

Assignment - 1

- 1) What is Data Warehousing? What are the differences between OLTP and OLAP? Why have a separate datawarehouse

Ans. Data Warehousing provides architectures and tools for business executives to systematically organize, understand and use their data to make strategic decisions.

→ It refers to a database that is maintained separately from an organization's operational databases.

→ According to William H., "A data warehouse is a subject-oriented, integrated, time-variant and nonvolatile collection of data in support of management's decision making process"

The key features:-

i) Subject-oriented: A datawarehouse is organized around major subjects, such as customer, supplier, product and sales.

→ Rather than concentrating on the day-to-day operations and transaction processing of an organization, a data warehouse focuses on the modeling and analysis of data for decision makers.

ii) Integrated: A data warehouse is usually constructed by integrating multiple heterogeneous sources such as relational databases, flat files and on-line transaction records.

iii) Time-Variant: Data are stored to provide information from a historical perspective.

iv) NonVolatile: A datawarehouse is always a physically separate store of data transformed from the application data found in the operational environment.

→ Due to this separation, a data warehouse does not require transaction processing, recovery and concurrency control mechanisms.

Differences between OLAP and OLTP

OLTP	OLAP
→ Day to Day processing	→ historical processing of information
of information	→ It is used for analyzing business to know about basic business
→ Run business	→ Used by executive
	→ used by DBA professionals
	→ Managers & Analysts
→ Focus is mainly on data in	→ Focus on information out
→ Provides primitive access	→ Provides summarized and highly detailed data and consolidated data
→ Tens of records are accessed	→ Number of records accessed
→ Database size max: Gigabytes	→ Database size min: terabytes

summarized levels, special data organization, access and implementation methods based on multidimensional views.

→ If we process OLAP queries on operational databases, performance degrades

→ Operational databases support concurrent processing of multiple transactions whereas OLAP only requires read-only access of data records for summarization and aggregation.

→ If we apply concurrent controlling mechanism on OLAP operations, it will decrease the performance of OLTP.

→ Decision Support requires historical data whereas operational databases don't maintain it.

→ To view of all these, we require a separate datawarehouse.

Q. Explain Multier datawarehouse Architecture? What is the role of OLAP Servers?

Ans: Datawarehouse adopt a three-tier architecture

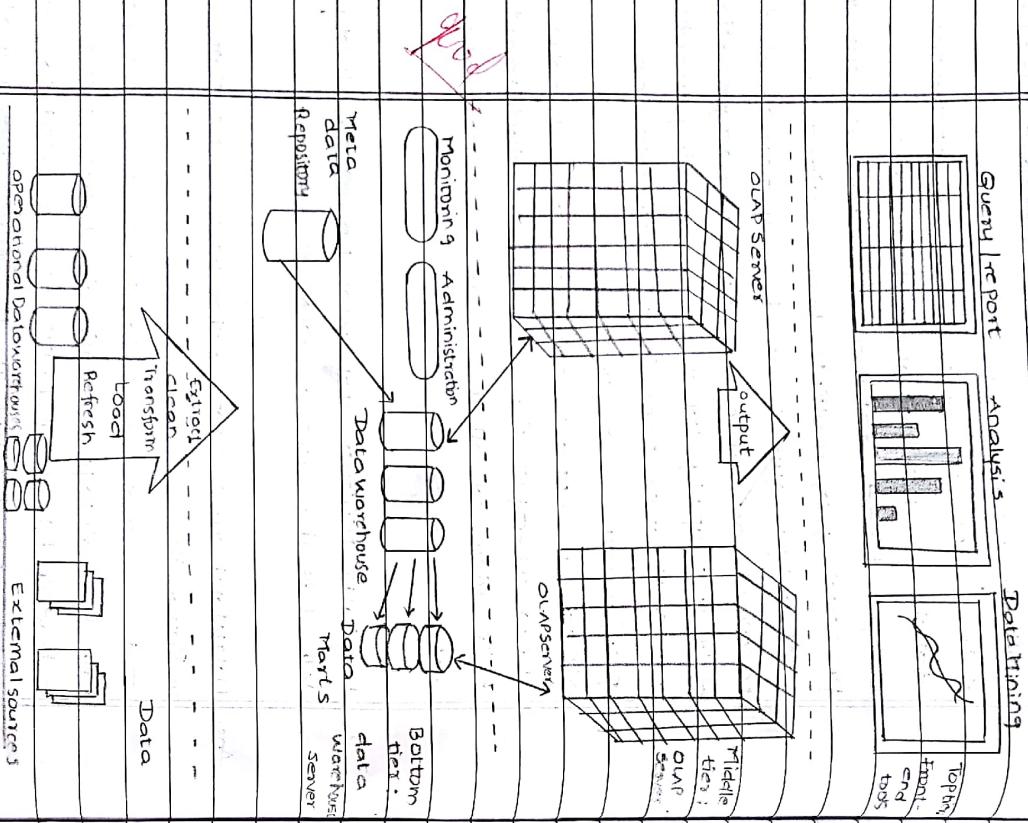
is Bottom tier:

→ It is a warehouse database server that is almost always a relational database system → Backend tools and utilities are used to feed data into the bottom tier from operational databases or other external sources

→ These tools and utilities perform data extract, loading and transformation as well as load and refresh functions to update the data warehouse.

→ The data are extracted using application program interfaces known as gateways.

A three-tier data warehousing architecture:



→ It is an OLAP server that is typically implemented using either a relational OLAP model that maps operations on multidimensional data to standard relational operations, or a MOLAP model, that is special-purpose server.

that directly implements multidimensional data and operations.

(3) Top Tier :

→ It is a front-end client layer.

→ It contains query and reporting tools, analysis tools, or data mining tools.

OLAP servers:

They present business users with multidimensional data from data warehouses or databases without concerning where data is stored.

Types of OLAP servers:

1) Relational OLAP (ROLAP) servers:

→ These are the intermediate servers that stand in between a relational back-end server and client front-end tools.

2) Dimensional OLAP (DOLAP) servers:

→ They use a relational or extended-relational DBMS to store and manage warehouse data and OLAP middle-ware to support missing pieces.

→ They include optimization for each DBMS backend, implementation of aggregation navigation logic and additional tools & services.

→ ROLAP technology tends to have greater scalability than DOLAP technology.

→ For example, DSS servers of Microstrategy employs ROLAP approach.

repository which stores information about the datawarehouse and its contents.

iii) Multidimensional (MOLAP) servers:

- These servers support multidimensional views of data through array-based multidimensional storage engines.
- They map multidimensional views directly to datacube array structures.
- The advantage of using a datacube is that it allows fast indexing to precomputed summaries.
- Sometime due to storage utilization data set may be sparse.
- For that case, Many MOLAP servers adopt a two level storage representation to handle dense and sparse data sets.
- whereas, sparse subcubes are stored as array structures whereas, dense subcubes employ compression technology for efficient storage utilization.

ii) Hybrid OLAP (HOLAP) servers:

- The hybrid OLAP approach combines ROLAP and MOLAP technology.
- It have greater scalability of ROLAP and the faster computation of MOLAP.
- ex: HOLAP server may allow larger volumes of detail data to be stored in a relational database while aggregations are kept in a separate MOLAP store.

Q. Explain in detail Datawarehouse Models.

Ans. From Architectural point of view, there are three datawarehouse models

- Enterprise warehouse
- DataMart
- Virtual Warehouse

→ An enterprise warehouse collects all the information about subject spanning the entire organization.

→ It provides corporate-wide data integration usually from one or more operational systems or external information providers & is cross-functional in scope.

→ It contains both detailed data and summarized data and can range in size from a few gigabytes to hundreds of gigabytes, terabytes or beyond.

→ It is implemented on traditional mainframe, computer super servers or parallel architecture platforms.

→ It requires extensive business modeling & may take years to design & build.

L1) Data Mart:

- A data mart contains a subset of corporate wide data that is of value to a specific group of users.
- It is confined to specific selected subjects.
- ex: Marketing-data mart is confined to customer, item, sales
- Data in data mart is summarized.

- DataMarts are implemented on low-cost departmental servers (unlike dataware based)
- Implementation cycle of datamart is weeks rather than months or years.
- If planning of data mart is not enterprise wide, it may involve complex integration in long run.
- Depending on source of data, data marts are of 2 types
 - Independent DataMart
 - Sources are from data as external info providers or from one or more data warehouses.
- Dependent Data Mart Sources are captured from enterprise
- or from data generated locally within a particular department geographic user
- Virtual Warehouse:**
- It is a set of views over operational databases.
- For efficient query processing, only some of the possible summary views may be materialized.
- A virtual warehouse is easy to build but requires excess capacity on operational database servers.
- If planning of data mart is not enterprise wide, it may involve complex integration in long run.

4. What is multidimensional data model? Explain with an example.

Ans: Data Warehouses and OLAP tools are based on a multidimensional data model.

→ This model views data in the form of a data cube.

→ A data cube allows data to be modeled and viewed in multiple dimensions.

→ It is defined by dimensions and facts.

Dimensions :-
~~~~~

→ They are the perspectives or entities with respect to which an organization wants to keep records.  
ex: AllElectronics may create a sales data warehouse in order to keep records of the store's sales with respect to the dimensions time, item, branch and location.

→ These dimensions allow the store to keep track of things like monthly sales of items and the branches and locations at which the items were sold.

Facts :-

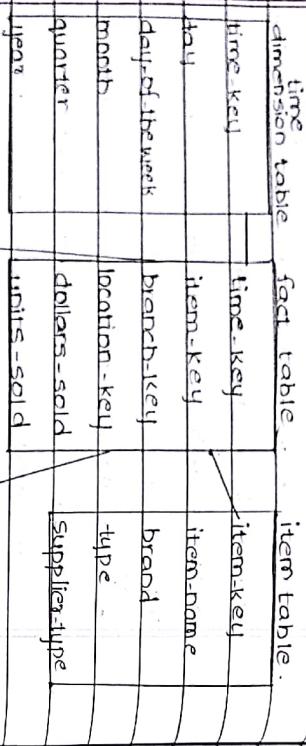
→ They are numerical measures.  
→ They are the quantities by which we want to analyze relationship between dimensions  
ex: facts for a sales data warehouse include dollars-sold (sales amount in dollars), units-sold (no. of units sold).

Star schema :- In star schema the data warehouse contains a large facttable (fact table) containing the bulk of the data with no redundancy.

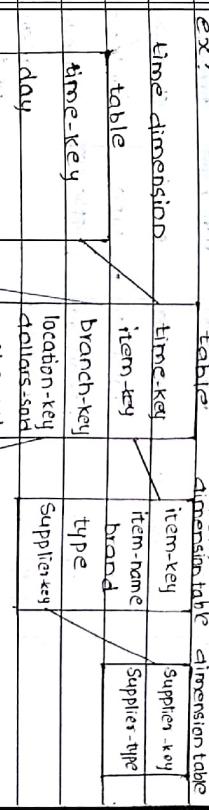
- ↳ a set of smaller attendant tables (dimension tables) one for each dimension.

Dimension tables) one for each dimension.

ex: sales



ex: sales fact table dimension table dimension table dimension table



→ The resulting schema graph forms a shape similar to a snowflake.

Snowflake :-  
→ The Snowflake schema is a variant of the star schema model.

- In snowflake, some dimension tables are normalized i.e., by further splitting the data into additional tables.

→ The resulting schema graph forms a shape similar to a snowflake.

In the above example, it contains one fact table sales, with keys of four dimensions along with 2 measures: dollars-sold and units-sold.

- And 4 dimension tables (time, item, branch, location)
- Sometimes the attributes in dimension table results in redundancy.

ex: location, location-key, street, city, province or state, country)

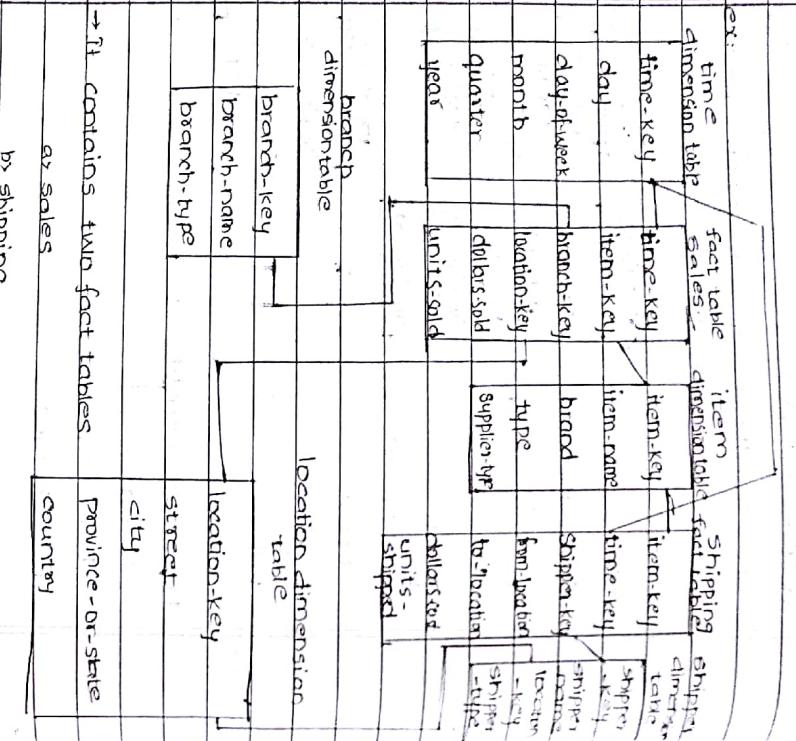
↳ (---, Vancouver, British Columbia, Canada)

→ C - -> Victoria, British Columbia, Canada)

### Fact constellation:

- It contains multiple fact tables to share dimension tables.

→ This schema is viewed as collection of stars and hence called **galaxy schema or fact constellation**.



### Data Cube Measures:

- A datacube measure is a numeric function that can be evaluated at each point in the datacube space.

→ A measure value is computed for a given point by aggregating the data corresponding to the respective dimension-value pairs defining the given point.

Measures are organised into 3 categories:

#### 1. Distributive:

→ If the result derived by applying the function to n aggregate values is the same as that derived by applying the function on all the data without partitioning, then the function can be computed in a distributed manner.

#### 2. example: counts

It can be computed for a data cube by first partitioning the cube into a set of subcubes, computing counts for each subcube and then summing up the counts obtained for each subcube. Subcube count is a distributive function.  
→ Distributive measures can be computed efficiently because they can be computed in a distributive manner.

→ Both sales and shipping shares the dimension tables item, time and location.

- Explain the following with an example
- Data cube measures
- Concept hierarchy

### 2. Algebraic:

→ It can be computed by an algebraic function with m arguments (where m is a bounded integer), each of which is obtained by applying a distributive aggregate function.  
ex: avg, can be computed by sum/counts functions.

→ A measure is algebraic if it is obtained by applying an algebraic aggregate function

3) Holistic:

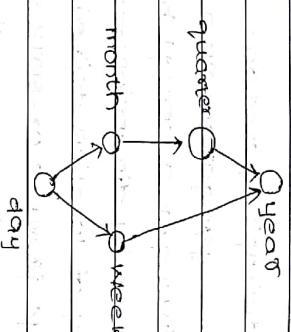
→ If there is no constant bound on the storage size needed to describe a subaggregate ex: median, mode, rank,

→ A measure is holistic if it is obtained by applying a holistic aggregate function.

5b) Concept hierarchy:

→ A concept hierarchy defines a sequence of mappings from a set of low-level concepts to higher-level, more general concepts.

→ Many concept hierarchies are implicit within the database schema.  
ex: Concept hierarchy of location:  
 $\text{street} < \text{city} < \text{province-or-state} < \text{country}$ .



→ All these attributes are related by state or province-or-state in total order.

→ A concept hierarchy that is a total or partial order among attributes in a database schema is called schema hierarchy.  
→ Concept hierarchies may be defined by discretizing or grouping values for a given dimension or attribute resulting in a set-grouping hierarchy.

If we consider values:

$\text{location} = \{\text{au}, \text{ca}, \text{usa}\}$

$\text{country} = \{\text{canada}, \text{usa}\}$

$\text{ProvinceOrState} = \{\text{British Columbia}, \text{Ontario}, \text{New York}, \text{Illinois}\}$

$\text{City} = \{\text{Vancouver}, \text{Victoria}, \text{Toronto}, \text{Ottawa}, \text{New York}, \text{Chicago}, \text{Buffalo}, \text{Milwaukee}\}$

To them, cities are mapped to province-or-state which are mapped to countries. These mapping

form a concept hierarchy for the dimension locator mapping a set of low-level concepts (i.e., cities) to higher-level, more general concepts (i.e., countries).

→ Alternatively, the attributes of a dimension may be organized in a partial order forming a lattice.

ex: for time dimension partial order is "day < month < quarter < week < year".

In the above example, interval  $(\$x \dots \$y)$  denotes the range from  $\$x$  (exclusive) to  $\$y$  (inclusive).

→ There may be more than one concept hierarchy for a given attribute or dimension, based on different user view points.

ex: user prefers to organize price by defining ranges for inexpensive, moderately-priced

and expensive.

→ Concept hierarchy allows data to be handled at varying levels of abstraction.

Q7 Explain different OLAP operations in the multidimensional data model?

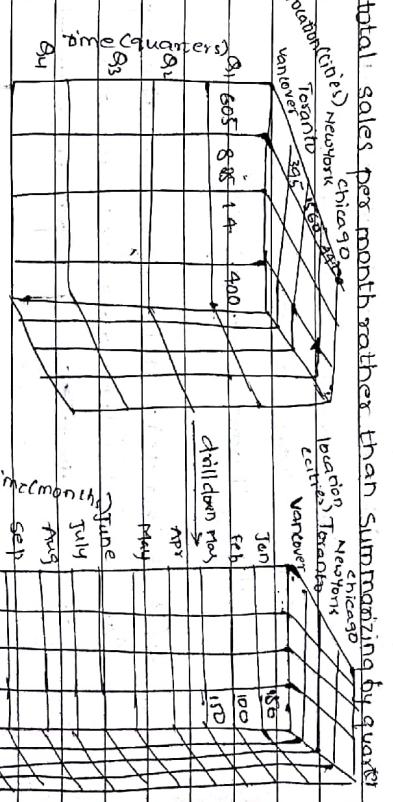
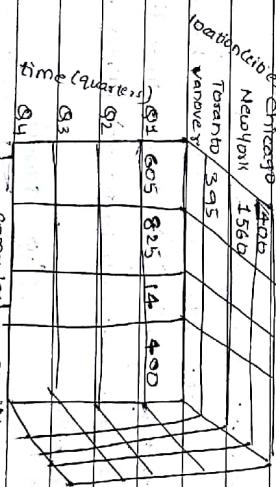
Ans. OLAP operations

i) Roll-up (drill-up) operation.

→ It performs aggregation on a data cube either by climbing up a concept hierarchy for a dimension or by dimension reduction.

ex: Hierarchy of location defined as total order "street < city < province-or-state < country"

The roll-up operation groups the data by city, country rather than grouping by city.



ii) Drill-down

→ It navigates from less detailed data to more detailed data.

→ It can be realized by stepping down a concept hierarchy for a dimension or introducing additional dimensions.

ex: Consider concept hierarchy for time

day < month < quarter < year

Drill down occurs by descending the time hierarchy from the level of quarter to more detailed level of month. So the resulting data cube details

total sales per month rather than summarizing by quarter

location (city) New York, Chicago, location Toronto, Vancouver

quarters Q1, Q2, Q3, Q4

time month

year

month

July

Aug

Sept

Oct

Nov

Dec

When rollup is performed by dimension reduction one or more dimensions are removed from the given cube.

→ For example: sales data cube contains only 2 dimensions location and time

On them Roll-up may be performed by removing the time dimension, resulting in an aggregation of the total sales by location rather than by location and by time.



## city bitmap index table

|    |   |   |                                                                                   |
|----|---|---|-----------------------------------------------------------------------------------|
| R1 | 1 | 0 | → bitmap indexing is advantageous compared to hash and tree indices               |
| R2 | 1 | 0 | hash and tree indices                                                             |
| R3 | 1 | 0 | → IL is especially useful for low-cardinality domains                             |
| R4 | 1 | 0 | because comparison, join and aggregation operations are reduced to bit arithmetic |
| R5 | 0 | 1 | it reduces processing time                                                        |
| R6 | 0 | 1 |                                                                                   |
| R7 | 0 | 1 |                                                                                   |
| R8 | 0 | 1 |                                                                                   |

→ Bit mapping leads to significant reductions in space and I/O since a string of characters can be represented by a single bit

## Join Indexing

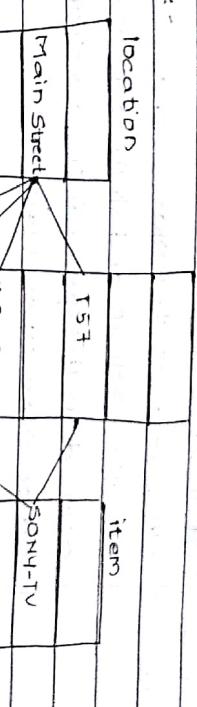
- It has gained popularity from its use in relational database query processing.
- Traditional indexing maps the value in a given column to a list of raw having that value

→ Whereas, join indexing requires the joinable rows of two relations from a relational database

- ex: if two relations R(BID, R) and S(B, SID) join on the attributes A and B, then the join index record contains the pair (BID, SID)
- where BID and SID are record identifiers from the R & S relations

- Join index records can identify joinable tuples without performing costly join operations.

→ Join indexing especially useful for maintaining the relationship between a foreign key and its matching primary keys from the joinable relation.



## Join index table for join index table for

| location    | sales-key | item    | sales-key |
|-------------|-----------|---------|-----------|
| Main Street | T57       | Sony-TV | T57       |
| Main Street | T238      | Sony-TV | T459      |
| Main Street | T884      |         |           |

## Join index table linking two dimensions

| location/item | sales   |
|---------------|---------|
| Main Street   | Sony-TV |
| Main Street   | T57     |

- 5) Explain the methods for efficient implementation of data warehouse systems.

- Ans: The purpose of materializing cuboids and constructing OLAP index structures is to speed up

query processing in data cubes.

Given materialized views, query processing should proceed as follows.

1) Determine which operations should be performed on the available cuboids:

→ This involves transforming any selection, projection, roll-up and drill-down operations specified in the query into corresponding SQL and/or OLAP operations.

→ For example, slicing and dicing a data cube may correspond to selection and/or projection operations on a materialized cuboid.

2) Determine to which materialized cuboids the relevant operations should be applied.

→ This involves identifying all of the materialized cuboids that may potentially be used to answer the query, pruning the above set using knowledge of "dominance" relationships among the cuboids, estimating the costs of using the remaining materialized cuboids and selecting the cuboid with the least cost.

ex: Determine which materialized cuboid(s)

should be selected for OLAP operation.

→ Let the query to be processed be on with the condition "year = 2004" and there are 4 materialized cuboids available

1) year, item-name, city

2) year, brand, country

3) year, brand, province-or-state

4) item-name, province-or-state where year = 2004

which should be selected to process the query?

1) finer granularity data cannot be generated from coarser-granularity data.  
→ ② cannot be used because country is more general concept than province-or-state cuboid.

2, 3 and 4 can be used to process the query because → they have the same set (or) a superset of the dimensions in the query.

3) When we compare costs, cuboid 4 would cost the most because both item-name and city are at a lower level than the brand and province-or-state concepts specified in the query.

→ If there are not many year values associated with items in the cube, but there are several item-names for each brand, then cuboid ③ will be smaller than ④ and thus ③ should be chosen to process the query. But if efficient indices are available for cuboid ④, then ④ may be a better choice.

Q1 what is datamining? what are the motivating challenges and origins of data mining.

Ans. Data Mining is the process of automatically discovering useful information in large data repositories.

→ Data mining techniques are deployed to scour large databases in order to find new and useful patterns that might otherwise remain unknown.

→ They also provide capabilities to predict the outcome of a future observation, such as predicting a newly arrived customer will spend more than \$100 at a department store.

### Motivational challenges of data mining:

The specific challenges that motivated the development of data mining are:

#### 1. Scalability:

- Because of advances in data generation and collection, data sets with sizes of gigabytes, terabytes or even petabytes are becoming common.
- If data mining algorithms are to handle this massive data sets, then they should be scalable.
- Many data mining algorithms employ special search strategies to handle exponential search problems.
- To access individual records efficiently scalability requires efficient implementation of novel data structure.
- Scalability can be improved by using sampling, developing parallel & distributed algorithms.

#### 2. Heterogeneous and Complex Data:

- Traditional analysis methods deal with data sets containing attributes of same type, either continuous or categorical
- As ~~now-a-days~~<sup>need for</sup> data-mining has grown in many fields, this led to techniques that handle heterogeneous attributes.
- Recently complex data objects emerged like web pages, containing semi-structured text and hyperlinks and ~~DNA~~<sup>DNA</sup> data with sequential and three-dimensional structure.
- Techniques developed for mining such complex objects should take into consideration relationships in the data, such as temporal and spatial autocorrelation, graph connectivity and parent-child relationships between the elements in semi-structured text and XML documents.

#### 3. High Dimensionality:

- Now-a-days datasets have many attributes (hundreds or thousands) rather than a handful.
- Datasets with temporal or spatial components also tend to have high dimensionality.

ex: Consider a dataset that contains measurement of temperature at various location.

→ If the temperature measurements are taken repeatedly for an extended period, the no. of dimensions increase in proportion to the no. of measurements taken.

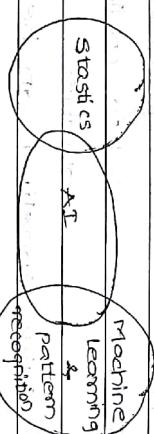
→ Traditional techniques developed for low dimensions, may not work for them.

→ Some Data Analysis algorithms, the computation complexity increases rapidly, as the dimensionality increases.

#### 4. Data Ownership and Distribution:

- The data needed for an analysis is not stored in one location or owned by one organization.

- But it is geographically distributed among regions belonging to multiple entities.
- This requires the development of distributed data mining techniques.
- The challenges faced by distributed data mining algorithms include:
- i> How to reduce the amount of communication needed to perform the distributed computation.
  - ii> How to effectively consolidate the data mining results obtained from multiple sources.
  - iii> How to address data security issues.
- Non Traditional Analysis:-
- The traditional statistical approach is based on hypothesize and test paradigm.
  - In other words, a hypothesis is proposed and an experiment is designed to gather the data and then the data is analyzed with respect to the hypothesis.
  - It is highly labor intensive.
  - Current data analysis tasks often require the generation and evaluation of thousands of hypotheses.
  - Consequently, the development of some data mining techniques has been motivated by the desire to automate the process of hypothesis generation & evaluation.
  - Origins of data mining.
  - The data mining is drawn upon ideas of sampling, estimation and hypothesis testing from statistics.



[Distributed Technology, Parallel computing, Distributed]

Explain different data mining tasks in detail with examples:

Ans: Data mining tasks are divided into 2 categories:

i) Predictive Tasks:-

→ The objective of these tasks is to predict the value of a particular attribute based on the values of other attributes.

→ The attribute to be predicted is known as target attribute, for making the prediction used

are known as the explanatory or independent variables.

(Q) Search algorithms, modelling techniques and learning theories from AI, pattern recognition and machine learning.

→ Data mining has been quick to adapt ideas from other areas like

Like Database Systems are needed to provide support for efficient storage, indexing and query processing.

ii> Parallel computing, are important in addressing the massive size of some data sets.

iii> Distributed techniques can help address the issue of size and are essential when the data cannot be gathered in one location.

### Descriptive tasks:

Here, the objective is to derive patterns that summarize the underlying relationships in data.

→ Descriptive data mining tasks are often exploratory in nature and frequently require postprocessing techniques to validate and explain the results.

### Predictive modelling:

→ It refers to the task of building a model for the target variable as a function of the explanatory variables.

→ There are 2 types of predictive modelling tasks:

→ Classification: used for discrete target variable

→ Regression: used for continuous target variables

→ Predictive modeling can be used to identify customers that will respond to a marketing campaign, predict disturbances in the Earth's ecosystem, or judge whether a patient has a particular disease based on results of a medical test.

ex: predicting the type of flower:

→ Consider the task of predicting a species of flower based on the characteristics of flower.

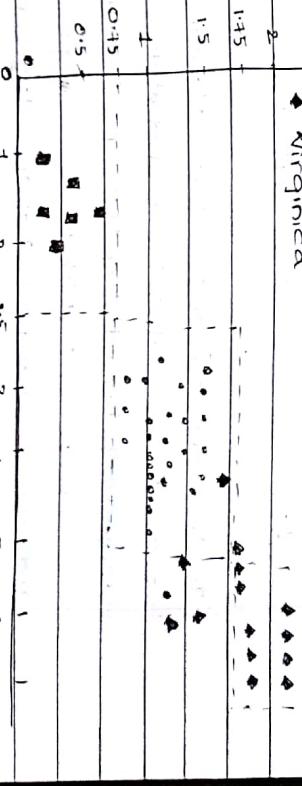
→ Consider classifying an Iris flower whether it belongs to one of following three Iris

Species: Setosa, Versicolour or Virginica

For this classification, we need a data set that contains characteristics of various flowers

of these 3 species.

→ In addition to the species, it contains additional attributes (Sepal width, Sepal length, Petal length and Petal width).



Based on these categories of petal width (low [0, 0.75], medium [0.75, 1.75], high [1.75, 2.5])

and petal length ([0, 0.25], [0.25, 0.5], [0.5, 0.75], [0.75, 1.0]), the following rules can be derived.

→ Petal width low & petal length → Setosa

→ Petal width medium & petal length medium → Versicolour

→ Petal width high & petal length high → Virginica

### Associative Analysis:

→ It is used to discover patterns that describe strongly associated features in the data.

→ The discovered patterns are typically represented in the form of implications (rules), or feature subsets.

ex: Market Basket Analysis:

Transaction-ID      Items

1                      \$Bread, Butter, Diaper, Milk,

2                      \$Coffee, Sugar, Cookies, Salmon

3                      \$Bread, Butter, Coffee, Diapers, Milk, Eggs



er: Credit Card Fraud Detection)

→ A credit card company records the transactions made by every credit card holder, along with personal information such as credit limit, age, annual income and address.

→ When a new transaction arrives, it is compared against the profile of the user.

→ If the characteristics of the transaction are very different from the previously created profile, then the transaction is flagged as potentially fraudulent.

Q1] What is dataset? Explain different types of dataset.

Ans: A data set can often be viewed as a collection of data objects where data objects can be record, point, vector, pattern, event, case, sample, observation or entity.

The different types of data set are:-

1> Record Data:

→ Data set is a collection of records, each of which consists of a fixed set of data fields.

→ Record data is stored either in flat files or in relational databases.

2> Different types of record data:

→ Transaction or Market Basket Data.

→ It is a special type of record data, where each record involves a set of items.

Ex: Consider a grocery store:

The set of products purchased by a customer during one shopping trip constitutes a transaction.

while the individual products that were purchased are the items.  
This type of data is called market basket data.

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | Yes    | Single        | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 3   | No     | Single        | 70.5    | No        |
| 4   | No     | Married       | 100.15  | No        |

| Tid | Refund | MaritalStatus | Taxable | Defaulted |
|-----|--------|---------------|---------|-----------|
| 1   | No     | Married       | 100.15  | No        |
| 2   | No     | Married       | 12.5K   | No        |

### 3) The Sparse Data Matrix:

→ It is a special case of a data matrix in which the attributes are of same type & are asymmetric i.e., only non-zero values are important.

ex: Transaction data is an example of a sparse data matrix that has only 0-entries

|            | team | coach | play | ball | score | game | win |
|------------|------|-------|------|------|-------|------|-----|
| Document 1 | 3    | 0     | 5    | 0    | 2     | 6    | 0   |
| Document 2 | 0    | 7     | 0    | 2    | 1     | 0    | 0   |

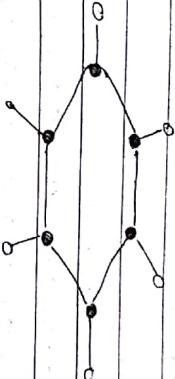
### Graph-Based Data:

→ Graph is a powerful representation for data specific cases

↳ The graph captures relationship among data objects  
↳ The data objects themselves are graphs

Data with Relationships among Objects:  
→ The relationships among objects frequently convey important information  
→ In them, data objects are mapped to nodes of the graph, while the relationships among objects are captured by the links between objects & link properties such as directions & weight  
ex: Webpages on worldwide web.

Data with objects that are graphs: If objects have structure, that is objects contain subobjects that have relationships then such objects are frequently represented



### Sequential Data:

→ It also referred to as temporal data, it is an extension of record data where each record has a time associated with it

ex: Movie Customers Items purchased  
t1 c1 A,B  
t2 c2 A,C  
t3 c3 C,D  
  
Customer Time & items purchased  
c1 C1:A,B) C2:C,D)  
c2 (t3: A,D)  
c3 (t2:A,C),  
  
Sequence Data:

→ It consists of a data set, that is sequence of individual entries such as sequence of words or letters.

→ It is similar to sequential data, but there no time stamps instead there are positions in an

as graphs.  
ex: Structure of chemical compounds can be represented as graphs.

ordered sequence  
e.g. characteristics of plants & animals  
can be represented in the form of sequence  
of molecules that are present as genes

Genetic Information  
Genetic sequence  
concrete sequence  
→ It is human genetic code formulation  
DNA is constructed as a sequence

### Q1 Explain Data Preprocessing Tasks.

Time series data:  
→ It is a special type of sequential data  
in which each record is a time series i.e.  
a series of measurements taken over time  
e.g. financial data set might contain  
objects that are time series of the daily  
prices of various stocks

- Ans. Data is often collected for unspecified applications
- Data may have quality problems that need to be addressed before applying a data mining technique
    - Noise and outliers
    - Missing values
    - Duplicate data
  - Preprocessing may be needed to make data more suitable for data mining

Data Preprocessing tasks:  
Combining two or more objects into a single object

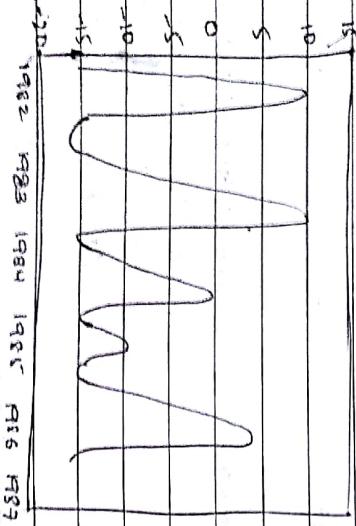
| Transaction ID | Item    | Store Location | Date     | Price   |
|----------------|---------|----------------|----------|---------|
| 101123         | watch   | chicago        | 09/06/04 | \$25.99 |
| 101123         | Battery | chicago        | 09/06/04 | \$5.99  |
| 101124         | shoes   | minneapolis    | 09/06/04 | \$39.99 |

### Spatial Data..

Some objects have spatial attributes  
such as position or areas as well as  
other types of attributes

An example of spatial data is a weather data that is collected for a variety of geographical locations

- Spatial data examples are the sciences
- engineering data sets that are the result of measurements or model output taken at regularly or irregularly distributed points or 2 or 3 dimensional grid or mesh



aggregation is the process of eliminating attributes eg

- reducing the possible values of data from 365 days to 12 months.
- This type of aggregation is online Analytics

processing

Purpose of Aggregation:

i. Data Reduction

- reduce the no. of attributes or objects
- ii. High-level view of the data
- easier to discover patterns
- iii. More "stable" data

- Aggregated data tends to have less variability

Sampling:

- If is a technique employed for selecting a subset of the data
- It is used in data mining for the following reasons

- is It may be too expensive or too time consuming to process all data
- ii. To measure a classifier's performance

- the data may be divided in a training set and a test
- iii. To obtain a better balance between class distribution

Sampling Techniques:

Simple Random Sampling

- Every sample of size  $n$  has the same chance of being selected.
- Perfect random sampling is difficult to achieve in practice

→ Use Random Numbers.

There are 2 variations in it

- is Sampling without replacement
- ii. A selected item cannot be selected again - removed from the full data set once selected

is Sampling with replacement

- Items can be picked up more than once for the sample - not removed from the full dataset once selected.

Stratified Sampling:

- Split the data into several partitions, then draw random samples from each partition.
- Each partition may correspond to each of the possible classes in the data.
- The no. of items selected from each partition is proportional to partition size.

## Dimensionality Reduction:

- Determining dimensions that are important for modeling

→ Reason for dimensionality reduction

- ↳ Many data mining algorithms work better if the dimensionality of data is lower

- ↳ Allows the data to be more easily visualized

- ↳ If dimensionality reduction eliminates irrelevant features or reduces noise, the quality of results may improve
- ↳ can lead to a more understandable "way"

### Curse of dimensionality:

- Data analysis becomes significantly harder as the dimensionality of the data increases.

### Dimensionality reduction using Principal Component Analysis:

- Some of the most common approaches for dimensionality reduction, particularly for continuous data, use a linear or non-linear projection of data from a higher dimensional space to a lower dimensional space.

## Dimensionality Reduction by Feature Subsets

selection.

- The reduction of dimensionality by selecting new attributes that are a subset of old is known as feature subset selection or feature selection

### Redundant features:

- They contain no information that is useful for data mining tasks at hand
- ↳ Student ID numbers would be irrelevant to the task of predicting their GPA

### Irrelevant features:

- If there are features present, they can reduce classification accuracy & the quality of clusters that are found

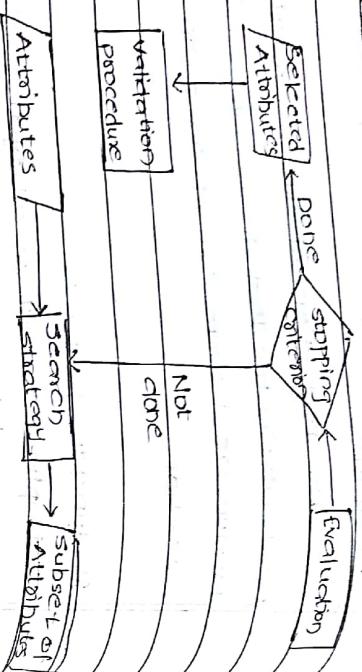
There are 3 standard approaches to feature selection:

#### → Embedded Approaches: feature selection occurs naturally as part of the data mining algorithm

- Filter Approaches: features are selected before the data mining algorithm is run
- Wrapper Approaches: use the target data mining algorithm as a black box to find the best subset of attributes

- PCA is a linear Algebra technique for continuous attributes that find new attributes that are linear combinations of original attributes
- are orthogonal to each other
- capture the maximum amount of variance in data

## Architecture for Feature subset selection



## Feature Creation:

Sometimes a small number of new attributes can capture the important information in a data set much more efficiently than the original attributes.

- Also, the no. of new attributes can be often smaller than the no. of original attributes hence, we get benefits of dimensionality reduction.

### → Three general methodologies

#### i) Feature Extraction:

→ One approach to dimensionality reduction is feature extraction, which is creation of a new, smaller set of features from the original set of features.

#### ii) Mapping the Data to a New Space

Sometimes, a totally different view of the data can reveal important & interesting features.

ex: Applying Fourier transformation to data to detect time series patterns.

## Feature construction

→ Features have the necessary information but not in the form necessary for the data mining algorithm.

→ In this case, one or more new features constructed out of original features may be useful.

ex: There are 2 attributes that record values

and mass of a set of objects

→ Suppose there exists a classification model based on material of which the objects are constructed.

→ Then a density feature constructed from the original 2 feature would help classification.

## Discretization and Binarization

→ Discretization is the process of converting a continuous attribute to a discrete attribute.

→ common example: rounding off real nos. to integers

→ Some data mining algorithms require that the data be in the form of categorical or binary attributes

→ Thus it is often necessary to convert continuous attributes to categorical

attributes and/or binary attributes

→ Its pretty straightforward to convert categorical attributes to discrete or binary attributes

### Variable Transformation:

→ It refers to a transformation that is applied to all values of an attribute i.e. for each object, the transformation is applied to the value of the attribute for that object.

→ There are 2 important types of attribute transformations.

↳ Simple function Transformations

→  $x' = \log x, \text{ex. } \sqrt{x} \text{ etc.}$

↳ Standardization or Normalization

$$\text{std. dev}(x) = \sqrt{\frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2}$$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n (x_1 + x_2 + \dots + x_n)$$

$$\sigma_x = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2}$$

b) Explain the similarity & dissimilarity between simple attributes, for the following vectors  $x$  &  $y$ . Calculate the indicated similarity

or distance measure.

$x = [1, 1, 1, 1]$ ,  $y = [2, 2, 2, 2]$

$$Sxy = \frac{1}{3} \sum_{k=1}^4 (x_k - 1)(y_k - 2)$$

$$= \frac{1}{3} [(1-1)(2-2) + (1-1)(2-2) + (1-1)(2-2) + (1-1)(2-2)]$$

$$= 0$$

a) cosine b) correlation.  
c) Euclidean d) L1 Norm, e) L2 Norm

$$\cos(x, y) = 0$$

a) cosine.

$$x \cdot y = 1 \cdot 2 + 1 \cdot 2 + 1 \cdot 2 + 1 \cdot 2$$

$$= 2 + 2 + 2 + 2 \Rightarrow 8$$

$$\|x\| = \sqrt{1+1+1+1} = \sqrt{4} = 2$$

$$\|y\| = \sqrt{1+1+1+1} = \sqrt{4} = 2$$

$$\cos(x, y) = \frac{8}{2 \times 4} = 1$$

b) correlation.

$$\bar{x} = 4/4 = 1$$

$$\bar{y} = 8/4 = 2$$

$$\text{std. dev}(x) = \sqrt{\frac{1}{3} \sum_{k=1}^3 (x_k - \bar{x})^2}$$

c) Euclidean.

d) L1 Norm

e) L2 Norm

$$P_1 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$P_2 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$P_3 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$P_4 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$Sx = \sqrt{\frac{1}{15} \sum_{k=1}^5 (x_k - 0.667)^2}$$

$$\begin{array}{l} P_1 \\ P_2 \\ P_3 \\ P_4 \end{array} = \sqrt{\frac{1}{5} [(0.82688) - 0.667]^2} = 0.4066.$$

$$\text{a)} \quad \begin{array}{c} P_1 \\ P_2 \\ P_3 \\ P_4 \end{array} = \begin{array}{ccccc} 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{array}$$

$$Sy = \sqrt{\frac{1}{5} \sum_{k=1}^5 (y_k - 0.667)^2} = 0.4066$$

$$\begin{array}{l} P_1 \\ P_2 \\ P_3 \\ P_4 \end{array} = \begin{array}{ccccc} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{array}$$

$$\text{c)} \quad \begin{array}{c} P_1 \\ P_2 \\ P_3 \\ P_4 \end{array} = \begin{array}{ccccc} 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{array}$$

$$\text{d)} \quad \begin{array}{c} P_1 \\ P_2 \\ P_3 \\ P_4 \end{array} = \begin{array}{ccccc} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{array}$$

$$\text{Sxy} = \frac{1}{5} \sum_{k=1}^5 (x_k - 0.667)(y_k - 0.667) = 0.06668$$

$$\text{Correlation} = \frac{0.06668}{(0.4066)(0.4066)} = 0.40333$$

c) Taccord

$$M_{100} = 1 \quad T = M_{11} / (M_{100} + M_{01})$$

$$M_{01} = 1 \quad = 3/3$$

$$M_{10} = 1 \quad = 1$$

$$M_{11} = 3$$

a) Cosine

$$x \cdot y = 1 \cdot 1 + 1 \cdot 1 = 3$$

$$\|x\| = \sqrt{1+1+1+1} = 2$$

$$\|y\| = \sqrt{1+1+1+1} = 2$$

$$\cos(x,y) = \frac{3}{2 \times 2} = \frac{3}{4} = 0.75$$

b) Correlation

$$\bar{x} = 4/4 = 0.667$$

$$\bar{y} = 4/4 = 0.667$$