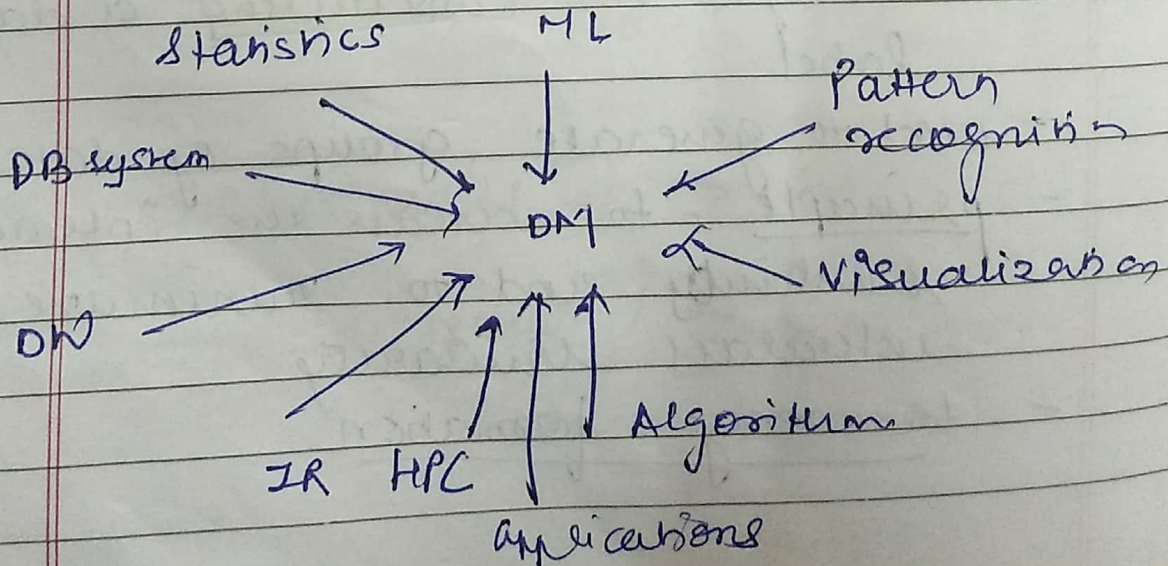
- used to generate groups of data
- principle - to increase the interclass similarity and to minimize the interclass similarity
- taxonomy formation

⑤ Outlier Analysis

- data set may contain data entries that do not comply with the general behaviour / modelling of data. These data objects are called outliers. / noise / exception
- In fraud detection, outlier is an interesting pattern than the regularly occurring pattern.
- anomaly mining
- cluster analysis can be used for the purpose wherein an object distorts from any cluster is an outlier

Classification of DM systems

- Domains that influence the development of DM methods.
application-driven field



①

Statistics

- study collection, analyze, interpret or explain of present data
 - statistical model - mathematical form describing behaviour of objects in a target class
eg: classification & data characterization
- Statistical model can be built
- Also, a DM task can be built on top of a statistical model.

②

ML

- how systems can learn on the basis of data.
 - automatically recognize patterns & make intelligent decisions based on it
- Supervised Learning :- synonym for classification.
- Learning from the labelled example of the training dataset.
- Unsupervised Learning :- clustering
- training data / ip examples are not class labelled.
 - grouping on general feature but semantic meaning not drawn.

Semi-supervised - Both labelled & unlabelled.

labelled - as labels used to classify
unlabelled - refine boundaries.

Active Learning :- users are asked to label the examples so as to keep users active in the process

(B) DB Systems & DW

DB systems - creation, maintenance & use of db for org. & end users

- data storage
- indexing
- accessing
- query proc.

DM can be used to extend the capability of conventional DB.
- analysis

DW & DM - DW integrates data originating from multiple sources & various time frames & consolidates in a multidimensional space.

(A) IR

science of searching document/info. in a document.

doc \rightarrow text / multimedia & may over ~~internet~~ web.

IR assumes (1) unstructured data

(2) queries formed from keyword

Text Document | bag of words

Language Model - function that generates

these bag of words.

Similarity sim^m of document is similarity sim^m corresponding lang. model.

Topic

text doc may have 1 / more topics. IR & DM help find major topic in a collection of document & for each doc, major topic is involved.

Applications

(1) Business Intelligence

(2) Web Search Engines

uses (1) BI :- helps acquire better understand of customers, the market, supply & resources & the competitors.

BI due to DM provides historic, current & predictive views of business operation

customer feedback on similar products, effective market analysis, strengths & weakness of competitors, smart decision

OLAP (Online Analytical Processing) is BI relies on DW & multidimensional data mining

CRM - customer relation mgmt
classifying

↑

characterization - general feature understanding.

Predictive analysis

↑

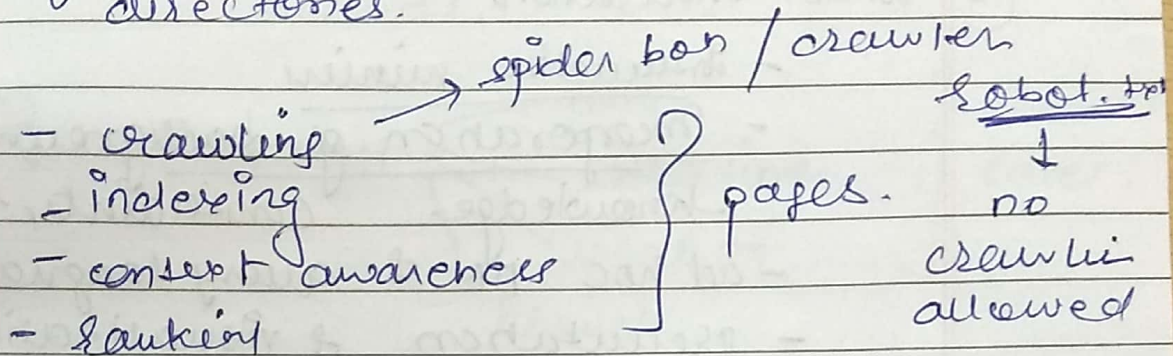
classification & prediction

(2)

Web Search Engines

↳ specialized computer server that searches for info on the web.

- search results returned as hits (list)
- web page, images
- search engines, working algorithmically or times search & return results from data available public db / directories.



- scaling DM algo to cloud.
- context awareness.
- previous user query based search, few times - queries.

* Major Issues in DM

① Mining Methodology

- ① - Mining various new kinds of knowledge - tasks
- ② - Mining in Multidimensional dataspaces.
- ③ - DM - interdisciplinary.
IR + NLP
Big Mining + BI
- ④ - Business power of discovery in networked environment

- Page _____
- Handling uncertainty, noise or inconsistency of data
 - Pattern evaluation / pattern consistency
 - guided mining

② User Interaction

- Interactive mining
- Incorporation of background knowledge, constraints, rules
- ad hoc DM & query languages
- presentation & visualization

③ efficiency & scalability-

① algorithms.

② parallel, dis, incremental

③

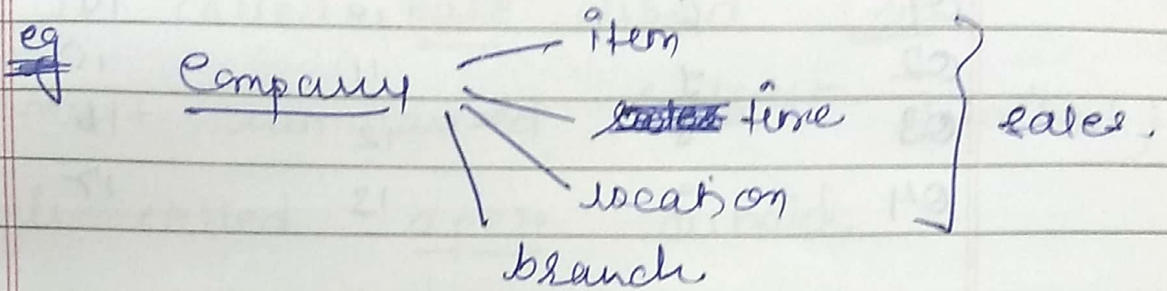
④ diverse db

⑤ DM & Soc.

Multidimensional Data Modelling

Data cube - to be viewed & analysed in various/ dimensions multiple.

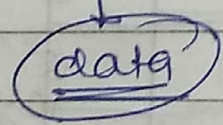
dimension - perspective / entity wrt which an org. wants to keep records.



each dimension has a table associated to it, dimension table.

eg item :- item-id
item-name
item-type
item-brand.

these multidimensional data is based around a single theme.
i.e. sales, theme is represented by fact table.
key of dimension
measures.



Quantity to be analyzed through which, we analyse relationship betⁿ the dimensions.

eg → item key, time key, loc key, branch key, unit sold, dollar, etc.

time & item & view location for Chicago
 NY, Toronto, Vancouver,

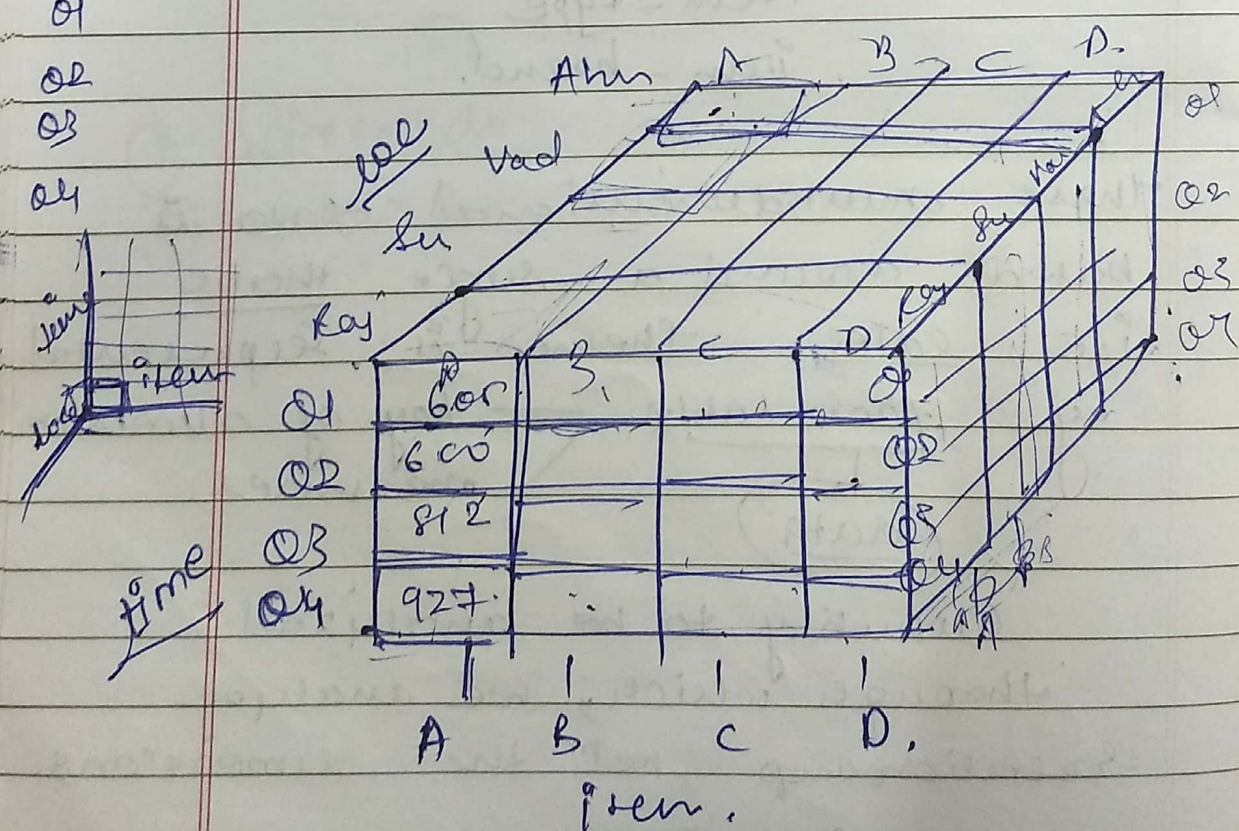
2D location = Ahmedabad.
~~Vancouver~~

Sales data

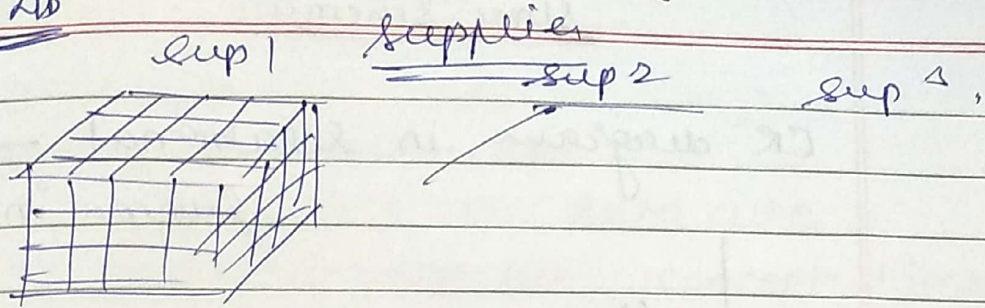
Time	home	Comp	phone	see
01	6	10	14	18
02	7	11	15	19
03	8	12	16	20
04	9	13	17	21

3D view Ahm, Vad, Surat, Raj

Time	Ahmedabad	Vad.	Surat	Raj
Time	A	B	C	D
01				
02				
03				
04				

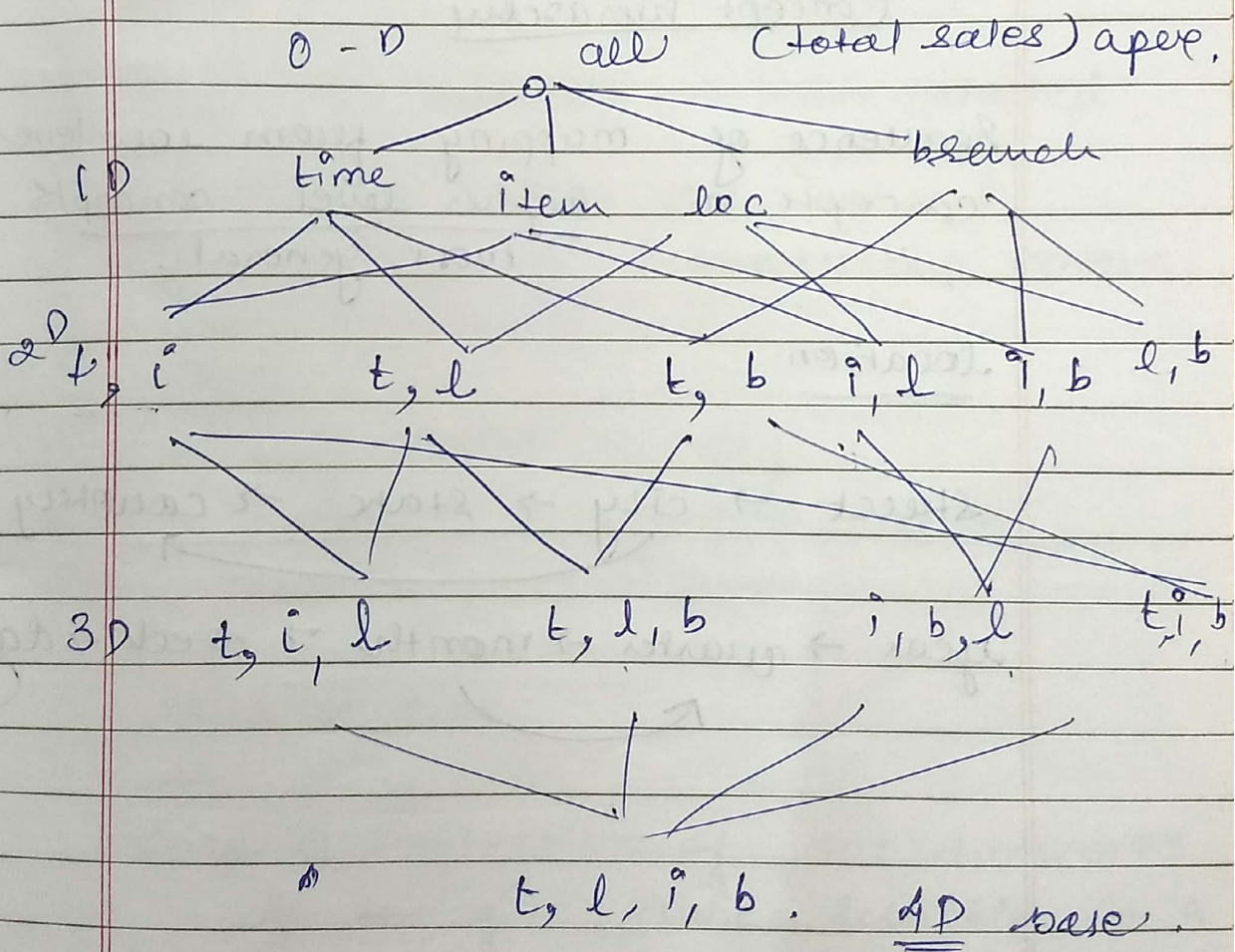


AD



least summarized (most detailed)
is called base cuboid, AP

most summarized (~~most~~ detailed)
is called apex cuboid.



Star Schema

ER diagram in relational → Star Schema in DWH

↳ Star

↳ Snowflake

↳ constellation, galaxy.

Concept hierarchy

Sequence of mapping from low level concepts to higher level concepts more general.

location

street → city → state → country

year → quarter → month → week → day

OLAP operations

roll up → drill up.

aggregation on data cube

climbing up the concept hierarchy

attribute reduction.

~~state~~ ^{city} state.
city → country

one/more dim removed.

drill down → roll ~~down~~ ^{up} down.

less detailed → more detailed.

stepping down the concept hierarchy, or introducing more attributes.

country → city.

month → ~~quarter~~ week.

quarter → month.

slice & dice.

slice → selection of one dimension of the give cube resulting in a subcube.

	A	B	C	D
Ah	□	□	□	□
Va	□	□	□	□
Su	□	□		
Ra	□	□		

time = 01 01

dice - selection of 2/more dimensions

loc = ahm

quantity = Q1

Ahm

A B C D

Q1

- - - -

pivot - rotate

alternative data representation

item
loc

others - drill across more than 1 table
drill through relational SQL

- each dimension has a table associated with it - dimension table

eg item
item_key
item_brand
item_type
item_name

eg location
location_key
city
province/state
country

3 - central theme - fact is represented by fact table

↓

consists of keys of all dimensions & the measures with the relationship amongst these dimensions needs to be analyzed.

fact table

sales
item-key
loc-key
time-key
<u>branch-key</u>
sold-units
earned-profit

cr table

sol

NON-NRM

DW facilitates higher level executives to effectively analyze, organize & understand data to make strategic decisions.

DW is a huge data repository maintained separately from the org. operational database.

It is a subject-oriented, integrated, time-variant, non volatile collection of data in support of mgmt's decision making process.

Subject oriented - DW is organized around major subjects like customers, sales, products, supplier - rather than processing the data on a day to day basis the data is analyzed to be used for decision making.

Integrated - usually concerned merging multiple heterogeneous data sources eg relational db, flat files, online transaction logs etc. Data cleaning & integratⁿ done before data merger into the DW.

es Time-variant - past 5-10 years of data (historic) perspective. Every key structure consists of an implicit/explicit time key element.

Nonvolatile - Always physically separate from ^{app} data in operational env. It doesn't realize transaction processing, recovery, concurrency control, it only realizes loading & access.

Managers / Executives / Analysts
- use warehouse to quickly obtain an overview of data & to make sound decision based on info in warehouse.

- eg:
- ① Inc customer focus (analyzing buyer patterns)
 - ② Repositioning products to inc sales
 - ③ Analyzing op & looking for profit
 - ④ CRM.

Heterogeneous data integration -

wrapper/integrators built on top of multiple/het. databases. when query posted to a client site,

a meta dictionary used to translate into query appropriate for each individual heterogeneous db involved. These queries are mapped & sent to local query processors. Results obtained are integrated into a global answer set.

query-driven approach

update driven integrated in advance & stored in DW.

OLAP vs OLTP

online transaction processing - day to day activities of an org like inventory, purchase, manufacture, payroll, banking, registrarⁿ & accounting.

online analytical processing -

data analysis & decision making, over persistent data

Why DW?

classmate

Date _____

Page _____

① Major reason for separate DB of DW is performance

operation db is made/tuned on basis of known tasks and workload like indexing, hashing with p-key, searching for particular records, & to optimize "canned" queries

↓
common & stored in

a way that they can easily be executed without to see-enter details of query each time.

DW queries are complex. Includes computation of large data groups at summarized level & may req. use of special data organization, access, index methods based on MD view.

② Operational db has multiple transactions. Concurrency control & recovery mechanism required

OLAP query has read only access of data for sum, aggregation

(B) Diff structure, content, uses

Decision support needs historic data,
 ODB do not maintain historic data

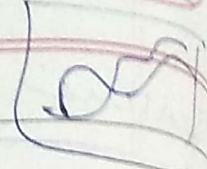
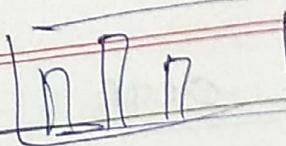
Decision support requires consolidated
 (aggregatⁿ & consolidatⁿ) from
 her sources.

ODB contain only detailed record data
 of transaction.

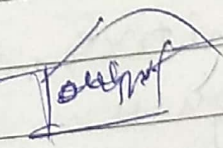
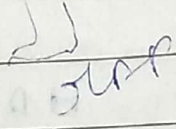
feature	OLTP	OLAP
characteristic	operational proc.	informational proc.
orientation	transactional	analysis
user	clerk, DBA, DB professional	mg, executive analyst ⁿ
function	day to day op.	long term info.
DB design	ER based	star, snowflake, constellation
data	current, up to date	historic, accuracy maintained over time

feature	OLTP	OLAP
summarization	primitive, higher detailed	summarized, consolidated
view	detailed, flat schema	M.D
unit of work	short transaction	complete query
access	read/write	mostly read
focus	data	info over
operation	index / hash	lots of actions
no. of records accessed	less	millions
no. of users	thousands	hundreds.
DB size	GB to \uparrow GB	\geq TB
priority	\uparrow performance, \uparrow availability	\uparrow flexibility, end user autonomy
metric	latency t/p	over t/p response time.

top tier

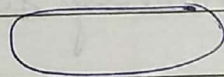


middle tier

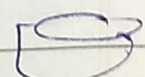
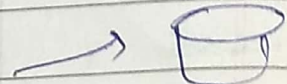
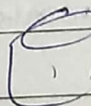


bottom tier

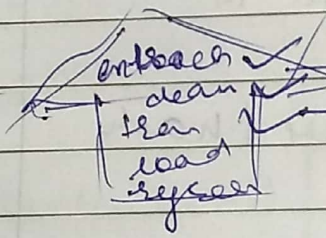
metadata ref



DW



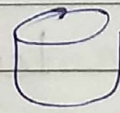
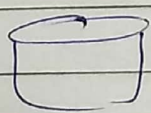
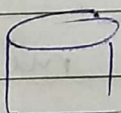
mainten



unified format

Back end

Backend tools & utilities



data

Database

External Info sources

customer profile provided by external consultants

gateway

ODBC
JDBC

SQL code

Metadata repository - contains information about the data warehouse.

① DW structure - schema, view, dimension, hierarchy & derived data defⁿ, data mart locatⁿ & content

② operational metadata - data lineage (history of migrated data & sequence of transf. applied to it)

- (DW statistics, error report)

- data - archived/active

dw server ③ Algorithms used for data summarizatiⁿ.

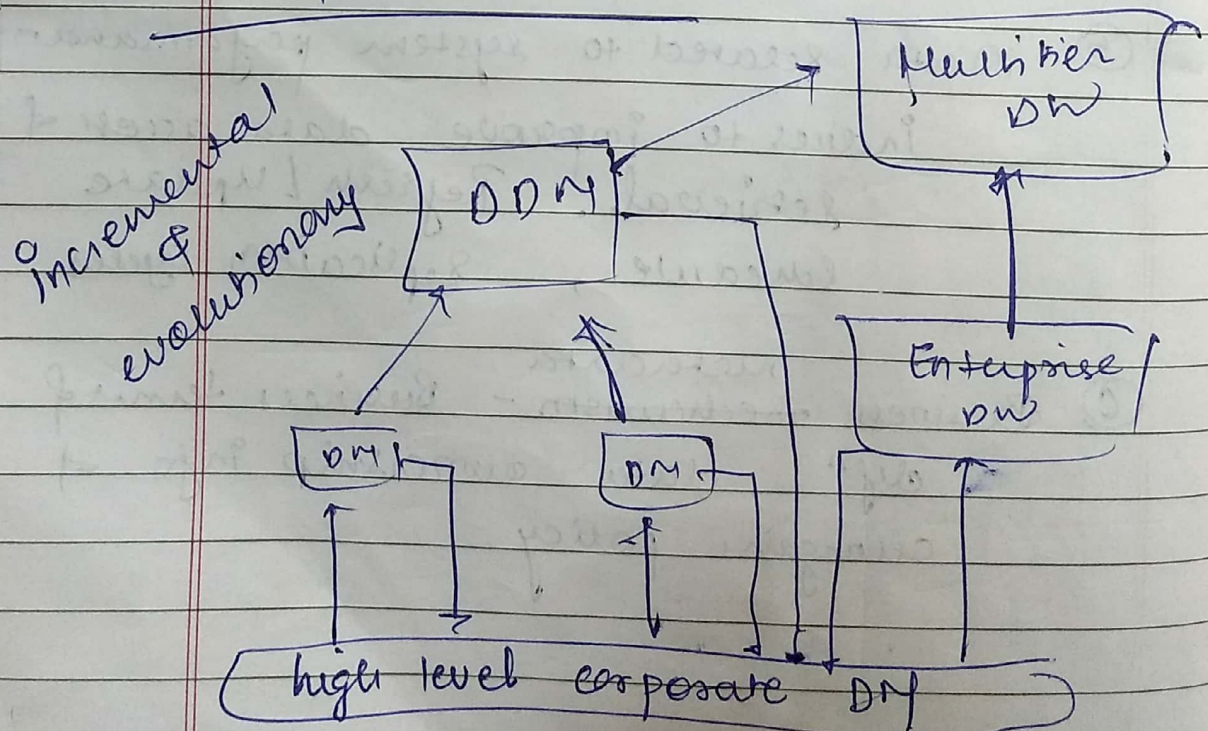
④ Mapping from operational env to DW - source db, gateway, e, c, + rules

⑤ Data related to system performance indices to improve data access & retrieval. Refresh / update schedule, replication cycles

⑥ Business metadata mechanism - Business terms & defⁿ, data ownership info & charging policy.

- ① Extraction - gathers information from multiple heterogeneous sources.
- ② Cleaning - detects errors & rectifies them.
- ③ Transformation - converts data from host format to warehouse format
- ④ Load - sort/summarize/consolidate, computes views, checks integrity & partitions.
- ⑤ Refresh - propagates updates from data sources to dw.

DW, DM & VDW



DW development approach

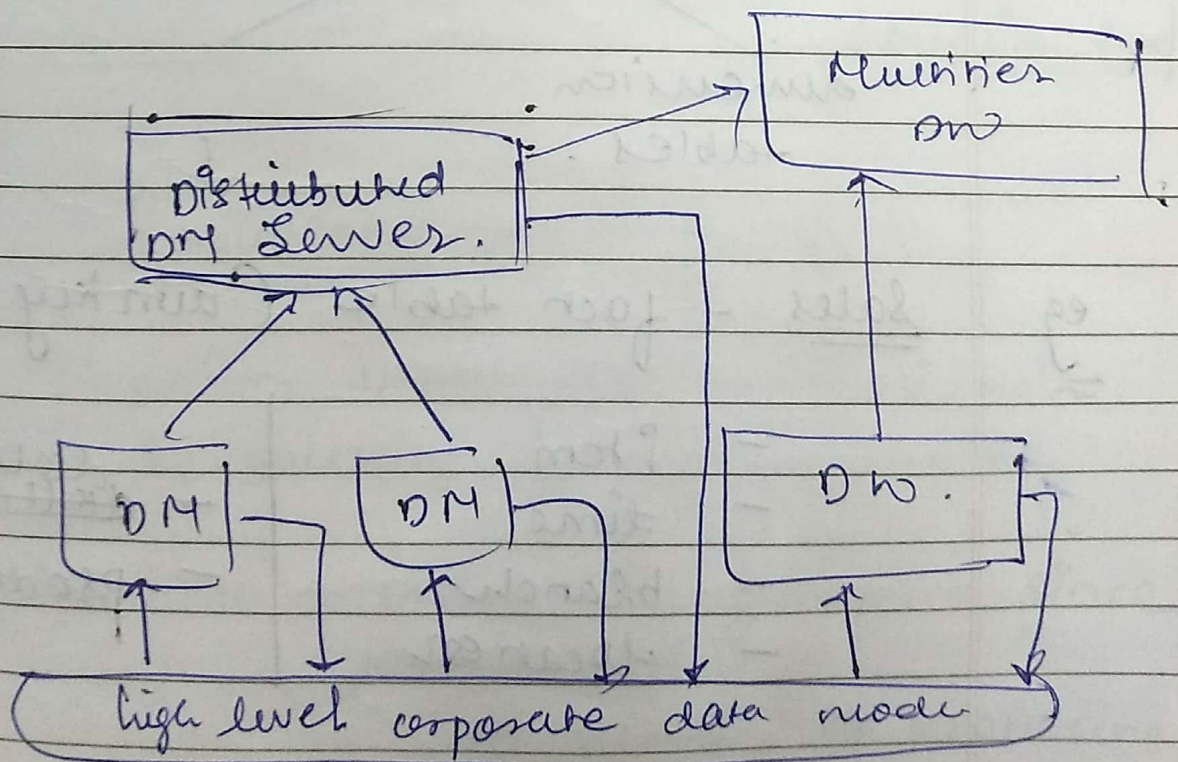
- Incremental & Evolutionary.

First - high level corporate data model is defined in a month or two, - provides integrated view of data among diff subjects & usages. (Starware model defined later)

Second - Independent Data Marts are prepared in parallel to datawarehouse on the basis of model.

Third - Distributed DM is constructed to integrate diff DM via hub servers.

Finally - Multitier datawarehouse is constructed ~~where~~

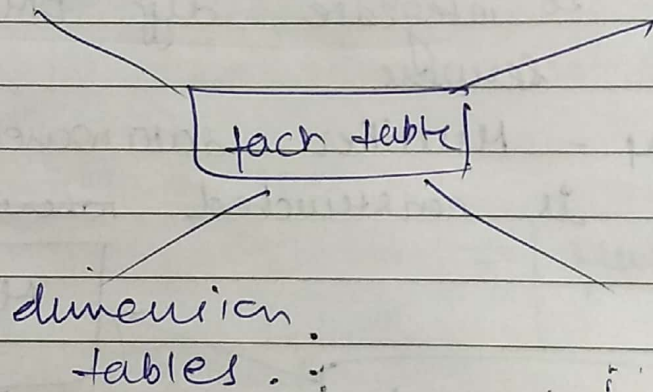


Star schema

① large central table - fact table
↓
bulk of data
with no redundancy
central
theme table

② set of smaller attendant tables
(dimension table) one for each
dimension.

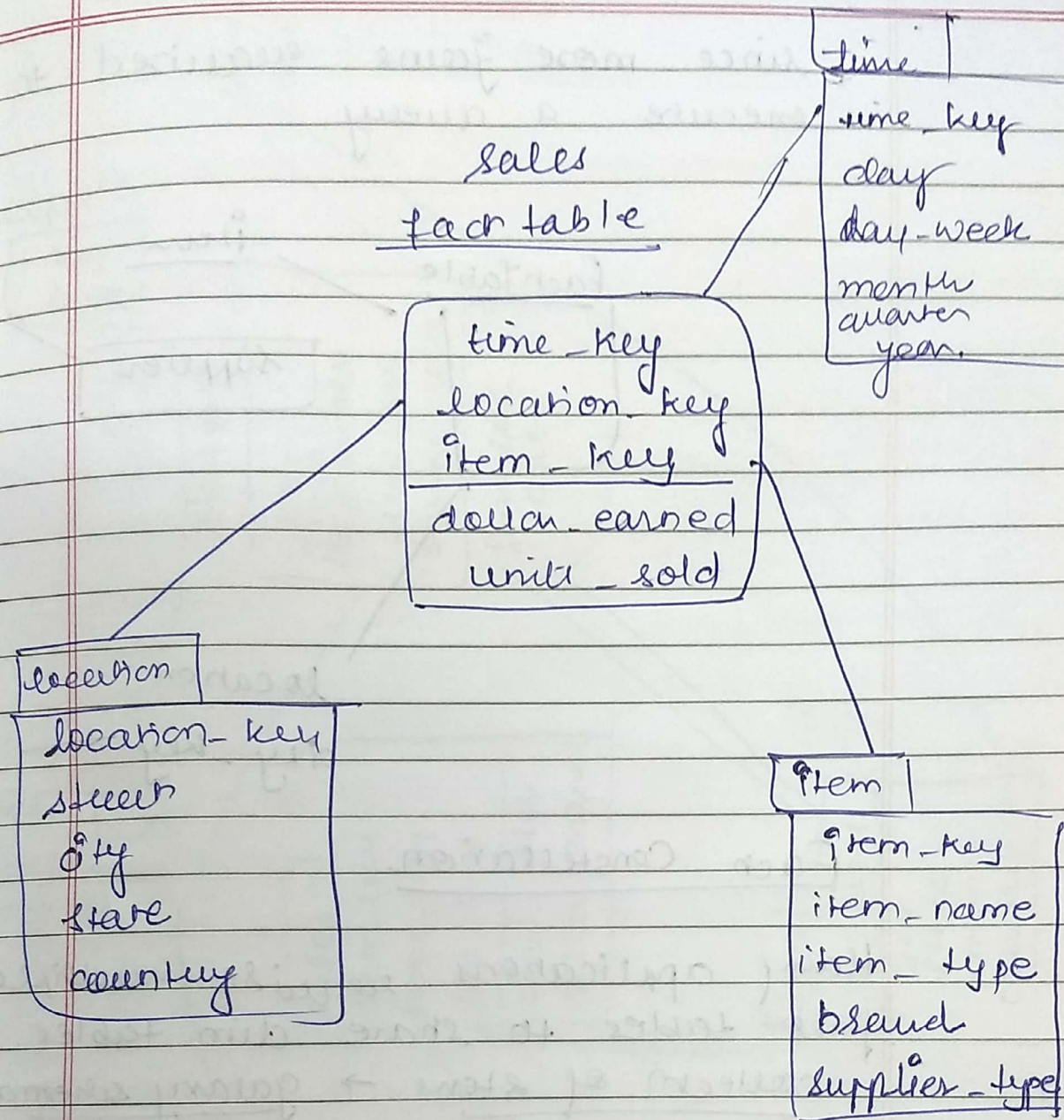
③ Graph resembles a starburst
pattern radiating
from
central
object



eg Sales - fact table (dim key + m measures)

- item	- dollar earned
- time	- dollar sold
- branch	- product sold
- location	

level
loc
ste
o
s
o



Incoflate Schema

earnes)

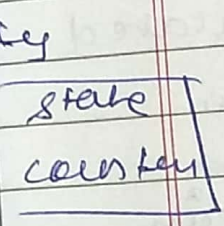
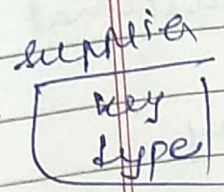
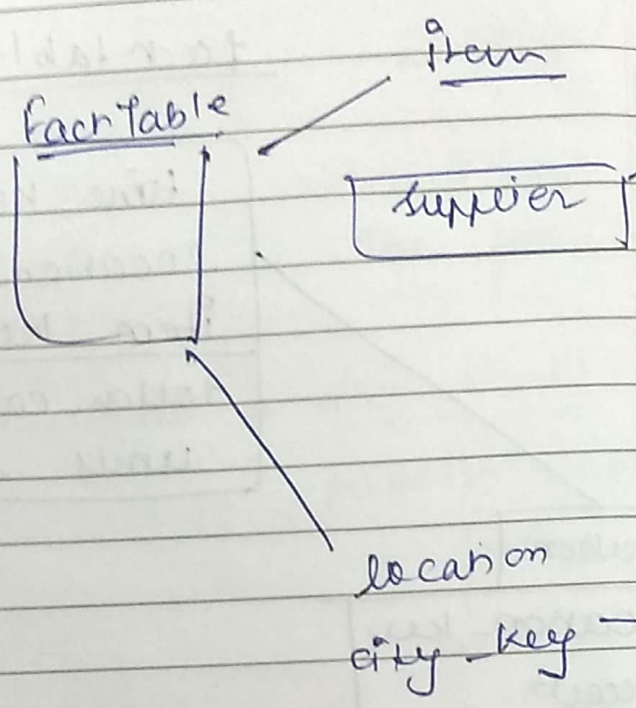
Dimension tables are normalized.
by splitting into more tables.

at
old

difference - saves storage space since
normalized.

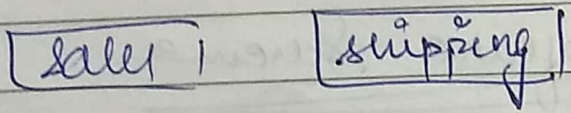
reduces effectiveness of browsing

since more joins required to execute a query.



Fact Constellation

Many applications require multiple fact tables to share dim tables.
collection of stars → galaxy schema
fact constellation



Item
 Item name
 brand
 supplier

state
 country

supplier
 key
 type

time
 time key
 date
 day of week
 month
 quarter
 Year

sales
 item-key
 time-key
 loc-key
 unit-sold
 dollar-amount

supplier
 item-key
 loc-key
 supplier-key
 from loc
 to loc

loc
 loc-key
 city
 state
 province
 country

state-key
 loc-key

DW Implementation

classmate

Date _____

Page _____

Rolap

Molap.

Relational

Multidimensional

Storage &
fetch

main
dw

proprietary
MDDBS

Data
form

relational
tables

MD arrays
of data cubes

volume

large.

limited
summaries

tech.

complex
queries

precalculated

data
cubes

Spase matrix
tech.

view

MD dynamically
shared

static

MD
stored.

access

slow

faster

efficient

① which OLAP op?

② materializ em

nto

aper

mesh

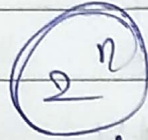
full
partial.

base (mesh

specific,

least

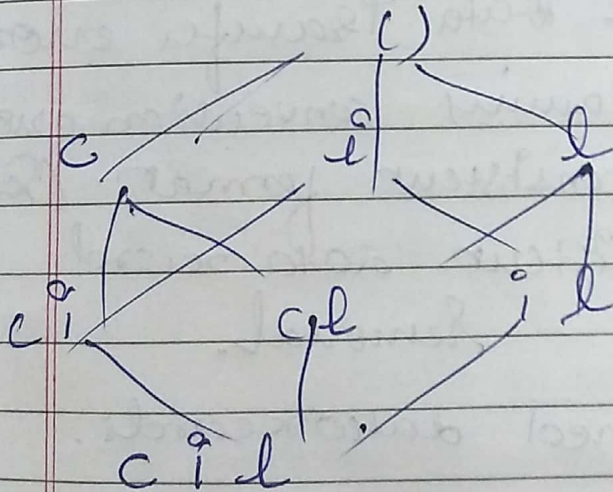
summarized



cube of dimension

group by

compute cube, group by city,
item.



Preprocessing

① Data Quality

- Accuracy
- Completeness
- Consistency
- Timeliness
- Believability
- Interpretability.

Inaccuracy - } machine / system fault
Incompleteness - } human entry fault
Inconsistency - } disguised missing values

- ✓ Data Transfer errors
- ✓ Naming convention error
- ✓ Inconsistent format (Date)
- ✓ Inconsistent data record removal.
- ✓ Inferred data records.

Address not present in db.

↳ db accurate for SEO (IT) department.

↳ db not accurate for Sales / marketing dept.

Timeliness - Hardly report

if data not submitted in time it fails timeliness.

Believability - if errors in data have earlier occurred in failures, the next time using the same dataset is not trusted by users.

Interpretability - how easy it is to understand data.

Accounting codes in a db are not understood by the Sales Department.

Major Task

- Data Cleaning - fixing missing values

- smoothing noisy data
- removing outliers
- resolving inconsistencies

- Data Integration -

- id, customer id
- naming convention to be resolved
- remove redundancies.

- Data Reduction -

Reduction in volume of data producing same analytical results.

- Dimensionality Reduction -
 - "compressed" representation
 - Principal component analysis
 - attribute subset selection
- Numerosity Reduction -
 - Data replaced by alternative, smaller rep. parametric model, eg. regression / log-linear models
 - nonparametric
 - histograms,
 - clustering,
 - sampling,
 - data aggregation.

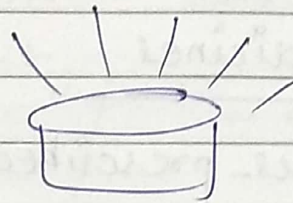
Normalization,
Discretization,
Concept hierarchy generation -

~~data~~ transforming data to a format that is such that all attributes fall in a range such that they don't outweigh each other

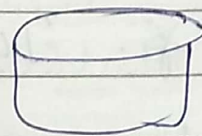
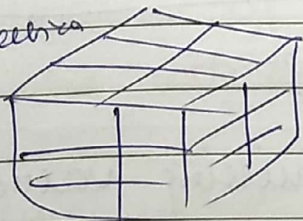
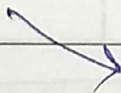
age \rightarrow raw values

\rightarrow adult, young, senior.

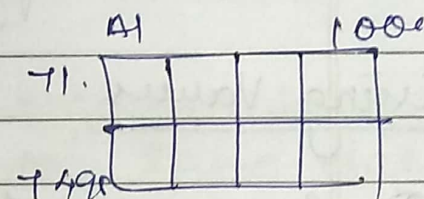
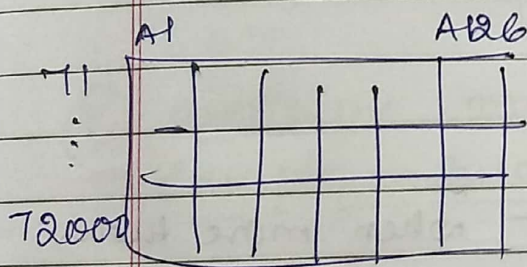
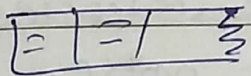
cleaning



Integration



Integration.



Reduction

Transformation

- 2, 32, 100, 59, 48 \rightarrow - 0.2, 0.3
 1.00, 0.5,
 - 0.3

hospital

- location
- doctor
- patient
- medicines

medicines prescribed
charges.

chemist shop

- location.
- chemist
- patient
- medicine

medicines. ~~prescribed~~ ^{supplied}
profit-earned.

*. Data Cleaning

- ↳ Dealing with missing values
- ↳ correcting inconsistencies
- ↳ smoothing out noise
- ↳ identifying outliers.

Missing Values

① Ignore the tuple :- when more than one attr. missing. otherwise, we lose out on other attr that were contributors to patterns.

②. fill values manually :- time consuming & not feasible for large datasets.

(2) Using global constant to fill values. -

Replace all missing values by a global constant.

eg. "Unknown" / "-x"

might also generate patterns
↳ not foolproof

(4) Measure of central tendency for the attribute → middle value.

symmetric

data distribution

asymmetric

data distribution

mean

eg. income

median

(5) attribute mean / median for a sample belonging to same class.

eg: customers in same credit risk → same mean

(8) Most probable value

- regression
- inference based
- decision tree based

Null values

Noisy Data

random error / variance in the measured variable

numeric attribute - noisy data
smoothness

① Binning - sorted data values

bucket
bms

- neighborhood
- boundaries
- central tendency

local smoothing

equal frequency - bin size

① smoothing by mean.

② smoothing by bin medians

③ smoothing by bin boundaries

4 8 15 21 21 24 25 28 34

equal freq

4	8	15
21	21	24
25	28	34
mean - 9 9 9		
22	22	22
29	29	29

11	18
34	38
3	3

12

(equal width)

0-12

13-26

27-34

(4) (8)

15

24

21

24

20

28

34

median - 8 8 8
21 21 21
25 28 34

boundaries - 4 8
15 24 24 29
28 34

(Regression)

Linear Regression involves finding the "best line" that fits two attributes so that one attr can be used to predict the other.

logistic / multiple Reg → to find more attributes.

$$y = B_0 + B_1 x$$

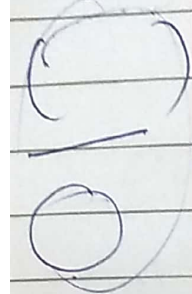
intercept

slope

$$B_0 = \text{line } \text{mean}(y) - B_1 \cdot \text{mean}(x)$$

$b_0 \approx \text{mean}(x) \cdot b_1$

$$b_1 = \frac{\text{mean}(x) \cdot x - \text{mean}(x)^2}{\text{mean}(x) \cdot y - \text{mean}(y) \cdot \text{mean}(x)}$$



Outliers

Clustering

Cleaning Process

① Identify Discrepancies
detection

- poor design forms
- data decay
- deliberate error
- human entry
- inconsistency
- instrument

How to identify?

use data about data - metadata

- ✓ what data type?
- ✓ what an acceptable value?
- ✓ mean?
- ✓ median?
- ✓ variance?

range

Field Overloadin - 31 out of 32 bits used for storage.

1 bit reserved.

3 rules for analysis

- ① unique ✓
- ② null ✓
- ③ consecutive ✓

Data Subbing / Data Auditing tools

↓
 domain knowledge
 (parsers)
 (fuzzy)

↓
 rules of
relationships
correlation
outliers

Data Transfer

Data Migration Tools

Inconsistency
Removal

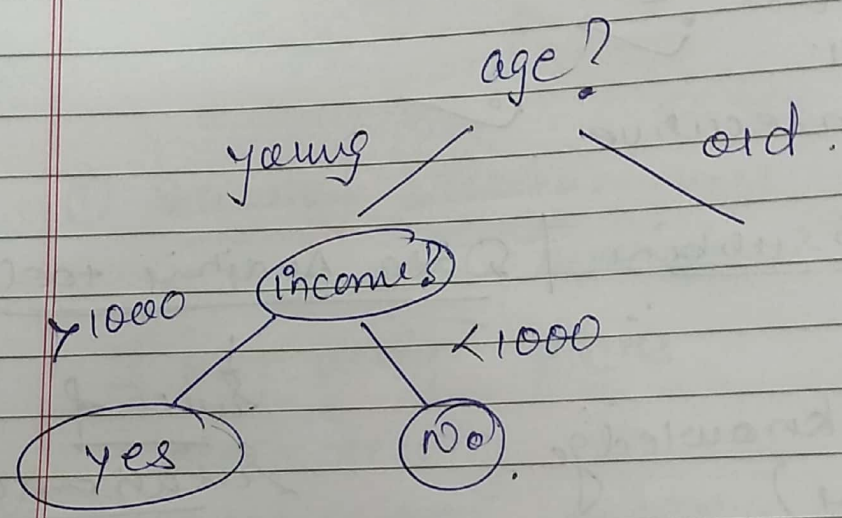
ETL

etl used

Potter's Wheel

①

			label
A	1500	young	✓
B	1600	old	X
C	2000	young	✓



$$\text{income}(x, "10000") \wedge \text{age}(x, "young")$$

②

$$\rightarrow \text{credit risk}(x) \rightarrow \text{loan}(x, "granted")$$

x - mean x
 y - mean y

$$y = (B_0) + (B_1) \cdot x$$

$$B_1 \Rightarrow \sum (\text{mean } x \cdot x - \text{mean } y \cdot y)$$

$$\sum (x - \text{mean } x)^2$$

$$B_0 = \bar{y} - B_1 \cdot \bar{x}$$

$$B_1 = \frac{\sum (x - \bar{x}) \cdot (y - \bar{y})}{\sum (x - \bar{x})^2}$$

$$B_0 = \bar{y} - B_1 \cdot \bar{x}$$

$$y = B_0 + B_1 \cdot x$$

Correlation

$$X^2 = \frac{c}{\sum_{i=1}^c \sum_{j=1}^r} \frac{e_{ij}^2}{e_{ij}}$$

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

-1 +0 +1

classmate

Date _____
Page _____

	male	female	tot
fiction	250 (90)	200 (360)	450
non fiction	50 (210)	1000 (840)	1050
tot	300	1200	1500

$$e_{mf} = \frac{2}{300} \times 450$$

$$= 90$$

$$e_{mnf} = \frac{210}{300} \times 1050$$

$$= 210$$

$$e_{ff} = \frac{30}{1500} \times 1200$$

$$= 360$$

$$e_{fnf} = \frac{90}{1500} \times 1050$$

$$= 840$$

co-relations =

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

①

②

③

④

⑤

⑥

⑦

⑧

$$\sum_{i=1}^n \frac{(a_i - \bar{A})(b - \bar{B})}{n \bar{A} \bar{B}}$$

Test

① Which one of them provides better/efficient query processing?

(a) Star (b) Schema

② Name 5 OLAP operations

③ ~~Write 2 diff~~ Which ~~form of~~ data does DB & DW star?

④ Name 5 DM functionalities.

⑤ Diff betⁿ supervised & unsupervised learn?

⑥ What is freq. itemset mining?

⑦ Diff betⁿ single dim. associatⁿ rule & MD asso. rule with eg.

⑧ Which is most summarized?

840) ² + (200-360) ² apex cube / base cube

840.

⑨ 4 properties of DW

(10) Design DW, DM & VDW.

Integratⁿ.

* tuple Duplication DV conflicts

(1) weiger

(2) hotels

(3) gradins.

* Data value conflict tuple duplication.

denormalized.

* Redundancy & correlation

derived attribute

covariates

Data
integration

Nominal Data

classmate

Date _____
Page _____

① χ^2 correlation test for Nominal Data

	male	female	total
fiction	250 (90)	200 (360)	450
non fiction	50 (210)	1000 (840)	1050
total	300	1200	1500.

χ^2 - square

Are gender and preferred-reading correlated?

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r (o_{ij} - e_{ij})^2 \rightarrow \text{expected frequency.}$$

\downarrow
 e_{ij}
 observed frequency
 of count

$$e_{ij} = \frac{\text{count}(A = a_i) \times \text{count}(B = b_j)}{\text{total}(n)}$$

A has c values $a_1 \dots a_c$ (cols)
 B has r values $b_1 \dots b_r$ (rows)

Contingency table

count $(A = a_i)$ no of tuples having value a_i for A

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

1500 people

Date _____
Page _____

- gender

- preferred - reading

fiction/non fiction

observed joint summaries in table

expected :-

$$e_{11} = \frac{\text{count}(M) \times \text{count}(F)}{n}$$

$$= \frac{300 \times 450}{1500}$$

$$= 90$$

(data distribution)

$$\chi^2 = \frac{(250 - 90)^2}{90} + \frac{(60 - 210)^2}{210}$$

$$+ \frac{(200 - 360)^2}{360}$$

$$+ \frac{(1000 - 840)^2}{840}$$

$$= \underline{\underline{507.93}}$$

$$\text{Degree of freedom} = (C - 1) \times (R - 1)$$
$$= 1 \times 1$$

$$= 1 \rightarrow 0.001$$

significance level α

10.828

upper percentage point of χ^2 .

Correlation coefficient for numeric

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

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

$$-1 < r_{A,B} < 1$$

$r_{AB} > 0 \Rightarrow$ +vely correlated.
A \uparrow with B

\nexists can prove that this is redundancy \nexists

$r_{AB} = 0 \Rightarrow$ A & B independent

$r_{AB} < 0 \Rightarrow$ A & B are negatively correlated. They discourage each other.

Covariance

mean - expected.

$$E(A) = \frac{\sum_{i=1}^n a_i}{n}$$

$$E(B) = \frac{\sum_{i=1}^n b_i}{n}$$

$$\text{covariance} = \frac{E(A - \bar{A})(B - \bar{B})}{n}$$

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

$$\rho_{ab} = \frac{\text{COV } ab}{\sigma_A \sigma_B}$$

$$\boxed{\text{COV } ab = E(A \cdot B) - \bar{A}\bar{B}}$$

if one of A is larger than \bar{A} then one B \uparrow \bar{B} . positive

one larger & one below \rightarrow negative

Independent, $E(A \cdot B) = E(A)E(B)$

time	A allelec	B nighted
t1	6	20
t2	5	10
t3	4	14
t4	3	5
t5	2	5

$$E(A) = \frac{6 + 5 + 4 + 3 + 2}{5} = \frac{20}{5} = 4$$

$$E(B) = \frac{20 + 10 + 14 + 5 + 5}{5} = \frac{54}{5} = 10.8$$

$$COV = \frac{120 + 80 + 56 + 15 + 10. - 4 \times 10.8}{5}$$

$$= 7$$

positive covariance

①. Tuple Duplication

- Duplication arises due to denormalized tables (done to improve performance)
- Inaccurate data / updates / Deletion of records

②. Data Value Consistency & Reso.

- Metric | British Imperial
- A to F+
- Hotel prices

Entity Identification Problem

Schema Integration
Subject Matching

format

ETP

attribute

customer id & car id.

- Structure ✓
- Metadata ✓

$$E(A \cdot B) - \mu_{AB}$$

Date _____

Page _____

$$\sqrt{\frac{\sum(x-\bar{x})^2}{n}}$$

Covariance

t1	6	20
t2	5	10
t3	4	14
t4	3	5
t5	2	5

$$E(A) = \frac{6+5+4+3+2}{5} = 4$$

$$E(B) = \frac{20+10+14+5+5}{5} = 10.8$$

$$\text{Covariance} = \frac{6 \times 20 + 5 \times 10 + 4 \times 14 + 3 \times 5 + 2 \times 5}{5}$$

$$- (4)(10.8)$$

$$= 50.2 - 43.2$$

$$= 7$$

+ve covariance \rightarrow related.

-ve covariance \rightarrow ^{negatively} ~~not~~ related.

0 covariance \rightarrow not related.

Data Reduction

- to obtain reduced representation of data set that is much smaller in volume but more efficient to have mining performed over it.

(maintain integrity of original dataset)

① Sampling

① Simple Random Sample without Replacement (SRSWOR)

of size s

$s < \underline{N}_{total}$ (D)

each tuple = $\frac{1}{N}$

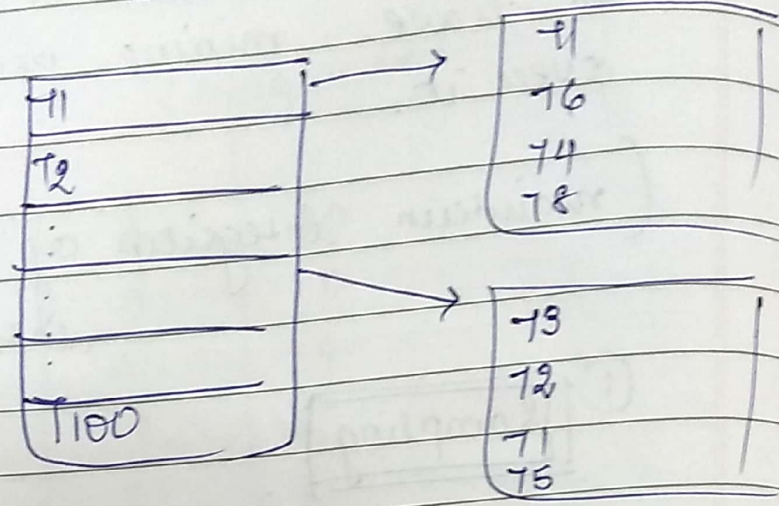
② Simple Random Sample with Replacement of size s (SRRWR)

each time a tuple is drawn from D, recorded & replaced, to be drawn again.

③ Cluster Sample \rightarrow If tuples in D are grouped into M

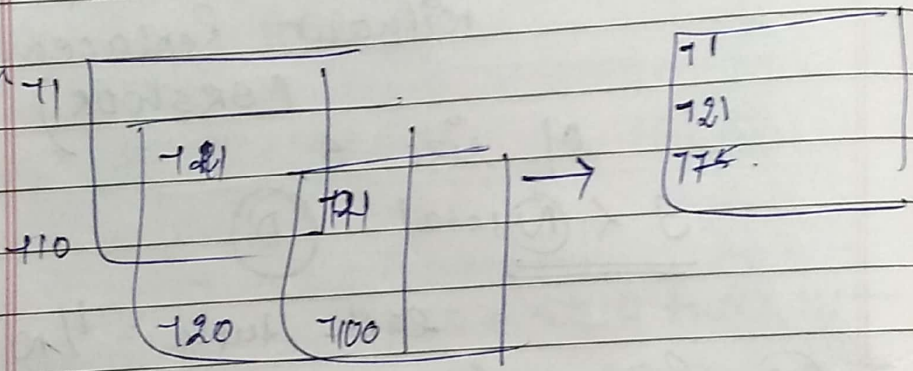
disjoint clusters then an SRS of s clusters can be obtained when $S \times M$

SRSWOR



SRSWOR

cluster

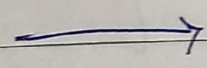


stratified sampling

according to age

cluster attribute

- 71 young
- 72 young
- 760 young
- 712 middle-aged
- 716 middle-aged
- 788 middle-aged
- 711 senior
- 778 senior



- 738 youth
- 791 youth
- 7101 middle aged
- 721 senior

(52) / (17)

Advantages:-

efficiency

Disadvantages:-

Sampling errors

bias / discrimination

↳ Iterative Procedure

sample size ↑ or ↓

size of sample depends upon the Population.

②

* Histograms

- uses Binning to approximate a data distribution.

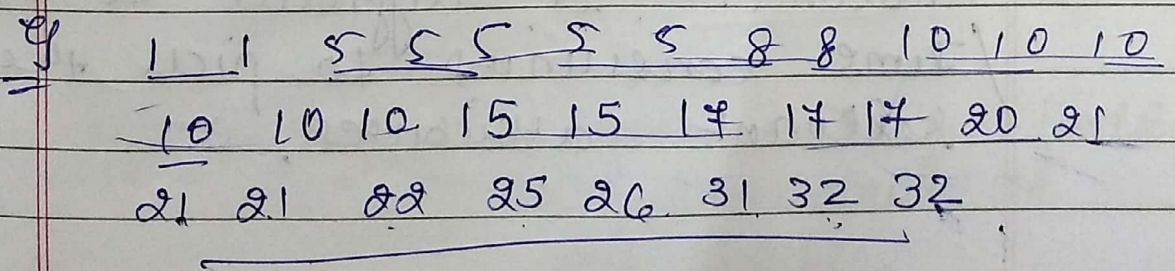
Histogram for attribute A, partitioning the ^{data} distribution of A into disjoint subsets → buckets / bins.

one element in bucket → singleton bucket.

single attribute

52 / 17

ed



size = ∞ frequency 5 element
width level
range

(B)

~~Attribute Selection~~

Attribute Subser Selection

- there might be some attributes in the datasets that might be irrelevant to the mining task

eg - to classify customer on the basis of their liking or not liking ^{to purchase} a popular new CD or A company when notified about sale

age / music-taste are relevant
address / phone no are all irrelevant

When data behaviour is not well known, it is difficult to pick / time consuming to pick the relevant attributes.

Attribute Subset Selection reduces the dataset size by removing the irrelevant / redundant attributes.

goal: to find minimum set of attributes such that the distribution of data classes is as close as possible to the original distribution obtained with all attributes

for n attributes, 2^n subsets.

If n increases & no. of classes increase, search for optimal subset becomes exhaustive.

∴ Greedy Methods are used

while looking through attributes they choose the one that appears to be the best choice at the moment

↳ make locally optimal decision to reach globally optimal solⁿ.

Test - statistical test used for finding if attr^s dependent on one another.

① Stepwise forward selection:-

Initial set

{ A_1 , ^{empty} A_2 , ..., A_6 }

process starts with an empty set

Initial reduced set

{ A_1 }

{ A_1 , A_3 }

{ A_1 , A_3 , A_6 }

Best original attribute is defined and added

Best from remaining is kept on being added

② Stepwise backward elimination

A_1 , ..., A_6

A_1 , A_4 , A_5 , A_6

A_1 , A_4 , A_6

starts with

full set & removes

attributes one

by one

③ Combination at each step selects best & worst both

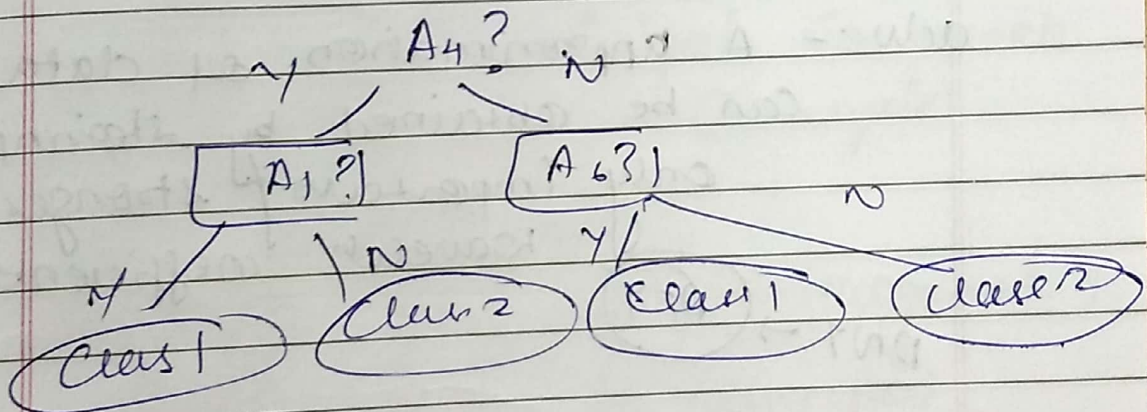
④ Decision Tree Induction

ID3, CART, C4.5

a flowchart where internal nodes
 all test on attribute & each
 branch is outcome of test
 internal leaves are classes

at each stage, best attribute is
chosen for partitioning.

Stopping criteria - threshold



[A1 A4 A6]

- ① Information gain
 - ② Gini Index
 - ③ gain ratio
- } max ✓

② Wavelet Transform

DWT Discrete Wavelet Transform

linear signal processing technique that when applied to X (data vector) transforms it to a numerically different X' of wavelet coefficients

& vectors are same length
Here, n -dimensional data

$$X = (x_1, \dots, x_n)$$

adv:- A approximation of data can be obtained by storing only important strongest wavelet coefficients.

DWT \rightarrow (DFT)

\rightarrow Length L of i/p data vector must be 2^n .
pad necessary as (L, m)

\rightarrow each transform — data smoothing / with avg
— weighed diff to
blur out
detected features

→ these 2 functions are applied to pair of data points in X than x_i to all pair x_{2t}, x_{2t+1}

Results in 2 datasets of $1/2, 1/2$
) (lower freq. version of i/p)

→ two func^{ns} recursively applied to data set in previous loop.

Selected values in previous itera^{ns} are designated as low-level transforms.

②. PCA principal Component Analysis

① normalization

② eigenvector (orthogonal)

③ the principal

eigen value associated with these vectors

magnitude

higher significance kept

ng