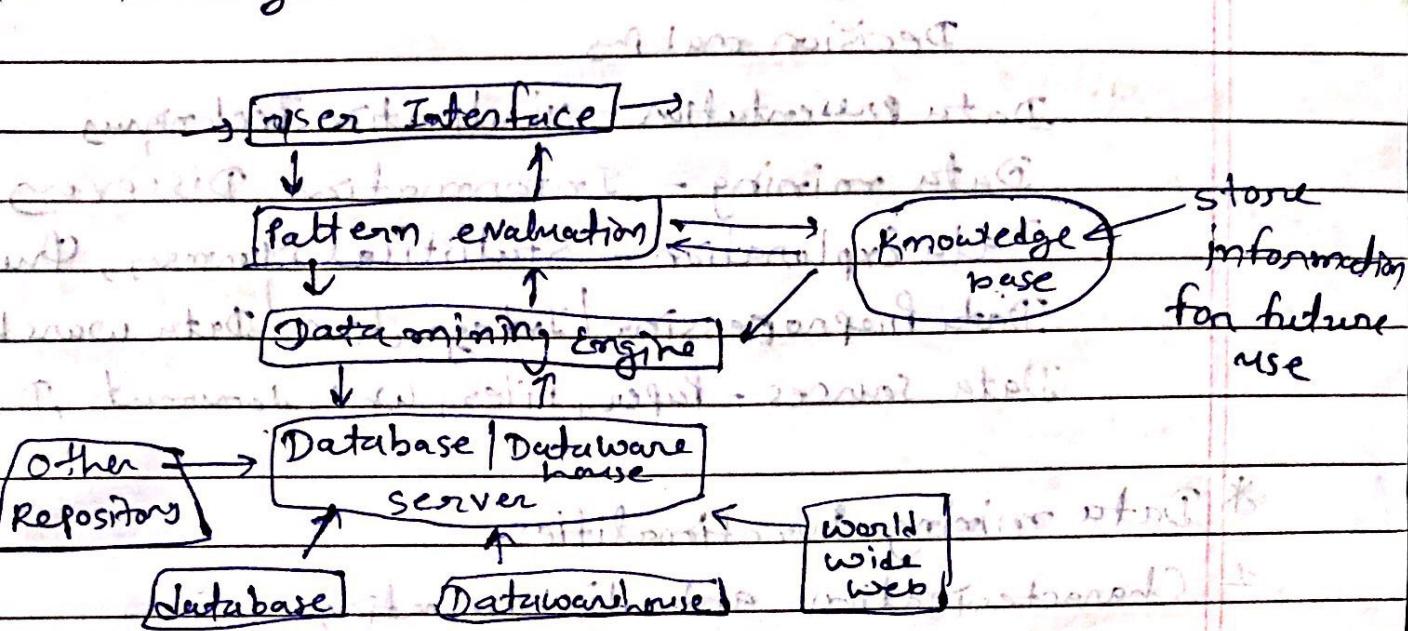


Data mining refers to extracting knowledge from large amount of data



- Data mining - Automated analysis of massive datasets
- Extraction of interesting - non-trivial, implicit, previously unknown, potentially useful patterns or knowledge from huge amount of data.
- Knowledge discovery in databases, data archaeology, data dredging, information harvesting, business intelligence

* Web Mining Framework

- ① Data Cleaning
- ② Data Integration from multiple source
- ③ Warehousing the data
- ④ Data cube construction
- ⑤ Data selection for datamining
- ⑥ Data mining
- ⑦ Presentation of the mining result
- ⑧ Pattern & Knowledge to be used or stored into knowledge base

प्रश्न अभी वृत्तालये स भगवान् जयतीह साक्षात् ॥

Business Intelligence

Decision making

Data presentation - visualization Techniques

Data mining - Information Discovery

Data Exploration - Statistical, Summary, Query, Reporting

Data Preprocessing / Integration / Data warehouse

Data sources - Paper, Files, Web documents, Database system

* Data mining Functionalities

1 Characterization and discrimination

generalization, summarize and contrast data characterization

2 Patterns, association, correlation

3 Classification and prediction

4 Cluster analysis - group data from new classes

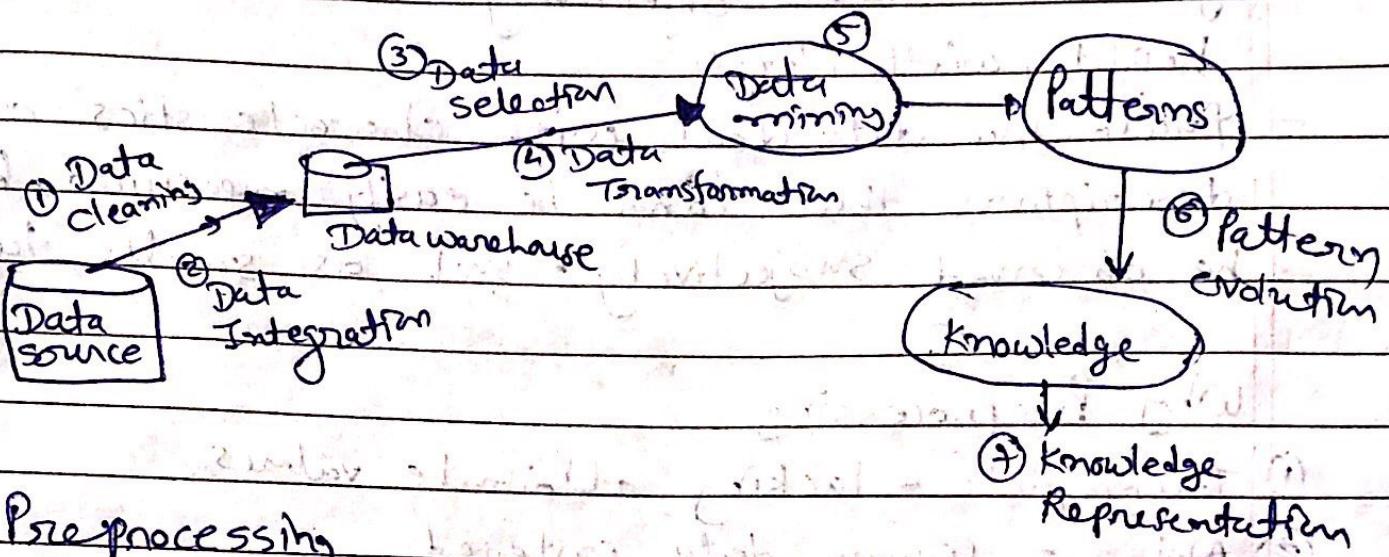
5 Outlier analysis - useful in fraud detection

6 Trend and evolution analysis

Regression Analysis, Periodicity, Similarity based Analysis

7 Statistical Analysis

Knowledge Discovery from Data



* Preprocessing

Types of Attributes

- Nominal - Ex ID number, eye color, zip codes
only name is not sufficient to predict
- Ordinal - Ex ranking, grades, height in cm
comparative value
- Interval - Ex calendar dates, temperature in celsius
range of value
- Ratio - Ex temperature

Computation can perform

* Properties of attributes / values

- (1) Distinctness: $= 1 \neq n$
- (2) Order $< >$
- (3) Addition: $+ , -$
- (4) Multiplication: $* , /$

- Quantitative - data deals with numbers and things you can measure objectively: dimensions such as height, width, length
- Qualitative data deals with characteristics and descriptors that can't be easily measured, but can be observed subjectively - such as smells, tastes, colors

Why Preprocessing

- ① Incomplete - lacking attribute values
- ② Noisy - wrong data entered
- ③ Inconsistent - discrepancies in codes → ms 2003 file not open

* Data Quality

- ① Completeness - As long as the data meets the expectation then the data is considered complete
- ② Consistency - data across all systems reflects the same information and are in sync with each other.
- ③ Conformity - Data follow some standard definitions like data type, size and format → mobile number - 10 digits
- ④ Accuracy - degree to which data is correct
- ⑤ Integrity - validity of data. all data can be traced & connected to other data
- ⑥ Timeliness - information is available when it is expected and needed

पृष्ठा 11 वृत्तालय से भगवान् ज्यतीह साक्षात् ॥

* Data Cleaning

- Parsing - put data in proper format ex Name MiddleName
- Correcting - give correct information
- Standardizing - Data convert into one standard ex ${}^{\circ}\text{C}$
- matching - Remove same tuples
- Consolidation - make one unit from two dependent units
- Dealing with missing data - Data is not present
- Dealing with Incorrect & noisy data

* Dealing with missing values

- Ignore the tuple - use when low priority, tuple is empty
- Filling Manually
- Using global constants - "Null", 'N/A'
- Use attribute mean / median
- mean: normal distribution
median: skewed distribution (Asymmetric)
- Use of determining Algorithms (decision tree, clustering)

* Missing Data Reasons

- ① equipment malfunction
- ② misunderstanding
- ③ inconsistent with ^{other} record data and thus deleted
- ④ Data may not be considered important at the entry time
- ⑤ not register history or changes of the data

॥ अपने सभी वृत्तालयों से भगवान् जयतीह साक्षात् ॥

- * Noisy Data - error or variance in measured Variable
 - Due to ① Instrument ② Data entry problems
 - ③ Data Transmission ④ Technology limitation
 - ⑤ Inconsistency in naming

* Dealing with Noisy Data

- ① Binning
 - ① equal partitioned Bin - divide Data into equal parts
 - ② Bin mean - replace Data with mean
 - ③ Bin median - replace Data with median
 - ④ Bin boundaries - replace data who value is near to either of boundary
 - equvi width : Binsize
 - equvi depth : Bin-number

* Data Integration

- Give unified view
- sources: multiple Databases, data cubes, flat files
- It is use to make Data warehouse which is use in Analysis

* Schema Integration

- Problem: customer-id in one database, c-id in second database. Is both tuple having number? is issue in Integration
- ॥ वृत्तालये स भगवान् जपतीह साक्षात् ॥

* Entity identification Problem

- identify real world entity like Person mit = Neel
- * Detecting and Resolving data value conflicts
- for the same real world entity, attribute values from different sources are different.

* Redundancy in Integration

- (1) Object identification - same attribute may have different names in different database

Ex db:mit , db2: mitkumar

- (2) Derivable Data: One attribute can be derived from other Ex age can be derived from Date of birth
- Redundant attributes may be able to be detected by correlation analysis and covariance analysis

$$\text{cov}(A, B) = E((A - \bar{A})(B - \bar{B}))$$

$$= E(AB) - \bar{A}\bar{B}$$

$$-1 \leq \text{cov}(A, B) \leq 1$$

$\text{cov}(A, B) > 0$ then A, B are +vely correlated

$\text{cov}(A, B) = 0$; no correlation

$\text{cov}(A, B) < 0$ " " -vely Correlated

* Data Integration approaches

- (1) Tight Coupling - Data combined from different sources into a single physical location through the process of Extraction, Transformation and Loading

② Loose Coupling -

Data only remains with the actual source database
which contains only those fields which are required.

* Data Reduction

- Dataset much smaller but produce same output
- Complex data set may take very long time to Analyse on dataset so reduction is required.

①

Dimension Reduction

Curse of Dimension - When dataset increase, value becomes increasingly sparse.

- Density and distance between points is much more so clustering, and analysis become less meaningful.
- Possible combination of subspaces will grow exponentially.

② Dimensionality Reduction

- Avoid the curse of dimensionality
- help eliminate irrelevant features and reduce noise

- Reduce time & space required

③ Attribute subset selection

- Redundant attributes - duplicate info (temperature in °C, F°)
- irrelevant attributes

④ Data Compression

- string compression (lossy & lossless)
- audio-video compression

॥ एति शब्दोऽप्येति वृत्तालये स भगवान् जयतीह साक्षात् ॥



LIC

भारतीय जीवन बीमा निगम
LIFE INSURANCE CORPORATION OF INDIA



Parallel Systems

- In parallel systems many operation performed parallel, in which the computational steps are performed sequentially.

A coarse-grain parallel machine have less number of processor

बिपीनकुमार अम. पटेल - massively parallel or fine-grain parallel machine uses thousands of small processor
संसद्य - अधिकारी के लिये हेतु प्रबंधक वर्ष

Bipinkumar M. Patel
Member of the Zonal Manager's Club for Agents
828 83A

→ two main measures of performance of a database system
 ① throughput - number of tasks that can be completed in a given time interval

② Response Time - amount of time it takes to complete a single task from the time it is submitted

* Speedup - running a given task in less time by increasing the degree of parallelism

Scallop - Handling larger task by increasing the degree of parallelism

→ Execution time of a task on the larger machine is T_L

" " " " " smaller " " " T_S

$$\text{Speedup} = \frac{T_S}{T_L}$$

→ Linear speedup if the speedup is "N" when the larger system has N times the resources than smaller system

→ If speedup is less than "N" then it's called sublinear speedup

→ Let Φ be a task, Φ_N is "N" times bigger than Φ

Φ on Given Machine M_S is T_S

Φ_N on a parallel machine M_L is T_L

Linear scallop if task Φ if $T_L = T_S$

Sublinear scallop if $T_L > T_S$

→ Batch scallop - the size of database increases, and the tasks are large jobs whose run time depⁿ on the size of database

कार्यालय : परबड़ी के पास, पो. समरखा.

निवास : सरदार पार्क सोसायटी, पो. समरखा, ता. जी.आणंद-388360.

फोन : (ग) (02692) 256673 (मो) 98242 32755.

Offl. : Nr. Parabadi, Po. Samarkha.

Road : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. : (ग) (02692) 256673 (मो) 98242 32755. E-mail : bipinpatellic@yahoo.com

→ Transaction scallop - rate at which transactions are submitted to the database increase and the size of the database increase proportionally to the transaction rate.

→ A machine that scales up linearly may perform worse than a machine that scales less than linearly, simply because the latter machine is much faster to start off with.

* factors diminish speedup and scalup ↴

- ① startup costs - cost with initiating a single process.
 - In parallel operation thousands of processes, startup-time may overshadow the actual processing time, speedup - effect
- ② Interference - process executing parallel often access shared resources, a slowdown process (systembus, disk, locks), speedup effect depends on ~~the~~ slowest step task.
 - Ex 100 is divided into 10 task → If one task is of 20 size then output: 5 rather than: 10.
- ③ skew - breaking down single task into number of parallel steps,

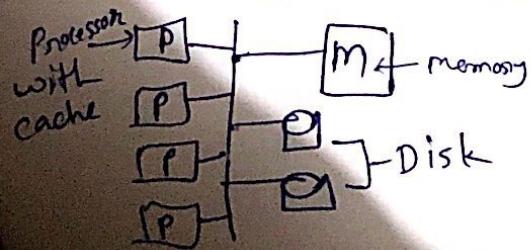
* Interconnection Networks

- ① Bus - suitable for small number of processors
 - Bus send and receive data from a single communication channel
 - Not work well to scale and increasing parallelism
- ② Mesh - two dimension each node connected to Four adjacent node
 - 3D each node connected to 6 adjacent nodes
 - Node which not connected direct communicate via other nodes
 - Scales better with increasing parallelism
- ③ Hypercube - each node is far from each other maximum $\log(n)$ links.
 - where as in mesh components may be $\approx (m-1)$ links
 - Thus, communication delay in mesh is High than hyper cube away

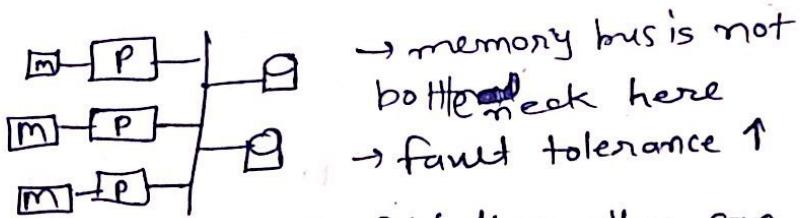
* Parallel Database Architectures

① Shared memory

- extremely efficient communication between processors
- One processor can send message to other using memory controller
- Maximum 32 or 64 processor due to bus interconnection
- Adding more CPU not work because they wait for work and access memory
- Maintaining Cache coherency become overhead with more CPU



② Shared Disk - all disks directly via interconnection network, attach with processor



बिपीनकुमार अम. पटेल

सरपाय - अधिकारीजी के सिवे हेतु प्रतिष्ठित अन्व

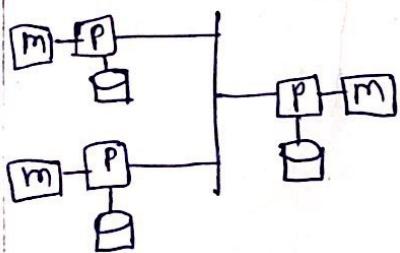
Bipinkumar M. Patel
Member of the Zonal Manager's Club for Agents
82883A

→ If one CPU fail then other can do that task due to disk is shared. - RAID architecture is used in Disk → Prob^m - scalability

→ Disk subsystem is now a bottleneck; when database makes a large number of accesses to disks

→ Large number of CPU can use but communication between CPU's slow

③ Shared Nothing + One processor may communicate with others by a high-speed interconnection Network

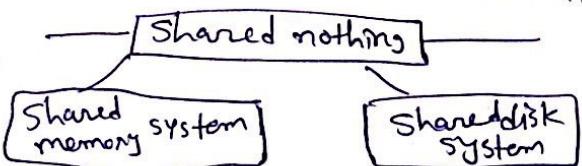


→ For I/O, no interconnection Network required so it is easily scalable

→ Drawback - costs of communication from local disk

→ software interaction acquire at both ends

④ Hierarchical + Combination of above three



→ Reduce the complexity of programming such systems have yielded "Distributed Virtual memory" - logically divide single memory between systems

- Non Uniform Memory Architecture - One memory divide different amount

* Distributed Systems +

→ The computers in distributed system are referred to by a number of different names, such as sites or nodes.

→ A local transaction is one that access data only from sites where the transaction was initiated

कार्यालय : परबड़ी के पास, पो. समरखा.

निवास : सरदार पार्क सोसायटी, पो. समरखा, ता. जी.आणंद-388360.

फोन : (नि) (02692) 256673 (मो) 98242 32755.

→ A global transaction, one that either access data in a site different from the one at which the transaction was initiated, or access data in different sites

Offi. : Nr. Parabadi, Po. Samarkha.

Resi. : Sardar Park Society, Po. Samarkha, Ta. Dist. Anand-388360.

Ph. : (R) (02692) 256673 (M) 98242 32755. E-mail : bipinpatellic@yahoo.com

Advantages of Distributed

- ① Sharing data - One site may be able to access the data residing at other sites.
- ② Autonomy - each site is able to retain a degree of control over data that are stored locally. In distribution system, there is a global database administrator responsible for entire system. Each local database administrator have local autonomy
- ③ Availability - If one site fails in system, the remaining sites may be able to continue operating. If data items are replicated in several sites, then 24×7 file available

- Easily find failure in system, Although recovery from failure in distributed system is difficult

* Two phase Protocol use in distributed system r

For example bank have four center in one city and two person simultaneously use one account then according 2PC if all branch are ready then and then either of person can commit his action. ($SOI^M \rightarrow$ Atomicity)

→ For concurrency control - locking system is implemented

Disadvantage of Distributed

- ① Software - development cost
- ② greater potential for bugs
- ③ Increased processing overhead

* Network Types r

- local area Networks - CPU allocated in same geographical Area
 - wide area network - CPU allocated distributed over Area
 - Distributed operating system require
- * LAN r numerous small computers are beneficial than one ~~one~~ single large system. It is used in office. channel is twisted pair, coaxial cable, fiber optics, wireless
- Storage Area Network is between disk and LAN. This system gives highly availability and scalability. RAID organization is used
- * WAN r
- WAN, have significant latency - speed of light delay, queuing delay
- Discontinuous Connection - Mobile Internet
- Continuous connection - Colored Internet