# DATA MINING

**Mahesh T R**
**Vikram Singh**

# Data Mining

# Data Mining

Mahesh T R
Vikram Singh

**BOOKS ARCADE**
KRISHNA NAGAR, DELHI

# Data Mining

Mahesh T R
Vikram Singh

# BOOKS ARCADE

# CONTENTS

# CHAPTER 1

## INTRODUCTION TO DATA MINING

Mahesh T R

Associate Professor, Department of Computer Science Engineering, Faculty of Engineering and Technology,
JAIN (Deemed-to-be University), Karnataka – 562112
Email Id- t.mahesh@jainuniversity.ac.in

The term "data mining" describes the process of obtaining information from vast volumes of data. Actually, the phrase is a misnomer. Thus, knowledge mining, which places a focus on mining from enormous volumes of data, would have been a more fitting term for data mining.It is a computer process that combines techniques from artificial intelligence, machine learning, statistics, and database systems to find patterns in massive data sets.

The main objective of data mining is to take information from a data collection and organize it so that it may be used in other ways. The fundamental elements of data mining are

- A. Automatic pattern recognition
- B. forecasting probable outcomes
- C. creation of knowledge that can be used
- D. Concentrate on big data and databases

**The Purpose of Data Analysis**

The term "data mining" comes from the analogy between mining a mountain for a vein of precious mineral and looking for useful business information in a vast database, such as locating related goods in terabytes of store scanner data. In both cases, finding the value requires either combing through a vast quantity of information or carefully probing it. Given large, high-quality datasets, data mining technologies may provide new business prospects by offering these features. Automated trend and behavior prediction. The technique of discovering predictive information in huge datasets is automated by data mining. Questions that previously needed significant manual analysis may now be swiftly and immediately addressed from the data. Targeted marketing is a prevalent illustration of a predictive issue. The targets that are most likely to maximize the return on investment for future mailings are identified using data from previous promotional mailings. Predicting bankruptcy and other types of default as well as identifying demographic subgroups who are likely to react similarly to certain situations are additional predictive difficulties. Automated pattern recognition of undiscovered patterns in a single step, data mining technologies scan databases to find previously undetected patterns. Analyzing retail sales data to find apparently unrelated goods that are often bought together is an example of pattern discovery. Detecting fraudulent credit card transactions and locating abnormal data that can indicate data entry keying mistakes are two further pattern detection issues.

**Data mining tasks**

Six categories of frequent tasks are involved in data mining:

Anomaly detection, also known as outlier/change/deviation detection, is the process of finding uncommon data records that may be intriguing or data problems that need further research.The process of association rule learning, also known as dependency modelling, looks for connections between variables. For instance, a supermarket may compile

information on client shopping patterns. The supermarket may find out which goods are usually purchased together using association rule learning and utilize this knowledge for marketing. This is also known as market basket analysis.Finding groupings and structures in the data that are somewhat "similar" without employing pre-existing data structures—is the problem of clustering.Classification is the process of applying established structure to fresh data. An email software could try tocategorise a message as "genuine" or "spam," for instance.Regression looks for a function that models the data with the least amount of error.Summarization, which includes report preparation and visualization, gives the data set a more condensed form.

## Data mining architecture

The key parts of a typical data mining system can include the following (Figure 1.1).



**Figure 1.1: Data mining architecture**

Knowledge Base: This is the body of subject-matter expertise that is used to direct searches or gauge how intriguing patterns that emerge are. Concept hierarchies, which are utilized to arrange attributes or attribute values into various degrees of abstraction, might be a part of this knowledge.It may also include information like user beliefs, which may be used to judge a pattern's interest level depending on how surprising it is. Metadata and extra interestingness restrictions or thresholds are other instances of domain knowledge (e.g., describing data from multiple heterogeneous sources).

## Data Mining Engine:

This component of the system is crucial for data mining and should include a collection of functional modules for operations including characterization, association and correlation analysis, classification, prediction, cluster analysis, outlier analysis, and evolution analysis.The pattern evaluation module interacts with the data mining modules and often uses interestingness criteria to direct search efforts toward intriguing patterns. To exclude patterns that have been found, it could apply interestingness levels. Alternatively, depending on how the data mining technique is implemented, the pattern assessment module may be coupled with the mining module. Pushing the assessment of pattern interestingness as far into the mining process as feasible is highly advised for effective data mining in order to limit the search to just the interesting patterns.

This module acts as a conduit for user and data mining system communication. It enables user interaction with the system through specification of a data mining task or query, input of data to help narrow the search, and execution of exploratory data mining based on intermediate data mining results. The user may also examine database and data warehouse schemas or data structures, assess mined patterns, and see the patterns in various ways using this component. The data mining process involves extracting different models, summaries, and derived values from a given set of data. The following phases make up the overall experimental process as it is applied to data-mining issues:

## 1. Outline the issue and create the theory

The majority of data-based modelling investigations are carried out inside a specific application area.Therefore, developing a relevant issue statement often requires domain-specific expertise and knowledge. Unfortunately, a clear issue definition is often sacrificed in favour of a focus on the data-mining approach in many application studies. A modeller often defines a collection of variables for the unidentified dependence in this stage, along with, if practical, a generalised version of this dependency as an initial hypothesis. At this point, many hypothese may have been developed for a single issue. The first phase requires the competence of both a data-mining model and an application area. In actuality, it often entails tight collaboration between the application expert and the data mining specialist. This collaboration continues throughout the whole data-mining process in effective data-mining applications; it does not end in the beginning.

## 2. Gather the data

How the data are created and gathered is the focus of this stage. There are often two different options. The first method is known as a planned experiment and occurs when the data generating process is managed by a professional (modeller). The second choice is the observational strategy, which is used when the expert cannot affect the process of gathering data. The majority of data-mining applications presuppose an observational framework, namely random data production. The sample distribution is often partly and implicitly provided throughout the data collecting process or is wholly unknown after the data are gathered. However, it is crucial to comprehend how data gathering influences its theoretical distribution since this kind of a priori information may be highly helpful for modelling and, eventually, for the findings' ultimate interpretation. Additionally, it's crucial to confirm that the data used for estimating a model and the data used subsequently for testing and implementing a model are from the same, same sampling distribution. If not, the predicted model cannot be properly applied to the findings in a final application.

## 3. Data preprocessing

Data are often "collected" in the observational environment from the existing databases, data warehouses, and data marts. At least two typical activities are generally involved in data preprocessing:

**Outlier identification (and elimination):** An outlier is a data result that is out of the ordinary and inconsistent with the majority of observations. Outliers often occur from measurement mistakes, coding problems, and recording errors. Occasionally, outliers are also just naturally anomalous data. Such unrepresentative samples have a significant impact on the final model. There are two ways to deal with outliers: either a. find them at the preprocessing stage and eventually delete them, or b. create robust modelling techniques that are insensitive to outliers.

**Scaling, encoding, and feature selection:** Data preparation involves a number of processes, including variable scaling and various encoding techniques. For instance, a feature with a range of [0, 1] and another with a range of [100, 1000] may not have the same weights in the applied approach and will thus have a distinct impact on the ultimate outcomes of data mining.For further study, it is advised to scale them and equalise the weight of both attributes. Additionally, application-specific encoding techniques often accomplish dimensionality reduction by offering fewer useful characteristics for later data modelling.

These two categories of preprocessing jobs are only two examples from a broad range of preprocessing activities that may be performed throughout a data-mining process.The various stages of data mining should not be fully separated from the data pretreatment procedures. Together, all actions in the data-mining process might create new and better data sets for next iterations. In general, a good preprocessing strategy incorporates a priori information in the form of application-specific scaling and encoding to offer an appropriate representation for a data-mining tool.

### 4. Calculate the model:

The primary goal in this phase is choosing and using the suitable data-mining technology. This procedure is not simple; in fact, implementation is often based on a number of models, and choosing the appropriate model is a separate operation. The fundamentals of learning and discovering from data are laid forth. Later, certain methods used to carry out a successful learning from data process and to create a suitable model.

### 5. Evaluate the model and make judgments:

Data-mining models should typically aid in decision-making. Therefore, in order for such models to be helpful, they must be interpretable, since people are unlikely to base their judgments on intricate "black-box" models. Keep in mind that the aims of the model's accuracy and the accuracy of its interpretation are somewhat at odds with one another. Simple models are often easier to understand, but they are also less precise. High dimensional models are supposed to help modern data-mining techniques provide very precise findings. The issue of understanding these models, which is likewise crucial, is treated as a distinct process with particular methods for validating the outcomes. Numerous pages of numerical results are not what a user wants. He cannot summaries, interpret, or apply them for sound decision-making because he does not comprehend them.

Systems for Data Mining Classification: The following factors may be used to classify data mining systems:

1. Technology Database Statistics
2. Visualization of Machine Learning Information Science
3. Various Disciplines

Additional classification criteria include:

a) classification based on the kind of datasets that were mined
b) classification based on the kind of knowledge that was mined
c) classification based on the kind of strategies being used
d) Adapted classification methods for various applications
e) classification based on the kind of datasets that were mined

Depending on the kind of databases mined, we may categorise the data mining system. Data models, data kinds, and other factors may be used to categorise database systems.

Additionally, the data mining system may be categorised appropriately. For instance, if the database is categorised based on the data model, the mining system may be relational, transactional, object-relational, or data warehouse-based.

**Classification based on the kind of knowledge:**

The kind of knowledge that was mined may be used to categorise the data mining system. It implies that data mining systems are categorised according to features like:

   i.     Characterization
   ii.    Discrimination
   iii.   Analysis of Association and Correlation
   iv.   ClassificationsPrediction
   v.    Analysis of Clustering Outliers
   vi.   Analysis of Evolution

**Classification based on the kind of strategies being used**

The data mining system may be categorized based on the types of methods used. We may categories these tactics based on the amount of user engagement required or the analytic methodologies used.

**Adapted classification methods for various applications**

The application-adapted data mining system may be categorized. These include the following:

**Key Problems in Data Mining:**

Mining several types of information from databases. - Different users have different needs. Additionally, many users may be interested in various types of information. Data mining must thus cover a wide variety of knowledge finding tasks.

Interactive knowledge mining at different levels of abstraction. - Because it enables users to narrow their search for patterns and provide and modify data mining requests depending on returned findings, the data mining process must be interactive.

Background information may be included to help with the discovery process and to communicate the patterns that are found. Background information may be utilised to represent the patterns found not just succinctly but also at several levels of abstraction.Ad hoc data mining and query languages for data mining. Ad hoc mining job description languages should be connected with data warehouse query languages and enhanced for effective and adaptable data mining.Presenting and visualising the findings of data mining. The patterns must be articulated in high level languages and visual representations after they have been found. The users should have no trouble understanding these representations.

Managing noisy or incomplete data - It is necessary to use data cleaning techniques that can manage noisy or incomplete data while mining the data's regularities. The accuracy of the patterns found will be low if there are no data cleansing measures.Evaluation of the pattern: The problem's interest level is considered. The patterns found should be intriguing since they either reflect existing information or don't provide anything new.Efficiency and scalability of data mining algorithms - Data mining algorithms need to be efficient and scalable in order to efficiently extract information from enormous amounts of data in databases.

Algorithms for parallel, distributed, and incremental mining. The development of parallel and distributed data mining algorithms is driven by issues like the enormous size of databases, the

broad dissemination of data, and the complexity of data mining techniques. These algorithms split the data so that it may be processed in parallel moving forward. The results from the partitions are then combined. By using incremental algorithms, databases are updated without having to re-mine the data.

**Knowledge Discovery in Databases (KDD)**

Some individuals approach data mining in the same way as knowledge discovery, while others see it as a crucial phase in the process. The following is a list of the stages in the knowledge discovery process:

   i.   Data cleaning is the process of removing noise and erroneous data.
   ii.  Data integration is the process of combining data from many sources.
   iii. Data Selection: In this stage, the database is searched for information relevant to the analytical job.
   iv.  Data Transformation: Using summary or aggregation processes, data are converted or consolidated into mining-ready forms in this stage.
   v.   Intelligent techniques are used in this stage to extract data patterns, which is known as data mining.
   vi.  Data patterns are assessed in this stage, which is called pattern evaluation.
   vii. Knowledge Presentation – Knowledge is presented at this stage.

**Data Warehouse:** To assist management's decision-making process, a data warehouse is a subject-oriented, integrated, time-variant, and non-volatile collection of data.

**Topic-Oriented:** An analysis of a certain subject area may be done using a data warehouse. An example of a specific topic is "sales."

**Integrated:** A data warehouse combines information from several sources. For instance, sources A and B may have many methods for identifying a product, but a data warehouse will only have one method for doing so.

**Time-Variant:** A data warehouse houses old data. Over instance, one may access data from a data warehouse for the last three, six, twelve, or even more years. In contrast, a transactions system often just keeps the most current data. For instance, a transaction system could only have a client's most current address, but a data warehouse might include all addresses connected to a certain customer.

**Non-volatile:** Data that has been stored in a data warehouse won't be altered. In a data warehouse, historical data should never be changed.

**Data Warehouse Design Methodologies:**

A data warehouse may be constructed using a top-down methodology, a bottom-up methodology, or a hybrid of the two.The general design and planning are where the top-down method begins. When the technology is established and well-known and the business issues that need to be resolved are well-defined and understood, it may be helpful.Experiments and prototypes are the first steps in the bottom-up method. In the early stages of business modelling and technological development, this is helpful. It enables a company to advance at a much lower cost and to assess the advantages of the technology before making important investments.

In the combined method, an organisation may maintain the quick implementation and opportunistic application of the bottom-up strategy while also taking advantage of the planned and strategic nature of the top-down approach.

*The stages in the warehouse design process are as follows:*

Select a business process to model, such as sales, account management, inventories, orders, invoices, shipments, or the general ledger. A data warehouse model should be used if the business process is organisational and contains several complicated object collections. A data mart model should be used, nevertheless, if the process is departmental and concentrates on the analysis of a single kind of business process.Select the business process's direction of grain. The basic, atomic level of information that must be represented in the fact table for this operation is called the grain; examples include individual transactions, individual daily snapshots, and so forth.

Select the dimensions that will be applied to each entry in the fact table. Time, item, customer, supplier, warehouse, transaction type, and status are examples of typical dimensions.Select the metrics that will be used to fill each fact table entry. Common measurements include numerical additive amounts like dollars and units sold.

**A Three Tier Data Warehouse Architecture:**



**Figure 1.2: Three Tier Data Warehouse Architecture**

Layer-1: A warehouse database server, which is nearly typically a relational database system, makes up the lowest tier. Data from operational databases or other external sources (such as customer profile information given by external consultants) is fed into the bottom tier using back-end tools and utilities. To keep the data warehouse up to date, these tools and utilities also execute data extraction, cleansing, and transformation (such as merging comparable data from many sources into a single format). Utilizing application programme interfaces called gateways, the data are extracted. The underlying DBMS has a gateway that enables client applications to create SQL code that may be run at a server (Figure 1.2).

Microsoft's ODBC (Open Database Connection) and OLEDB (Open Linking and Embedding for Databases) as well as JDBC (Java Database Connection) are examples of gateways.A metadata repository, which houses details on the data warehouse and its contents, is also included at this layer.

A relational OLAP (ROLAP) model or a multidimensional OLAP is commonly used to create an OLAP server in the intermediate tier.An enhanced relational DBMS called an OLAP model translates operations on multidimensional data to conventional relational operations. A server designed specifically for multidimensional data and activities, or a multidimensional OLAP (MOLAP) model.

Tier 3: The top tier consists of a front-end client layer that includes query and reporting tools, analytical tools, and/or data mining capabilities (such as trend analysis, prediction, and so on).

**There are three different data warehouse models.**

**1. Enterprise warehouse:** An enterprise warehouse compiles all of the data on topics across the whole company.It offers cross-functional data integration for the whole company, often from one or more operational systems or outside information sources.

Its size may vary from a few gigabytes to hundreds of gigabytes, terabytes, or beyond, and it often includes both comprehensive and summary data.Supercomputers, parallel architecture platforms, or conventional mainframes may all be used to construct a corporate data warehouse. It may take years to develop and construct and requires substantial business modelling.

**2. Data mart:** For a particular user group, a portion of corporate-wide data is stored in a data mart. The focus is limited to a few carefully chosen topics. A marketing data mart, for instance, may limit its contents to customers, items, and sales. Data marts often include summary versions of their data.

Data marts are often deployed on Windows- or UNIX-based, low-cost departmental servers. Instead of being measured in months or years, the deployment cycle of a data mart is more likely to be measured in weeks. However, if its design and planning weren't enterprise-wide, it may eventually entail difficult integration.Data marts may be characterized as independent or dependent depending on the source of the data. Independent data marts get their data from a variety of operating systems, outside information sources, locally created data inside a certain department or region, or from both. Enterprise data warehouses serve as the primary source of dependent data marts.

**3. Virtual warehouse:** An operational database is seen via a collection of virtual warehouses. Only a small portion of the potential summary views may materialise for effective query processing.

Although creating a virtual warehouse is simple, doing so calls for more space on active database servers.

**Metadata Archive:**

*Data about data are called metadata.*

Metadata are the data that define warehouse objects when utilised in a data warehouse. For the data names and definitions of the specific warehouse, metadata are established. In order to date any extracted data, identify the source of the extracted data, and add any missing fields that were added by data cleaning or integration procedures, additional metadata is produced and recorded.

**Metadata repository is as follows:**

An explanation of the data warehouse's structure, including the data warehouse schema, view, dimensions, hierarchies, and definitions of derived data, as well as the locations and contents of data marts.Data lineage (history of transferred data and the order in which changes were performed to it), currency of data (current, archived, or deleted), and monitoring data are all examples of operational metadata (warehouse usage statistics, error reports, and audit trails). Measure and dimension defining techniques, data on granularity, divisions, topic areas,

aggregation, summarizing, and preset queries and reports are some of the algorithms used for summarization. Source databases and their contents, gateway descriptions, data partitions, data extraction, cleaning, transformation rules and defaults, data refresh and purging rules, and security are all included in the mapping from the operational environment to the data warehouse (user authorization and access control).Data about system performance, including indices, profiles, and schedules for refresh, update, and replication cycles. These data also contain indices and profiles that enhance data access and retrieval performance.Data ownership information, billing policies, and company terminology and definitions are all examples of business metadata.

**OLAP (Online analytical Processing):**

Online analytical processing, or OLAP, is a method for quickly responding to multi-dimensional analytical (MDA) questions.The wider field of business intelligence, which also includes relational databases, report authoring, and data mining, includes OLAP as one of its subcategories.Users of OLAP technologies may interactively evaluate multidimensional data from many angles.

**Three fundamental analytical activities make up OLAP:**

**Drill-Down and Consolidation (Roll-Up)**

**Making Slices and Dice**

Data that may be amassed and calculated in one or more dimensions are combined during consolidation. For instance, all sales offices are gathered in one place at the sales division or sales department to forecast sales trends.

Users may browse through the information using the drill-down method. Users may, for example, see the sales of certain goods that make up a region's sales.

A feature known as "slicing and dicing" allows users to remove (slice) a specified set of data from an OLAP cube and inspect (dice) the slices from various angles.

**OLAP types:**

Relational OLAP (ROLAP) is an OLAP technique that works with relational databases directly. New relational tables are constructed to retain the aggregated data when the base data and dimension tables are saved. It is dependent on a unique schema design.To mimic the slicing and dicing capabilities of conventional OLAP, this technology depends on modifying the data contained in the relational database. Slice and dice operations are essentially identical to adding a "WHERE" clause to a SQL query.

ROLAP tools query the regular relational database and its tables to get the data needed to answer the question rather than using pre-calculated data cubes.

Due to the methodology's lack of restrictions on the contents of a cube, ROLAP tools allow for the asking of any query. ROLAP may also drill down to the database's most basic level of information.

Multidimensional OLAP (MOLAP): MOLAP is frequently referred to as simply OLAP and is the "classic" type of OLAP.Instead of storing this data in a relational database, MOLAP saves it in an efficient multi-dimensional array storage. Processing is thus required, which is the pre-computation and storing of data in the cube.A data cube, which is a pre-calculated data collection, is often used by MOLAP tools.The data cube includes every response that

might be provided to a certain set of questions.MOLAP tools may swiftly write back data into the data set and have a very quick reaction time.

### Hybrid OLAP (HOLAP):

A database will split data between relational and specialised storage, although there is no clear industry consensus on what exactly qualifies as hybrid OLAP.For certain vendors, a HOLAP database could, for the greater amounts of comprehensive data, employ relational tables and, for the lower quantities of more-aggregate or less-detailed data, use specialised storage, at least in part.By integrating the strengths of the two methodologies, HOLAP tackles the drawbacks of MOLAP and ROLAP.

*Both relational data sources and pre-calculated cubes may be used with HOLAP tools.*

### Data Mining Applications

Companies in the retail, communication, financial, and marketing industries, for example, employ data mining extensively. Identify the influence on sales, customer satisfaction, and company profitability of the pricing, consumer preferences, and product positioning. data analysis helps a shop to create items and promotions that will assist the business attract customers by using point-of-sale data of their purchases the client. Data mining is extensively utilized in the following fields:

### Healthcare Data Mining:

The potential for data mining in healthcare to enhance the healthcare system is quite high. For better understanding and to identify, it makes use of data and analytics excellent practices that will raise the quality of medical treatment while cutting expenses. Data mining techniques used by analysts include machine learning, statistics, soft computing, data visualization, and multidimensional databases. Patient forecasting may be done using data mining for each category. The processes make sure that patients get intense care when and where it is needed. Data mining helps with healthcare as well insurance companies to spot fraud and misuse.

### Using data mining to analyze market baskets:

A modelling technique based on a hypothesis is market basket analysis. You are more likely to buy other things if you purchase a certain category of goods. Another collection of items. The shop may be able to comprehend a customer's purchasing habits using this strategy. This information might help merchant in recognizing the needs of the consumer and adjusting the stores layout appropriately. By comparing various analytical methods, Results between different shops and clients from various demographic groups may be compared.

### Education and Data Mining

Education data mining is a freshly developed discipline that focuses on creating methods for discovering information from the data produced from environments that are teaching. EDM goals are acknowledged as confirming students' future learning behavior and researching the effects of support for education and encouragement of scientific study. Using data mining, a company can make accurate judgments and forecast the outcomes of the student. With the findings, the institution may focus on what to teach and how to educate.

### Manufacturing engineering and data mining

The finest resource a manufacturing organization has is knowledge. Data mining technologies might be advantageous to uncover patterns in a complicated production procedure. When

building systems at the system level, data mining may be utilized to determine how product architecture and product clients' portfolio, data, and informational demands. It may also be used to predict the time, cost, and expectations for product development the extra duties.

**CRM (Customer Relationship Management) Data Mining:**

Getting and keeping customers, as well as fostering their loyalty, are the main goals of customer relationship management (CRM) customer-focused tactics. A company organization must gather and evaluate data to have a good connection with its Customers data. The gathered data may be utilized for analytics using data mining methods.

**Using data mining to identify fraud:**

Frauds cause billions of dollars in losses. Traditional fraud detection techniques are quite complex and time-consuming. Data Mining produces informative patterns and transforms data into knowledge.

The data of every user should be protected by an ideal fraud detection system. Users the records used in supervised techniques are divided into two categories: fraudulent and non-fraudulent. This information is used to build a model, and the method is designed to determine whether or not the document is false.

**Lie detection via data mining:**

It's not a huge thing to catch a criminal, but it might be quite difficult to get the truth out of him. Data may be used by law enforcement using mining methods to look into crimes, keep an eye on any terrorist communications, etc.

This method likewise uses text mining, and it examines data, which is often unstructured text, for significant patterns. Comparing the data gathered during the earlier investigations and a lie detection model is created.

**Financial Banking Data Mining**:

With each new transaction, the financial system's digitalization is expected to produce vast amounts of data.

The data analysis method may aid bankers in the resolution of commercial banking and financial issues by finding patterns, fatalities, and correlations in Business data and market expenses that managers or executives are not immediately aware of because the data volume is too vast or are produced on the screen too quickly by professionals. These statistics may be used by the management for more effective targeting, recruiting, retaining, segmenting, and keep a lucrative client.

**Benefits from Data Mining**

A. Data mining makes ensuring a business is gathering and analyzing trustworthy data. It is often a stricter, more organized procedure that formally recognizes a problem, collects information about the issue, and makes an effort to come up with a solution. Consequently, data mining increases the profitability, effectiveness, or operational strength of a company.

B. Applications for data mining might appear quite varied, but the general method can be used to practically any new or old application. Almost any business issue that depends on data may be researched and examined, and almost any form of data can be collected. Data mining may be used to qualify evidence.

Data mining's ultimate objective is to analyze unstructured data to see whether there is any coherence or connection between the data. With the help of data mining, a business may use the information it already has to produce value that would not otherwise be possible. Otherwise not be very obvious. Despite the fact that data models might be complicated, they can also provide amazing outcomes, such hidden patterns and distinct approaches.

--------------------

# CHAPTER 2

## ANALYSIS USING PREDICTIVE DATA MINING

Gadug Sudhansu
Assistant Professor, Department of Computer Science Engineering,
Faculty of Engineering and Technology, JAIN (Deemed-to-be University), Karnataka – 562112
Email Id- s.gadug@jainuniversity.ac.in

Each of the following data mining techniques addresses a different set of business concerns and offers special insights into each of them. However, being aware of the type of business issue you are dealing with will also assist you in selecting the strategy that will work the best and produce the best results. The following are the two basic types of data mining:

- Analysis using Descriptive Data Mining
- Analysis using Predictive Data Mining

Predictive data-mining analysis, as the name suggests, focuses on information that may be used to predict future events in the business world. The following four categories further break down predictive data mining:

- Analysis of Classification
- Analysis of Regression
- Serious Analysis of Time
- Analysis of Prediction

**Classification Analysis in Data Mining**

A key method that contributes to the categorization of data in data mining is classification and must make choices throughout the classification process in order to combine the data and provide the standards for categorizing the data sets. Choosing the input variables is the initial stage in categorization. A variety of acknowledgements and data instances are also necessary for classification.Data mining relies heavily on the classification approach, which categorizes data and predicts group membership for data instances. Make choices throughout the classification process in order to combine the data and provide the standards for categorizing the data sets. Choosing the input variables for a data mining system is the first step in classifying it. A variety of acknowledgements and data instances are also necessary for classification.

➢ **Lifecycle of Data Classification**

The data categorization life cycle provides a great framework for managing the data flow into an organization. Businesses must take each level of compliance and data security into consideration and these are able to carry it out at every step, from creation to deletion, with the use of data categorization. The following phases make up the data life cycle, including:

**Origin:** Sensitive data is produced in a variety of forms, including emails, Word, Excel, Google documents, social media, and websites.

**Role-based practice:** Role-based security limits are applied to all sensitive data by labelling it according to internal protection policies and agreement guidelines.

**Storage:** This is where keep the data have gotten, along with encryption and access restrictions.

**Sharing:** Data is continuously shared between agents, customers, and coworkers via a variety of platforms and devices.

**Archive:** Data is finally archived in the storage systems of an industry in this step.

**Publishing:** Data may reach consumers via publication. The dashboards are then available for viewing and downloading.

**Work of Classification**

The two phases of the data classification process are:

1. Constructing a classifier or model
2. Classifiers Used for Classification

**Construction of the Classifier or Model**

The learning phase or stage is now in progress.

The classifier is constructed in this stage by the classification algorithms.The training set, which consists of database tuples and the corresponding class labels, is used to create the classifier.A category or class is the term used to describe each tuple that makes up the training set. These tuples are also known as data points, samples, or objects.

**Using a Classifier to Classify**

The classifier is used in this stage to do classification. Here, test data are used to calculate the classification rule accuracy estimates. If the accuracy is deemed satisfactory, the categorization rules may be applied to the new data tuples.

➢ **Data classification difficulties**

**Labeling Critical Resources and Assets**

This makes it easier for security experts to identify which data is essential. This in turn motivates the expenditure of money to safeguard the specified assets. Mislabeling is one of the most frequent errors made while labelling data. In order to mitigate this risk, the process should be administered by solid, seasoned, and reliable specialists who should coordinate with all levels of senior management and important process owners to guarantee adequate resource allocation for security.

**Specify and categories**

In the majority of circumstances, securing information that ranks significantly higher on the data categorization scale should get much more resources than protecting information that ranks far lower. One of the many effects of incorrect data labelling is the waste of resources. In order to identify a data categorization system that suits the scale of their organization and their protection demands, security experts should make use of industry best practices.

**Privilege Control**

The idea of "privilege management" always brings to mind the idea of segregation of responsibilities, or the idea of dividing up the tasks involved in a particular process such that no one individual may carry out the full process without supervision (e.g., segregating the task of adding payees to the payroll system and printing checks). When segregation of roles is

not in place within an organization, the risk of fraud and data loss is significantly raised. Every sector is susceptible to the danger of division of tasks, and privilege management should carefully take access to data into account.

**Uphold Compliance**

Compliance and audits are often seen in different ways. Security experts might perceive one form of light as a tool for a third party to assess measures in place around their data. A means by which others may express their judgement on the efficiency of the controls in place regarding data that is often compared to a standard. A system that may have been created by the professional being audited may be criticized by outsiders using the second light, the hostile type. Companies that use the former are better positioned to safeguard their data. As an example, a data breach is more expensive than the perceived inconvenience that bi-annual or yearly audit(s) bring, both financially and to the brand. Data-related federal rules and regulations are intended to benefit, not harm, the organizations that must adhere by them. They provide a uniform approach to safeguarding resources and the data that keeps the metaphorical engine running. As indicated, performing audits of the data protection procedures in place is the best approach to ensure compliance, with special attention paid to essential data at the higher levels of the categorization system. To assist guarantee thorough coverage of the systems or processes being audited, qualified individuals with knowledge of the sector and relevant regulations should be included in the audit process.

**Continuity of Operations and Disaster Recovery**

Safety of the workforce and the continuance of corporate operations come first when a calamity strikes. When attempting to restore crucial activities for the success of a firm, information security experts don't want to be concerned about the protection of crucial data. The idea that the most important data should be secured and kept in separate places, both conceptually and physically, is a key component of data categorization. To assist prevent data loss or leakage in the case of a catastrophe, a formal, tested, and documented Business Continuity and Disaster Recovery Plan should be in place.

**Awareness**

Every level of the organization has to get awareness training on their roles and responsibilities in a business continuity and disaster recovery plan as well as how to manage different circumstances. Although it might be difficult for businesses to justify investing in a robust security awareness program, the advantages have been proved to exceed the disadvantages. It offers the advantage of lowering the possibility that individuals may be responsible for crucial data breaches or fall prey to social engineering.

**Classification of Data Mining Systems**

In addition to this, a data mining system may be categorized according to the kind of

- Mining databases
- Information gleaned
- Methods used
- Application modifications

**Classification Using the Mined Databases**

A data mining system may be categorized based on the types of databases it mines. Data models, data kinds, and other factors may be used to categories database systems.

Additionally, the data mining system may be categorized appropriately. A relational, transactional, object-relational, or data warehouse mining system, for instance, if categories a database according to the data model.

**Using the kind of Knowledge Mined to Classify**

A data mining system may be categorized based on the kind of knowledge it extracts. It implies that the features of the data mining system are used to categories the system.

A. Characterization
B. Discrimination
C. Analysis of Association and Correlation
D. Classification
E. Prediction
F. Analysis of Outliers
G. Analysis of Evolution
H. Using Techniques to Determine Classification

A data mining system may be categorized based on the kind of approaches it employs and may categories these tactics based on the amount of user engagement required or the analytic methodologies used.

**Adapted Classification Based on Applications:**

A data mining system may be categorized based on the applications it has been used for. Here are some applications:

i. Finance
ii. Telecommunications
iii. DNA
iv. Stock exchanges
v. E-mail

**Regression Analysis**

A data mining method known as regression is used to forecast the numerical values of a given data collection. Regression may be used, for instance, to forecast the price of a product or service or other factors. Additionally, it is employed in many different sectors for trend research, financial forecasting, and business and marketing behavior.

**Regression:**

A supervised machine learning method known as regression is used to forecast any characteristic with a continuous value. Any business organization may investigate the connections between the goal variable and the predictor variable using regression. It is a very important tool for analyzing data that may be utilized for time series modelling and financial forecasting. Regression is the process of fitting a line or a curve to a large number of data points. The gap between the data points and the remedy ends up being the smallest because of how things work out. Regressions of the linear and logistic kind are the most common. In addition, a variety of additional forms of regression may be used, depending on how well they work on a particular data set.

Regression can forecast all dependent data sets when independent variables are represented in their expression, and the trend is known for a limited time. Although there are certain limitations and presumptions, such as the independence of the variables and their innate

normal distributions, regression offers a useful method for predicting variables. Assume, for instance, that two variables, A and B, are taken into account and that their combined distribution is a bivariate distribution by nature. In such situation, these two variables may not only be connected but also independent. To be employed, the marginal distributions of A and B must be derived. To make sure the regression is appropriate, the data must be thoroughly examined and subjected to a number of preparatory tests before conducting regression analysis. In certain circumstances, non-parametric testing are available.

**Different Regressions**

Regression comes in a variety of forms, including the following:

Finding the best line to fit two traits (or variables) such that one attribute may be used to predict the other is known as linear regression. A development of linear regression called multiple linear regression involves more than two features and the record being fitted to a multidimensional space.

The least squares approach is used to implement the best fit line in linear regression, which reduces the sum of the squares of the deviations from each data point to the regression line. Because certain deviations are squared, the positive and negative deviations have not cancelled.

1. **Polynomial Regression:** The regression equation is referred to be a polynomial equation if the power of the independent variable is greater than 1.
2. **Logistic Regression** The logistic regression approach is used when the dependent variable is binary in character, such as 0 and 1, true or false, success or failure. As a result, the goal value (Y), which is often utilised for classification-based tasks, has a range of 0 to 1. It does not demand for all independent and dependent variables to have a linear relationship, in contrast to linear regression.
3. **Ridge Regression:** Ridge regression describes a method for computing multiple regression data with multicollinearity issues. The continuance of a linear connection between two different variables is multicollinearity.

Least Absolute Shrinkage and Selection Operator, often known as LASSO, stands for Lasso Regression. Shrinkage is used in the linear regression technique known as lasso regression. A central point, often known as the mean, is shrunk towards in Lasso regression. Compared to other regression methods, the lasso method is best suited for straightforward and sparse models with a variety of parameters. Regression techniques like this one work well for models that suffer from multicollinearity.

**Data mining: Regression vs. Classification**

Regression and classification both use similar principles. Classification and regression are two significant prediction issues in data mining. You should be able to predict outputs from fresh data if you provide training inputs and outputs and develop a function that links the two. The sole difference between the two is that Classification's outputs are discrete whereas Regression's outputs are not.

Certain terminology, such as "Logistic Regression," may be used to describe either a Classification approach or a Regression method. As a consequence, it might be challenging for the user to choose when to use the Classification and Regression Method in Data Mining, as below mention Table 2.1.

**Table 2.1 shows the classification of refression**

| Regression | Classification |
|---|---|
| In Regression, the predicted data types are ordered. | The predicted data is often unordered in classification. |
| Regression may be of two different types: linear and non-linear. | The two forms of classification are binary classifiers and multi-class classifiers. |
| The calculations for the Regression procedure are done using the Root Mean Square Error. | In the classification process, calculations are mostly made by evaluating efficiency. |
| Examples of regression include regression trees and linear regression. | A classification example is the decision tree. |

**Time Series Analysis**

A particular method of examining a set of data points gathered over a period of time is called a "time series analysis." Instead of merely capturing the data points sporadically or arbitrarily, time series analyzers capture the data points at regular intervals over a predetermined length of time.

But this kind of study involves more than just gathering data over time. To put it another way, time is a key variable since it both reveals how the data changes throughout the duration of the data points and the outcomes. It offers an extra source of data as well as a predetermined sequence of data dependencies.To maintain consistency and dependability, time series analysis often needs a lot of data. A large data collection guarantees that your analysis can sift through erratic data and that your sample size is representative. Additionally, it guarantees that any trends or patterns are not outliers and can take seasonal variation into consideration. Time series data may also be utilised for forecasting, which is the process of making predictions about the future based on the past.

**Time-Series Movements by Category:**

**Trend movements over the long term:**

The broad direction of a time series' movement over an extended period of time. It demonstrates the overall propensity of the data to rise or fall over an extended period of time.

**Cycles or variants on cycles:**

Oscillations that are long-term along a trend line or curve. Think about business cycles. The oscillation time for this movement is more than a year.

**Seasonal alterations or movements:**

Patterns that a time series seems to follow throughout comparable months of next years that are almost similar. If the data are collected hourly, daily, weekly, or monthly, this variance will be visible in the time series.

**Irregular or erratic motions:**

These changes occur unexpectedly, without control, and without warning. They are merely random or irregular variations, not regular variations.

**Uses for time series analysis:**

**Time Series in the Business and Financial Sector**

    A. The majority of business, investment, and financial choices are based on projections of expected financial market developments and demand.

    B. The dynamic and influential behavior of financial markets may be explained via time series analysis and forecasting. An expert can foresee the necessary forecasts via the analysis of financial data for critical financial applications in a variety of domains, including risk evolution, option pricing & trading, portfolio design, etc.

    C. Time series analysis, which may be used to forecast interest rates, foreign exchange risk, stock market volatility, and many other things, has evolved into an integral aspect of financial analysis. Financial forecasting is used by policymakers and business professionals to decide on production, purchasing, market sustainability, resource allocation, etc.

**Medical-related time series**

A data-driven industry, medicine has developed and is still making significant advances in time series analysis of human knowledge.

**Research case**

    A. Consider the scenario where time series data are combined with the medical technique CBR (case-based reasoning) and data mining. These synergies are crucial as the pre-processing for feature mining from time series data and may be helpful to examine the progression of patients over time.

    B. In the field of medicine, it is crucial to look at how behavior changes over time rather than drawing conclusions based just on the time series' absolute values. The typical demonstration of linking time series with case-based monitoring is to identify heart rate variability in conjunction with respiration based on the sensor values.

    C. However, time series in the context of the epidemiology domain have only lately and slowly evolved since methods to time series analysis necessitate recordkeeping systems so that data should be linked through time and gathered accurately at regular intervals.

    D. Healthcare applications utilizing time series analysis have produced significant prognostication for the industry as well as for individual patients' health diagnoses after the government has installed enough scientific devices to collect excellent and long temporal data.

**Medical Equipment**

Time series analysis has entered medicine with the development of technologies such ECGs, which were first developed in 1901, are used to diagnose cardiac disorders by capturing the electrical pulses that go through the heart.

Developed in 1924, the electroencephalogram (EEG) measures electrical activity and brain impulses. Medical professionals now have additional chances to use time series for medical diagnostics because to these advancements.

As a consequence of the development of wearable sensors and smart electronic healthcare equipment, people may now take routine measures automatically and with little input, leading to a reliable collection of longitudinal medical data for both ill and healthy people.

**Series Time in Astronomy**

Different fields of astronomy and astrophysics are among the present and modern applications where time series plays a key role, Astronomical specialists are skilled in time series for calibrating devices and researching things of their interest since astronomy, being specialized in its field, heavily depends on graphing objects, trajectories, and exact measurements.

i.   Time series data has a long history in the astronomy field; for instance, sunspot time series were documented in China in 800 BC, making sunspot data collecting as well-recorded natural occurrences. Time series data have an essential influence on understanding and quantifying everything about the cosmos. Similarly, time series analysis was used in earlier eras.

ii.  The search for variable stars used to estimate stellar distances, and to comprehend the process through which the cosmos changes over time, one must examine transient occurrences like supernovae.

iii. These systems, which rely on the wavelengths and light intensities of light to transmit time series data in real time, enable astronomers to see phenomena as they happen.

iv.  Astroinformatics and astrostatistics are new fields of study that have emerged in recent decades as a result of data-driven astronomy; these paradigms integrate key fields including statistics, data mining, machine learning, and artificial intelligence. Here, time series analysis would play a role in the quick detection and classification of astronomical objects as well as the independent characterization of unique events.

**Time Series in Weather Prediction**

i.   Aristotle, a Greek philosopher, conducted study on weather events in antiquity with the goal of determining the origins and consequences of weather variations. Later, scientists began to collect weather-related data, recording it on an hourly or daily basis and storing it in various places, using the instrument "barometer" to calculate the status of atmospheric conditions.

ii.  Newspapers started printing personalized weather predictions throughout time, and as technology developed, forecasts eventually went beyond just basic weather conditions.

iii. Many countries have set up tens of thousands of weather forecasting stations all around the globe in order to undertake atmospheric observations using computer techniques for quick compilations.

iv.  These stations are outfitted with highly functioning equipment and are linked to one another in order to gather meteorological data at various areas and predict weather conditions at all times as needed.

**Business Development Time Series**

As the process examines prior data trends, time series forecasting assists organizations in making wise business choices. It may be helpful in predicting future possibilities and occurrences in the following ways.

a) **Reliability:** Time series forecasting is very trustworthy when the data has a wide range of time intervals in the form of many observations across a longer time horizon. By using data observations at varied time periods, it delivers illuminating information.

b) **Growth:** Time series is the best asset to use while evaluating endogenous as well as overall financial performance and growth. Endogenous growth is essentially the improvement of internal human capital inside firms that leads to economic

development. Time series forecasting, for instance, may be used to analyze the effects of any policy variable.

c) **Estimating trends**: It may be done using time series approaches. For instance, these methods look at data observations to see when sales of a certain product are increasing or decreasing.

d) **Seasonal patterns:** Variations in recorded data points may reveal seasonal patterns and variations, which serve as the foundation for data forecasting. The information gathered is important for markets whose goods vary seasonally and helps businesses manage their product development and delivery needs.

## Prediction in Data Mining

Prediction is the second approach to use data mining. It is often used to find various facts. Similar to categorization, the data set's behavior preserves inputs and comparable numerical output values. The algorithm (division) obtains the model or a prediction that is consistent with the dataset's behavior. A numerical output should be detected by the model when the new information is provided. Despite the categorization, neither the class name nor the notes are present in this method. The current valued action or command value is estimated by the model. Most often, regression (or growth) is employed for prediction. Predicting the cost of a home based on variables like the number of apartments, the size of the location, and other factors is an example of prediction. An organization has the authority to determine how much cash was exchanged during a discussion.

A. Identifies items in a datasheet that are unknown or missing.
B. The categorization model was developed to forecast the results.
C. Is independent of class label.

To support hypotheses about potential future volumes of transactions, the predictive data mining method, for instance, may utilize algorithm-based tools to search through a client database and look at previous transactions. In other words, the data may aid in predicting future business events, enabling company executives to make appropriate plans. In other words, the data may aid in predicting future business events, enabling company executives to make appropriate plans.

## Techniques

### A number of linear regressions

This technique is used to a dataset to forecast the response variable based on the predictor variable or to investigate the link between a response and predictor variable, such as the correlation between student test scores and demographic data like income and parental education.

### K Nearest Neighbor

This prediction approach splits a training dataset into groups of k observations using a Euclidean Distance metric to assess similarity between "neighbours," much as the classification method with the same name above. The response value for each participant in the validation set is predicted using these groupings.

### Tree of Regression

A regression tree is a kind of decision tree that is used to approximate real-valued functions rather than for classification purposes. Like other regression algorithms, XLMiner makes the assumption that there are one or more input (predictor) variables and one or more output

(response) variables. The variable for the output is numerical. Input variables might be a combination of continuous and categorical variables when using the generic regression tree construction approach. When each decision node in the tree performs a test on the value of an input variable, a decision tree is created. The expected output variable values are included in the tree's terminal nodes.

**Network neural**

Artificial neural networks are based on how the human brain functions and is organised. These networks analyse one record at a time and "learn" by contrasting their initial, essentially arbitrary prediction of the record with the actual value of the response variable that is known. In order to change the network's algorithm a second time, errors from the original forecast of the first data are put back into the system. This keeps happening over and over again.

**Predictive Analytics Techniques**

Based on the data acquired by predictive analytics, several businesses have used these tactics to boost turnover, meet objectives, and increase profits.

1. Coordinating supply and demand.
2. Fraud avoidance.
3. Establishing long-term inventory.
4. Client satisfaction
5. Choosing the right pricing point to maximize profit.

**Prediction problems**

The hardest task is getting the data ready for prediction. The preparation of data involves the following tasks:

**Data cleaning:** Reducing noise and handling missing values are two aspects of data cleaning. Using smoothing methods, noise is reduced, and the issue of missing data is resolved by substituting the missing value with the value that occurs most often for that feature.

**Relevance Analysis:** The irrelevant qualities may also be contained in the database. In order to ascertain if two qualities are related, the correlation analysis approach is performed.Any of the techniques described below may be used for data transformation and reduction.

**Normalization:** The data are transformed via normalization. By narrowing the range of all values for a specific characteristic, a procedure known as normalization is used. Normalization is carried out when neural networks or other techniques requiring measurements are used in the learning process.

**Generalization:** The information may also be changed by being combined with a bigger concept. Hierarchies are a useful notion for this.

**Comparison of Methods of Classification and Prediction**

Here are the standards used to compare classification and prediction techniques:

**Accuracy:** Classifier ability is referred to as the classifier's accuracy. It accurately predicts the class label, and predictor accuracy describes how effectively a particular predictor can make an educated estimate about the value of a predicted characteristic for fresh data.

**Speed:** This is the amount of computing required to create and use the classifier or predictor.

**Robustness:** This term describes a classifier's or predictor's capacity to provide accurate predictions from noisy input data.

**Scalability:** The capacity to efficiently build a classifier or predictor in the presence of a vast quantity of data.

**-------------------**

# CHAPTER 3

# ANALYSIS USING DESCRIPTIVE DATA MINING

Vanitha K

Assistant Professor, Department of Computer Science Engineering, Faculty of Engineering and Technology, JAIN (Deemed-to-be University), Karnataka – 562112
Email Id- k.vanitha@jainuniversity.ac.in

The primary objective of jobs involving descriptive data mining is to condense or transform the data into relevant information. The following four kinds may be used to further categories the descriptive data-mining tasks:

    i.     Analysis of clusters
    ii.    Conclusion analysis
    iii.   Analysis of association rules
    iv.   Analysis of sequence discovery
    v.    Analysis of Cluster

Data mining use this technique to create pertinent item clusters with related features. Most people have a problem with classification, but if they completely understand how both of these systems actually work, they won't have any issues. In contrast to categorization, which places objects into predetermined classifications, clustering places things in classes that it specifies. Consider the following example to help you understand it better:

Consider yourself at a library with a wide selection of books. Your current objective is to arrange these books so that users may quickly find books on any particular topic. Consequently, in this situation, clustering may be used to put related books on the same shelf, after which the shelves can be given titles or classes that make sense. So whenever a reader is looking for books on a particular topic, they only need to visit that shelf. As a result, he won't need to search through the entire library to find the book he wants to read.

**Data mining clustering**

A collection of various data items is categorised as comparable things by clustering. A single group is a collection of data. Based on how closely related the data are, cluster analysis divides data sets into several groupings. Following the grouping of the data into different categories, each group is given a label. The classifying process aids in adjusting to change.

Therefore, if we were to describe clustering in data mining, we might state that the process essentially entails grouping related abstract objects into groups. Cluster analysis is the method used to divide them into these categories and store them there.

**Data mining using Cluster Analysis**

In data mining, cluster analysis refers to the process of identifying groupings of items that are similar to one another within the group but distinct from the objects inside other groups. Data sets are split into groups or classes during the clustering process in data analytics depending on how similar the data sets are. Following that, the data types associated with each of these classes are labelled. You may better comprehend the analysis by going through the clustering in data mining example.

**The Uses of Data Mining Group Analysis**

Data clustering analysis has several applications, including image processing, data analysis, pattern identification, market research, and a lot more. Companies may find new categories in their customer database by using data clustering. According on the trends of their customers' purchases, they may also divide the current client base into different categories. Data classification might also be done in accordance with purchase trends.

In the study of biology, taxonomy, or the categorization of creatures using cluster analysis, is quite frequent. Our grasp of some of the most often seen intrinsic structures of certain populations or species may be gained by understanding how clustering can assist locate and group species with comparable genetic traits and capabilities. Data mining clustering is used to identify areas. Lands that are similar to one another are found in the earth observation database.A group of homes in the city is characterised based on location, value, and housing type. By categorising the files on the internet, clustering in data mining aids in the finding of information. Applications for detection also make advantage of it. By analysing the pattern of fraud, clustering in data mining makes it simple to find credit card fraud. Cluster analysis may be used as a method to assist someone obtain understanding of the data clusters if they wished to notice the features of each data cluster.

It aids in comprehending each cluster's attributes. It is possible to comprehend the distribution of the data, and it serves as a tool for data mining.

**Clustering in Data Mining Requirements**

**Interpretability**

Clustering should provide useable, comprehendible, and interpretable results. In data analytics, clustering is mostly used to put chaotic data in groups according to how similar their characteristics are.

**Aids in handling corrupted data**

The data is often chaotic and unorganised. As a result, it cannot be swiftly examined, which is why the clustering of data in data mining is so important. By grouping together related data elements, grouping may give the data some structure.The data specialist finds it easier to digest the data and learn new things as a result. It is significantly simpler to analyse data that has already been clustered and labelled than it is to analyse unstructured data. Additionally, it reduces space for mistake.

**High Dimensional**

Both data with large dimensions and tiny amounts of data may be handled via data clustering. Any dimension of data must be able to be handled by the clustering methods used in data mining.

**Clusters of attribute shape are found**

Data mining clustering algorithms should be able to recognise clusters of any form. These techniques shouldn't be constrained to merely locate compact, spherical clusters.

*Various classifications are not regarded as cluster analyses.*

**1. Graph Partitioning -** Cluster analysis is not the kind of categorization in which regions are not the same and are only grouped together based on mutual synergy and significance.

**2. Search results -** In this sort of categorization, groupings are made based on the information provided by other sources. It doesn't qualify as a cluster analysis.

**3. Simple Segmentation -** Cluster Analysis does not include the division of names into distinct registration groups based on last names.

**4. Supervised Classification -** Cluster analysis cannot be stated to be a sort of classification that uses label information since cluster analysis includes grouping based on pattern.

**Clustering techniques for data mining**

The following is an explanation of the many clustering techniques used in data mining (Figure 3.1):



**Figure 3.1: Clustering techniques for data mining**

Data mining clustering (CLUSTER ANALYSIS)

- A. A method based on partitioning
- B. Method based on density
- C. Method based on centroid
- D. A Hierarchical Approach
- E. Grid-Based Approach
- F. Method Based on Models

**A method based on partitioning**

Data is partitioned into several subgroups using the method.Assume for the moment that the database's n items and data are divided up into a division using the partitioning procedure. As a result, question n will be used to represent each segment.

**Dividend Clustering**

This shows that there must be exactly one item in each group and that each object must belong to precisely one group.

**A method based on density**

Based on the dense population of data set members, these algorithms generate clusters at certain locations. For group members in clusters, it aggregates some range concept to a density standard level. Such methods are less effective in identifying the group's Surface regions.

**A method based on centroid**

In this kind of OS grouping mechanism, practically every cluster is referenced by a vector of values. Each item belongs to the group with the least amount of value variation when compared to other groups. The most important algorithmic issue of this kind is the need for set limits on the number of groups. This strategy is popular for solving optimization issues and is the most closely related to the topic of identification.

**Fourth-Level Method**

A given collection of data items will result in a hierarchical decomposition thanks to the approach. We may categorise hierarchical approaches based on how the hierarchical breakdown is created. The following procedure is shown.

**Approaches: Agglomerative and Divisive**

Button-up Approach is another name for the aggregative approach. Here, we start with each item that makes up a distinct group. It keeps fusing pieces or groupings that are close together.

The Top-Down Approach is another name for the Divisive Approach. All of the items in the same cluster are where we start. This approach is rigorous in that it cannot be reversed after a fusion or division has taken place.

**Grid-Based Approach**

Instead of splitting the data into a grid, grid-based approaches operate in the object space. Based on the features of the data, the grid is separated. This approach makes managing non-numeric data simple. The grid's division is unaffected by the order of the data. A grid-based paradigm has the significant benefit of speedier execution.

**The following are some benefits of hierarchical clustering.**

**Model-Based Approach**

This technique makes use of a hypothetical probability distribution-based model. The clusters are found using this approach, which clusters the density function. It displays the geographic distribution of the data points.

**Clustering in Data Mining: Application**

In many domains, including biology, plant and animal classification by attributes, and marketing, clustering may be helpful. It can be used to find clients with similar behaviour in a given customer record. The clustering analysis is widely utilised in a variety of applications, including market research, pattern recognition, data and image processing. Clustering may assist advertising in identifying various client segments. And the purchasing habits of their customers may be identified.

It helps biologists classify genes with similar functions and gain insight into population-inherent structures by establishing plant and animal taxonomies. Finding sections of the land with comparable uses is also made simpler by clustering in an earth observation database. Groups of homes and flats may be distinguished by home type, price, and location. Finding information on the internet is also made easier by the grouping of materials. As a data mining function, the cluster analysis is a tool for obtaining understanding of the distribution of data and observing the characteristics of each cluster.

**Conclusion Analysis**

A group (or set) of data is stored via the Summarization analysis in a more condensed and understandable format. A simple illustration can enable us to grasp it: In order to generate graphs or determine averages from a particular collection (or group) of data, you could have utilized summarization. One of the most well-known and accessible types of data mining is this one.

**Analysis of Association Rules**

It might be viewed as a strategy that, in general, can help us find some fascinating connections (dependency modelling) between different variables in large datasets. By employing this method to find specific concealed patterns in the data, we may be able to identify the elements in the data. It also helps to identify the coexistence of numerous variables that frequently occur together in the dataset. It is common practice to analyses and forecast customer behavior using association rules. It is also actually recommended in the retail industry analysis. This approach is also used to create product clusters, organize store layout, create catalogues, and analyze the data from shopping carts. IT programmers can create software that can employ machine learning by using association rules. In other words, this data mining strategy helps uncover connections amongst two or more elements. This technique uncovers the hidden pattern in the data set. Large volumes of data are analysed using association rule mining to uncover intriguing linkages and relationships. This rule displays the number of times an itemset appears in a transaction. A market-based analysis serves as a common illustration. One of the most important methods used by big organizations to demonstrate correlations between goods is market-based analysis. It enables merchants to discover connections between the products that customers usually purchase together.

**Analysis of Sequence Discovery**

Finding fascinating patterns in data is the fundamental goal of sequence discovery analysis, which uses a subjective or objective criterion for how attractive a pattern is. The issue at hand is frequently identifying frequent sequential patterns in connection to a frequency support measure. Some people frequently confuse the Sequence discovery analysis and the Time series analysis with time series since they both incorporate close-by, order-dependent observations. However, if people take a closer look at both of them, they may easily avoid the misconception because Sequence discovery analysis employs discrete values or data, whereas the Time series analysis method involves numerical data.

**--------------------------**

# CHAPTER 4

# DATA MINING ARCHITECTURE

Chandramma
Assistant Professor, Department of Computer Science Engineering,
Faculty of Engineering and Technology, JAIN (Deemed-to-be University), Karnataka – 562112
Email Id- r.chandramma@jainuniversity.ac.in

A data source, data mining engine, data warehouse server, the pattern assessment module, graphical user interface, and knowledge base are the key elements of data mining systems (Figure 4.1).

**Data sources:** It include databases, the World Wide Web (WWW), and data warehouses. These sources may provide data in the form of spreadsheets, plain text, or other types of media like pictures or videos. One of the largest sources of data is the WWW.

**Database Server:** The database server houses the real, processed data. According to the user's request, it does data retrieval tasks.

**Data Mining Engine:** One of the key elements of the data mining architecture is the data mining engine, which executes various data mining operations including association, classification, characterization, clustering, prediction, etc.

**Modules for Pattern Evaluation:** These modules are in charge of spotting interesting patterns in data, and sometimes they also work with database servers to fulfil user requests.

**Graphic User Interface:** Because the user cannot completely comprehend the intricacy of the data mining process, a graphical user interface enables efficient user-data mining system communication.

**Knowledge Base:** A crucial component of the data mining engine, Knowledge Base is very helpful in directing the search for outcome patterns. The knowledge base may sometimes provide input to data mining tools. The information in this knowledge base could come from user experiences. The knowledge base's goal is to improve the reliability and accuracy of the outcome.
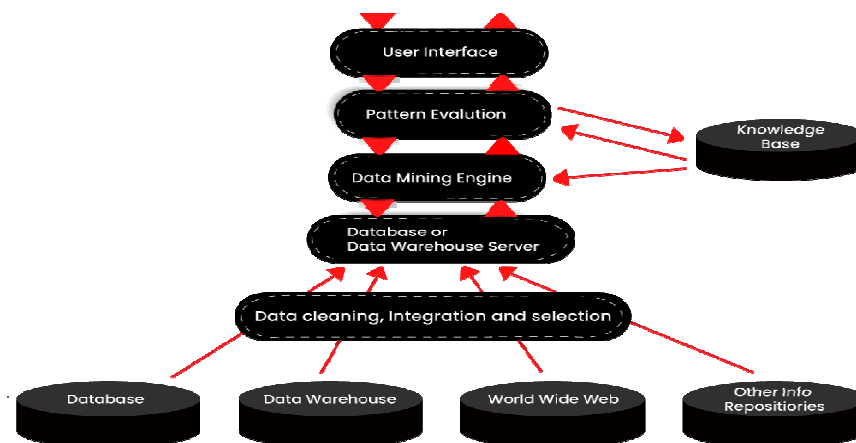


**Figure 4.1: Data Mining Architecture**

### Data Mining Architecture Types

Through commercial transactions, data is gathered and stored in relational database systems. These business procedures were also developed to provide analytical reports. Users in the corporate world should decide on it. Therefore, a decoupled data mining system should be implemented.

There are four potential architectural solutions to this query:

### Data Mining With No Coupling

The data mining system under this design does not make use of any database features. Data from specific data sources are retrieved through a no-coupling data mining method. A database is not used in any way by the no-coupling data mining architecture. That already organizes, stores, accesses, and retrieves data extremely well. No-coupling architecture is regarded as a subpar design for systems that use data mining. But it is used in simple data mining procedures.

### Mining Loose Coupling Data

A database is used in this design by the data mining system to retrieve data. Data mining system extracts information from a database in a loosely coupled data mining architecture. Additionally, it saves the outcome in such systems. The architecture of data mining is for memory-based data mining systems. That does not necessarily need great performance and scalability.

### Mining of Semi-Tight Couplings

The data mining system makes use of a number of data warehousing system capabilities in semi-tight coupling. To carry out certain data mining operations, that is. This comprises aggregation, indexing, and sorting. For improved performance, an interim result from this may be saved in a database.

### Data Mining With Tight Coupling

A data warehouse is regarded as an information retrieval component in tight coupling. Data mining jobs are carried out using the whole of a database's or data warehouse's features. This design delivers system scalability, high performance, and integrated information.

The tight-coupling data mining architecture has three levels:

**Info Layer:** Data layer may be referred to as a database or data warehouse system. All data sources communicate with this layer.

The data layer stores the findings of data mining. In order to show to the end user, we may do it via reports or another kind of visualization.

**An Application Layer for Data Mining:** It is used to get information out of a database. Here, some kind of transformation process is required. Data must be transformed into the proper format to do this. The next step is to process the data using different data mining methods.

**The First Layer:** It offers the user-friendly and straightforward user interface. To engage with a data mining system is to do that. The user is shown data mining results in visualization form at the front-end layer.

**Various Data Mining Elements**

1. Data mining is a crucial technique that allows for the extraction of unknown but possibly relevant information from a vast quantity of data. A data mining system structure is made up of the many components that make up the data mining process. The following are some of the main elements of data mining:
2. One or more databases, data warehouses, spreadsheets, or other forms of data repositories make up an information repository. Techniques for data integration and cleansing may be used to the data.
3. Database or data warehouse server: Based on the user's data mining request, the database or data warehouse server is responsible for retrieving the relevant data.
4. Knowledge base: This is the body of subject-matter expertise that may direct a search or gauge how intriguing a design will be as a consequence.
5. The key component of the data mining system is the data mining engine, which consists of a number of functional modules for tasks including characterization, association and correlation analysis, classification, prediction, cluster analysis, outlier analysis, and evolution analysis.
6. In order to narrow the search to interesting designs, the pattern assessment module often uses interestingness metrics and interacts with the data mining framework.

**Clustering in Data Mining**

**Data Mining Employs Clustering:**

Due to the wide range of applications for clustering analysis, this issue has been developing in data mining. The popularity of these methods must be influenced by the introduction of several data clustering tools in recent years and their extensive usage in a wide variety of applications, including image processing, computational biology, mobile communication, medicine, and economics.

The fact that the techniques for clustering data cannot be standardized is their fundamental drawback. With certain data sets, the advanced method could provide the greatest results, but with others, it might not work at all or perform badly. No notable progress has been made so far despite several attempts to standardize algorithms that can function adequately in all circumstances.

There have already been several clustering tools suggested. Each method has benefits and drawbacks, however, and isn't applicable in all cases.

**Scalability:**

Scalability in clustering indicates that as increase the number of data items, the time it takes to accomplish clustering should roughly scale to the algorithm's complexity order. For instance, know that K-means clustering takes on time, where n is the total number of objects in the data.

The time needed to cluster data items should roughly rise ten times if the number of data objects is multiplied by 10. It implies that there ought to be a linear connection. If that is not the case, then there is a problem with how are doing it.

Data must be scalable; else, cannot get the desired outcome.

**Interpretability:**

Clustering results need to be understandable, useful, and able to be interpreted.

**Finding clusters with the shape attribute:**

The clustering method should be able to locate clusters of any form. They shouldn't be restricted to merely measures of distance, which have a propensity to find a spherical cluster of tiny sizes.

**Adaptability to many sorts of attributes:**

Any kind of data, including numeric interval data, binary data, and categorical data, should be able to be processed by algorithms.

**The capacity to handle noisy data**

Data that is erratic, incomplete, or false may be found in databases. There are few algorithms that are sensitive to such input and may produce clusters of low quality.

**Superior dimensionality**

In addition to handling high-dimensional data, clustering technologies should also be able to handle low-dimensional data.

**Cluster analysis applications**

 A. It has a broad range of uses, including data analysis, pattern recognition, and image processing.
 B. Marketing professionals may define their client groups using their purchase behaviours and use this information to identify separate groups within their consumer base.
 C. It may be used to biology by determining the taxonomies of plants and animals and discovering genes with similar properties.

By categorizing web-based materials, it also aids in the finding of information.

**Clustering Approaches:**

According to the following categories, it may be categorized

 a) Model-Based techniques
 b) A hierarchy of methods
 c) Method Based on Constraints
 d) Grid-Based Technique
 e) Partitioning Technique
 f) Density-Based Approach

**Partitioning Technique**

Consider a scenario in which the partitioning technique creates 'k' partitions of data from a database of 'n' objects. Every partition will stand in for a cluster and k n. The data will be divided into k groups that meet the conditions listed below.

 1. There is at least one item in each group.
 2. Exactly one group must contain each item.
 3. Remedies to keep in mind
 4. The partitioning technique will provide an initial partitioning for k partitions, for example.
 5. The partitioning is then improved by transferring items from one group to another using the iterative relocation approach.

**Model-based techniques**

To determine which model fits the data the best, a model is proposed for each cluster in this procedure. By grouping the density function, this approach finds the clusters. It displays the data points' spatial distribution. With the use of basic statistics and accounting for outliers or noise, this approach also offers a way to automatically calculate the number of clusters. As a result, reliable clustering approaches are produced.

**Limitation-based Approach**

By incorporating user- or application-oriented requirements, this approach does clustering. User expectations or anticipated clustering results' desirable qualities are examples of constraints and have a direct line of contact with the clustering process thanks to constraints. Users or application requirements may specify constraints.

**Hierarchical Techniques**

The specified collection of data items is divided into a hierarchy using this procedure. According on how the hierarchical breakdown is created, hierarchical approaches may be categorized. There are two methods used here.

**Agglomerative Method**

The bottom-up strategy is another name for this one. In this, each item forms a different group at the beginning. The nearby items or groupings keep becoming merged together. This continues until either the termination condition is satisfied or all of the groups are combined into one.

**Disparaging Method**

The top-down strategy is another name for this one. In this case, begin with every item in the same cluster. A cluster is broken up into smaller clusters throughout the continuous iteration. Up until all objects in one cluster are present or the termination condition is met. This procedure is rigorous in that once merging or splitting has been carried out, it cannot be undone.

**Methods to Enhance the Quality of Hierarchical Clustering**

Here are the two methods utilized to raise the caliber of hierarchical clustering:

i.   At each hierarchical partitioning, carefully analyze the object links.
ii.  Integrate hierarchical agglomeration by creating micro-clusters of items using a hierarchical agglomerative algorithm, followed by macro-clustering on the micro-clusters.

**Density-based Approach**

This approach is grounded on the idea of density. The fundamental principle is to keep expanding a given cluster as long as the density in the area reaches a certain threshold, which means that for each data point within a particular cluster, the radius of that cluster must include at least a certain number of points.

**Grid-based Technique**

In this, the things are arranged in a grid. The object space is quantized into a grid-like structure comprising a limited number of cells.

**Advantages**

    A. Quick processing time is this method's main benefit.
    B. It is solely reliant on the quantity of cells present in each dimension of the quantized space.

**Clustering requirements for data mining**

The following are some justifications for clustering's significance in data mining.

**Scalability:** In order to operate on big datasets, clustering algorithms must be extremely scalable.

**Capacity to handle various qualities**: Algorithms should be able to handle many types of data, including category, binary, and numerical data.

**Finding clusters using the shape attribute:** The method should be able to find clusters of any form and should not be constrained by distance measurements.

**Interpretability**: The findings must be complete, applicable, and easy to understand.

**High dimensionality**: The method should be able to deal with high-dimensional space as well as low-dimensional data.

**Data Normalization in Data Mining**

When dealing with qualities on multiple scales, normalization is usually necessary; otherwise, the efficacy of a significant and equally important attribute may be diluted since other attributes have values on a greater scale. This means that while executing data mining activities, it may result in bad data models when numerous characteristics exist but have values on various scales. So they are normalized to bring all the attributes on the same scale. In order to put all the qualities on the same scale, they are normalized.

**Requirement of Normalization in Data mining:**

The fundamental purpose of data normalization is to reduce or eliminate duplicate data. Data duplication is a serious problem. This is due to the fact that maintaining similar data in several locations while storing it in relational databases is becoming more difficult. Data normalization is a helpful process in data mining since it makes it possible to get the following benefits:

    A. Applying data mining methods to a collection of normalized data is much simpler.
    B. Data mining methods used on a collection of normalized data provide more precise and efficient results.
    C. The extraction of data from databases becomes much quicker when the data has been standardized.
    D. Data that has been normalized may be analyzed using more specialized techniques.

**Advantages of Data Mining Normalization:**

Data mining normalization techniques are beneficial because they make it possible to achieve the following advantages:

    A. It is significantly simpler to use normalization techniques in data mining on a set of normalized data.
    B. When data mining normalization techniques are used on a set of normalized data, the results are more precise and efficient.

C. After the data has been standardized, database data extraction is substantially quicker.
D. More specific data analysis techniques might be used to normalize data.

**Data mining Normalization Techniques:**

The main three data normalizing methods used in data mining—Z-score normalization, min-max normalization, and decimal scaling normalization—will be covered in this article. The following list of data mining normalization techniques includes.

**Normalization of Z-score:**

One Data Mining normalization technique that quantifies how far a data point deviates from the mean is the Z-Score value. It figures out if the standard deviations are higher or lower than the mean. It might range from -3 to +3 standard deviations. For data analysis that calls for comparing a value to a mean (average) value, such as test or survey results, Z-score normalization procedures are advantageous.

**Normalization to Min and Max:**

The difference between 500 and 1000000 or the difference between 0.5 and 1 is easier to understand. When the range between the lowest and highest numbers is less, the data is easier to grasp. Using the min-max normalization approach, a dataset is transformed into a scale from 0 to 1. In this data normalization process, the original data is linearly altered. The formula below is used to alter each value after retrieving the lowest and maximum values from the data.

The following formula is used:

$$\frac{(v - \min X)}{(\max X - \min X)} * (\text{new\_max } X - \text{new\_min } X) + \text{new\_min } X$$

Where:

A. The attribute data is X.
B. The absolute values of X's minimum and maximum are denoted by Min (X) and Max (X).
C. v' is the updated value of every item in the data.
D. v represents the previous value of each data item.
E. The maximum and minimum values of the range are new max (X) and new min(X).

**Normalization of Decimal Scaling:**

Decimal scaling is another method of normalization used in data mining. To operate, an integer is rounded to the next decimal place.

By moving the decimal point of the integers, it normalizes the data. Using this method, normalize the data by dividing each data value by the biggest absolute value of the data. The algorithm below normalizes the data value vi to vi'.

$$V' = \frac{k}{10^z}$$

Where:

- After decimal scaling, v' is the updated value.
- k denotes the value of the attribute.
- Integer z now controls the movement of the decimal point.

**Benefits of using data mining normalization:**

    a.  The use of data mining methods is made simpler
    b.  The efficiency of data mining techniques increases.
    c.  The information is transformed into a form that can be understood by everyone.
    d.  Faster databases may be used to extract data
    e.  The data may be analyzed in a certain way.

A technique for arranging data across several linked databases is called data normalization. It enables tables to be changed in order to get rid of duplicate data and undesirable traits including insertion, update, and deletion anomalies. Data mining normalization methods are multi-stage processes that organize data into tables while deleting redundant information from relational databases. It is essential because, if the dataset is not normalized while being vast and full of useful features, one of the features can outperform the others. Data mining normalization techniques are used to resolve this problem.

**--------------------**

# CHAPTER 5

# DATA PREPROCESSING

Mohammed Zabeeulla A N
Assistant Professor, Department of Computer Science Engineering,
Faculty of Engineering and Technology, JAIN (Deemed-to-be University), Karnataka – 562112
Email Id- z.mohammed@jainuniversity.ac.in

One of the most important internal aspects of the well-known knowledge innovation from the data processor is data preprocessing in data mining speeches. Data that was instantly extracted from the source may likely have mistaken, inconsistencies, or most importantly be unwilling to be taken into consideration for a data mining approach. The industry's frightening numerical statistics, current science, calls, and commercial applications to the need for more difficult activities are examined. It is possible to change the unfavorable into possible during data preparation. Data preparation includes the removal of noisy components from the information, data reduction procedures, and methods for reducing the complexity of the data.

**Needs:**

Getting good results from the deep learning and machine learning models requires the information to be organized in the right way. Data in an identifier format is only sometimes used by certain deep learning and machine learning models. The information set should be organized such that many deep learning and machine learning algorithms are performed in one information group, with the best one being chosen, as an extra step in the analysis and data preparation process.

**Cleaning Data**

There may be a lot of useless information and gaps in the data. Data cleansing is done to handle this portion. It entails dealing with erroneous data, noisy data, etc.

**Data Absent:**

This problem occurs when there are gaps in the data. It may be dealt with in a number of ways. Among them are:

**Leave the tuples out:** This method only works when a tuple has numerous missing values and our dataset is rather sizable.

**Complete the blank values:** There are many methods to do this job. You may opt to manually fill in the missing values, use the attribute mean, or use the value that is most likely.

**Unclean Data:**

Data that is noisy has no significance and cannot be understood by computers. It may be produced as a result of poor data gathering, incorrect data input, etc. Following are some options for handling it:

**Binning Technique:** To smooth data, this procedure applies to sorted data. The whole set of data is separated into equal-sized pieces before the work is finished using a variety of

techniques. Each section is dealt with independently. To finish the operation, one may utilise boundary values or replace all the data in a segment with its mean.

**Regression:** By fitting the data to a regression function, the data in this case may be made smooth. The regression utilized may be linear (containing one independent variable) or multiple regression (having multiple independent variables).

**Clustering:** This strategy creates a cluster from related data. The outliers may not be seen or they could be outside of the clusters.

**Transformation of data**

This method is used to change the data into formats that are suited for the mining process. This entails the following:

**Normalization:** To scale the data values inside a certain range, this is done (-1.0 to 1.0 or 0.0 to 1.0)

**Selecting an attribute:** With the assistance of the provided set of characteristics, additional attributes are created using this method to aid the mining process.

**Discretization:** By doing this, interval levels or conceptual levels are used to substitute the raw values of numerical attributes.

**Generation of Concept Hierarchies:** Here, qualities are transferred from a lower level of the hierarchy to a higher one. As an example, the attribute "city" may be changed to "country."

**Data compression**

Considering that a process called data mining is used to manage vast amounts of data. Analyses in these situations grew more difficult when dealing with large amounts of data. We use data reduction approach to get rid of this. It attempts to lower the price of data storage and analysis while increasing storage efficiency. The many stages to data reduction are:

A. **Aggregation of Data Cubes:** For the purpose of building the data cube, an aggregation process is conducted to the data.
B. **Choose an attribute subset**: Only the most important characteristics should be utilized; the rest may be ignored. The level of significance and the attribute's p-value may be used to execute attribute selection. An attribute whose p-value is higher than the level of significance can be eliminated.
C. **Reduced Numbers:** This makes it possible to save data models rather of whole data sets, such as regression models.
D. **Reduced Dimensionality:** Data size is decreased using encoding techniques, which may be either lossy or lossless. Such reductions are referred to as lossless reductions if the original data can be recovered after being compressed; otherwise, they are referred to as lossy reductions. Wavelet transformations and PCA are the two most effective dimensionality reduction techniques (Principal Component Analysis)

**Data Processing Types**

Based on the data source and the procedures the processing unit takes to produce an output, there are several kinds of data processing. Processing raw data is not something that can be done in a one-size-fits-all manner.

Data is gathered and processed in batches when using this method of data processing. For big volumes of data, it is used. Take the payroll system, for instance.

1. **Processing by a single user:** This is often done by a single individual for their own use. Even tiny offices may use this method.
2. **Many programming processing:**

This method enables the Central Processing Unit to store and run multiple programmes at once (CPU). Within a single computer system, data is divided into frames and processed using two or more CPUs. It also goes by the name of parallel processing. Additionally, the various programming approaches improve each computer's general productivity. Forecasting the weather is a nice example of multiple programming processes.

**Real-time processing:**

This method makes it possible for the user to communicate directly with the computer system. This method makes data processing simpler. This method, which was created specifically to carry out one job, is sometimes referred to as the direct mode or the interactive mode method. It's a kind of ongoing online processing that never stops. For instance, using an ATM to withdraw cash.

**Online processing:**

This method makes it easier to enter and use data immediately; it does not store or gather data before processing it. The method was created to decrease data input mistakes since it checks the data throughout and makes sure that only corrected data is entered. This method is often used to online applications. Take barcode scanning as an example.

**Time-sharing Processing:**

This kind of online data processing enables several users to share an online computer system's resources. When immediate results are required, this method is used. Moreover, this method is time-based, as the name would imply. Some of the main benefits of time-sharing processing include the following:

a. Multiple users may be serviced at once.
b. The processing times for all users are almost equal.
c. Interaction with already running applications is conceivable.

**Distributed processing:**

It is a specialized data processing method in which a network of computers is created by connecting several computers (placed remotely) to a single host computer. A high-speed communication network keeps all of these computers linked together. The master database, however, is maintained and monitored by the central computer system. This makes computer-to-computer communication easier.

**Data Reduction in Data Mining**

**Data Reduction**

Data reduction is the process of lowering the amount of storage space necessary. Data minimization may save expenses and improve storage effectiveness. Storage providers often use the words raw capacity and effective capacity, which relate to data after reduction, to characterize storage capacity. Data reduction may be done in a variety of ways. Data deduplication, compression, and single-instance storage are the three basic forms. On storage systems, data deduplication, sometimes referred to as data dedupe, removes redundant data segments. When a request is made to access a redundant segment, it only saves that segment once and utilizes that copy. Compared to single-instance storage, data deduplication is more

precise. When files are found, such as email attachments delivered to many recipients, single-instance storage only keeps one copy of the file. Similar to dedupe, single-instance storage points to the one saved copy in lieu of Similar to dedupe, single-instance storage points to the one saved copy in replacement of duplicates.

**Data Reduction in Data Mining**

A method for obtaining information from a large database is called data mining. Data analysis and mining are both impractical and infeasible when working with big volumes of data since they require a long time to complete. Data reduction strategies maintain data integrity while lowering the amount of data. One way to make the original data smaller so that it may be represented in a much lower size is via data reduction. Data reduction methods are used to create a reduced version of the dataset that is much lower in volume while maintaining the integrity of the original data. Lowering the data yields the same analytical results while increasing the efficiency of the data mining process. Data minimization has no impact on the result of data mining. In other words, the outcomes of data mining before and after data reduction are the same or very similar. Compacting information is the aim of data reduction. When the quantity of data is less, it is simpler to apply complex and computationally costly algorithms. The number of rows (records) or columns (columns) (dimensions) in the data may be lowered.

**Approaches to Data Reduction:**

The data reduction technique may result in a condensed description of the original data that is substantially lower in number but maintains the original data's quality. The following are some approaches or strategies for data reduction in data mining, including:

**Diminished Dimensions:**

If we find data that is just tangentially relevant to our investigation, we use the characteristic required for it. Dimensionality reduction reduces the amount of the original data by removing features from the data set under consideration. By deleting useless or redundant attributes, it reduces the size of the data. These three methods for lowering dimensionality are provided. Wavelet Transform: Assume that a data vector A is converted by the wavelet transform into a numerically distinct data vector A' such that both the A and A' vectors have the same length. Data may be compressed by maintaining the least possible portion of the strongest wavelet coefficients. The Wavelet transform is useful for data cubes, sparse data, and skewed data.

Assume that to analyze a data set with n attributes using principal component analysis. The data collection may be described using the k unique tuples with n characteristics that the major component analysis identifies. This may lead to dimensionality reduction, where the original data is cast on a much smaller space. On sparse or skewed data, principal component analysis may be used.

**Reduced Numerosity:**

The numerosity reduction shrinks the original data to a smaller size and conveys it in a much more condensed manner. Parametric and non-parametric numerosity reduction methods are also available.

- **Parametric numerosity:** It reduction saves just the data parameters; the original data is not kept. One method for minimizing parametric numerosity is to use the regression and log-linear methodology.

- **Non-Parametric:** A non-parametric numerosity reduction approach has no model. Regardless of the magnitude of the data, the non-parametric strategy delivers a more uniform reduction, although it may not do as much data reduction as the parametric technique. A minimum of four non-parametric data reduction methods include data compression, histograms, clustering, sampling, and data cube aggregation.

### Aggregation of Data Cubes

Data is compressed using this technique into a more palatable format. Data Cube Aggregation is a multidimensional aggregation that reduces the amount of data while still accurately representing the original data set. Take into account the following instance: All Electronics' quarterly sales data from 2018 to 2022. Simply sum the quarterly sales for each year to get the annual sale for that period. Aggregation provides you with the necessary data, which is much lower in size, and we accomplish data reduction using this strategy without sacrificing any data.

### Compression of data

The act of changing, encoding, or otherwise manipulating data so that it takes up less space is known as data compression. Data compression produces a compact representation of information by minimizing duplication and encoding data in binary form. Data that can be successfully retrieved from its compressed form is referred to as lossless compression. On the other side, lossy compression happens when the original form cannot be recovered from the compressed version. Dimensionality and numerosity reduction techniques are also used for data compression.

### Operation of Discretization

A method for splitting continuous natural characteristics into data with intervals is data discretization. We substitute some of the constant values of the attributes with labels with very small intervals. This suggests that mining results are provided in an understandable and condensed way.

### Data Transformation in Data Mining:

Prior to data mining, data transformation is a crucial data pretreatment technique that must be used to produce patterns that are simpler to interpret. . In two steps of the data pipeline for data analytics projects, data can be modified. Data transformation is the middle phase of an ETL (extract, transform, and load) process, which is commonly used by businesses with on-premises data warehouses. The majority of businesses now increase their computer resources with latency estimated in seconds or minutes by using cloud-based data warehouses. Organizations can load raw data directly into the data warehouse and perform preload transformations at query time thanks to the scalability of the cloud platform.

It's challenging to track or comprehend raw data. Because of this, it needs to be normalized before any information can be extracted from it. The process of transforming raw data into a format that makes it easier to conduct data mining and recover strategic information is known as data transformation. In order to change the data into the right form, data transformation techniques also include data cleansing and data reduction. The process of altering the format, structure, or values of data is known as data transformation. Data can be modified at two points in the data pipeline for initiatives including data analytics. Data transformation serves as the middle phase in an ETL (extract, transform, load) process, which is commonly used by businesses with on-premises data warehouses.

**Data Transformation Types**

In order to simplify the data mining process, the following are some of the most popular types of data transformation processes:

A. **Bucketing/Binning:**It is the process of dividing data into several ranges known as buckets. It improves data organisation and reduces the possibility of insignificant observational mistakes. This sort of data transformation in data preprocessing places numerical data into multiple buckets and then applies various thresholds to turn it into categorical data.

B. **Format Revision:** Processing various data kinds in a set is one of the main challenges in data mining. The format revision kind of data transformation is used to address this problem. Through the conversion of all data into a standardised format, this process standardises data.

C. **Data Splitting:** Another significant data transformation in the data preparation procedure is data splitting. For training, testing, or experimental purposes, it divides or splits information from a data column into several columns.

**Methods for Data Transformation**

Before analysis or storage in a data warehouse, there are a number of data transformation techniques that can help organise and clean up the data. Here are a few of the more popular techniques:

**Smoothing:** This is a data processing technique that purges the dataset of distorted or nonsensical data. To find particular patterns or trends, it also finds slight changes to the data.

**Aggregation:** Data aggregation is the process of gathering unprocessed data from several sources and storing it in one location for accurate analysis and reporting. This method is essential if your company collects a lot of data.

**Discretization:** In order to increase efficiency and facilitate analysis, this data transformation approach creates interval labels in continuous data. Decision tree techniques are used in the process to reduce a large dataset into a small set of categorical data.

**Generalization:** Generalization transforms low-level attributes to high-level ones via concept hierarchies, producing a distinct data snapshot.

**Attribute Building:** By constructing new attributes from a set of existing ones, this technique enables the organisation of a dataset.

**Normalization:** In order to improve the effectiveness of extraction and data mining applications, normalisation changes the data such that the attributes remain within a given range.

**Manipulation:** Manipulation is the process of transforming data to improve its readability and organisation. Tools for data manipulation assist in finding patterns in the data and transforming it into a form that may be used to produce insight.

**Data transformation benefits**

Businesses can benefit from data transformation in a number of ways. Here are a few of the key benefits of data transformation, including:

**Better Scheduling:** Data that has been transformed is simpler to use for both people and computers.

**Improve data quality:** Poor data comes with a lot of expenses and hazards. Data transformation can assist your company in getting rid of inconsistent data and other quality problems like missing values.

**Query more quickly:** Standardized and stored at a source location, enabling rapid and simple retrieval of altered data.

**Improvements in Data Management:** Data is continually being produced by businesses from a growing number of sources. It might be difficult to manage and comprehend metadata if there are inconsistencies in it. Your metadata is refined during data translation, making it simpler to manage and comprehend.

**Greater Use of Data:** Despite the fact that corporations may continually acquire data, much of it is left unanalyzed. By standardizing and improving the usability of your data, transformation makes it simpler to get the most out of it.

**Data transformation drawbacks**

Even though data transformation has many advantages, it might be difficult to change data successfully because of things like:

A. Costly data transformation is a possibility. The cost is determined by the particular infrastructure, software, and processing-related instruments employed. Licensed, computing resources, and recruiting the required employees are a few examples of expenses.
B. Processes for transforming data can be resource-intensive. It can slow down other activities to do transformations in an on-premises data warehouse after data has been loaded or to perform transformations on data before feeding it into apps. The transformations can be completed after loading a cloud-based data warehouse because the platform can expand to handle increased demand.
C. Carelessness and a lack of knowledge might cause issues during change. Lacking the necessary subject matter expertise, data analysts are less likely to spot inaccurate data since they are less familiar with the gamut of accurate and allowed values.
D. Transformations that don't meet an organization's needs can be carried out. For one application, a business might alter information into a certain format, only to change it back to its original format for another application.

There are several methods for preparing data. To eliminate noise and fix discrepancies in the data, data cleaning may be used. Data integration is the process of combining data from several sources into one cohesive data storage, such a data warehouse. It is possible to apply data transformations such as normalisation.

For instance, normalisation may increase the precision and effectiveness of distance-based mining algorithms. Aggregation, the removal of duplicate characteristics, or clustering are a few examples of how data reduction might lower the quantity of the data. These methods may be combined; they are not exclusive of one another.

To repair inaccurate data, for instance, data cleaning may require translating all entries for a date field into a standard format. When data processing techniques are used before mining, the overall quality of the patterns mined and/or the time needed for the actual mining may both be significantly improved.

## Process the Data

Assume you are a manager at AllElectronics and have been given the task of reviewing the data collected by the firm on the sales at your branch. You instantly started working on this project. You carefully go through the database and data warehouse of the business to find and choose the features or dimensions that will be used in your study, such as the item, price, and units sold. Alas! Several of the properties for different tuples don't have any recorded values, as you can see. You want to know whether each item you bought was marketed as being on sale for your analysis, but you find that this information has not been recorded. Additionally, users of your database system have noted mistakes, odd numbers, and consistency issues with some of the data stored during transactions. In other words, the data you want to analyse using data mining techniques are inconsistent (e.g., have discrepancies in the department codes used to categorise items), noisy (have errors or outlier values that deviate from the expected), and incomplete (lack attribute values or certain attributes of interest, or contain only aggregate data). Greetings from the real world!

Large real-world databases and data warehouses often have inconsistent, noisy, and incomplete data. Numerous factors may lead to incomplete data. It's possible that some attributes, including customer details for sales transaction data, are not always accessible. It's possible that some information wasn't entered since it wasn't deemed crucial at the time. Relevant information may not be captured as a result of a misunderstanding or faulty technology. Data that didn't match what had previously been recorded may have been removed. Additionally, it's possible that the history of the data or any updates to it were missed. It may be necessary to infer missing data, especially for tuples with missing values for certain properties.

There are several potential causes of noisy data (having incorrect attribute values). The tools used to acquire the data might be flawed. Data input mistakes might have been caused by either humans or computers. Data transfer errors might also happen. Technology could be constrained by things like a small buffer size for synchronised data flow and consumption. Additionally, mismatched naming standards, data codes, or input field formats, such as date, might lead to incorrect data.

### Data cleansing is also necessary for duplicate tuples.

By adding missing values, reducing noise in the data, locating or eliminating outliers, and resolving discrepancies, data cleaning processes attempt to "clean" the data. Users are less inclined to accept the findings of any data mining if they think the data are unreliable. Additionally, inaccurate data might confuse the mining process and provide misleading results. Although the majority of mining processes contain some methods for handling imperfect or noisy data, they are not always reliable. Instead, they can focus on preventing the data from being overfit to the function that is being modelled. This would include data integration—the combining of several databases, data cubes, or files. The names of certain properties that describe a particular notion, however, may vary throughout databases, leading to inconsistencies and redundancy. For instance, the property used to identify customers can be called cust id in one data store and customer id in another. For attribute values, there may also be naming discrepancies. A first name could be recorded as "Bill" in one database, "William" in another, and "B." in a third, all while having the same spelling. Additionally, you believe that certain characteristics may be deduced from others (e.g., annual revenue). The knowledge discovery process may be slowed down or complicated by having a lot of duplicate material. It is obvious that actions must be done to assist prevent redundancies during data integration in addition to data cleansing. Normally, while preparing the data for a

data warehouse, data cleaning and integration are carried out as a preparatory stage. To find and eliminate redundancies that may have arisen through data integration, further data cleaning might be done.

Returning to your data, let's imagine you've chosen to analyse it using a distance-based mining approach like neural networks, nearest-neighbor classifiers, or clustering. These techniques provide better results if the data being examined have been normalised, or scaled to a certain range like [0.0, 1.0]. Examples of characteristics in your client data are age and yearly pay. Age often has significantly smaller values than the yearly income characteristic. As a result, if the variables are not normalised, yearly income distance measurements will often be greater than age distance measurements. Additionally, obtaining aggregate data on sales by customer region—information that is not included in any precomputed data cube in your data warehouse—would be helpful for your research. You quickly notice that further data pretreatment steps like normalisation and aggregation are data transformation activities that would help the mining process succeed. About data integration and data transformation.

You ponder your findings even more as you say, "Hmmm." "The data collection I've chosen for study is HUGE, which will inevitably cause the mining process to sluggish. Data reduction creates a significantly smaller, reduced version of the original data collection that yet yields essentially the same analytical conclusions. Numerous methods exist for data reduction. These include data aggregation (for example, creating a data cube), attribute subset selection (for example, eliminating pointless attributes through correlation analysis), dimensionality reduction (for example, using encoding schemes like minimum length encoding or wavelets), and numerosity reduction (for example, "replacing" the data by alternative, smaller representations like clusters or parametric models). Concept hierarchies may be used to "reduce" data by replacing low-level ideas, such as city for client location, with higher-level concepts, such as region or province or state. The ideas are arranged into different degrees of abstraction using a concept hierarchy. Data reduction techniques like data discretization are highly helpful for automatically creating idea hierarchies out of numerical data.



**Figure 5.1: Data Processing**

The phases of data preparation are summarized in Figure 5.1. Remember that the categories listed above are not exclusive of one another. For instance, both data reduction and data cleaning may be considered to constitute the elimination of superfluous data. In conclusion, real-world data are often inaccurate, lacking, and inconsistent. By enhancing the quality of the data, data preparation methods may also aid to increase the precision and efficacy of the following mining operation. Preparing data for analysis is a crucial phase in the knowledge

discovery process since good judgments need good data. Making decisions may benefit greatly by finding data abnormalities, fixing them quickly, and analyzing less data.

**Describing and Summarizing Data**

A comprehensive understanding of your data is crucial for effective data preparation. You may use descriptive data summarising methods to pinpoint your data's usual characteristics and highlight the data values that need to be regarded as outliers or noise. Thus, before delving into the specifics of data pretreatment approaches, we first discuss the fundamental ideas of descriptive data summarization. Users would want to learn about data characteristics related both central tendency and dispersion of the data for several data preparation jobs. Measures of data dispersion include quartiles, interquartile range (IQR), and variance, while measures of central tendency include mean, median, mode, and middle. Understanding the distribution of the data is greatly aided by these descriptive statistics. The statistical literature has extensively examined these metrics. We need to investigate how they may be calculated effectively in huge datasets from the perspective of data mining. The concepts of distributive measure, algebraic measure, and holistic measure must be specifically introduced. Choosing an effective implementation for a measure depends on knowing what sort of measure it is that we are working with.

**Cleaning of Data**

Real-world data often include gaps, are noisy, and are inconsistent. Data cleaning (or data cleansing) procedures make an effort to complete blanks in the data, level out noise while identifying outliers, and fix discrepancies. You will learn the fundamental techniques for data cleansing in this part.

**Absence of Values**

Consider that you must examine the sales and client information from AllElectronics. You see that a large number of tuples lack recorded values for a number of fields, including customer income.

**The following techniques:**

Disregard the tuple: When the class label is absent, this is often done (assuming classification is a part of the mining operation). Unless the tuple has multiple characteristics with missing values, this approach is not particularly successful.

It is particularly bad when the variance in the percentage of missing data for each characteristic is large.Fill in the missing value manually: Generally, this method takes a lot of time and may not be practical for big data sets with several missing values.

**Fill in the blank with a global constant:** In lieu of each missing attribute value, use a single constant, such as "Unknown" or. If missing values are substituted with, for example, "Unknown," the mining software can erroneously believe that they constitute an interesting notion since they all share the value of "Unknown." Thus, although being straightforward, this approach is not infallible.

**If a value is absent, use the attribute mean to fill it in**. For instance, assume that AllElectronics customers make an average annual income of $56,000. Use this value to fill in the income value that is missing.

**For all samples in the same class as the supplied tuple**, use the attribute mean:

The average income value for clients who fall into the same credit risk category as the supplied tuple, for instance, should be used to fill in the missing value when categorising clients based on their credit risk.

**Fill in the missing value with the most likely value**: this may be done using decision tree induction, regression, or inference-based techniques that use the Bayesian formalism. You may build a decision tree to forecast the revenue numbers that are absent by utilising the other customer variables in your data collection, for instance.The data are skewed by methods that was entered may not be accurate. However, the sixth method is a well-liked tactic. It makes the maximum use of the available data to forecast missing values when compared to the other approaches. There is a higher likelihood that the connections between income and the other characteristics are maintained if the estimate of the missing value for income takes the values of the other attributes into account.

It's vital to remember that a missing number cannot always indicate that the data is incorrect! For instance, applicants could be required to provide their driver's licence number while requesting a credit card. Naturally, candidates without a driver's licence may leave this box empty. Respondents should be allowed to indicate values like "not applicable" on forms. Each attribute should ideally contain one or more restrictions on the null condition. The guidelines may outline whether null values are permitted as well as how to manage or alter them. If information will be entered into a field at a later stage of the business process, it is also possible to purposefully leave it empty. As a result, although we may make every effort to clean the data after it has been collected, a proper database architecture and data input process should assist to reduce the amount of mistakes or missing values in the first place.

**Noisy Data Binning**

Binning techniques smooth sorted data values by looking at their "neighborhood," or the values nearby. The values from the sort are divided up into a number. Of containers or "buckets." Binning techniques accomplish local smoothing since they look at the surrounding variables. Binning methods. In this illustration, the price data are sorted first and then divided into equal-frequency bins of size 3. (i.e., each bin contains three values). When smoothing by bin means, the mean value of the bin is used to replace each value in the bin. As an example, the average of the numbers 4, 8, and 15 in Bin 1 is 9. Consequently, the number 9 is used to replace each of the bin's original values. Similarly, smoothing by bin medians, which substitutes the bin median for each bin value, may be used. The minimum and maximum values in a certain bin are known as the bin borders when smoothing by bin boundaries is used. The nearest boundary value is then used to replace each bin value. In general, the smoothing effect increases with increasing breadth. Alternative bin configurations include equal-width bins, where each bin's interval range of values is fixed.

Regression: By fitting the data to a function, such as using regression, data may be smoothed. Finding the "best" line to fit two traits (or variables) is the goal of linear regression, which enables one attribute to predict the other. A variation of linear regression known as multiple linear regression involves more than two features and fits the data to a multidimensional surface. Clustering: When comparable numbers are grouped together into "clusters," outliers may be found. Values that are outliers are ones that don't fit into the collection of clusters.

Numerous discretization-based data reduction techniques also include data smoothing. The number of different values for each characteristic is decreased, for instance, by the binning approaches discussed above. For logic-based data mining techniques like decision tree induction, which continually perform value comparisons on sorted data, this serves as a type of data reduction. Concept hierarchies are a data discretization technique that may also be

used to smooth out the data. The quantity of data values that must be processed by the mining process is decreased by a concept hierarchy for price, for example, which may map actual price values into low, moderately priced, and costly categories.

**Cleaning Data as a Process**

Inaccurate data is a result of missing values, noise, and inconsistencies. We have examined methods for managing missing data and smoothing data up to this point. "However, cleaning up data is a major task. Discrepancy detection is the initial stage in the process of cleansing data. Several causes, such as poorly designed data input forms with several alternative sections, human mistake in data entry, purposeful errors (example: respondents who did not wish to provide personal information), and data degradation, may lead to discrepancies (e.g., outdated addresses). Inconsistent data representations and the inconsistent application of codes may also result in discrepancies. System flaws and mistakes in instrumentation equipment that record data are other sources of differences. When data are (inadequately) utilised for goals other than those for which they were originally intended, errors may also result. Inconsistencies resulting from data integration might also exist (e.g., where a given attribute can have different names in different databases).

Use whatever prior knowledge you may have about the characteristics of the data as a starting point. Metadata is a term used to describe this information or "data about data. Examples of probable outliers include values that are greater than two standard deviations from the mean for a certain characteristic. You may develop your own scripts in this phase or utilise some of the tools we'll cover later on. You could discover noise, outliers, and strange values from this that need examination. Additionally, the data should be evaluated for unique, consecutive, and null rules. Each value of the specified property must be distinct from all other values for that attribute, according to a unique rule. According to a sequential rule, all values for the property must be distinct and there cannot be any missing values between the lowest and highest values (e.g., as in check numbers). When a value for a certain attribute is unavailable, for example, a null rule defines the usage of blank spaces, question marks, special characters, or other strings to signify this and how these values should be handled. There are a few possible explanations for missing values, including the following:

The person who was initially asked to provide a value for the attribute refuses and/or discovers that the information requested is irrelevant (for example, a license-number attribute left blank by nondrivers);

The data entry person is unaware of the proper value. The value is to be provided by a later step of the process. The null rule should describe how to represent the null condition, for instance, by storing a blank for character attributes or zero for numerical attributes, or by using any other conventions that may be in use.

The stage of discrepancy identification may be aided by a variety of various commercial technologies. Simple domain knowledge, such as an understanding of postal addresses and spell-checking, is used by data scrubbing systems to identify problems and rectify the data. When cleaning data from various sources, these technologies use parsing and fuzzy matching methods. Data auditing tools look for inconsistencies by examining the data to identify patterns and correlations, then identifying patterns that are violated by the data. They are variations on tools for data mining. To uncover correlations, for instance, they may use statistical analysis; to locate outliers, they would use clustering.

It may be possible to manually fix certain data discrepancies utilising other sources. For instance, a paper trail might be used to fix mistakes made during data input. However, the

majority of mistakes call for data transformations. The second stage of the data cleansing procedure is this. To put it another way, after we identify differences, we usually need to define and use a number of transformations to fix them.The data transformation stage may be aided by commercial technologies. Simple modifications may be provided using data migration tools, such as changing the string "gender" to "sex." Through a graphical user interface, ETL (extraction, transformation, and loading) tools enable users to define transformations (GUI). For this stage of the data cleaning process, we often also choose to build custom scripts since these solutions generally only handle a few number of transforms.

Iterations are used in the two-step procedure of discrepancy identification and data transformation (to fix discrepancies). But it takes a long time and is prone to mistakes. More differences could be introduced by certain changes. Some nested inconsistencies may not be discovered until after others have been corrected. For instance, a mistake like "20004" in a year field could not be discovered until all date values have been converted to the same format. Transformations are often carried out in batches while the user waits and receives no feedback. The user can only go back and make sure that no additional anomalies have been unintentionally produced once the change is finished. Before the user is pleased, usually many iterations are needed. Any tuples that cannot be handled automatically by a particular transformation are normally written to a file without being provided any information as to why they failed. As a consequence, there isn't much interaction during the whole data cleansing process.

Increased interaction is emphasised in new data cleansing methodologies. For instance, Potter's Wheel is a freely accessible data cleaning tool that combines discrepancy detection and transformatio. By creating and debugging each individual transformation using a spreadsheet-like interface, users may progressively assemble a collection of transformations. Graphs or examples may be used to illustrate the transformations. On the records that are now visible on the screen, results are shown instantly. The user has the option to "delete" any modifications that led to further mistakes by choosing to reverse them. On the most recent converted view of the data, the programme automatically does discrepancy checking in the background. As differences are discovered, users may progressively build and improve transformations, resulting in more effective and efficient data cleansing.

The creation of declarative languages for the definition of data transformation operators is another strategy for enhancing interaction in data cleaning. Such work focuses on the design of robust SQL extensions and methods that allow users to effectively articulate data cleaning criteria. It is crucial to continually updating the metadata to reflect new information as it is learned about the data. As a result, data cleaning for next iterations of the same data store will go more quickly.

### Data Transformation and Integration

Data integration, or the combining of data from many data silos, is often required for data mining. It can also be necessary to change the data into formats suitable for mining. Both data integration and data transformation are discussed in this section.

### Data Integration

Your data analysis assignment is probably going to require data integration, which compiles data from several sources into a cohesive data store, similar to data warehousing. Multiple databases, data cubes, and flat files are a few examples of these sources. Several concerns need to be taken into account while integrating data. Integration of schemas and object

matching might be challenging. Such information may be utilised to prevent mistakes during the integration of the schema.

**Transformation of Data**

Data transformation involves transforming or consolidating the data into mining-ready formats. These processes may be used in data transformation:

Smoothing is a technique that helps to reduce data noise. These methods consist of grouping, regression, and binning. Aggregate is the process of applying summary or aggregation procedures on data. To calculate monthly and yearly totals, for instance, the daily sales data may be combined. This process is usually used to build a data cube for the examination of data at various granularities.

Using idea hierarchies to substitute low-level or "primitive" (raw) data with higher-level concepts to generalise the data. For instance, categorical characteristics like street may be applied to more abstract ideas like city or nation.

Similar to how age values may be translated into higher-level notions like youth, middle age, and seniority. Normalization is the process of scaling attribute data to fit inside a narrow range, such as 0.0 to 1.0 or 1.0 to 1.0.The technique of creating new attributes from the set of existing attributes in order to aid in the mining process is known as attribute creation (or feature construction).

**Data compression**

Consider that you have chosen some data for analysis from the All Electronics data warehouse. The data set will probably be enormous! Large-scale complex data analysis and mining may be time-consuming, rendering such study impractical or unfeasible. Data reduction methods may be used to produce a condensed version of the data set that is substantially lower in size while closely preserving the original data's integrity. In other words, mining on the smaller data set should be more effective while still yielding essentially the same analytical conclusions.

The following are some data reduction techniques:

1. Data cube aggregation, which entails applying aggregating processes to the data prior to building a data cube.

2. Attribute subset selection, which allows for the detection and elimination of redundant, weakly related, or irrelevant characteristics or dimensions.

3. Dimensionality reduction, which minimises the size of the data set using encoding techniques.

4. Data replacement or estimation using alternative, more manageable data representations, such as parametric models (which require only store the model parameters in lieu of the real data) or nonparametric techniques like clustering, sampling, and the application of histograms.

5. Discretization and idea hierarchy creation, where ranges or higher conceptual levels are used to substitute raw data values for attributes. The automated creation of idea hierarchies benefits greatly from data discretization, a kind of numerosity reduction. Being able to mine data at various levels of abstraction, discretization and idea hierarchy development are effective techniques for data mining.

**Aggregation of Data Cubes**

Assume you have gathered the information necessary for your analysis. These numbers represent the quarterly sales for AllElectronics during the years 2002 to 2004. However, you are more concerned with the yearly sales (total for the year) than the quarterly total. As a consequence, the data may be combined such that the final data summarises the total sales annually rather than quarterly. The amount of the final data set is reduced without any essential information being lost for the analysis process. Data warehousing book goes into great length on data cubes. Here, we quickly present a few ideas. Multidimensional aggregated data is stored in data cubes.

Each cell has an aggregate data value that, in multidimensional space, corresponds to the data point. Each attribute may have a hierarchy, enabling the examination of data at several levels of abstraction. For instance, a branch hierarchy may enable the division of branches into regions according to their addresses. Data cubes provide quick access to precomputed, summarised data, which is advantageous for both data mining and online analytical processing.

The basic cuboid is the cube formed at the lowest level of abstraction. The base cuboid should represent a specific entity of interest, such as a customer or a salesperson. In other words, the analysis should be able to employ the lowest level. The apex cuboid is a cube at the greatest degree of abstraction. A data cube may alternatively refer to a lattice of cuboids since data cubes designed for different levels of abstraction are sometimes referred to as cuboids. The amount of the generated data is further decreased with each higher degree of abstraction. The smallest cuboid that is appropriate for the particular job should be chosen when responding to requests for data mining.

**Clustering**

Techniques for clustering treat tuples of data as objects. In order to make things inside a cluster "similar" to one another and "dissimilar" to objects in other clusters, they divide the objects into groups or clusters. According to a distance function, similarity is often described in terms of how "near" the items are to one another in space. The greatest separation between any two items in a cluster, or its diameter, may be thought of as the cluster's "quality." The average distance between each cluster item and the cluster centroid, which is used to denote the "average object" or average location in space for the cluster, is known as centroid distance. The centroid of each cluster is shown with a "+", which depicts a 2-D plot of customer data with regard to customer locations in a city. There are three distinct data clusters.

The cluster representations of the data are employed in data reduction to replace the real data. The kind of the data will determine how successful this strategy is. Data that can be grouped into separate clusters performs far better than smeared data. Multidimensional index trees are generally utilised in database systems to provide quick data access. They may also be used to cluster data at several resolutions for hierarchical data reduction. This may be used to provide rough responses to inquiries. With the root node representing the whole space, an index tree recursively divides the multidimensional space for a given collection of data items. These balanced trees often have internal and leaf nodes. The space that each parent node represents is collectively represented by the child nodes that each parent node holds keys and pointers to. Pointers to the data tuples they represent (or the actual tuples) are included in each leaf node.

In order to hold data at different degrees of resolution or abstraction, an index tree may also contain aggregate and detail data. It offers a hierarchy of clusterings of the data collection, with a name for each cluster that describes the data it contains. An index tree may be seen as a hierarchical histogram if we treat each child of a parent node as a bucket.

-----------------------

# CHAPTER 6

# CLASSIFICATION AND PREDICTION

Santhosh S

Assistant Professor, Department of Computer Science Engineering, Faculty of Engineering and Technology,
JAIN (Deemed-to-be University), Karnataka – 562112
Email Id- s.santhosh@jainuniversity.ac.in

To generate models representing significant data classes or to forecast future data trends, two types of data analysis are classification and prediction.Prediction uses continuous-valued functions and categorical (discrete, unordered) labels.For instance, based on a person's salary and line of work, we may create a classification model to classify bank loan applications as safe or dangerous, or a prediction model to forecast how much a prospective consumer would spend on computer equipment.Instead of predicting a categorical label, a predictor that predicts a continuous-valued function, or ordered value, is created.The statistical technique of regression analysis is most often used to numerical prediction.Researchers in machine learning, pattern recognition, and statistics have suggested a wide variety of categorization and prediction techniques.

The majority of algorithms are memory-resident and usually assume a minimal amount of data. On top of this work, more recent data mining research has developed scalable classification and prediction methods that can manage enormous disk-resident data sets.

Classification and Prediction Issues:

**Data preparation for classification and prediction:**

To aid increase the precision, efficacy, and scalability of the classification or prediction process, the data may go through the following preparation processes.

Data cleaning: This is the preparation of data to eliminate or decrease noise (by using smoothing methods) and the handling of missing values (for example, by replacing a missing value with the most frequent value for that attribute or with the most likely value based on statistics).

This phase may assist clear up uncertainty during learning, even if the majority of classification algorithms incorporate some capabilities for addressing noisy or missing data.

Analysis of Relevance:Many of the data's properties can be duplicated.To determine if any two supplied qualities are statistically connected, correlation analysis may be performed.

For instance, a high correlation between the traits A1 and A2 can indicate that one of the two doesn't need to be included in future investigation.Also possible are irrelevant properties in a database. In these situations, attribute subset selection may be used to identify a smaller collection of attributes such that the probability distribution of the resultant data classes closely resembles the original distribution obtained using all attributes.

As a result, characteristics that do not contribute to the classification or prediction job may be found via relevance analysis, which takes the form of correlation analysis and attribute subset selection.

This kind of study may increase the scalability and effectiveness of categorization.

## Data Reduction and Transformation

When neural networks or techniques requiring distance measurements are utilised in the learning stage, the data may be modified through normalisation.

All values for a specific property are scaled during normalisation such that they all fit within a narrow range, such as -1 to +1 or 0 to 1.

The information may also be changed by applying it to ideas at a higher level. Concept hierarchies might be used for this. In particular, ongoing valued characteristics may benefit from this.For instance, discrete ranges, such as low, medium, and high, may be generalised from numerical values for the characteristic income. Similar to this, category characteristics like street may be applied to more abstract ideas like city.Numerous additional approaches, including wavelet transformation, principal components analysis, and discretization methods including binning, histogram analysis, and clustering, may also be used to compress data.

## Comparing Classification and Prediction Techniques

Accuracy: A classifier's accuracy is measured by how well it can predict the class label of brand-new or previously undiscovered data (i.e., tuples without class label information).

A predictor's accuracy measures how effectively it can forecast the value of the predicted characteristic for brand-new or previously unobserved data.Speed is a term used to describe the computing expenses associated with creating and using the specified classifier or predictor.

Robustness is the capacity of a classifier or predictor to provide accurate predictions in the presence of noisy or missing value data.Scalability is the capacity to build the classifier or predictor effectively in the presence of vast volumes of data.Interpretability: This describes the degree of comprehension and perception that the classifier or predictor offers.

Since interpretability is subjective, evaluating it might be more challenging.

Decision tree induction is the process of learning decision trees using training tuples that have been classified.

A decision tree is a tree structure that resembles a flowchart and where each internal node represents a test on an attribute.Each branch is an example of a test result.A class label is stored in each leaf node.The root node is the topmost node in a tree (Figure 6.1).
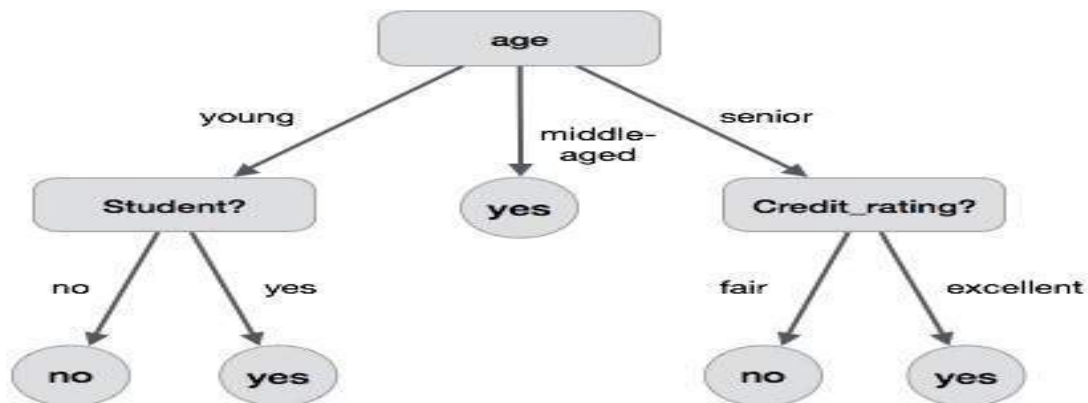


**Figure 6.1: Comparing Classification and Prediction Techniques**

Decision tree classifiers are ideal for exploratory knowledge discovery since they may be built without the need for topic expertise or parameter tuning.

a. High-dimensional data may be handled via decision trees.
b. Their use of a tree-like representation of learned information makes it simple for people to understand and process.
c. Decision tree induction's learning and classification phases are easy and quick.
d. Decision tree classifiers are often accurate.
e. Many application fields, including medical, manufacturing and production, financial analysis, astronomy, and molecular biology, have employed decision tree induction methods for categorization.

## Induction from a Decision Tree

A supervised learning technique called a decision tree is used in data mining for techniques of classification and regression. We can use this tree to aid in decision-making. The classification or regression models are produced as a tree structure by the decision tree. It divides a data set into smaller subgroups while while slowly building the decision tree. The decision nodes and leaf nodes make up the final tree. There are at least two branches on a decision node. The categorization or choice is shown in the leaf nodes. We are limited in our ability to divide leaf nodes further. The root node is the topmost decision node in a tree and is connected to the best predictor. Categorical and numerical data may both be handled by decision trees.

Key elements:

Entropy: Entropy is a typical unit of impurity measurement. The decision tree assesses the impurity or randomness of data sets.

## Induction from a Decision Tree

Information Gain: After the dataset is divided, information gain is the decrease in entropy. Another name for it is entropy reduction. Finding the qualities that provide the most data gain is the foundation of decision tree construction.

Induction from a Decision Tree

A decision tree is essentially a flowchart diagram with terminal nodes that represent choices. The dataset may be measured for entropy to determine how to segment it until all of the data is in the same class.

Decision trees are valuable

1. It allows us to carefully consider all of the potential effects of a choice.
2. It gives us a framework to evaluate the worth of outcomes and the likelihood that they will materialise.
3. It enables us to base our judgements on the most accurate predictions and available facts.

In other terms, a decision tree is a hierarchical tree structure that may be used to design a series of simple decision rules to divide a large collection of records into more manageable sets of the class. An algorithm for dividing a large diverse population into smaller, more homogenous, or mutually exclusive groupings is called a decision tree model. The classes must be of a qualitative kind, such as categorical, ordinal, or binary. In contrast, the characteristics of the classes may be any variables from nominal, ordinal, binary, and

quantitative values. In essence, a decision tree generates a set of rules that may be used to identify the class using the characteristics and class data that are provided. A hierarchy of segments inside a segment is created as a consequence of the application of one rule after another. Each section of the hierarchy is known as a node, and the hierarchy is known as a tree. The individuals from the succeeding groups are more and more similar to one another with each successive division. Recursive partitioning is the name for the algorithm used to create a decision tree. The formula is called CART (Classification and Regression Trees)

Consider the factory provided as an example.

**Induction from a Decision Tree**

A healthy economy has a 0.6 (60% chance), which results in a $8 million profit; a poor economy has a 0.4 (40% chance), which results in a $6 million profit. Expanding factor costs $3 million.The likelihood of a good economy is 0.6(60%), which results in a profit of $4 million, while the probability of a poor economy is 0.4, which results in a profit of $2 million.

Based on the available facts, the management teams must make a data-driven choice about whether to grow or not.

(0.6 * 8 + 0.4 * 6) - 3 = $4.2M is the net expand.

(0.6*4 + 0.4*2) Net Not Expand - 0 = $3M

$4.2M is more than $3M, hence the plant has to be enlarged.

Decision tree algorithm: Although it may seem complicated, the decision tree algorithm's fundamental approaches are as follows:

The three parameters that make up the algorithm are D, attribute list, and attribute selection method.D is often referred to as a data partition.D is initially made up of all of the training tuples and their associated class levels (input training data).A list of characteristics defining the tuples is included in the parameter attribute list.The attribute that "best" separates the provided tuples into their respective classes is specified by the attribute selection technique, or attribute selection method.Applying an attribute selection measure is what the attribute selection method procedure does.A decision tree does not need information to be scaled, which is an advantage.Additionally, missing values in the data have little to no impact on how a decision tree is constructed.The technical team and stakeholders can easily and automatically understand a decision tree model.Decision trees need less effort to prepare the data during pre-processing than other methods.Data normalisation is not necessary when using a decision tree.

**Bayesian Classifiers**

The relationship between the attribute set and the class variable is often non-deterministic in nature. In other words, even if a test record's attribute set is similar to some of the training samples, its class label cannot be inferred with confidence. These conditions could develop as a result of noisy data or the existence of certain perplexing elements that affect categorization, but they are not taken into account in the research. Consider the problem of determining, for instance, whether a person is likely to get liver disease based on eating habits and productivity. Even while most individuals who eat well and exercise regularly have a lower likelihood of developing liver disease, this may still be the case for certain people for other reasons. For instance, as a result of consuming high-calorie street cuisine and

abusing alcohol. The leaning problem may become vulnerable as a result of analyses that determine if a person's diet is healthy or whether their exercise effectiveness is adequate.

The Bayes theorem is used in Bayesian classification to forecast the occurrence of any event. Statistical classifiers with Bayesian probability understandings are known as Bayesian classifiers. The theory describes how a probability may be used to indicate a degree of belief.

Thomas Bayes is credited with creating the Bayes theorem, which is a method that utilises evidence to derive limits on an unknown parameter. Thomas Bayes was the first to apply conditional probability in this way.

**The Naive Bayes algorithm**

A classification strategy known as naive Bayes algorithms is based on the application of Bayes' theorem with the firm belief that all predictors are independent of one another. Simply put, the presumption is that a feature's inclusion in a class is unrelated to any other feature's presence in the same class. A phone could be deemed smart if it has a touch screen, internet access, an excellent camera, etc. Even while each of these traits is reliant on the others, they each individually increase the likelihood that the phone is a smart phone.

Finding the posterior probabilities, or the likelihood of a label given certain observable characteristics, P(L | features), is the core goal of Bayesian classification. We may describe this quantitatively using the Bayes theorem as follows:

$$P(L|features) \; = \; \frac{P(L)P(features|L)}{P(features)}$$

The prior probability of the class is P (L).

The likelihood, or P(features|L), measures the likelihood of a predictor in a given class.

P(features) is the predictor's prior probability.Using the Scikit Learn Python library, we can design a Naive Bayes model. Scikit Learn is the most helpful library for this. Under the Scikit Learn Python package, we have the three kinds of Naive Bayes models listed below:

**Naïve Gaussian Bayes**

The premise that the data from each label is obtained from a straightforward Gaussian distribution makes it the most basic Nave Bayes classifier.

**Multiple Naive Bayes**

Multinomial Nave Bayes is another helpful Nave Bayes classifier in which the features are assumed to come from a straightforward Multinomial distribution. Such naive Bayes models are best suited for features that reflect discrete counts.

**Naïve Bernoulli Bayes**

Bernoulli Naive Bayes is a crucial model since it assumes that features are binary (0s and 1s). An application of Bernoulli Naive Bayes may be text categorization using the "bag of words" model.A multilayer feed-forward neural network: A multilayer feed-forward neural network is learned on using the backpropagation technique.

A collection of weights are learned iteratively to predict the class label of tuples.An input layer, one or more hidden layers, and an output layer make up a multilayer feed-forward neural network.The properties that were assessed for each training tuple correspond to the

inputs to the network. The units that make up the input layer are concurrently given the inputs. These inputs go through the input layer before being concurrently weighted and routed to a subsequent layer known as the hidden layer.The hidden layer units' outputs may be fed into the input of a subsequent hidden layer, and so on. The number of layers that are concealed is arbitrary.

Units in the output layer receive the weighted outputs of the final hidden layer as input and emit the network's prediction for specific tuples.

**Back propagation-based classification**

A neural network learning algorithm is back propagation.A neural network is made up of linked input/output units with weights assigned to each connection.In order to be able to anticipate the right class label of the input tuples, the network learns throughout the learning phase by modifying the weights.Due of the links between units, neural network learning is also known as connectionist learning.Since training neural networks takes a long period, they are better suited to situations where this is possible.Backpropagation learns by processing a training set of tuples repeatedly and comparing each tuple's predicted value with the known goal value.

The training tuple's known class label, or a continuous value, may serve as the target value in classification issues (for prediction).The weights are adjusted for each training tuple in order to reduce the mean squared error between the network's predicted value and the actual target value. Backpropagation is the term for the process of making changes in a direction that is "backwards," i.e., from the output layer down via each hidden layer to the first hidden layer.The weights will ultimately converge, but this is not a certainty, and the learning process will come to an end.

**Advantages:**

  a. They have a high threshold for noisy data, and they can categories patterns they haven't been taught to recognize.
  b. They may be used if you are unsure about the connections between characteristics and classes.
  c. Unlike the majority of decision tree algorithms, they work well with continuous-valued inputs and outputs.
  d. On a variety of real-world data sets, including handwritten character recognition, pathology and laboratory medicine, and teaching a computer to speak English text, they have proven effective.
  e. Since neural network algorithms are naturally parallel, they may be run more quickly by using parallelization methods.

**Algorithm K-Nearest Neighbor (KNN)**

One of the simplest machine learning algorithms, based on the supervised learning method, is K-Nearest Neighbour. The K-NN method makes the assumption that the new case and the existing cases are comparable, and it places the new instance in the category that is most like the existing categories. A new data point is classified using the K-NN algorithm based on similarity after all the existing data has been stored. This implies that utilising the K-NN method, fresh data may be quickly and accurately sorted into a suitable category.

Although the K-NN approach is most often employed for classification issues, it may also be utilised for regression. Since K-NN is a non-parametric technique, it makes no assumptions about the underlying data.It is also known as a lazy learner algorithm since it saves the

training dataset rather than learning from it immediately. Instead, it uses the dataset to execute an action when classifying data.The KNN method simply saves the information during the training phase, and when it receives new data, it categorises it into a category that is quite similar to the new data.

Consider the following scenario: We have a photograph of a creature that resembles both cats and dogs, but we are unsure of its identity. Therefore, because the KNN method is based on a similarity metric, we may utilise it for this identification. Our KNN model will look for similarities between the new data set's characteristics and those in the photos of cats and dogs, and based on those similarities, it will classify the new data set as either cat- or dog-related.

**K-NN at work**

The following method may be used to describe how the K-NN works:

    A.  Step 1: Decide on the neighbours' K-numbers.
    B.  Calculate the Euclidean distance between K neighbours in step two.
    C.  Step 3: Based on the determined Euclidean distance, choose the K closest neighbours.
    D.  Step 4: Count the number of data points in each category among these k neighbours.
    E.  Step 5: Assign the fresh data points to the category where the neighbour count is highest.
    F.  Step 6: Our model is complete.

**Deciding on K's value for the K-NN algorithm**

The following are some things to keep in mind while choosing K's value in the K-NN algorithm:

The ideal value for "K" cannot be determined in a specific fashion, thus we must experiment with several numbers to discover the one that works best. K is best represented by the number 5.A relatively small number of K, such K=1 or K=2, might be noisy and cause outlier effects in the model.Although K should have large values, there may be some issues.

**Benefits of the KNN Algorithm**

1. It is easy to put into practise.
2. It can withstand noisy training data.
3. If there is a lot of training data, it could work better.
4. K's value must always be determined, and sometimes that may be difficult.
5. The high computing cost is caused by the need to determine the separation between each data point for each training sample.

**Fuzzy collection of methods**

The drawback of rule-based categorization systems is that they need abrupt cutoffs for continuous characteristics. Consider the following rule for approving client credit applications, for instance. In essence, the regulation states that requests from clients who have had a job for two years or more and who have a significant salary (i.e., of at least $50,000) were approved:

If (years employed) >= 2 and (income) > $50,000, credit is "accepted."

If a client has worked for at least two years and her salary is, let's say, $50,000, she will be given credit; nevertheless, if it is $49,000, she will not. It may seem unreasonable to use such strict thresholding. Instead, fuzzy logic may be added to the system to enable the definition of

"fuzzy" thresholds or limits. Fuzzy logic employs truth-values between 0.0 and 1.0 to express the degree of membership that a specific value has in a given category rather than having a clear boundary between categories or sets. Therefore, using fuzzy logic, we may express the idea that a $49k income is somewhat high but not as high as a $50k income.

Systems for data mining that do categorization may benefit from fuzzy logic. The benefit of working at a high degree of abstraction is provided. The following are typical applications of fuzzy logic in rule-based systems:

Fuzzy values are created from attribute values. Calculating the fuzzy membership or truth values is done by fuzzy logic systems, which often include graphical tools to help users.

More than one fuzzy rule could be applicable for a certain new sample. Each relevant rule casts a vote for a category's inclusion. The truth-values for each anticipated category are typically included.The system's return value was created by adding the amounts acquired above.

Depending on how complicated the fuzzy membership graphs are, this method may require multiplying by the mean truth-value of each category and weighting each category according to its truth sum. For categorization, fuzzy logic systems have been employed in many fields, including finance and healthcare.

**Prognosis**

Statistical regression modelling approaches may be used to simulate the prediction of continuous values. For instance, we would wish to create a model to forecast a new product's sales potential given its pricing or the pay of college graduates with 10 years of work experience. Many issues may be resolved using linear regression, and many more can be resolved by transforming the variables in order to turn a nonlinear issue into a linear one. We are unable to address regression in great depth due to space constraints. Instead, this section offers a simple overview of the subject. You will be acquainted with the concepts of linear, multiple, and nonlinear regression as well as extended linear models at the conclusion of this lesson. Regression issues may be resolved using a variety of software programmes. Examples include SAS, SPSS, and S-plus (all available at http://www.mathspfr.com), among others.

**Multiple and Linear Regression**

Data are modelled using a straight line in linear regression. The simplest kind of regression is linear regression. Bivariate linear regression assumes that the variance of the random variable Y (referred to as the response variable) is constant and models it as a linear function of the random variable X (referred to as the predictor variable), with the formula Y=. Regression coefficients a and b specify the Y-intercept and slope of the line, respectively. The least-squares approach, which reduces the error between the real data and the line's estimate, may be used to solve for these coefficients.

**Nonlinear Regression**

By adding polynomial terms to the fundamental linear model, polynomial regression may be represented. We may convert the nonlinear model by transforming the variables. We may turn a nonlinear model into a linear one that can be solved using the least squares approach by applying transformations to the variables.

Some models (like the sum of exponential factors, for instance) are intractably nonlinear and cannot be transformed into linear models. In certain situations, it could be feasible to derive least square estimates by doing in-depth calculations on more difficult equations.

## A different regression model

Continuous-valued functions may be modelled using linear regression. Due in great part to its simplicity, it is frequently utilised. Theoretically, generalised linear models serve as the basis for using linear regression to describe categorical response variables. Contrary to linear regression, where the variance of Y is constant, in generalised linear models, the variance of the response variable V is a function of the mean value of V. Logistic regression and the Poisson linear function of a collection of predictor variables are two common varieties of extended linear models. Poisson regression is widely used to describe count data since they usually have a Poisson distribution.

Discrete multidimensional probability distributions are roughly represented by log-linear models. They might be used to calculate the probability value corresponding to each data cube cell. Assuming for the sake of illustration that we have data for the variables city, item, year, and sales, we can see that because the log-linear technique requires all attributes to be categorical, continuous-valued attributes (like sales) must first be discretized. Based on the 2-d cuboids for city and item, city and year, city and sales, and the 3-D cuboids for item, year, and sales, the approach may then be used to estimate the probability of each cell in the 4-D base cuboid for the provided characteristics. This enables the construction of higher-order data cubes from lower-order ones via an iterative method. The method scaled very effectively to accommodate several dimensions. The log-linear model is helpful for data compression (because the smaller-order cuboids combined often require less space than the base cuboid) and data smoothing changes than cell estimates in the base cuboid, in addition to prediction.

## Accuracy of Classifiers

Estimating classifier accuracy is crucial because it enables one to assess how well a particular classifier will categorise untrained data—data on which the classifier has not been trained—in the future. For instance, if historical sales data are used to train a classifier to forecast customer buying behaviour, we would need an assessment of how well the classifier can anticipate future customer purchase behaviour.

## Classifier Accuracy Estimation

Due to the learning algorithm's (or model's) excessive reliance on the data, using training data to create a classifier and then estimating its accuracy might lead to falsely optimistic estimations. Based on randomly chosen partitions of the input data, holdout and cross-validation are two popular methods for evaluating classifier accuracy. In the holdout approach, the input data are randomly partitioned into two independent sets, a training set and a test set. The training set typically receives two thirds of the data, while the test set receives the last third. The classifier is created using the training set, and its accuracy is calculated using the test set. Since just a subset of the original data was utilised to create the classifier, the estimate is pessimistic. The holdout approach is modified by random subsampling, in which it is applied k times. The average of the accuracies gained from each iteration is used to calculate the overall accuracy estimate.

In k-fold cross-validation, the original data are randomly divided into k roughly equal-sized, mutually exclusive subsets, or "folds," S1, S2,..., Sk. Iterative training and testing is done k times. The subset S is put aside in iteration I to serve as the test set, while the other subsets are utilised collectively to train the classifier. In other words, the first iteration's classifier was trained on subsets S1,.., Sk, and tested on Si; the second iteration's classifier was trained on subsets S2,..,Sk, and tested on Si; and so on.

The accuracy estimate is calculated by dividing the total number of correctly classified samples throughout all k iterations by the entire sample size of the starting data. In stratified cross-validation, the folds are stratified such that the samples in each fold have a class distribution that is roughly similar to the distribution in the original data.

Bootstrapping, which evenly samples the training examples with replacement, and leave-one-out, which is k-fold cross-validation with k set to 5, the number of beginning samples, are further techniques for assessing classifier accuracy**.**

**---------------------**

# CHAPTER 7

# CLUSTER ANALYSIS

H S Shreenidhi
Assistant Professor, Department of Computer Science Engineering,
Faculty of Engineering and Technology, JAIN (Deemed-to-be University), Karnataka – 562112
Email Id- hs.shreenidhi@jainuniversity.ac.in

Large volumes of data are analysed using association rule mining to uncover intriguing linkages and relationships. This rule displays the number of times an itemset appears in a transaction. A market-based analysis serves as a common illustration.One of the most important methods used by big organisations to demonstrate correlations between goods is market-based analysis.It enables merchants to discover connections between the products that customers usually purchase together.

**Incremental clustering and disregard for the input record order:**

Some clustering algorithms must create a new clustering from scratch when newly added data (such as database updates) cannot be included into the current clustering structures Scratch. There are clustering methods that are sensitive to the input data's orderAccording on the sequence in which the input items are presented, such an algorithm may produce significantly different clusterings when given a set of data objects.

It is crucial to create algorithms that are insensitive to input order as well as incremental clustering techniques.High dimensionality refers to the presence of several dimensions or characteristics in a database or data warehouse. Numerous clustering algorithms excel at working with data that has two to three dimensions or less. For up to three dimensions, human eyes are capable of making accurate clustering quality judgements. Given that such data might be sparse and extremely skewed, it can be difficult to locate groups of data items in high-dimensional space. Constraint-based clustering: In real-world applications, clustering may be required to be performed under a variety of limitations.

Let's say your task is to choose the sites for a city's new automated teller machines (ATMs). You may cluster homes to make this decision while taking into account limitations like the city's rivers and transportation networks, as well as the kind and amount of consumers in each cluster. Finding data clusters that fulfil certain requirements and exhibit effective clustering behaviour is a difficult challenge.

Users anticipate that clustering findings will be understandable, intelligible, and useful. In other words, clustering may need to be connected to certain semantic applications and interpretations. It's critical to investigate how an application aim may affect the choice of clustering characteristics and methodologies.

**Principal Clustering Techniques:**

    A. Partitioning Techniques
    B. Hierarchical Approaches
    C. Model-Based Methods
    D. Density-Based Methods
    E. Grid-Based Methods

**Methods for Partitioning:**

A partitioning technique divides the data into k divisions, each of which represents a cluster and where k = n. In other words, it divides the data into k groups that together meet the conditions listed below. At least one item must be present in each group, and each object must belong to precisely one group.An initial partitioning is made using a partitioning technique. After that, an iterative relocation method is used to move items from one group to another in an effort to optimise the partitioning.The typical standard for a successful partitioning is that items belonging to the same cluster are near or connected to one another, whilst those belonging to other clusters are separated or quite dissimilar.

**Hierarchical Methods**: A hierarchical technique divides a given collection of data items into a hierarchy. Depending on how the hierarchical breakdown is created, a hierarchical technique may be categorised as either agglomerative or divisive. The bottom-up strategy, also known as the agglomerative approach, begins with each item creating its own group. Up until all of the groups are merged into one or until a termination condition is met, it gradually combines the items or groups that are near to one another.

Starting with all of the items in the same cluster is where the dividing method, also known as the top-down approach, begins. A cluster is divided into smaller clusters with each iteration, until all objects are in one cluster at the end, or until a termination condition is met.The drawback of hierarchical approaches is that once a step (merge or split) is completed, it cannot be undone. This rigidity is advantageous since it reduces processing costs by relieving the need to consider an infinite number of options.

There are two methods for enhancing the effectiveness of hierarchical clustering:

Perform a thorough examination of object "linkages" at each hierarchical partitioning, as in Chameleon, or combine hierarchical agglomeration and other strategies by first using a hierarchical algorithm to group objects into microclusters, and then performing macroclustering on the microclusters using a different clustering technique, like iterative relocation.Density-based methods: The majority of partitioning techniques group items based on their proximity to one another. Such algorithms struggle to detect clusters of any forms and can only locate spherical-shaped clusters.

On the basis of the idea of density, several clustering techniques have been created. For each data point within a given cluster, the neighbourhood of a specified radius must include at least a certain number of points. This is because the overall goal is to keep increasing the given cluster as long as the neighbourhood density above some threshold. This technique may be used to remove noise (outliers) and find clusters of any form.The usual density-based approaches DBSCAN and its extension OPTICS build clusters in accordance with a density-based connectivity analysis. Using an examination of the value distributions of density functions, the DENCLUE approach group's objects into clusters.

**Grid-Based Methods**

The object space is quantized using grid-based approaches into a limited number of cells that make up a grid structure.On the grid structure, or on the quantized space, all clustering operations are carried out. This method's key benefit is its quick processing time, which is often independent of the quantity of data items and only reliant on the number of cells in each dimension of the quantized space. A common illustration of a grid-based approach is STING. Wave Cluster uses both grid-based and density-based wavelet processing for clustering analysis.

**Model-Based Approaches:** Model-based methods posit a model for each cluster and determine which model best fits the data.

By creating a density function that depicts the geographical distribution of the data points, a model-based approach may be able to find clusters. It also results in a method for automatically calculating the number of clusters based on common statistics, accounting for "noise" or outliers, and producing reliable clustering techniques. Clustering High-Dimensional Data and Constraint-Based Clustering are two of the 4.3 tasks in data mining.

**High-Dimensional Data Clustering:**

Because many applications call for the study of objects with a lot of attributes or dimensions, it is a job that is especially crucial in cluster analysis.For instance, written documents may include tens of thousands of phrases or keywords, and DNA microarray data may provide details on the levels of expression of tens of thousands of genes under hundreds of different circumstances.

The curse of dimensionality makes it difficult to cluster high-dimensional data.Some dimensions may not be important. The data get sparser as the number of dimensions rises, rendering the average density of points across the data likely to be low and the distances between pairs of points useless. As a result, a new clustering technique has to be created for high-dimensional data. Two well-known subspace clustering techniques, CLIQUE and PROCLUS, look for clusters inside subspaces of the data rather than over the full data space. Another clustering process called frequent pattern-based clustering extracts separate frequent patterns from subsets of often occurring dimensions. It groups things and creates meaningful clusters using these patterns.

**Clustering with Constraints:**

It is a clustering method that employs application- or user-specified restrictions to produce clustering. A constraint is an excellent tool for expressing to the clustering process the user's expectations or describing the characteristics of the intended clustering results. Different types of limitations may be provided, either by the user or in accordance with the needs of the programme. With impediments present and clustering under user-specified limitations, spatial clustering is used. Additionally, pairwise restrictions are used in semi-supervised clustering to enhance the quality of the final grouping.

**Conventional Partitioning Techniques:**

The most popular and widely used partitioning techniques are the k-Means and k-Medoids approaches.

**Technique Based on Centroids: K-Means Analysis:**

The k-means method divides a collection of n items into k clusters based on the input parameter k, resulting in high intracluster similarity but low interclustersimilarity.The mean value of the items in a cluster, which may be thought of as the cluster's centroid or centre of gravity, is used to quantify cluster similarity.

The following is how the k-means algorithm works.

First, it chooses k objects at random, each of which at this point represents the cluster mean or centre. Based on the separation between each surviving item and the cluster mean, each object is then allocated to the cluster to which it is most comparable. The new mean for each cluster is then calculated. The criteria function is iterated via this method until convergence.

**Methods for Hierarchical Clustering:**

A hierarchical clustering approach groups data items into clusters in a tree-like structure. A pure hierarchical clustering method's capacity to modify after a merge or split decision has been made detracts from its quality. This means that the approach cannot go back and change a specific merge or split decision if it subsequently turns out to have been a bad option.

Depending on whether the hierarchical breakdown is generated top-down or bottom-up, hierarchical clustering approaches may be further categorised as either agglomerative or divisive.

Agglomerative hierarchical clustering: This bottom-up approach begins by grouping every item into its own cluster before combining these small, initial clusters into bigger and larger ones until every object is in a single cluster or until certain termination requirements are met. The majority of hierarchical clustering techniques fall under this heading. Only their interpretations of intercluster similarity vary.

Divisive hierarchical clustering: This top-down method starts with every item in a single cluster, which is the opposite of agglomerative hierarchical clustering.

It breaks the cluster into smaller and smaller fragments until each item creates a cluster on its own or until it fulfils certain termination requirements, such as obtaining the appropriate number of clusters or having cluster diameters that are all within a particular range.

**Cluster analysis with constraints:**

Finding clusters that adhere to user-specified preferences or requirements is known as constraint-based clustering. Constraint-based clustering may use a variety of strategies depending on the kind of constraints.

**There are several types of restrictions**

Restrictions on certain items: We might define restrictions on the things that will be grouped. For instance, in a real estate application, one would want to geographically cluster just those expensive homes that cost more than a million dollars. The collection of items to be clustered is constrained by this restriction. Preprocessing makes it simple to manage, and the issue then becomes a case of unconstrained clustering.

The user may choose to provide a preferred range for each clustering parameter. Restrictions on the use of clustering parameters Typically, clustering parameters are quite particular to the chosen clustering technique. Examples of parameters are e, the radius, and the minimum number of points in the DBSCAN method, or k, the required number of clusters in a k-means algorithm. These user-specified parameters typically only affect the algorithm itself, even though they may have a significant impact on the clustering results.

As a result, their processing and fine tuning are often not regarded as a kind of constraint-based clustering. We may define various distance or similarity functions for certain properties of the items to be clustered, or different distance measures for particular pairings of objects. Constraints on distance or similarity functions. We may utilise various weighting algorithms for height, body weight, age, and skill level when grouping athletes, for instance. Although the mining results will probably vary as a consequence, the clustering process itself may not change. However, in certain circumstances, such modifications might make evaluating the distance function challenging, particularly if the clustering procedure is closely related to them.

Constraints set by the user on the attributes of each cluster: The user may choose to provide preferred qualities for the generated clusters, which may have a significant impact on the clustering process. Using a weak kind of supervision, semi-supervised clustering based on partial supervision may greatly enhance the quality of unsupervised clustering. Pairwise restrictions, or pairings of items tagged as belonging to the same or separate clusters, may be used to achieve this. Semi-supervised clustering is the term used to describe a limited clustering procedure. Outlier Analysis: Some data items deviate from the model or behaviour of the data as a whole. Outliers are those data items that are significantly different from or inconsistent with the rest of the data. Many data mining methods aim to either completely remove or significantly reduce the impact of outliers. However, because one person's noise might be another person's signal, this could lead to the loss of crucial hidden information. For example, in the context of fraud detection, where outliers may signal fraudulent behaviour, the outliers may be of special importance. Outlier mining, also known as outlier discovery and analysis, is a fascinating data mining challenge. It may be used to the detection of fraud, for instance, by identifying anomalous use of telecommunications or credit card services. Additionally, it is helpful in medical analysis for recognising unexpected reactions to different medical treatments as well as in targeted marketing for identifying the purchasing patterns of clients with exceptionally high or extremely low incomes. The following is a description of outlier mining: Find the top k items that are noticeably different, exceptional, or inconsistent from the rest of the data given a collection of n data points or objects and k, the predicted number of outliers. Determining whether data in a given data collection may be perceived as inconsistent and finding an effective way to mine the outliers that have been identified as inconsistent are the two main components of the outliermining challenge.

Statistical Distribution-Based Outlier Identification, Distance-Based Outlier Detection, Density-Based Local Outlier Detection, and Deviation-Based Outlier Detection are a few examples of outlier detection techniques.

**Outlier detection using statistical distributions:**

The statistical distribution-based method to outlier identification makes the assumption that the provided data set has a distribution or probability model (such as a normal or Poisson distribution), and then uses a discordancy test to identify outliers with regard to the model. The properties of the data set, as well as distribution factors like mean and variance and the anticipated number of outliers, must be understood in order to apply the test.

An analysis of statistical discordance looks at two hypotheses:

1. A practical theory
2. One alternate theory

**Distance-Based Outlier Detection:**

The concept of distance-based outliers was developed to overcome the primary restrictions imposed by statistical methodologies. If at least a portion, pct, of the objects in D are at a distance greater than dmin from object o, then object is a distance-based (DB)outlier with parameters pct and dmin, or a DB(pct;dmin)-outlier. Alternatively, we might conceive of distance-based outliers as those things that do not have enough neighbours, where neighbours are determined based on distance from the supplied item, rather than depending on statistical tests. Distance-based outlier identification generalises the principles of discordancy testing for diverse standard distributions in contrast to statistically-based techniques. The extra computation that may be involved in choosing discordancy tests and fitting the observed distribution into a standard distribution is avoided by distance-based outlier identification. It

can be shown for various discordancy tests that if an object, o, is an outlier based on the test provided, then o is also a DB(pct, dmin)-outlier for any appropriately determined pct and dmin. Assuming a normal distribution, if objects that deviate by three standard deviations or more are regarded as outliers, then a DB(0.9988, 0.13s) outlier may be used to extend this concept.

For mining distance-based outliers, many effective algorithms have been created.a method based on indexes

The index-based technique searches for neighbors of each item within a radius of dmin using multidimensional indexing structures, such as R-trees or k-d trees, given a data collection. Let Mbe the maximum number of objects in an outlier's dmin-neighborhood. Therefore, it is evident that object to is not an outlier after M+1 neighbors of object o are discovered. The worst-case complexity of this approach is O(n2k), where n is the number of items in the data collection and k is the number of dimensions.

Even though the procedure of creating an index might be computationally costly, our complexity assessment solely considers the search time. Nested-loop Method: This approach uses nested loops to reduce the amount of I/O operations while maintaining the same computational complexity as the index-based algorithm. It splits the data set into many logical blocks and the memory buffer space in half.I/O efficiency may be achieved by carefully selecting the sequence in which the blocks are put into either half.

**Detection of Local Outliers Based on Density:**

The global or overall distribution of the provided collection of data points, D, is a prerequisite for both statistical and distance-based outlier identification. Data are often not dispersed equally, however. When examining data with distinctly diverse density distributions, these approaches have trouble.

**Deviation-Based Outlier Identification:** To identify unusual objects, deviation-based outlier detection does not need statistical tests or distance-based measurements. Instead, it looks at the primary properties of the items in a group to identify outliers. Outliers are objects that "deviate" from this description. Therefore, in this technique, outliers are often referred to as deviations. We examine two methods for deviation-based outlier identification in this section. The first compares each item in a collection one by one, whereas the second uses an OLAPdata cube strategy.

-------------------

# CHAPTER 8

# ASSOCIATION PATTERN MINING: ADVANCED CONCEPTS

Vikram Singh

Assistant Professor, School of Computer & System Sciences, Jaipur National University, Jaipur, India,
Email Id- vikram@jnujaipur.ac.in

It is challenging to apply association pattern mining algorithms' big output for task-specific activities since they often uncover a large number of patterns. One explanation for this is because the great majority of correlations found can be irrelevant or redundant for a given application. The following advanced techniques are covered in order to make association pattern mining more application-sensitive:

**Summarization:** Association pattern mining often produces a lot of data. A smaller collection of identified item sets is much simpler to comprehend and integrate for an end user. Numerous summarizing techniques will be covered in this chapter, including discovering maximum item sets, closed item sets, and no redundant rules.

**Querying:** Users may want to query a large number of item sets when they are accessible to get more concise summaries. Several specific, query-friendly summarization techniques will be covered in this chapter. The plan is to use a two-phase method that preprocesses the data to provide a summary. Then, this summarization is questioned.

**Including limitations:** In many practical situations, it may be desirable to include application-specific constraints in the process of creating the item set. Contrary to a two-phase "preprocess once query-many" method, a constraint-based algorithm permits the use of considerably lower support levels for mining, even if it may not always yield online results.

**Pattern synthesis**

Regular itemset mining algorithms often find a lot of trends. The output's vastness makes it difficult for consumers to comprehend the findings and draw conclusions that are useful. The great majority of the created patterns often include duplicate information, which is a significant finding. This is so that all subsets of a frequent itemset are likewise frequent, which is ensured by the downward closure feature. In frequent pattern mining, there are several types of compact representations that maintain varying degrees of information about the actual collection of frequent patterns and their support values. Maximal frequent itemsets, closed frequent itemsets, and various approximations are the most well-known representations. The amount of information lost in the summary representation varies across these representations. Closed representations are completely lossless in terms of the membership and support of itemsets. Maximal representations are lossy in terms of the membership of itemsets but lossless in terms of the support. In application-driven circumstances, approximate condensed representations often provide the most useful practical choice while being lossy with regard to both.

**Maximal Patterns**

The idea of maximum item sets was briefly examined. The definition of maximum itemsets is repeated for convenience. With the subsetting technique, all the itemsets may be obtained from the maximum itemsets, but their support values cannot. Because they do not preserve

information about the support values, maximum itemsets are lossy. The idea of closed itemset mining is used to offer a lossless representation in terms of the support values.

Using any common itemset mining technique to locate all itemsets would be a simple approach to find all the maximum itemsets. Then, in a postprocessing step, only the maximum ones may be kept by looking at itemsets in decreasing order of length and eliminating appropriate subsets. Until all itemsets have been either evaluated or eliminated, this procedure is repeated. The maximum itemsets are those that were left in place at termination.The problem with this strategy is that it is ineffective. The number of maximum frequent itemsets may be orders of magnitude fewer than the number of frequent itemsets when the itemsets are extremely lengthy. Designing methods that may directly reduce portions of the search space of patterns during frequent itemset discovery may make sense in such circumstances. With the idea of lookaheads, the majority of tree-enumeration techniques may be altered to reduce the search area for patterns.

## Closed Arrangements

A closed pattern, also known as a closed itemset, may be simply defined as follows:

If none of its supersets have precisely the same support count as X, then an itemset X isclosed. Closed frequent pattern mining methods need itemsets that are both frequent and closed. Take into account a closed itemset X and the collection of closed itemsets S(X) that are subsets of X and have the same support as X. A closed frequent itemset mining technique will only yield X as the item set from S(X). You may refer to the itemsets in S(X) as the equi-support subsets of X.

## A significant finding is as follows:

Observation Let S(X) be the equi-support subsets of a closed itemset called X. The set of transactions T (Y) containing Y is the same for any itemset Y S(X).Additionally, no itemset Z exists outside of S(X) such that the set of transactions in T(Z) is identical to T. (X).The downward closed characteristic of frequent itemsets leads to this discovery. The set of transactions T (Y) is always a superset of T for each suitable subset Y of X. (X). T (X) and T (Y) are the same, but, if the support values of X and Y are the same.

Furthermore, the support of Z X must be the same as that of X if any itemset Z S(X) provides T (Z) = T (X). Z X must be a legitimate superset of X because Z is not a subset of X. If X is assumed to be closed, this would result in a contradiction.

It's critical to realise that any nonredundant counting detail required with regard to each itemset in S is encoded in the itemset X. (X). The single representative itemset is sufficient since every itemset in S(X) describes the same collection of transactions.

## Approximation in Terms of Itemsets

It is also possible to define the approximation in terms of itemsets in a variety of ways, and is a clustering-related term. Conceptually, the aim is to choose representatives J = J1... Jk from the clusters formed from the set of frequent itemsetscalF. The concept of "approximate sets" is likewise dependent on a distance function since clusters are always created with regard to a distance function Dist(X, Y) between itemsets X and Y.

## Pattern Querying

Despite the fact that the compression strategy offers a succinct overview of the common itemsets, there may be circumstances when users want to query the patterns with certain

properties. The relevant sets of patterns for an application are provided by the query replies. The whole collection of patterns is often substantially larger than this relevant selection. Here are a few instances:

    A. List all frequently occurring patterns that include X and have a minimum support of minsup.
    B. List all association rules that include X and have a minimal amount of support (minsup) and confidence (minconf).

One option is to thoroughly scan every common itemset and report the ones that adhere to the user-specified restrictions. However, when there are many recurring patterns, this is quite inefficient. When searching for intriguing subsets of patterns, there are two groups of strategies that are typically employed:

**Preprocess-once query-many paradigm**: The first strategy entails mining all of the itemsets with little help and arranging them into a lattice- or hierarchical data structure. The first stage only has to be completed once offline, therefore there may be enough computing resources available. To increase the amount of patterns maintained in the initial phase, a low degree of support is used. With the summary generated in the first step, several questions may be answered in the second part.

**Constraint-based pattern mining:** In this method, the miner is explicitly given the user-specified restrictions. Even while such a method may take longer for each query, it enables pattern mining at far lower support levels than the first method could. This is due to the fact that the restrictions may decrease the pattern sizes in the itemset discovery algorithm's intermediate phases and can, as a result, allow the discovery of patterns at considerably lower support values than a (unconstrained) preprocessing phase.

**Preprocess-once-ask-many Paradigm**

For the purpose of handling basic queries, this specific paradigm is particularly useful. In these situations, the trick is to first identify all of the common patterns at a relatively low support value. It is then possible to assemble the generated itemsets into a data structure for querying. The itemset lattice, which may be thought of as a graph data structure for querying, is the most basic kind of data structure. However, data structures may also be used to query itemsets (Figure 8.1).
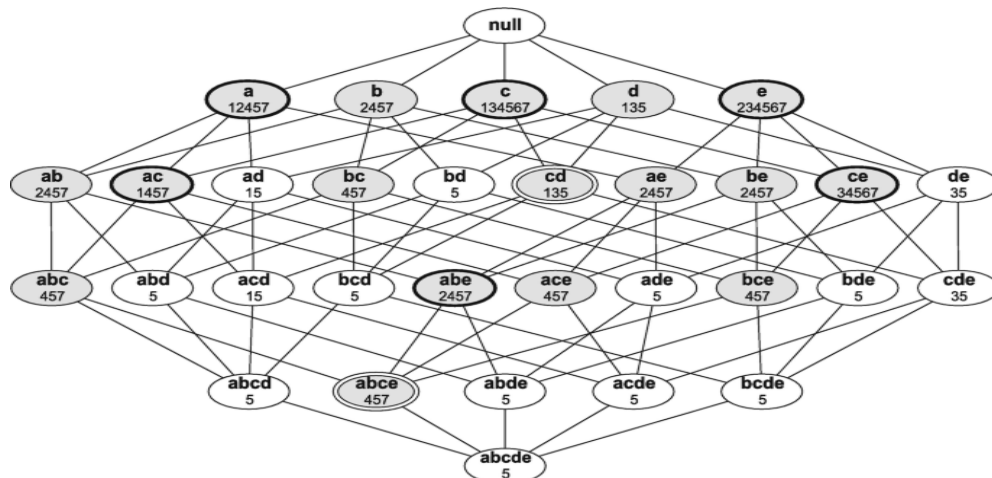


**Figure 8.1: Preprocess-once-ask-many Paradigm**

Adapted from the bag-of-words representation literature for information retrieval.

**Making Use of the Itemset Lattice**

The space of itemsets may be described as a lattice, as was mentioned in the preceding chapter.

Itemsets above the dashed boundary are numerous, and itemsets below the border are rare; this is a convenience replication of Figure 8.1 from the previous chapter.Theitemsets are mined at the lowest level of support s feasible in the preprocess-once query-many paradigm so that a significant percentage of the lattice (graph) of itemsets may be kept in main memory.

Running time is not a top priority at this level since it is a preprocessing step. For effective traversal, the lattice's edges are represented as pointers. A hash table also links theitemsets to the graph's nodes.A variety of crucial characteristics of the lattice, such downward closure, make it possible to identify nonredundant association rules and patterns.

This framework is capable of successfully responding to a wide range of support minsup questions. Here are a few instances:

1. One uses the hash table to map to the itemset X to find all itemsets that include set X at a certain level of minsup. The required supersets of X are then found and reported after traversing the lattice. By utilising a traversal in the other direction, a similar strategy can be used to find all the common itemsets present in X.

2. By identifying nodes that do not have edges to their immediate supersets at the user-specified minimum support level minsup, it is feasible to find maximum itemsets immediately during the traversal.

3. By moving through the lattice structure both uphill and downward from X for a certain number of steps, it is feasible to locate nodes within a defined hamming distance of X and a given minimum support (Figure 8.2).



**Figure 8.2: Making Use of the Item set Lattice**

With the use of this method, no redundant rules may also be determined. For any itemset Y Y, for instance, the confidence and support of the rule X Y cannot be larger than that of the rule X Y. The rule X Y is thus unnecessary in relation to the rule X Y.

Strict redundancy is what this is known as. Additionally, the rule I Y Y is redundant with regard to the rule I Y Y solely in terms of the confidence for any item set I. Simple redundancy is what is meant by this. In terms of both simple redundancy and stringent redundancy, the lattice structure offers an effective method for locating such nonredundant rules. For precise search techniques on how to locate such rules, the reader is directed to the bibliographic notes.

**Using Data Structures to Support Querying**

The usage of disk-resident representations for querying is preferred in certain circumstances. The memory-based lattice traversal technique is probably ineffective in these circumstances. The inverted index and the signature table are the two most often used data structures. The main disadvantage of employing these data structures is that, unlike the lattice structure, they do not provide an orderly exploration of the collection of common patterns.

Both transactions and itemsets may utilise the data structures. Some of these data structures, such as signature tables, work especially well for itemsets, however, since they explicitly take use of correlations across itemsets for effective indexing. Keep in mind that correlations are more important for itemsets than for raw transactions. Below, both of these data formats are briefly explained.

The inverted index is a kind of data structure that is used for locating sparse set-valued data, such as the text that is represented as a bag of words. Frequent patterns may be effectively retrieved using an inverted index since they are also sparse sets drawn across a much wider universe of things.

An individual itemset-id is given to each itemset. A hash function may be used to create this quickly. The tid used to identify transactions and this itemset-id are identical. The itemset-id may be used to index a secondary data structure where the itemsets themselves are kept. The same hash algorithm that was used to generate the itemset-id may be used to produce a hash table as this secondary data structure. There is a list for each item in the inverted list. A list of itemset-ids is referenced by each item.

Particularly helpful for inclusion searches over tiny groups of things is the inverted representation.

Consider a search for all sets of things that include the tiny set of objects X. On the disc are kept the inverted lists for each item in X. These lists' intersection is identified. The itemsets are not provided; just the relevant itemset-ids are. The pertinent itemsets may be retrieved from disc and reported, if needed. The recovered itemset-ids must be used to get access to the secondary data structure on disc in order to do this. The inverted data structure has an extra burden since it could need to make a random disc request. Such a method may not be workable for massive query replies.

Inverted lists work well for inclusion searches over small sets of items but fall short for similarity queries over larger itemsets. The inverted index has the drawback of treating each item separately and failing to take advantage of the strong correlations among the items in the itemset. Additionally, retrieving whole itemsets rather just their ids is more of a challenge. The signature table is the preferred data structure in these circumstances.

Applications for Putting Associations to Work numerous real-world situations may benefit from association pattern mining, which has many applications. These applications will be briefly discussed in this section.

**Connection with Other Data Mining Issues**

Classification, grouping, and outlier identification are a few more data mining issues that the association model is closely connected to. These data mining issues can be successfully solved via association patterns. These connections will be briefly examined in this section.

The chapters on these various data mining issues also cover many of the pertinent algorithms.

**Application to Classification, Section**

The issue of association pattern mining and categorization are closely connected. Association-rule mining and rule-based classifiers are closely related concepts. A basic summary of these sorts of classifiers is given below;

X is the antecedent and Y is the consequent in the rule X Y. In associative classification, the antecedent comprises the feature variables, and the consequent Y is a single item that corresponds to the class variable. The training data is mined for these rules. The conventional support and confidence indicators are often not used to decide the rules. Instead, it is necessary to identify the regulations that discriminate the most against the various classifications. Think about an itemset X and the classes c1 and c2, for instance. It makes intuitive sense that the itemset X is discriminative between the two classes if there is a significant difference in the absolute confidence levels of the two rules, X c1 and X c2. Therefore, such discriminatory regulations need to be decided by the mining process.It's interesting to note that the association structure may be used to the classification issue rather well with just a few relatively minor modifications. The Classification Based on Associations (CBA) framework is an illustration of one such classifier.

**Application to Clustering**

Since association patterns identify strongly connected subsets of variables, they may be used to identify dense areas in quantitative data following discretization.On top of the modified data, association patterns are found. These areas' overlapping data points are referred to be subspace clusters. Naturally, this method reports clusters that have a lot of overlap with one another. The groups that are produced, however, fit the data's dense areas and provide important details about the underlying clusters.

**Applications to Outlier Detection**

Outliers in market basket data have also been identified using association pattern mining.The important concept is that most association patterns in the data "cover" the majority of the transactions, whereas the outliers are defined as transactions that are not. When the matching association pattern is present in a transaction, it is said to be covered by that association pattern. This method is very helpful when the data is high dimensional and conventional distance-based techniques are difficult to use. Such a strategy is especially successful since transaction data is naturally high dimensional.

**Market Basket Analysis**

Market basket analysis is the model issue for which the association rule mining problem was initially put up. The goal of this challenge is to identify guidelines for consumer purchasing behaviour. For a retailer, being aware of these restrictions may be quite helpfulFor instance, if an association rule indicates that the sale of beer implies the sale of diapers, a retailer may utilise this knowledge to enhance the positioning of their products on the shelves and their marketing strategies. For market basket analysis, unusual or surprising rules in particular are the most instructive. Such choices are the main focus of many of the conventional and alternative methods for market basket analysis.

**Collaboration and Recommendations Filtering**

The two aforementioned applications have a common ground in the broader issue of collaborative filtering and recommendation analysis. With collaborative filtering, customers are given suggestions based on the purchasing habits of other users who share their

characteristics. Localized pattern mining is very helpful in this situation. The goal of localised pattern mining is to identify patterns within certain segments after clustering the data into these segments. The patterns from each segment often provide a sharper picture of the patterns among similar consumers and are more resilient to noise from the worldwide data dispersion. A specific pattern for movie titles, such as "Gladiator, Nero, Julius Caesar," may not have enough support globally, for instance, in a system that recommends movies. However, such a pattern could have enough support among buyers who have similar interests in historical dramas. Applications like collaborative filtering use this method. The necessity to concurrently identify the clustered segments and the association rules makes the localised pattern mining issue significantly more difficult.

**Web Log Evaluation**

The study of web logs is a typical application for pattern mining techniques. For instance, the collection of pages visited during a session is identical to a market-basket data set of transactions. Frequent visit to a group of Web sites reveals interesting connections in user behaviour with regard to certain Web pages. Site managers may use these information to enhance the website's architecture. For instance, it could be beneficial to build a connection between two Web sites if they are regularly browsed together during a session but are not already connected. The most complex types of Web log analysis often go beyond the set-wise structure of frequent itemset mining and focus on the temporal features of logs.

**Bioinformatics**

Numerous cutting-edge bioinformatics technologies, such mass spectrometry and microarray technologies, enable the gathering of various sorts of extremely high-dimensional data sets. Gene-expression data, which can be written as a $n$ $d$ matrix with a significantly high number of columns $d$ compared to normal market basket applications, is a famous example of this kind of data. The number of columns in a microarray application may reach 100,000 with relative ease. Numerous applications for finding regular patterns in these data sets include the identification of important biological characteristics that are encoded in them. Long pattern mining techniques, such maximum and closed pattern mining, are highly helpful in these situations.

**Additional Uses for Complex Data Types**

Algorithms for frequent pattern mining have been expanded to include more intricate data types as temporal, geographic, and graph data. These sophisticated data types are covered in distinct chapters in this book. Here is a quick explanation of some more sophisticated applications:

**Analyzing temporal web logs:** The analytical method is substantially enhanced by the incorporation of temporal data from Web logs. When future events can be anticipated from the existing pattern of events, for instance, specific patterns of accesses that show up often in the logs may be utilised to create event prediction models.

**Patterns of spatial co-location:** Patterns of spatial co-location provide light on the spatial connections between various people. To such situations, frequent pattern mining algorithms have been extended.

**Applications of chemical and biological graphs:** The identification of structural patterns sheds light on the characteristics of many real-world situations, such as chemical and biological components. Models for classification are also made using these patterns.

**Software bug analysis:** Call graphs are often used to depict the architecture of computer systems. The study of the call graphs' common patterns and significant departures from these patterns reveals the flaws in the underlying programme.

-----------------------

# CHAPTER 9

# MINING OBJECT, SPATIAL, MULTIMEDIA

Vikram Singh

Assistant Professor, School of Computer & System Sciences, Jaipur National University, Jaipur, India,
Email Id- vikram@jnujaipur.ac.in

Scientific research and engineering design are two examples of sophisticated, data-intensive applications that need the storage, access, and analysis of complex but generally organised data items.These objects can't be represented in data relations as straightforward, consistently organised records (i.e., tuples). Object-relational and object-oriented database systems were designed and developed in response to such application needs. Both types of systems are concerned with the effective storage and access of significant quantities of disk-based complicated structured data items. These systems classify a large collection of complex data items into groups called classes, which are then arranged into hierarchies called class/subclass groups. An object-identifier, a set of attributes that may contain complex data structures, set- or list-valued data, class composition hierarchies, and multimedia data, and a set of methods that specify the computational procedures or rules specific to the object class are all associated with each object in a class. The efficient indexing, storing, accessing, and manipulating of complex objects in object-relational and object-oriented database systems has been the subject of substantial study in the area of database systems. Numerous publications on database systems, particularly those on object-oriented and object-relational database systems, explore the technologies addressing these problems.

The systematic analysis and mining of such data is a step beyond the storage and access of massively scaled, complex object data. Building multidimensional data warehouses for complex object data and performing online analytical processing (OLAP) in such data warehouses are two of the key objectives in this, along with developing efficient and scalable techniques for knowledge mining from object databases and/or data warehouses.Since these data make up the majority of the new types of complex data objects, the second duty is mostly addressed by the mining of certain types of data (such as geographical, temporal, sequence, graph- or tree-structured, text, and multimedia data). In this chapter, much like in Chapters 8 and 9, we continue to examine techniques for mining complicated data. Therefore, the major topics of discussion in this part will be how to build object data warehouses and conduct OLAP analysis on data warehouses for such data.The restriction on the permitted data types for dimensions and measures is a significant drawback of many commercial data warehouse and OLAP systems for multidimensional database analysis. The majority of data cube implementations limit measurements to basic, aggregated values and constrain dimensions to non-numeric data. This section discusses how to 77nalyzing77 complicated structured objects for OLAP and mining in object databases, as well as how to build object cubes for OLAP and multidimensional data analysis for complex objects. It is crucial to research how each component of such databases may be 77nalyzing77n, as well as how the 77nalyzing77n data can be utilized for multidimensional data analysis and data mining, in order to ease generalization and induction in object-relational and object-oriented databases.

## Application of Structured Data in General

The ability to store, read, and represent complex structure-valued data, such as set- and list-valued data and data with nested structures, is a key characteristic of object-relational and object-oriented databases.An attribute with a set value may be either homogeneous or

heterogeneous in kind. Set-valued data can typically be 78nalyzing78n in one of two ways: (1) by generalizing each value in the set to its corresponding higher-level concept, or (2) by determining the general behavior of the set, such as its size, its types or value ranges, its weighted average for numerical data, or the major clusters it forms. Additionally, generalization may be carried out by using several generalization operators to investigate distinct generalization routes. A heterogeneous set is the end result of generalization in this instance.

### Spatial Combination and Approximation

### Generalizing multimedia data

Another major method of generalization is by aggregation and approximation. They are particularly helpful when trying to generalize qualities with vast sets of values, intricate structures, or multimedia or geographical data.Consider spatial data as an example. According to land use, we would want to group together certain geographic sites into clusters that may be commercial, residential, industrial, or agricultural. A group of geographic locations must often be combined for such 78nalyzing78n78on using spatial operations, such as spatial union or geographical clustering techniques. The methods of aggregation and approximation are crucial for this kind of 78nalyzing78n78on. In a spatial merge, it is required to ignore certain dispersed regions with various kinds if they are unnecessary to the research and combine the regions of comparable types within the same general class as well as calculate the total areas, average density, or other aggregate functions. Spatial aggregation and approximation may also be used as data 78nalyzing78n78on operators by other spatial operators, such as spatial-union, spatial-overlapping, and spatial-intersection (which may need the fusion of dispersed tiny areas into large, clustered regions).

Class/Subclass Hierarchies and Object Identifiers Generalization. Even after the data has been structurally 78nalyzing78n, it does not alter. The 78nalyzing78n78on of an object, however, may be carried out by referring to its related hierarchy since objects in an object-oriented database are classified into classes, which are structured into class/subclass hierarchies. A generic object identification may therefore be described as follows. The object's identification is first 78nalyzing78n to the lowest subclass identifier to which the object belongs. By moving up the class/subclass hierarchy, the identification of this subclass may then be 78nalyzing78n to a higherlevel class/subclass identifier. By moving up the related class/subclass hierarchy, a class or subclass may be 78nalyzing78n to its corresponding superclass.

Some characteristics or methods of an object class are not explicitly stated in the class but are inherited from higher-level classes of the object because object-oriented databases are arranged into class/subclass hierarchies. Numerous inheritance is a feature of several object-oriented database systems, allowing properties to come from multiple superclasses when the class/subclass "hierarchy" is arranged in the form of a lattice. The object-oriented database's query processing may be used to determine an object's inherited attributes. It is superfluous from the perspective of data 78nalyzing78n78on to differentiate between data stored inside the class and data inherited from its superclass. The data mining procedure will consider the inherited data in the same way as the data contained in the object class and conduct 78nalyzing78n78on as appropriate as long as the set of relevant data are gathered via query processing.

Methods play a key role in object-oriented databases. Additionally, objects may inherit from them. The use of techniques may be used to derive several behavioral characteristics of things. It is hard to generalize a method since a method is often described by a computational

function or by a set of deduction rules. On the data obtained from the use of the approach, however, generalization may be done. In other words, generalization may be done on the data after the collection of task-relevant data has been obtained by using the approach.

## Class composition hierarchy generalization

An object's attribute may be made up of or described by another object, some of whose attributes may then be made up of or described by additional objects, constituting a hierarchy of class composition. It is possible to think about 79nalyzing79n79onon a class composition hierarchy as 79nalyzing79n79onon a collection of layered structured data (which may be limitless if the layering is recursive).

In theory, a reference to a composite object might go down the associated class composition hierarchy via a lengthy chain of references. The semantic connection between the original object and the referred composite item is often weaker the longer the series of references travelled. As an example, the attribute automobiles owned by an object class student may relate to a different object class car, which could have the attribute auto dealer, which might refer to the manager and kids of the dealer. It goes without saying that it is doubtful that there are any intriguing general patterns between a student and the kids of her auto dealer boss. The descriptive attribute values and methods of a class of objects should thus be 79nalyzing79n, with only limited reference to its closely related component through its closely related links in the class composition tree. In other words, 79nalyzing79n79on should be done on the objects in the class composition hierarchy that are semantically related to the class or classes that are currently being focused, not on those that have only distant and relatively weak semantic linkages, in order to uncover interesting knowledge.

## Building and Mining of Object Cubes

Data 79nalyzing79n79on and multidimensional analysis are applied to classes of objects rather than specific items in an object database. The main challenge is how to make the 79nalyzing79n79on processes work together among various attributes and methods in the class, since a group of objects in a class may share multiple characteristics and methods and the 79nalyzing79n79on of each property and method may apply a series of 79nalyzing79n79on operators (es).

Consider the idea that applying a series of class-based 79nalyzing79n79on operators to various characteristics may be considered as a generalization-based data mining process. Generalization may go on until the final class only comprises a few 79nalyzing79n objects that can be succinctly and broadly described as a rule. Examining each attribute (or dimension), 79nalyzing79n79 each attribute to simple-valued data, and building a multidimensional data cube—referred to as an object cube—can be used to efficiently accomplish the 79nalyzing79n79on of multidimensional attributes of a complex object class. After being built, an object cube may be subjected to multidimensional analysis and data mining in a way similar to relational data cubes.

Be aware that 79nalyzing79n79 a collection of values to single-valued data is not necessarily desirable from the perspective of an application. Think about the attribute keyword, which may include a number of terms that describe a book. It is not very logical to reduce this group of keywords to a single value. It is challenging to create an object cube having the term dimension in this situation. When addressing the creation of spatial data cubes in the next part, we will touch on certain advancements made in this regard. The development of methods for processing set-valued data in the creation of object cubes and object-based multidimensional analysis, however, continues to be a difficult research problem.

**Divide-and-Conquer Generalization-Based Mining of Plan Databases**

We look at a scenario of mining essential patterns of successful activities in a plan database using a divide-and-conquer technique to demonstrate how 80nalyzing80n80on may be crucial in mining big datasets.

A varied series of acts make up a plan. A vast collection of plans is known as a planbase, or simply a plan database. The process of extracting important patterns or information from a planbase is known as plan mining. Plan mining may be used to detect important patterns from action sequences in the repair of autos or to determine the trip patterns of business 80nalyzing in an air flight database. Plan mining is distinct from sequential pattern mining, which involves thoroughly 80nalyzing a large number of regularly recurring sequences. Plan mining, on the other hand, is the process of extracting significant or relevant generalized (sequential) patterns from a plan base. The plan mining method may be improved in a number of ways. For instance, a minimal support threshold may be used to gauge the degree of generalization and make sure that a pattern covers a significant number of examples, just as it is done in association rule mining. Plan mining may also investigate other operators, such less than.

Other versions involve mining sequence patterns including multidimensional characteristics, such as the patterns involving both airport size and location, or extracting connections from subsequences. Prior to looking at the combined sequence patterns, such dimension-combined mining also necessitates the high-level generalization of each dimension.

**Spatial Data Analysis**

Maps, processed remote sensing or medical imaging data, VLSI chip layout information, and other vast quantities of space-related data are all kept in spatial databases. There are several characteristics that set spatial databases different from relational databases. They typically contain topological and/or distance information, which is 80nalyzing by complex, multidimensional spatial indexing structures. To access them, one must use methods for accessing spatial data, which frequently call for techniques for spatial reasoning, geometric computation, and spatial knowledge representation. The process of extracting information, geographical connections, or other intriguing patterns from spatial databases but not expressly recorded there is known as spatial data mining. Such mining necessitates the combination of geographical database technology with data mining. It may be used to comprehend spatial data, find correlations between spatial and nonspatial data, build spatial knowledge bases, rearrange spatial databases, and improve spatial searches. Geographic information systems, geomarketing, remote sensing, picture database exploration, medical imaging, navigation, traffic management, environmental research, and many more fields that employ spatial data are predicted to benefit greatly from its 80nalyzing80n. Due to the vast volume of spatial data, the variety of geographic data types, and the difficulty of spatial access methods, the investigation of effective spatial data mining approaches is a critical task.

A well-liked method for examining geographic data and 80nalyzing geographical data is statistical spatial data analysis. While the phrase "spatial statistics" is often connected with discrete space, the term "geostatistics" is frequently associated with continuous geographic space. One often assumes statistical independence between various pieces of data when using a statistical model to analysenonspatial data. In contrast to traditional data sets, spatially distributed data does not exhibit this independence because, in reality, spatial objects are frequently related or, to be more precise, spatially co-located, meaning that the closer two objects are to one another, the more likely it is that they will share similar characteristics. For instance, locations that are geographically adjacent to one another are likely to have

comparable natural resource, climatic, temperature, and economic conditions. "Everything is connected to everything else, but adjacent things are more related than distant things," is even thought of as the fundamental rule of geography. Geographical autocorrelation is a term used to describe this kind of tight spatial interdependency. Spatial statistical modelling techniques have been successfully created based on this idea. With a focus on efficiency, scalability, collaboration with database and data warehouse systems, enhanced user involvement, and the discovery of new forms of information, spatial data mining will further develop spatial statistical analysis techniques and expand them for massive volumes of spatial data.

## Construction of a spatial data cube and spatial OLAP

By combining combine geographical data with relational data to create a data warehouse that supports spatial data mining. A subject-oriented, integrated, time-variant, nonvolatile collection of both spatial and nonspatial data is known as a spatial data warehouse, which is used to enable spatial data mining and decision-making processes using spatial data.The creation and usage of geographic data warehouses present a number of difficult problems. The integration of geographical data from many sources and systems presents the first difficulty. Different industrial companies and governmental organisations often store spatial data using different data formats. Data formats vary by vendor (e.g., ESRI, MapInfo, Intergraph) as well as by structure (e.g., raster- vs. vector-based geographic data, object-oriented vs. relational models, various spatial storage and indexing structures). There has been a significant amount of work done on the interchange and integration of heterogeneous geographic data, which has prepared the way for the creation of spatial data warehouses and the integration of spatial data. The implementation of quick and adaptable online analytical processing in geographical data warehouses is the second barrier. The three model of the star schema is effective because it offers a clear and ordered warehouse structure and makes OLAP operations easier, is a preferred option for modelling geographical data warehouses. However, both dimensions and metrics could have spatial components in a spatial warehouse.

## A spatial data cube has three different kinds of dimensions:

Only nonspatial data are present in a nonspatial dimension. For the warehouse in Example 10.5, nonspatial dimensions for temperature and precipitation may be created since they both include nonspatial data with nonspatialgeneralisations (such "hot" for temperature and "wet" for precipitation).A dimension that has spatial data at the most fundamental level but transitions to nonspatial data at a particular high level is known as a spatial-to-nonspatial dimension. For the United States map, for instance, the spatial dimension city communicates geographic information. Let's imagine the string "pacific northwest" serves as the dimension's generalised geographical representation of, say, Seattle. Despite the fact that "pacific northwest" is a geographical idea, it is not represented spatially (since, in our example, it is a string). It serves as a nonspatial dimension as a result.A spatial-to-spatial dimension is one whose high-level generalised data are all spatial at the primitive level. In a spatial data cube, we separate two categories of measures:Only numerical data are present in a numerical measure. The monthly income of an area, for instance, may be one measure in a geographic data warehouse such that a roll-up could calculate the overall revenue by year, by county, and so on. A group of references to spatial objects are included in a spatial measure. Only numerical measurements and nonspatial dimensions are included in a nonspatial data cube. A spatial data cube's OLAP operations, such as drilling and pivoting, may be implemented similarly to those for non-spatial data cubes if the cube has spatial dimensions but no spatial measurements.Regarding the calculation of spatial measurements in the creation of spatial data cubes, at least three options are available:

Do not execute precomputation of spatial measures in the spatial data cube; instead, collect and save the associated spatial object pointers. This may be done by keeping a pointer to a group of spatial object pointers in the associated cube cell and, if needed, calling up the collection to execute a spatial merge (or other calculation) on the related spatial objects. This method is a good choice if only spatial display is needed (i.e., no actual spatial merge needs to be done), if only a small number of regions need to be merged in each pointer collection (so that the on-line merge is not very expensive), or if the on-line spatial merge computation is quick (recently, some efficient spatial merge methods have been developed for fast spatial OLAP).It is nevertheless advised to precompute part of the spatially related areas since OLAP results are often utilised for online spatial analysis and mining. This will speed up the analysis.Create a rough approximation of the spatial measurements in advance and save it in the spatial data cube. Under the presumption that it takes minimal storage capacity, this option is excellent for a rough view or crude assessment of the outcomes of a spatial mergingdata sets that are not part of the answer sets may be included in the exam, however itshould not permit a false-negative test that would rule out certain feasible solutions.Using a few rudimentary spatial evaluation algorithms, such as an MBR structure (which registers only two spatial points rather than a collection of complex polygons), we can first gather candidates that meet the minimum support threshold for mining spatial associations related to the spatial predicate close to. Then, we can evaluate the relaxed spatial predicate, g close to, which is a generalised close to covering a broader context that includes close to, touc, and more.

The MBRs that surround two spatial objects that are near together must also be close together, matching g close to. The opposite is not necessarily true, either; if the enclosing MBRs are near together, the two spatial objects may or might not be as well. Since only those that pass the quick test need to be further analysed using more costly spatial computing techniques, the MBR trimming is a false-positive testing tool for proximity.Only the patterns that are prevalent at the approximation level will need to be analysed by more intricate and precise, but more costly, spatial computing thanks to this preprocessing.

Along with mining spatial association rules, one can also want to find clusters of certain characteristics that regularly show together on a geospatial map. Mining spatial co-locations is basically the task at hand in this situation. One might think of finding spatial co-locations as a specific example of mining spatial relationships. However, fascinating traits are likely to coexist in close-by locations based on the spatial autocorrelation property. In light of this, exploring spatial co-location may be just what you need. Similar to what has been done for mining spatial association rules, efficient techniques for mining spatial co-locations may be created by investigating methodologies like Aprori and progressive refinement.

**Techniques for Spatial Clustering**

In a big, multidimensional data collection, spatial data clustering locates clusters, or heavily populated locations, according to some distance measurement.

**Classification of Space and Analysis of Space Trends**

In order to create categorization schemes that are relevant to certain geographical qualities, such as the vicinity of a district, highway, or river, spatial classification examines spatial objects. Spatial objects have several characteristics, such as having a university there, having interstates nearby, being close to a lake or the ocean, and so on. These characteristics may be utilised to identify intriguing categorization schemes and do relevance analysis.The identification of changes and trends along a geographical dimension is another topic that is addressed by spatial trend analysis. Trend analysis often identifies changes over time, such as

shifts in temporal patterns in time-series data. By replacing time with space, spatial trend analysis examines how nonspatial or spatial data tend to change with space. For instance, we may notice a pattern of changing economic conditions as we go farther from a city's heart or a tendency of changing temperature or flora as we travel further from an ocean. Regression and correlation analysis techniques are often used for these types of investigations, making use of geographical data structures and spatial access techniques.There are several applications where patterns change across time and space. For instance, there are both time and space-related traffic flows on roads and in urban areas.Time and space are very highly connected to weather patterns. The research into spatiotemporal data mining is still in its early stages, despite the existence of a few fascinating papers on spatial categorization and geographical trend analysis. It is necessary to investigate more approaches and uses for geographical categorization and trend analysis, particularly those involving time.

**Mine raster databases**

Systems for managing spatial databases often deal with vector data, which includes points, lines, polygons (regions), and the combinations of these shapes, such as networks or divisions. Maps, design graphs, and three-dimensional representations of the organisation of protein molecule chains are typical examples of this kind of data. However, a significant portion of space-related data—including satellite pictures, data from remote sensing, and tomography—is in digital raster (image) formats. Investigating data mining in raster or picture databases is crucial. The next section on the mining of multimedia data examines techniques for mining raster and picture data.

**Mining Multimedia Data**

A multimedia database system collects and maintains a large collection of multimedia data, including text, markups for text, images, graphics, voice, audio, video, and documents. In light of the widespread usage of audiovisual equipment, digital cameras, CD-ROMs, and the Internet, multimedia database systems are becoming more and more prevalent. The NASA EOS (Earth Observation System), different picture and audio-video databases, and Internet databases are examples of common multimedia database systems.Our examination of multimodal data mining in this part concentrates on picture data mining. The two sections that follow examine text data mining and web data mining. Here, we describe many techniques for mining multimedia data, such as multidimensional analysis, association mining, classification and prediction analysis, and similarity search in multimedia data.

**Similarity Search in Multimedia Data**

That is accurate. We focus on two primary families of multimedia indexing and retrieval systems for similarity searching in multimedia data:

(1) Content-based retrieval systems, which support retrieval based on the image content, such as colour histogram, texture, pattern, image topology, and the shape of objects and their layouts and locations within the image.

(2) Description-based retrieval systems, which build indices and perform object retrieval based on image descriptions, such as keywords, captions, size, and time of creation. If done manually, description-based retrieval is time-consuming.

Results from automation are often of low quality. Assigning keywords to photographs, for instance, may be a challenging and arbitrary operation. Because image-surrounded text information and Web linkage information can be used to extract the proper description and group images describing the same theme together, recent developments in Web-based image

clustering and classification methods have improved the quality of description-based Web image retrieval.Visual characteristics are used to index photos in content-based retrieval, which encourages object retrieval based on feature similarity, which is extremely desired in many applications.There are typically two types of queries in a content-based image retrieval system: image feature definition queries and imagesample-based queries.The photos that are comparable to the provided image sample are found via image-sample-based searches. The feature vector (or signature) that was derived from the sample is compared with feature vectors of previously extracted and indexed pictures in the image database in this search. Images that are similar to the example picture are returned as a result of this comparison. Image characteristics are specified or sketched in image feature specification queries, such as colour, texture, or form, which are then converted into feature vectors and compared to the feature vectors of the photos in the database. Numerous industries, including as e-commerce, television production, weather forecasting, and medical diagnostics, use content-based retrieval.

There are certain systems that enable both sample-based and image feature definition inquiries, such as QBIC (Query by Image Content). Systems exist that enable both content- and description-based retrieval.

Several methods for similarity-based retrieval in picture databases based on image signature have been put forward and researched:

Color histogram-based signature: In this method, independent of the picture's size or orientation, the signature of an image consists of colour histograms based on the colour composition of the image. Shape, picture topology, or texture information are not included in this technique. Thus, even if two pictures may be utterly unrelated conceptually yet have similar colour compositions and extremely diverse forms or textures, they may be seen as comparable.

**Multi-feature constructed signature:** Using this method, an image's signature is made up of a variety of characteristics, including the colour histogram, shape, picture topology, and texture. Images are indexed based on the information that was recorded with the extracted picture characteristics. Frequently, distinct distance functions for each characteristic may be constructed and then merged to get the final results. To look for pictures with these (similar) properties, multidimensional content-based search often employs one or a few probing features. Therefore, it may be used to look for photographs with comparable themes. The most often used strategy in practice is this one. Wavelet-based signature: This technique utilises an image's dominant wavelet coefficients to identify it. In a single integrated framework, wavelets collect information about shape, texture, and picture structure. 1 This boosts effectiveness and eliminates the need for several search primitives (unlike the second method above). This approach may, however, be unable to distinguish between photos containing identical items if the objects vary in position or size since it computes a single signature for a whole image.

**Wavelet-based signature with region-based granularity:** With this method, computation and comparison of signatures are done at the level of individual areas rather than the full picture. This is based on the fact that although portions in two photographs that are similar to one another may include identical regions, they may also be translated or scaled differently. As a result, the percentage of the area of the two pictures that are shared by matching pairings of areas from Q and T may be used to determine the similarity between the query image Q and the target image T. Such a region-based similarity search can locate photographs with related items, even when they have been scaled or translated.

**Analysis of Multidimensional Multimedia Data**

Multimedia data cubes may be generated and built similarly to conventional data cubes using relational data in order to enable the multidimensional analysis of big multimedia datasets. Additional dimensions and measurements for multimedia information, such as color, texture, and shape, may be included in a multimedia data cube. Examining Multimedia Miner, a prototype of a multimedia data mining system that augments the DB Miner system by managing multimedia data. The following describes how the sample database used to test the Multimedia Miner system was built. A feature descriptor and a layout descriptor are both included in each picture. Only the picture's descriptors are kept in the database; the actual image is not. The image file name, image URL, image type (such as gif, tiff, jpeg, mpeg, bmp, or avi), a list of all known Web pages referring to the image (i.e., parent URLs), a list of keywords, and a thumbnail used by the user interface for browsing images and videos are all included in the description information. Each visual attribute is represented by a collection of vectors in the feature descriptor. The three vectors are an MFC (Most Frequent Color) vector, an MFO (Most Frequent Orientation) vector, and a colour vector with the colour histogram quantized to 512 colours (8 8 8 for R, G, and B). The five most common colours and edge orientations are represented by five colour centroids and five edge orientation centroids, respectively, in the MFC and MFO. The edge orientations that are employed include 0, 22.5, 45, 67.5, 90, and so on. A colour layout vector and an edge layout vector are both included in the layout description. All photos, regardless of their original size, are given an 88 grid. The edge layout vector stores the number of edges for each orientation in each of the 64 cells, whereas the colour layout vector stores the most prevalent colour for each of the 64 cells. It is simple to generate grids with other dimensions, such as 44, 22, and 11.

Similar to HTML tags in Web pages, the Image Excavator part of Multimedia Miner leverages contextual information about images to extract keywords. It is feasible to build hierarchies of keywords mapped onto the directories where the picture was discovered by navigating online directory structures like the Yahoo! directory. In the multimedia data cube, these graphs serve as idea hierarchies for the dimension keyword.

Multimedia data cubes may come in a variety of sizes. These are a few instances: the Internet domain of the image or video, the Internet domain of pages referencing the image or video (parent URL), the keywords, a colour dimension, an edge-orientation dimension, and so forth. The size of the image or video in bytes; the width and height of the frames (or pictures), constituting two dimensions; the date the image or video was created (or last modified); the format type of the image or video; the frame sequence duration in seconds. For many numerical dimensions, concept hierarchies may be automatically formed. Predefined hierarchies may be utilised for various dimensions, such Internet domains or color.The creation of a multimedia data cube will make it easier to analyses multimedia material in a multidimensional way and mine several types of knowledge, such as summarization, comparison, categorization, association, and grouping. Figure 10.5 displays the Multimedia Miner Classifier module and its output.

An intriguing paradigm for the multidimensional analysis of multimedia data seems to be the multimedia data cube. With so many dimensions, it is important to keep in mind that implementing a data cube effectively might be challenging. In the case of multimedia data cubes, this dimensionality plague is very severe. Color, orientation, texture, keywords, and other characteristics may be represented as several dimensions in a multimedia data cube. Many of these traits, however, are set-oriented rather than single-valued.One picture, for instance, may be associated with a list of keywords. It could include a collection of things, each connected to a spectrum of colours. In the design of the data cube, there will be a

tremendous number of dimensions if we utilise each term as a dimension or each finely detailed colour as a dimension. However, failing to do so might result in the modelling of a picture at a relatively crude, constrained, and inaccurate scale. More study is required to determine how to create a multimedia data cube that might strike a balance between effectiveness and representational power. Multimedia data classification and prediction analysis.

For mining multimedia data, classification and predictive modelling have been utilised, especially in scientific fields like astronomy, seismology, and geoscience. In general, image analysis and pattern recognition may make use of every classification technique. In-depth statistical pattern analysis techniques are also well-liked for identifying subtle traits and creating superior models.

Example 10.8: Analysis of astronomical data using classification and prediction. We may build models for the identification of galaxies, stars, and other stellar objects based on characteristics like magnitudes, areas, intensities, picture moments, and orientation using sky photos that have been meticulously categorised by astronomers as the training set. The built models may then be compared to a huge number of sky photos captured by telescopes or spacecraft in order to discover new celestial bodies. Similar investigations have been effective in locating volcanoes on Venus.When mining picture data, data preparation is crucial and might involve feature extraction, data cleansing, and data transformation. The decomposition of pictures into eigenvectors or the use of probabilistic models to handle uncertainty are two strategies that may be examined in addition to basic pattern recognition techniques like edge detection and Hough transformations. Parallel and distributed processing are advantageous since picture data are sometimes in enormous amounts and may need significant processing resources. Image analysis and scientific data mining are closely related, and as a result, numerous image analysis and scientific data analysis methods may be used for image data mining classification and clustering.

The widespread usage of the World Wide Web has transformed it into a vast and rich collection of multimedia data. The Web has various photographs, pictures, albums, and video images in the form of online multimedia libraries, in addition to having these and other types of multimedia on almost every Web page. Such images and photographs, which are accompanied by text descriptions and placed at various Web page blocks or inside news or text pieces, may have a variety of functions, including acting as an integral part of the content, an advertising, or a suggestion of a different subject. Additionally, many Web sites are intricately connected to other Web pages. If utilised appropriately, such text, picture location, and Web linkage information may aid in the categorization and clustering of images on the Web or assist in understanding the contents of the text. Web data analysis now takes a significant turn toward data mining by effectively using the relative placements and connections among Web-based pictures, text, blocks of content, and page links.

**Association Mining in Multimedia Data**

Databases of images and videos may be mined for association rules involving multimedia elements. It is possible to see at least three categories:

associations between non-image content characteristics and image content This category includes rules that relate the image content to the term sky, such as "If at least 50% of the top section of the picture is blue, then it is likely to depict sky. Associations between picture components unconnected to spatial connections This category includes rules that relate to the contents of images, such as "If a picture has two blue squares, it probably also contains one red circle"associations between spatial relationship-related picture contents: This category

includes rules that relate spatial connections to the items in the picture, such as "If a red triangle is between two yellow squares, then it is probably a large oval-shaped object lies beneath."Wemay regard each picture as a transaction and look for recurring patterns among several photos to mine relationships among multimedia items.

There are several minute variations. First, there might be a lot of connections since a picture could include many objects, each with a variety of properties including colour, form, texture, keyword, and spatial placement. A feature may often be seen as being the same in two photographs at a certain resolution level but different at a finer resolution level. Promoting a progressive resolution refinement method is so crucial. In other words, while mining at a finer resolution level, we can only pay attention to patterns that have met the minimal support criterion after first mining frequently recurring patterns at a somewhat rough resolution level. This is due to the fact that patterns that are rare at a coarse resolution level cannot be rare at finer resolution levels. Without sacrificing the accuracy and comprehensiveness of the data mining findings, this multi-resolution mining technique significantly lowers the overall data mining cost.

In big multimedia datasets, this results in an effective way for mining frequent itemsets and relationships.Second, recurrence of the same items should not be disregarded in association analysis since a picture with several recurring objects is a crucial element in image analysis. For instance, an image with two golden circles is handled quite differently from one with only one. Contrast this with a transaction database, where it is often assumed that a person purchases milk whether they purchase one gallon or two. As a result, the concept of multimedia association as well as its metrics, including support and confidence, should be modified.Third, there are often significant spatial interactions between multimedia items, including above, below, between, close by, to the left of, and so on. In order to investigate object relationships and correlations, these qualities are particularly helpful. Interesting connections may be created by combining spatial linkages with other multimedia aspects that are content-based, such as colour, shape, texture, and keywords. As a result, topological spatial linkages and their characteristics as well as spatial data mining techniques become crucial for multimedia mining.

## Data mining for audio and video

In addition to still photos, a vast quantity of audiovisual data is becoming accessible in digital form, in digital archives, on the Web, in broadcast data streams, and in personal and professional databases. This sum is increasing quite quickly.Effective content-based retrieval and data mining techniques for audio and video data are in high demand. Examples of typical applications include finding and editing specific video clips in a TV studio, spotting suspicious people or scenes in surveillance footage, looking for specific events in a personal multimedia archive like MyLifeBits, identifying trends and outliers in weather radar recordings, and finding a specific melody or tune in your MP3 audio collection.

Industry and standardisation groups have made considerable advances in creating a set of standards for multimedia information description and compression, which will make it easier to record, search for, and analyse audio and video data. Examples of standard video compression techniques are MPEG-k (created by MPEG: Moving Picture Experts Group) and JPEG. A standard for describing the multimedia content data is the most current MPEG-7, officially known as "Multimedia Content Description Interface." It enables some degree of information meaning interpretation that may be sent to or accessed by a computer or device.The components that MPEG-7 standardizes support a wide variety of applications rather than being specifically designed for any one use. Three-dimensional models, still

images, audio, voice, video, graphics, and information on how these data elements are combined in the multimedia presentation are all included in the MPEG-7 audiovisual data description. The following components are standardized in MPEG-7 by the MPEG committee:

(1) A set of descriptors that each define the syntax and semantics of a feature, like colour, shape, texture, image topology, motion, or title;

(2) A set of descriptor schemes that each define the structure and semantics of the relationships between their components (descriptors or description schemes);

(3) A set of coding schemes for the descriptors;

(4) A description definition language (ddl) to specify schemes and descriptors. The retrieval of videos based on their content and video data mining are made much easier by such standardization.

Since there are too many photographs in a video clip and many of the adjacent ones could be very identical, it is impractical to approach the video clip as a series of individual still images and examine each image separately. It is preferable to approach each video clip as a collection of actions and events in time and to first temporarily segment them into video shots in order to capture the narrative or event structure of a film. A shot is a collection of frames or images where the video material does not suddenly shift from one frame to the next. In addition, the key frame in a video clip is the one that represents the shot best. The picture feature extraction and analysis techniques covered previously in the content-based image retrieval may be used to investigate each key frame. The order of the key frames will then be utilised to specify how the events in the video clip should play out. As a result, the fundamental tasks in video processing and mining are shot recognition and key frame extraction from video clips. The field of video data mining is continually developing. Before it is adopted as standard practise, there are still several research problems that need to be resolved. Important data mining tasks in this area include similarity-based preprocessing, compression, indexing and retrieval, information extraction, redundancy reduction, frequent pattern finding, classification, clustering, and trend and outlier detection.

---------------------------

# CHAPTER 10

# TEXT, AND WEB DATA

Dr. Nidhi Mathur

Associate Professor, School of Life & Basic Sciences, Jaipur National University, Jaipur, India,
Email Id- nidhi.mathur@jnujaipur.ac.in

Prior data mining research has mostly concentrated on structured data types including relational, transactional, and data warehouse data. However, in practice, a significant percentage of the information is saved in text databases (or document databases), which are enormous collections of documents from diverse sources, including books, digital libraries, e-mail messages, news stories, research papers, and Web sites. Due to the rising quantity of information that is accessible in electronic form, including electronic periodicals, different types of electronic documents, e-mail, and the World Wide Web (which may also be considered as a vast, interconnected, dynamic text database), text databases are expanding quickly. The majority of information is now electronically maintained in text databases in organisations like the government, business, and others. Since they are neither entirely unstructured nor entirely organised, the data kept in the majority of text databases are classified as semi-structured data. A document could, for instance, include a few structured elements like title, authors, publication date, category, and so on, but it might also have some totally unstructured text components like contents and abstract. In contemporary database research, there have been several studies on the modelling and use of semi-structured data. Additionally, methods for text indexing have been developed to handle unstructured materials as part of information retrieval approaches.The escalating volume of text data renders conventional information retrieval methods ineffective. Usually, only a tiny portion of the many documents that are accessible will be relevant to a certain user. It is challenging to create efficient queries for data analysis and information extraction without knowing what could be in the documents. Users need tools to compare various papers, rate the significance and applicability of the documents, or identify patterns and trends in a large number of documents. As a result, text mining is now a crucial and widely used subject in data mining.

## Information Retrieval and Analysis of Text Data

For many years, the discipline of information retrieval (IR) has advanced concurrently with database systems. Information retrieval is concerned with organising and obtaining information from a vast number of text-based documents, in contrast to database systems, which has concentrated on query and transaction processing of structured data. Since information retrieval and database systems deal with distinct types of data, several database system issues, such concurrency control, recovery, transaction management, and update, are often absent from information retrieval systems. Additionally, several frequent information retrieval issues, such as unstructured documents, rough keyword searches, and the idea of relevance, are often not encountered in conventional database systems.Information retrieval has a wide range of uses since text data is so prevalent. There are several information retrieval systems, including online document management systems, online library catalogue systems, and more recently created Web search engines.Finding relevant documents in a document collection based on a user's query—often a few keywords expressing an information requirement, but it might alternatively be an example of a relevant document—is a common information retrieval challenge. When a user has an immediate (i.e., short-term)

information requirement, like researching used cars to purchase, they take the initiative to "extract" the pertinent material out of the collection in a search issue like this. A retrieval system may also decide to "push" any recently added information item to a user if it is determined to be pertinent to the user's information need when a user has a long-term information need (for example, a researcher's interests). Information filtering is the term for this kind of information access procedure, and the accompanying technologies are often referred to as filtering systems or recommender systems.

**Methods for Retrieving Text**

Retrieval techniques may be broadly divided into two groups: They often either see the retrieval issue as a problem with document ranking or selection.The query is seen as establishing limitations for choosing relevant documents in document selection techniques. A common technique in this area is the Boolean retrieval model, in which a user enters a Boolean expression of keywords, such as "vehicle and repair shops," "tea or coffee," or "database systems but not Oracle," to represent a document as a collection of keywords. Such a Boolean query would be sent to the retrieval system, which would then return any documents that match the Boolean phrase. The Boolean retrieval approach often only works successfully when the user has extensive knowledge of the document collection and is able to build an effective query in this manner. This is because it is difficult to precisely specify a user's information demand using a Boolean query.The query is used by document ranking techniques to rank all documents according to relevancy. These techniques are more suited for exploratory queries and common users than document selection techniques. When a user submits a keyword query, the majority of current information retrieval systems return a ranked list of documents. There are a variety of ranking techniques that are based on a wide range of mathematical concepts, such as algebra, logic, probability, and statistics. All of these approaches share the common assumption that we can compare the keywords in a query with those in the documents and rate each document according to how well it fits the query. With a score calculated based on data like the frequency of terms in the document and the whole collection, the objective is to roughly estimate the level of relevance of a text.You should be aware that it is intrinsically difficult to determine with absolute certainty how relevant a group of keywords are to one another. For instance, it is difficult to put a number on how far apart data mining and data analysis are. So, verifying any retrieval strategy requires a thorough empirical study.The scope of this book certainly does not allow for a thorough explanation of all of these retrieval techniques. Next, we quickly go over the vector space model, which is the most often used strategy. Readers may consult information retrieval textbooks, which are included in the bibliographic notes, for information on other models. Although the vector space model is our primary emphasis, some of the procedures we explain are general to all approaches.

The vector space model's fundamental notion is as follows: In order to represent a document and a query as vectors in a high-dimensional space that correspond to all the keywords, we first calculate the similarity between the query vector and the document vector using the suitable similarity measure. The similarity values may then be used to document ranking. Most retrieval systems start with a preprocessing phase known as tokenization, which involves choosing keywords to represent documents. A text retrieval system often connects a stop list with a collection of documents in order to prevent indexing pointless terms. A group of words that are regarded "irrelevant" is known as a stop list. A, the, of, for, with, and other words that occur regularly are examples of stop words. Each document set's stop lists may be different. Database systems, for instance, may be a crucial term in a newspaper. However, it may be seen as a stop word in a collection of research papers given at a conference on

database systems.It's possible for many distinct words to have the same word stem. In order to retrieve just the common word stem for each set of words, a text retrieval system must recognize word clusters where the words are minor syntactic variations of one another. For instance, the terms drug, drugged, and drugs may all be considered as variations on the same word since they all have the same word stem, drug.

## Techniques for Text Indexing

There are a number of well-liked indexing methods for text retrieval, such as inverted indices and signature files.An inverted index is an index structure that keeps track of two hash-indexed or B+-tree-indexed tables: the document table and the term table. The document table is made up of a set of records for various documents, each of which has two fields: the document's doc id and the posting list, which is a list of terms that appear in the document and is sorted by some relevance metric.Each term entry in the term database has two fields: a term id and a posting list, where the posting list gives a list of document IDs where the term occurs.Finding all of the documents related with a given set of words or all of the terms associated with a given set of papers is simple with such structure. For instance, we may first locate a list of document IDs in the term table for each term, then intersect them to get the set of relevant documents, in order to identify all the documents related to a set of words. In business, inverted indices are often utilised. They are simple to use. The storage need might be extremely high due to the potential length of the posting lists. Although simple to use, they fall short when it comes to managing polysemy and synonymy, when two quite distinct words might have the same meaning (where an individual word may have many meanings).A signature record for each document in the database is kept in a file called a "signature file."

The size of each signature is set at b bits, which denote phrases. Following is a basic encoding technique. The initial value of each bit in a document signature is 0. If the phrase a bit represents occurs in the document, it is set to 1. If every bit set in signature S2 is likewise set in signature S1, then the signatures S1 and S2 match. Multiple terms may be mapped onto a single bit since there are often more terms than there are available bits. Because a document that fits the signature of a query does not always include the set of keywords in the query, such multiple-toone mappings increase the cost of the search. It is necessary to obtain, parse, stem, and verify the document. The list of terms may be improved by first doing frequency analysis, stemming, and filtering stop words, and then by encoding the list into bit representation using a hashing approach and overlaid coding technique. However, the main drawback of this strategy is that multiple-to-one mappings remain a difficulty.

## Techniques for Processing Queries

When a document collection's inverted index is built, a retrieval system may swiftly respond to a keyword query by determining which documents contain the query terms.To be more precise, we will keep score accumulators for each document and update them as we process each search word. We will get all of the papers that match each search keyword and raise their scores. In [WMB99], more advanced query processing methods are described.The system may learn from accessible samples of pertinent documents to enhance retrieval performance. Retrieval performance has been shown to be enhanced using a technique known as relevance feedback. When we lack instances that are as relevant, a machine might make the assumption that the top few papers in some early retrieval results are pertinent and can extract more related keywords to broaden a query. Pseudo-feedback, also known as blind feedback, is the practise of extracting helpful keywords from the most popular papers returned. Additionally, pseudo-feedback often results in enhanced retrieval performance.The fact that many current retrieval techniques rely on precise keyword matching is a significant

drawback. The intricacy of real languages, however, presents two significant challenges for keyword-based retrieval. First, there is the issue of synonymy, which arises when two words with the same or comparable meanings may have radically different surface forms. For instance, a user's query may include the phrase "automobile," yet a pertinent page might use the word "vehicle" instead. The second issue is polysemy, which occurs when the same term, such as mining or Java, may have many meanings.

## Text Dimensionality Reduction

We may create similarity-based indexes on text documents using the similarity metrics. Then, text-based queries may be expressed as vectors and searched for in a document collection by their closest neighbours. The number of words T and the number of documents D are often fairly big for any nontrivial document database, however. Due to the resultant frequency tables size T D and the issue of inefficient computing brought on by such large dimensionality. High dimensionality also produces relatively sparse vectors, which makes it more challenging to identify and take advantage of term correlations (e.g., synonymy).Dimensionality reduction methods like latent semantic indexing, probabilistic latent semantic analysis, and locality-preserving indexing may be utilised to solve these issues.

## Text Mining Methods

Based on the inputs used in the text mining system and the data mining tasks to be carried out, there are several techniques to text mining that may be categorised from various angles. According to the types of data they accept as input, the three main approaches are generally: (1) the keyword-based approach, which accepts a set of keywords or terms from the documents as input; (2) the tagging approach, which accepts a set of tags; and (3) the information-extraction approach, which accepts semantic information, such as events, facts, or entities uncovered by information extraction, as input. A straightforward keyword-based technique may only find links at a surface level, such as the rediscovery of compound nouns (such as "database" and "systems") or less significant co-occurring patterns (such as "terrorist" and "explosion"). It may not significantly deepen our grasp of the text. The tagging strategy may depend on tags retrieved either manually (which is expensive and impractical for huge collections of documents) or automatically (which may process a relatively limited number of tags and necessitates setting the categories beforehand). The information-extraction strategy is more sophisticated and may uncover some deep knowledge, but it requires semantic text analysis using machine learning and natural language processing techniques. This knowledge finding job is difficult.On the keywords, tags, or semantic data that have been collected, several text mining operations may be carried out. This list includes information extraction, association analysis, trend analysis, and document clustering. In the discussion that follows, we look at a number of these tasks.

## Analyses of Associations Based on Keywords

The term "keyword-based association analysis" is used here. This kind of study gathers groups of keywords or concepts that are commonly used together and then determines the associations or correlations between them.Like the majority of studies performed on text databases, association analysis begins with parsing, stemming, deleting stop words, and other text preprocessing techniques before using association mining algorithms. Each document in a document database may be thought of as a transaction, and a group of keywords in the document can be thought of as a group of items in the transaction. In other words, the database has the following structure: document id, collection of keywords.The issue of item association mining in transaction databases, where numerous intriguing approaches have

been devised, is therefore translated to the problem of keyword association mining in document databases. Be aware that a term or phrase might be formed from a group of often occurring, consecutive, or nearby keywords. The practise of association mining may be used to find noncom pound connections like as well as compound associations, such as domain-dependent words or phrases. Term-level association mining, as opposed to mining on individual words, refers to mining based on these relationships. In text analysis, word recognition and term-level association mining have two benefits:

(1) Terms and phrases are automatically tagged, eliminating the need for manual document tagging.

(2) The quantity of nonsensical results is significantly decreased, as is the mining algorithms' execution time.

Term-level mining may be used to discover associations between a group of identified terms and keywords with the use of this term and phrase recognition. Others may like to identify the largest collection of terms that appear together, while other users may prefer to find associations between pairs of keywords or terms from a given set of keywords or phrases. Therefore, ordinary association mining or max-pattern mining algorithms may be invoked depending on user mining needs.

## Analysis of Document Classification

Because there are so many documents available online, it is difficult yet necessary to be able to classify them automatically in order to make document retrieval and analysis easier. Automated document categorization is a crucial text mining problem. Document classification has been used for automated topic tagging (i.e., labelling documents), topic directory construction, Document writing style identification (which may help identify the authors of anonymous documents), and document hyperlink classification (i.e., categorising the goals of hyperlinks associated with a set of documents).As a general rule, do the following: First, a training set of pre-classified documents is selected. A classification system is subsequently developed by analysing the training set. Such a categorization technique often requires testing to be made more accurate. Other online publications may be classified using the resulting categorization method.This procedure seems to be comparable to relational data categorization. There is, however, a key distinction. Since each tuple is specified by a collection of attribute-value pairs, relational data are well-structured. For instance, the value "sunny" corresponds to the attribute weather outlook in the tuple "sunny, warm, dry, not windy, play tennis," the value "warm" relates to the attribute temperature, and so on. The attribute-value pairings with the highest discriminating power in predicting whether a person would play tennis are determined by the classification analysis. Document databases, on the other hand, are not organised using attribute-value pairs. In other words, a collection of keywords connected to a set of documents are not arranged according to a predetermined set of characteristics or dimensions. There might be thousands of dimensions in a collection of papers if we treat each unique keyword, phrase, or characteristic in the document as a separate dimension. As a result, popular relational data-oriented classification techniques like decision tree analysis may not be useful for categorising document databases.based on a few common categorization techniques that have been used to classify texts with effectiveness. These include association-based classification, support vector machines, Bayesian classification, feature selection techniques, and nearest-neighbor classification.In the vector-space concept, two documents are said to be comparable if their document vectors are similar. Based on the assumption that related documents should be given the same class label, this model drives the development of the k-nearest-neighbor classifier. All of the training

materials may be easily indexed and each one given a class designation. In order to extract the k documents from the training set that are most similar to a test document when it is submitted, we may use the test document as a query to the IR system. Based on the class label distribution of the test document's k closest neighbours, the class label of the test document may be ascertained. Such a distribution of class labels may also be improved, for example by using weighted counts rather than raw counts or by putting aside certain labelled documents for validation. This kind of classifier may attain accuracy similar to the best classifier by adjusting k and implementing the proposed modifications. However, the approach has a higher space and time cost compared to other types of classifiers since it requires nontrivial space to store (potentially redundant) training data and extra time for inverted index search.

The vector-space approach may give uncommon things a significant amount of weight notwithstanding the class distribution features. Such uncommon objects might make categorization inefficient. Let's look at a TF-IDF measure calculation example. Let's say that there are two terms, t1 and t2, in two courses, C1 and C2, each of which has 100 training materials. However, term t2 only appears in 20 documents in classC1 (i.e., 10% of the whole corpus), whereas term t1 appears in five documents in each class (i.e., 5% of the total corpus). Due to its rarity, term t1 will have a larger TF-IDF value, yet it is clear that t2 has more discriminative strength in this situation.The phrases in the training materials that are statistically uncorrelated with the class labels may be eliminated using a feature selection2 procedure. This will decrease the number of words needed for categorization, increasing accuracy and efficiency.The "cleaned" training papers that are produced after feature selection, which eliminates nonfeature words, may be utilised for efficient categorization. One of several well-liked methods that may be utilised for efficient document categorization is Bayesian classification. A Bayesian classifier first trains the model by calculating a generative document distribution $P(d|c)$ to each class c of document d, and then tests which class is most likely to generate the test document because document classification can be viewed as the calculation of the statistical distribution of documents in particular classes. Both techniques can effectively classify documents since they can handle high-dimensional data sets. The categorization of documents has also used other classification techniques. Support vector machines may be used to do efficient classification, for instance, if we represent classes as integers and create a direct mapping function from term space to the class variable.A approach for discriminative classification is also based on least-square linear regression.The last classification method we present is association-based classification, which categorises texts based on a collection of linked, common text patterns. Take note that frequent words are probably not good discriminators. As a result, only words with strong discriminative capabilities and low frequency will be chosen for document categorization. This is how such an association-based categorization technique works: First, using information retrieval and straightforward association analysis approaches, keywords and concepts may be recovered.

Second, keyword and term idea hierarchies may be found by leveraging existing term classes, like WordNet, or by depending on expert knowledge or specific keyword categorization methods. Class hierarchies may also be applied to documents in the training set. The next step is to utilise a word association mining technique to identify groups of related phrases that may be used to most effectively separate one class of texts from another. From this, a set of association rules for each document class is derived. These classification rules may be used to categorise fresh documents and can be arranged according to their discriminative strength and frequency of occurrence. An efficient association-based document classifier of this sort has been shown.The Web page linkage information may be utilised to further aid in the identification of document classes for categorization of Web documents. Analysis of

Document ClustersOne of the most important methods for unsupervised document organisation is document clustering. The clustering techniques may be used when texts are represented as term vectors. The document space, on the other hand, always has a very high dimensionality, perhaps from a few hundred to thousands. The documents should first be projected into a lowerdimensional subspace where the semantic organisation of the document space becomes obvious due to the curse of dimensionality. The conventional clustering techniques may then be used in the low-dimensional semantic space. The most well-known methods for doing this are spectral clustering, mixture model clustering, clustering using latent semantic indexing, and clustering using locality-preserving indexing. Here, we go through each of these approaches.

The standard clustering technique, such as k-means, is then used on the condensed document space in the spectral clustering approach after performing spectral embedding (dimensionality reduction) on the original data. The capacity of spectral clustering to handle extremely nonlinear data has just come to light (the data space has considerable curvature in every local location). It is possible to identify the document space's manifold structure thanks to its close ties to differential geometry. These spectral clustering techniques may employ the nonlinear embedding (dimensionality reduction), which is only determined on "training" data, which is a significant disadvantage. To learn the embedding, they must utilise all of the available data points. It is computationally difficult to learn such an embedding when the data collection is quite big. Due to this, spectral clustering cannot be used on huge data sets.The mixture model clustering approach incorporates multinomial component models often in order to model the text data using a mixture model. In the two processes of clustering, the model parameters are first estimated using the text data and any extra prior information, and then the clusters are inferred using the obtained model parameters. These techniques have the ability to cluster both words and documents at the same time, depending on how the mixture model is designed. Examples of such methods are Probabilistic Latent Semantic Analysis (PLSA) and Latent Dirichlet Allocation (LDA). Such clustering techniques may have the benefit of creating groupings that make it easier to compare texts.The linear dimensionality reduction techniques known as Latent Semantic Indexing (LSI) and Locality Preserving Indexing (LPI). In LSI and LPI, we may get the transformation vectors (embedding function). Such embedding functions are specified everywhere, allowing us to train the embedding function using a portion of the data before embedding the whole set of data into a low-dimensional space. Using utilising this technique, massive document data corpora may be handled by LSI and LPI clustering.In order to minimise the overall reconstruction error, LSI seeks to identify the best subspace approximation to the original document space, as was covered in the preceding section. In other words, LSI aims to find the most representative features for document representation rather than the most discriminative features. The ultimate purpose of clustering is to discriminate across papers with different meanings, hence LSI may not be the best method for doing so. LPI may have more discriminating power and seeks to understand the local geometrical structure. According to experiments, LPI is a better dimensionality reduction approach for clustering than LSI. When compared to LSI and LPI, the PLSI approach more clearly shows the latent semantic dimensions and is readily expanded to take into account any previous information or grouping preferences.

**The World Wide Web for Mining**

For news, ads, consumer information, financial management, education, government, e-commerce, and many other information services, the World Wide Web acts as a massive, widely dispersed, worldwide information service centre. Additionally, the Web offers abundant data mining sources due to its dynamic and comprehensive collection of hyperlink

information and Web page access and use statistics. But as the following facts show, the Web also presents significant obstacles to efficient resource and information discovery.The Web seems to be too big for data mining and data warehousing to be efficient. The Web is now hundreds of terabytes in size and continues to expand quickly. The majority of the information that many groups and organisations make available to the public is posted online. Setting up a data warehouse to copy, store, or integrate all of the data on the Web is quite difficult.

Web pages are far more complicated than any collection of conventional text documents. Web pages lack an overarching framework. They have a far wider range of writing styles and material than any collection of books or other conventional text-based documents. The Web is regarded as a large digital library, however the enormous number of papers in this library are not organised in any certain order of sorting.There isn't a category index, nor is there one for titles, authors, covers, tables of contents, or other information. Finding the information you want in such a library might be quite difficult!The Internet is a very dynamic source of information. The Web not only expands quickly, but its content is also updated often. The web regularly updates sites for news, financial markets, weather, sports, shopping, corporate ads, and many other things. Access records and linkage data are often updated as well.

The user populations served by the Web are quite diverse. More than 100 million workstations are presently connected to the Internet, and the number of users is continually growing quickly. The backgrounds, interests, and use goals of users might vary greatly. The structure of the information network may not be well understood by most users, and they may also be unaware of the high cost associated with a given search. By stumbling about in the "darkness" of the network, they may easily become lost or grow bored while making several access "hops" and anxiously waiting for some information. Only a tiny amount of the data on the Internet is either relevant or helpful. It is believed that 99% of Web users find 99% of the material on the Web to be worthless. Although it may not seem evident, it is true that a specific individual is often only interested in a very small fraction of the Web, whereas the rest of the Web includes material that the user is uninterested in and might saturate desired search results. These difficulties have sparked study on effective and efficient resource finding and use on the Internet.

There are various Web search engines with indexes. These do Web searches, index Web sites, and create and maintain sizable keyword-based indexes that assist in locating collections of Web pages that include certain terms. Using such search engines, a knowledgeable user could be able to retrieve documents rapidly by entering a list of closely related keywords and phrases. A straightforward keyword-based search engine, however, has significant drawbacks. First off, there are probably hundreds of thousands of papers on any given subject. This may result in a search engine returning a massive number of document entries, many of which are only tangentially related to the subject or may include low-quality content. Second, many highly relevant publications to a subject could not have keywords characterising them. This is known as the polysemy issue, which was covered in the text mining section above. For instance, the term "Java" might be used to describe the Java computer language, an Indonesian island, or freshly brewed coffee. Another example is that even the most well-known Web search engines like Google, Yahoo!, AltaVista, or America Online could not appear in a search for the term "search engine" if those services do not make that claim on their websites. This shows that finding Web resources online requires more than just a basic keyword-based Web search engine." Compared to Web searches using keywords, the more difficult work of "web mining" involves looking for web structures, ranking the significance of web contents, identifying the regularity and dynamism of web contents, and

mining web access patterns. Web mining, however, has the potential to significantly increase the strength of a web search engine since it can be used to identify authoritative web sites, categorise web documents, and clear up several ambiguities and nuances that come up in keyword-based web searches. Web content mining, web structure mining, and web use mining are the three broad categories into which jobs related to web mining may be divided. Web structures may instead be seen as a component of Web contents, allowing Web mining to be categorized into Web use mining and Web content mining.

## Web Mining of Multimedia Data

The Web offers a vast quantity of multimedia material in a variety of formats. These consist of audio, visual, graphic, and image content. Effective techniques for storing and retrieving such multimedia data are in increasing demand. The multimedia data on the Web have numerous distinct qualities as compared to general-purpose multimedia data mining. Web-based multimedia files are integrated into Web pages together with text and link data. The texts and links mentioned above may also be seen as aspects of the multimedia data. A Web page may be divided up into a number of semantic chunks using certain Web page layout mining algorithms (like VIPS).As a result, the multimedia data block may be seen as a whole. Finding and arranging multimedia blocks may be used to describe searching and organising multimedia material on the Web.

## Web document classification automatically

When Web pages are automatically categorized, a set of preset subject categories are used to give a class label to each document based on samples of previously classified papers. To create a categorization system for Web documents, for instance, the taxonomy from Yahoo! and the documents it is connected with may be utilised as training and test sets. Using categories from the same taxonomy, this system may then be used to categorise new Web content.

## Web Usage Analysis

Web use mining, which mines Weblog information to uncover user access patterns of Web sites, is another significant duty for Web mining in addition to mining Web contents and Web connection structures. Finding new clients for electronic commerce, enhancing the quality and delivery of Internet information services to end users, and enhancing Web server system efficiency may all be accomplished by analysing and investigating regularities in Weblog entries?Every time a Web page is accessed, a Web server typically records a (Web) log entry, also known as a Weblog post. It provides the date, the IP address from whence the request came, and the URL that was requested. A sizable quantity of Web access log data are being gathered for servers that support Web-based e-commerce. Popular websites could register Weblog records each day that are in the range of hundreds of gigabytes. Weblog databases provide extensive information on the dynamics of the Web. Therefore, it's crucial to create advanced weblog mining methods.We may take into account the following as we create strategies for Web use mining. First off, even if it is inspiring and thrilling to consider the many possible uses of Weblog file analysis, it is crucial to understand that the success of such applications relies on what and how much true and trustworthy information can be gleaned from the vast raw log data. It is often necessary to clean, compress, and modify raw Weblog data in order to extract and evaluate important and relevant information.In order to find the top N users, top N most frequently accessed Web pages, most frequently accessed time periods, and other information, a multidimensional view of the Weblog database can be built using the available URL, time, IP address, and Web page content information. This will help identify potential users, customers, markets, and other people.

Third, Weblog records may be used for data mining to uncover association patterns, sequential patterns, and Web accessing tendencies. It is often required to take extra steps in order to collect more user navigational data for Web access pattern mining in order to support thorough Weblog analysis. User surfing patterns of the Web pages in the Web server buffer are one example of this extra data.

Studies have been carried out using such Weblog files to examine system performance, enhance system architecture via Web caching, Web page prefetching, and Web page swapping, comprehend the nature of Web traffic, and comprehend user response and motivation. For instance, some research has suggested adaptable websites, which develop over time based on user access habits. Weblog analysis may also assist in creating user-specific Web services.Weblog information may be combined with web content and web linkage structure mining to aid in web page ranking, web document categorization, and the creation of a multilayered web knowledge base since weblog data offer information about what kind of people will access what types of web pages. Mining a user's interaction history and search context on the client side to extract helpful data for enhancing the ranking accuracy for the specific user is one especially intriguing use of Web usage mining. For instance, if a user enters the term "Java" into a search engine and then chooses "Java programming language" from the results to see, the system might assume that the user would find the presented snippet for this Web page interesting. Then, it may give "Java programming language" sites a higher ranking and keep "Java Island" pages off the search results page. As a result, since search is contextualised and tailored, its quality has increased.

----------------------

# CHAPTER 11

# MINING DATA STREAMS

Prerna Vyas

Assistant Professor, School of Computer & Systems Sciences, Jaipur National University, Jaipur, India,
Email Id- Prerna.vyas@jnujaipur.ac.in

Modern hardware technology has made it possible to capture data in new ways and at a faster pace than previously. For instance, a lot of daily activities like using a phone or credit card result in automatic data collecting. Similar to how wearable sensors and mobile devices have increased the amount of dynamically accessible data, so have new methods of data collection. These methods of data gathering make the crucial premise that data rapidly and continuously accrue throughout time. Data streams are the name given to these dynamic data collections.

The streaming paradigm makes the crucial premise that it is no longer feasible to retain all of the data due to resource limitations. Although such data may be archived utilising distributed "big data" frameworks, doing so incurs astronomical storage costs and limits the ability to interpret data in real-time. Such frameworks are often impractical due to excessive prices and other analytical factors. A different strategy is offered by the streaming framework, where real-time analysis may often be carried out using properly thought-out algorithms without a substantial investment in specialist equipment. The following are some examples of application domains that relate to streaming data:

**1. Transaction streams:** Purchases made by customers usually result in transaction streams. The information generated by using a credit card, a point-of-sale transaction at a grocery store, or an online purchase of an item are a few examples.

**Web click-streams:** A web click stream is the result of user activity on a well-known website. If the website is widely used, the pace of data creation could be high enough to call for a streaming strategy.

**Social streams:** Because of user engagement, online social networks like Twitter constantly produce enormous text streams. In most social networks, the number of actors grows superlinearly with the pace and volume of the stream.

**Network streams:** There are several traffic streams present in communication networks.

These streams are often searched for incursions, outliers, or other odd behaviour.

Due to the processing limitations brought on by the substantial amounts of continually incoming data, data streams pose a number of particular difficulties. The following restrictions, of which at least some are always present while others are only sometimes present, must be met by data streaming algorithms in particular:

**One-pass constraint**: It is expected that the data can only be processed once due to the continuous and quick generation of large amounts of data. All streaming models include this as a strict restriction. Rarely is it expected that the data would be preserved for further processing. For algorithmic advancement in streaming applications, this has important ramifications. Many data mining methods, in particular, are inherently repetitive and need repeated runs over the data. For use in the context of the streaming paradigm, such algorithms must be suitably adjusted.

**Concept drift:** The data may change over time in the majority of applications. Accordingly, other statistical features, including correlations between attributes, correlations between attributes and class labels, and cluster distributions, may alter with time.

Although it is virtually always true in real-world situations, not all algorithms will definitely make this assumption.

**Resource limitations:** A user may have very little influence over the external mechanism that generally generates the data stream. As a result, the user also has limited control over the stream's arrival pace. It may be challenging to carry out online processing continuously at peak times when arrival rates change over time. Tuples that can't be processed quickly in certain situations may need to be dropped. This practise is known as load shedding. Even though the streaming paradigm nearly always has resource limits, surprisingly few methods do.

**Massive-domain restrictions**: When an attribute's values are discrete, they may have a lot of different values. Consider a situation where pairwise communication analysis in an email network is needed. An email network with 108 users has in the neighbourhood of 1016 unique pairings of email addresses. The possibilities often reach the petabyte scale when described in terms of needed storage. It becomes very difficult to save even basic information like counts or the number of different stream items in such circumstances. As a result, a variety of specialised data structures have been developed for the synopsis generation of massive-domain data streams. Virtually all streaming techniques employ an online summary building approach throughout the mining process because to the enormous number of data streams. The main concept is to provide an online summary that may later be mined. Depending on the application, a synopsis may be created in a variety of ways. The sort of insights that may be gleaned from a summary is greatly influenced by its nature. Random samples, bloom filters, drawings, and different element-counting data structures are a few examples of synopsis structures. Additionally, you may use certain established data mining techniques, such clustering, to extract useful synopses from the data.

### Data Structures for Streams: A Synopsis

For various purposes, a broad range of synopsis data structures have been developed.

Two forms of synopsis data structures exist:

**Generic:** In this instance, the summary may be applied immediately to the majority of applications. A random selection of the data points is the sole such summary, albeit it cannot be utilised for specific tasks like distinct element counting. The technique of retaining a random sample from the data in the context of data streams is also known as reservoir sampling.

**Particular:** In this instance, the synopsis is made for a particular purpose, such frequent element counting or distinct element counting. The Flajolet-Martin data structure for distinct element counting and drawings for frequent element counting or moment computation are two examples of these data structures. Following is a discussion of several synopsis structure kinds.

### Sampling of a reservoir

One of the most adaptable techniques for stream summarization is sampling. Sampling's key benefit over other synoptic data structures is that it can be used to any situation. Almost any offline method may be used to process a sample of points that have been taken from the data.

Although it has limits for a limited number of applications like distinct-element counting, sampling should generally be thought of as the preferred approach in streaming circumstances. Reservoir sampling is the term used to describe the process used to keep a dynamic sample from the data in the context of data streams. A reservoir sample is the name given to the final sample.

For a straightforward issue like sampling, the streaming situation poses some intriguing difficulties. Because the complete stream cannot be stored on disc for sampling, there is a problem. Reservoir sampling aims to keep a dynamically updated sample of k points from a data stream constantly without explicitly saving the stream on disc at any particular moment. As a result, in order to preserve the sample, one must apply a set of effectively implementable actions for each new data point that enters the stream. In the static example, the likelihood that a data point will be included in the sample is given by the ratio k/n, where k is the sample size and n is the total number of data points.

The "data set" is not static in the streaming situation, and the value of n keeps rising over time. In addition, previously obtained data points that are not in the sample have been lost forever. As a result, the sampling technique functions with imperfect knowledge of the stream's historical context at any given time. In other words, two straightforward admission control choices must be made dynamically for each entering data item in the stream: The algorithm for reservoir sampling works as follows. The first k data points in the stream are always included in the reservoir for a reservoir of size k. The subsequent two admission control choices are then used for the nth data point of the entering stream.

- Add the nth incoming stream data item with probability k/n to the reservoir.
- If the newly arriving data point was inserted, then randomly evict one of the previous k data points from the reservoir to create space for the new point.

## Mining Patterns Frequently in Data Streams

In two distinct cases, the issue of frequent pattern mining in data streams is examined. The massive-domain scenario, which has a very high number of potential things, is the first conceivable outcome. Even the challenge of locating common products becomes challenging under such circumstances. Heavy hitters are another term for common things. The second situation is the typical one, when there are many things (but manageable numbers) that can fit in main memory.

Because the frequent counts may be directly preserved in an array in these circumstances, the frequent item issue is no longer as intriguing. In these circumstances, finding common patterns is of more relevance.

Because most common pattern mining methods need numerous iterations across the whole data collection, this is a challenging challenge. This is challenging due to the streaming scenario's one-pass restriction. The following will outline two distinct strategies. In the first of these methods, standard frequent pattern mining algorithms are combined with generic synopsis structures, while in the second, streaming versions of frequent pattern mining algorithms are created.

## Utilizing Synopsis Structures

The majority of streaming data mining issues, including frequent pattern mining, may be successfully solved using synopsis structures. Synopsis structures are especially appealing in the context of frequent pattern mining techniques because they allow for the use of a greater variety of algorithms or the incorporation of temporal decay.

**Sampling of a reservoir**

The most adaptable method for routine pattern mining in data streams is reservoir sampling.It may be used to frequent pattern mining or frequent item mining (in the massive-domain case). Using reservoir sampling is based on a straightforward principle:

- Keep a reservoir sample S taken from the data stream.
- Use a frequent pattern mining technique to uncover patterns in reservoir sample S.

As a function of the sample size S, it is able to obtain qualitative assurances on the frequently occurring patterns mined. The Chernoff bound may be used to estimate a pattern's likelihood of being a false positive. It is also feasible to have a guaranteed decrease in the amount of false negatives by using support standards that are somewhat lower. Pointers to these assurances may be found in the bibliographic notes. Because sampling and mining are neatly separated, reservoir sampling provides a number of flexibility benefits. On the memory-resident reservoir sample, almost any effective frequent pattern mining technique may be used.

Constrained pattern mining and intriguing pattern mining are two further variants of pattern mining algorithms that may be used. Conceptual sway is also rather simple to correct. When standard frequent pattern mining techniques are used with a decay-biased reservoir sample, the definition of the support is decay-weighted.

**Drawings**

Sketches may be used to identify common objects, however they are less useful for identifying frequent item groupings. The fundamental tenet is that drawings are often considerably better at properly approximating the relative numbers of more common things. This is due to the fact that every item's frequency estimate has an absolute constraint, meaning that the mistake relies more on the frequency of the stream as a whole than it does on the frequency of the individual item. As shown by Lemma 12.2.3, this is true. As a consequence, it is usually possible to estimate the frequencies of heavy hitters more precisely on a relative basis. The strong hitters may be identified using both the AMS drawing and the count-min sketch. There are links to several of these algorithms in the bibliographic notes.

**Lossy Counting Algorithm**

The frequent item or frequent item set counting methods may both be utilised with the lossy counting algorithm. According to the method, the stream is divided into segments S1–Si, where Si is the smallest segment and S1 is the largest. A user-defined tolerance on the necessary precision serves as the parameter.

The simpler issue of frequent item mining will be discussed first. The method keeps track of all the items' frequencies in an array and increases them when new things are added.

If there aren't many unique things, one may keep track of all the counts and just publish the ones that happen often. When there isn't enough overall space to keep track of all the various things' counts, a problem occurs. In these circumstances, infrequent items are discarded once the border of a segment Si is reached. Due to the fact that the bulk of the things in the stream are seldom used in reality, many items are removed as a consequence

**The grouping of data streams**

Because it may provide a concise summary of the data stream, the clustering issue is very important in the data stream situation. If a fine-grained clustering is applied, it is often

possible to replace reservoir sampling with a clustering of the data stream. Because of these factors, stream clustering often serves as a foundation for later applications, such streaming classification. Several typical stream clustering techniques will be covered in the sections that follow.

### STREAM Algorithm

The k-medians clustering approach serves as the foundation for the STREAM algorithm. The basic concept is to divide the stream into more manageable, memory-resident chunks. Thus, segments S1–Sr are created from the original data stream S. Maximum m data points arepresent in each section. A predetermined memory budget serves as the foundation for fixing the value of m.

Each segment Si can fit in main memory, allowing for the use of a more complicated clustering algorithm without having to worry about the one-pass restriction. For this aim, a range of other k-medians5 style algorithms may be used. In k-medians methods, each point in Si is given to the nearest representative from a set Y of k representatives from each chunk Si. The objective is to choose the representatives such that the sum of squared distances (SSQ) of the data points allocated to these representations is as little as possible.

### Algorithm CluStream

The clusters are dramatically altered over time by the idea drift in a changing data stream. The clusters from the previous day are considerably different from those from the previous month. The flexibility to choose the clusters based on one or more temporal horizons that are unknown at the start of the stream clustering process is desirable in many data mining applications. It is challenging to calculate clusters across various time horizons using traditional algorithms because stream data naturally imposes a one-pass limitation on the design of the methods.

The simultaneous preservation of the intermediate outcomes of clustering algorithms across all conceivable time horizons would be necessary for a direct adaptation of the STREAM algorithm to such a situation. With the growth of the data stream, the computing overhead of such a method rises and may quickly become a barrier for live deployment. Applying the clustering process with a two-stage technique, involving an online microclustering stage and an offline macro clustering stage, is a logical way to deal with this problem.

The stream is processed in real-time during the online micro clustering step in order to continually preserve the stream's summary yet comprehensive cluster data. Micro clusters are what they are known as. To provide the user a clearer grasp of the clusters across various time horizons and degrees of temporal granularity, the offline macro clustering step further condenses these precise clusters. In order to re-cluster these precise representations across user-specified time periods, it is necessary to preserve enough specific data in the micro clusters.

### Algorithm for Micro clustering

Any length of history from the present instant that is defined by the user may produce approximate clusters using the data stream clustering technique. This is accomplished by saving snapshots, or the micro clusters at certain points in the stream. The programm constantly keeps the most recent snapshot of the microclusters.Any time range may be utilised to extract microclusters using the additive property. This representation is subjected to the macroclustering phase.

**Period of the Pyramid**

To allow horizon-specific examination of the clusters, the microcluster data are routinely kept. The microclustering phase is when this upkeep is done. According to the recentness of the snapshot, several granularity levels for the microcluster snapshots are kept in this method. The ordering of snapshots may range from 1 to log(T), where T is the amount of time from the start of the stream in clock terms. According to the following guidelines, the order of a snapshot controls the amount of temporal granularity at which it is stored:

At time intervals of I where I is an integer, snapshots of the ith order are recorded. It is specifically saved when the clock value is precisely divisible by I for each snapshot in the $i^{th}$ order.

Only the latest l + 1 photographs of order I are kept in memory at any one moment.

**Stream Clustering in the Large-Domain**

The massive-domain situation is often seen in the stream context, as was previously stated. In many situations, working with a multidimensional data stream that draws each attribute's value from a wide range of potential values may be necessary. Because "concise" summaries of the clusters become substantially more space-intensive in these situations, stream analysis becomes much more challenging. This serves as the inspiration for a number of synopsis structures, including the Flajolet-Martin method, the count-min sketch, the bloom filter, and the AMS sketch.Due to the difficulties in keeping succinct statistics of the clusters, the data clustering issue also becomes more complex in the massive-domain scenario. Clustering massive-domain data streams has been made easier with the help of the current technique CSketch. In order to record the frequencies of attribute-value combinations in each cluster, this approach uses a count-min sketch. As a result, the number of clusters and the number of count-min drawings utilised are identical. The drawing is used as the representation for the (discrete) characteristics in the cluster in an online k-means style clustering. A dot product is calculated with regard to each cluster for each input data point.The calculation is carried out as shown below. The hash function hr() is applied to each attribute-value pair in the dimensions given a certain value of r. It is established what frequency the appropriate drawing cell has. For each of the d distinct dimensions, the combined frequencies of all the relevant drawing cells are calculated. The dot product is estimated in this way.

The minimum value across many hash functions (different values of r) is utilized to get a closer approximation. To eliminate bias against clusters containing plenty of data points, the dot product is divided by the aggregate frequency of the cluster's elements. Because the count-min sketch can compute the dot product precisely in a tiny area, this calculation can be done with accuracy. The cluster with which the data point has the biggest dot product is chosen as its home. The data for that specific cluster in the drawing are then updated. In terms of how data points are sequentially allocated to clusters, this method is similar to microclustering. It does not, however, put the merging and removal stages into practise. Additionally, for the preservation of cluster statistics, the sketch form is used rather than the microcluster representation. With regard to a clustering that has limitless space availability, theoretical guarantees on clustering quality may be shown. There are links to these findings in the bibliographic notes.

**Detection of Streaming Outliers**

Time-series data streams or multidimensional data contexts are often where the issue of streaming outlier identification occurs. In general, multidimensional data stream outlier

identification differs significantly from time series outlier detection. The latter scenario treats each time series as a separate entity, but for multidimensional data, temporal relationships are substantially weaker. The situation of the multidimensional stream is comparable to the analysis of static multidimensional outliers. The main change is that the analysis now includes a temporal component, but it is considerably weaker than it would be for time series data. Efficiency is a key consideration in the context of multidimensional data streams since it allows for fast identification of outliers. In the context of multidimensional data streams, there are two different types of outliers that might appear.

- One is based on the identification of outliers in certain data. A prime example of this kind of outlier is a first-ever news report on a particular subject. A novelty is another word for such an anomaly.
- The second is based on adjustments to the multidimensional data's overall patterns.

For instance, a sudden occurrence like a terrorist attack may cause a spike in news headlines on a certain subject. This is an outlier that has been combined based on a certain time range. Almost typically, the second kind of change point starts with a single first-type outlier. A single outlier of the first kind, however, may not necessarily turn into an aggregate change point. This and the idea of thought drift are closely connected. Although idea drift usually happens gradually, a sudden shift can be considered an outlier moment in time rather than an outlier data point.In this section, both types of outliers (or change points) will be covered.

**Individual Data Points as Outliers**

When the whole history of the data stream is utilized, the difficulty of identifying specific data points as outliers is strongly connected to the problem of unsupervised novelty detection. The issue of first tale detection is one that has received a lot of attention in the text domain. Such surprises often establish trends and may ultimately mix in with the standard data. However, it may not necessarily be a novelty when a single record is labelled an outlier in the context of a window of data points. By applying the related algorithms nearly directly to the window of data points in this context, proximity-based algorithms are especially simple to extend to the incremental situation.The streaming situation can be simply adapted for distance-based algorithms. The following changes are made to the original distance-based definition of outliers:It is quite simple to identify the outliers by calculating the score of each data point in the window when the complete window of data points can be kept in main memory. The addition and removal of data points from the window makes incremental preservation of the scores of data points more difficult. Additionally, certain algorithms, like LOF, demand that data like reachability distances be recalculated. The incremental scenario has been added to the LOF algorithm. The procedure involves the following two steps:

The statistics of the newly added data points, such as their LOF score and reachability distance, are calculated.The densities and reachability distances of the current locations in the window are updated along with their LOF ratings. In other words, many of the current data points' ratings need to be changed since the addition of a new data point has an impact on them. Only the location of the new data point is impacted, therefore not all scores need to be changed. The only LOF scores that are impacted when data points are erased are those in the immediate vicinity of the deleted point.Many of the aforementioned approaches are still rather costly in the context of the data stream since distance-based methods are well-known to be computationally expensive.Therefore, by using an online clustering-based method, the complexity of the outlier identification process may be considerably reduced.While it is typically not a good idea to use clustering-based approaches when there are few data points, streaming analysis makes an exception. In the context of a data stream, there are generally

enough data points available to keep the clusters at an extremely high degree of granularity. Unsupervised surprises are often linked to the development of new clusters in the context of a streaming clustering technique. When an incoming data point does not fall within a certain statistical radius of the data's existing clusters, the CluStream algorithm, for instance, specifically controls the generation of new clusters in the data stream. These data points can be regarded as outliers. As additional data points are added to the cluster at later phases of the algorithm, this is often the start of a new trend. Such data points may sometimes be associated with innovations, while in other instances they may be associated with patterns that were prevalent in the past but are no longer represented in the current clusters. Such data points are intriguing outliers in both scenarios. However, unless one is prepared to let the number of clusters in the stream to gradually expand, it is impossible to differentiate between these several categories of outliers.

## Outliers in Aggregate Change Points

The underlying data's abrupt shifts in aggregate local and worldwide patterns are often a sign of unusual occurrences. Numerous techniques also provide statistical means of calculating the magnitude of changes in the underlying data stream. Using the idea of velocity density is one technique to gauge notion drift. The goal of velocity density estimation is to create a velocity profile based on the data's density. The idea of kernel density estimation in static data sets is comparable to this.

## Classification for streaming

The influence of idea drift makes the task of streaming categorization more difficult. Using a reservoir sample to provide a succinct representation of the training data is one such method. It is possible to develop an offline model using this succinct formulation. Concept drift may be managed using a decay-based reservoir sample, if preferred. A benefit of this strategy is that any standard classification algorithm may be used since the difficulties posed by the streaming paradigm were already resolved at the sampling step. Additionally, some specialised techniques for streaming categorization have been put forward.

## Family VFDT

The concept of Hoeffding trees serves as the foundation for very quick decision trees (VFDT). The underlying notion is that a decision tree may be built using a sample of a very large data set and a properly planned methodology such that the resultant tree is very likely to be identical to what would have been produced using the original data set. The intermediate stages of the strategy are created with the Hoeffding bound in mind because it is utilised to estimate this probability. Because of this, these trees are known as Hoeffding trees.By developing the tree concurrently with stream arrival, the Hoeffding tree may be built gradually. The stream is assumed not to develop, therefore the set of points that have just arrived may be thought of as a sample of the whole stream.When there are enough tuples to calculate the accuracy of the appropriate split criterion, the higher levels of the tree are built sooner in the stream. Because statistics regarding lower level nodes can only be gathered after the higher level nodes have been built, the lower level nodes are erected later. As additional instances are added and the tree continues to expand, new layers of the tree are built. In order to conduct a split in the Hoeffding tree technique that is roughly equivalent to what would have been done with knowledge of the whole stream, it is necessary to quantify the moment at which statistically sufficient tuples have been gathered.As long as the same splits are utilised at each step, the same decision tree will be built on both the whole stream and the current stream sample. The approach's aim is to guarantee that the splits on the sample and the splits on the whole stream are the same. Consider the scenario where each

attribute6 is binary to make the explanation easier. If the same split attribute is chosen at each point, two algorithms will result in the exact same tree in this situation. Using a metric like the Gini index, the split attribute is chosen. Take a look at a specific node in the tree that was created using the original data and the identical node that was created using the sampled data.

## Supervised Micro cluster technique

In essence, the supervised micro cluster is a classification method that uses instances. It is assumed in this model that a training stream and a test stream are received concurrently throughout time. Concept drift makes it crucial to dynamically modify the model over time. The dominating class label among the top k closest neighbors is presented as the relevant result in the nearest-neighbor classification technique. Because of the growing size of the stream in the streaming scenario, it is challenging to effectively calculate the k closest neighbors for a certain test instance. However, a fixed-size summary of the data stream that doesn't become bigger as the stream progresses may be made using fine-grained micro clustering. The use of supervised micro clustering prevents the mixing of data points from various classes inside clusters. With just a few modest adjustments to the CluStream algorithm, such micro clusters are rather simple to maintain. The primary distinction is that throughout the cluster update process, data points are allocated to micro clusters that belong to the same class. As a result, labels relate to micro clusters rather than specific data points. The relevant label is stated to be the dominating label of the top k closest micro clusters. However, this does not take into consideration the modifications that must be made to the algorithm due to idea drift. Concept drift causes the trends in the stream to shift. In order to improve accuracy, it is more relevant to employ micro clusters from certain time ranges. Even while the most recent horizon may often be relevant, this may not always be the case if the stream's patterns abruptly switch back to prior trends. As a result, the validation stream is created from a portion of the training stream. To assess the correctness across various time periods, recent validation stream segments are used as test cases. The ideal horizon is chosen. Over this carefully chosen horizon, test cases are subjected to the k-nearest neighbor method.

## Ensemble Approach

Additionally, a reliable ensemble approach for classifying data streams was suggested. Since the approach may successfully take into account changes in the underlying data, it is also intended to address idea drift. Multiple classifiers are trained on each of the divided up portions of the data stream. The score on each of these chunks is used to get the final categorization score. Particularly, scores are obtained from successive sections of the data stream using ensembles of classification models, such as C4.5, RIPPER, and naïve Bayesian. The ensemble's classifiers are distributed according to their predicted classification accuracy in a time-evolving environment. Because the classifiers are constantly updated to enhance the accuracy for that section of the data stream, this guarantees that the technique may attain a greater level of accuracy. If the weights of all classifiers are allocated based on their predicted classification accuracy, it has been shown that an ensemble classifier delivers a lesser error than a single classifier.

## Classification for Massive-Domain Streaming

The multidimensional discrete properties seen in many streaming applications have a very large cardinality. Because of memory constraints, it becomes challenging to apply standard classifiers in these situations. These difficulties may be solved using the count-min sketch. For tracking frequently occurring r-combinations of items in the training data, where r is constrained above by a tiny integer k, each class is assigned a drawing. All potential r-combinations (for r k) are handled as pseudo-items that are added to the class sketch for each

incoming training data point. Different classes will have various pertinent pseudo-items that will appear at various frequency in the cells that correspond to their respective class drawings. The most discriminative cells in the various drawings may be found using these distinctions. In order to establish implicit rules connecting the pseudo-items to the various classes, numerous discriminative pseudo-items are required. Because they are held implicitly in the drawings rather than being explicitly realized, these rules are implicit. They are only ever obtained when a test instance is classified. Which pseudo-items match the mix of objects within them for a particular test case is determined. The class-specific drawings are used to get the statistics of the discriminative ones among them. The categorization of the test instance is then carried out utilizing these, generally speaking, in a rule-based classifier fashion.

-------------------------

# CHAPTER 12

# SOCIAL NETWORK ANALYSIS

Sachin Jain

Assistant Professor, School of Computer & Systems Sciences, Jaipur National University, Jaipur, India,
Email Id-sachin.jain@jnujaipur.ac.in

Humans have a natural propensity to connect with one another, which is an ingrained social need that existed long before the Web and Internet technology. Social connections in the past were accomplished by direct touch, postal correspondence, and telephone technology. In terms of the history of humanity, the last of these is likewise quite recent. However, the widespread use of the Internet and Web technology has created whole new opportunities for facilitating the seamless interaction of people who are dispersed geographically. The Web's visionary pioneers were aware of its enormous potential even in its early stages. But it took 10 years before the Web's entire social potential could be realised. Web-based social apps are still developing and producing a growing quantity of data nowadays. This data contains a wealth of knowledge on user preferences, connections, and social influences. So it makes sense to use this data to get analytical insights.

Although Facebook, LinkedIn, and other sizable online networks are often used to define social networks, they really make up a very tiny percentage of the interaction mechanisms made possible by the Web. In truth, the traditional study of social network analysis in the discipline of sociology came first, long before technologically enabled processes were widely used. Here are a few instances:

For more than a century, social networks have been thoroughly researched, but not from an online viewpoint. Because there were inadequate technical processes in place, data collecting was rather challenging in these situations. As a result, these research were often carried out using time-consuming and arduous manual data gathering techniques. An example of such an endeavour is Stanley Milgram's well-known six degrees of separation experiment from the 1960s, which tested if two randomly selected people on the globe could be linked by a chain of six connections via postal mail between participants. Such studies were sometimes difficult to carry out in a reliable manner since it was difficult to validate local mail forwarding. The findings have since been proven to be relevant to online social networks, where the ties between people are more readily quantified, despite the apparent shortcomings in the experimental context.

Indirect versions of social networks may be found in a variety of technical enablers, including telephones, email, and electronic chat messengers. These enhancers naturally foster social interaction by facilitating communication between various people.

Because they provide a high degree of user involvement, websites for sharing online media material, like Flickr, YouTube, or Delicious, may also be seen as indirect types of social networks. Additionally, social media platforms provide a variety of original methods for users to communicate with one another. Examples include tagging each other's photographs or creating blogs. Although the engagement in these situations is focused on a particular service, such content sharing, many basic social networking characteristics still hold true. These social networks are very valuable to mining applications. They are very rich in material, whether it be text, photographs, music, or video.

Particular interactions in professional groups may be used to build various social networks. Networks of citations and bibliographies are created within scientific groups. Because these networks are built on publications, they are also content-rich.

It is clear that these various network types highlight various aspects of social network analysis. Many of the basic issues raised in this chapter's discussion are relevant to these many circumstances, although in different contexts. Most common data mining issues, such clustering and classification, may also be used to social network research. Because networks are more complex than other types of data, it is also feasible to define a variety of more complicated problems, such as link prediction and social impact analysis.

## Social Networks: Foundations and Features

The social network is presumably formed as a graph with the formula $G = (N,A)$, where $N$ denotes the set of nodes and $A$ denotes the set of edges. Each member of the social network is symbolised by a node in $N$, also known as an actor. The links between the various actors are represented by the edges. These edges in a social network like Facebook correlate to friendship connections. Most of the time, these linkages are undirected, however certain "follower-based" social networks, like Twitter, may contain directed links. Unless otherwise stated, it will be assumed that the network $G = (N,A)$ is undirected by default.The nodes in $N$ could sometimes have content attached to them. This information might relate to remarks or other documents that members of social networks have uploaded. The social network is supposed to have n nodes and m edges. Some important characteristics of social networks will be described in the paragraphs that follow.

## Homophily

Homophily is a key characteristic of social networks that is used in a variety of contexts, including node categorization. Homophily's fundamental tenet is that nodes connecting to one another are more likely to share attributes. For instance, a person's Facebook friendship connections can be made up of former coworkers and classmates.Along with shared histories, the friendship connections may often point to shared interests between the two people. As a result, people who are connected may often have similar values, experiences, and interests.

## Triadic Closure and Clustering Coefficient

Triadic closure may be seen as the natural propensity of real-world networks to cluster. The following is the triadic closure principle:It is more probable that two people in a social network are already linked or will connect in the future if they have a mutual buddy.According to the triadic closure principle, the network's edge structure has a built-in correlation. This is a logical result of the fact that people who are related to one another are more likely to have similar backgrounds and encounter one other more often. Homophily and the idea of triadic closure are connected. A linked person is more likely to be connected to the same group of actors if their histories are similar, which also causes their attributes to be similar. While triadic closure may be thought of as the structural equivalent of homophily, homophily is often manifested in terms of the content qualities of node attributes. The network's clustering coefficient is intimately connected to the idea of triadic closure.

## Network Formation Dynamics

The formation of networks has an impact on a number of actual characteristics. Social networks and the World Wide Web are two examples of networks that are continually expanding as new nodes and edges are added. Intriguingly, the dynamic mechanisms through which networks from many fields develop exhibit a number of similar traits. The way in

which new edges and nodes are introduced to the network directly affects the network's final structure and the kind of mining that is most efficient. Consequently, the following will go through some typical characteristics of real-world networks:

**Preferential attachment:** As a network expands, a node is more likely to get additional edges the higher its degree. This is a logical result of the fact that people with plenty of connections often find it simpler to form new relationships.This is a logical result of the fact that people with plenty of connections often find it simpler to form new relationships. A model for the probability I in terms of the degree of node I is as follows: If I is the likelihood that a newly inserted node would connect itself to an existing node I in the network:

**The domain, such as a biological or social network**, from which the network is derived determines the value of the parameter. It is common to adopt the scale-free assumption in Web-centric domains. According to this presumption, 1 and the proportionality is linear. Scale-free networks are what these networks are known as. The Barabasi-Albert model is another name for this one. It is hypothesised that many networks, including the World Wide Web, social networks, and biological networks, are scale free, albeit this assumption is clearly just an approximation. In actuality, the scale-free assumption is not entirely compatible with many features of real networks.

**Small world property:** It is believed that the majority of genuine networks are "small world." This indicates that any pair of nodes' average route length is relatively short. In reality, Milgram's experiment from the 1960s suggested that any pair of nodes may be separated by around six. Typically, many models assume that the average route lengths increase as log (n(t)) for networks with n(t) nodes at time t. Even for extremely massive networks, this is a modest quantity. The typical route lengths of large-scale networks, such as Internet chat networks, are relatively short, according to recent research. Experimental evidence has shown that the dynamically variable diameters are even more constrained than the (modelled) log (n (t)) growth rate would imply. This is covered in more detail below.

**Densification nearly all real-world networks**, including the Web and social networks, grow over time more in terms of nodes and edges than they lose. In most cases, the effect of adding new edges outweighs the impact of adding new nodes. This suggests that as time passes, the networks progressively get denser, with the number of edges increasing super-linearly as the number of nodes does. The network displays the following densification power law if n(t) is the number of nodes and e(t) is the number of edges at time t.e(t) $\propto$ n(t) \s$\beta$ (19.3) (19.3)

The exponent has a number that ranges from 1 to 2. The network with the value of = 1 has no effect on the average degree of its nodes as a result of the network's expansion. A network with a value of = 2 has a total number of edges e (t) that, as n (t) grows, retains a constant percentage of the entire graph of n(t) nodes.

**Diminished diameters in the majority of real-world networks**, the average distances between nodes become shorter over time as the network gets denser. Contrary to what traditional theories predict, actual results show that the diameters do not grow as log (n(t)). The dominance of the addition of new edges over the inclusion of new nodes results in this surprising outcome. Keep in mind that the average distances between nodes would grow with time if the effect of adding new nodes were to predominate.

**Gigantic connected component:** A giant connected component appears as the network becomes denser over time. The concept of preferred attachment, which states that freshly arriving edges are more likely to attach themselves to the network's densely linked and high-degree nodes, is compatible with the creation of a massively connected component. This

feature also complicates network clustering techniques since, unless the algorithms are properly constructed, it generally results in imbalanced clusters.

The normal structure of online networks is significantly influenced by preferential attachment as well. A few number of extremely high-degree nodes, sometimes known as hubs, are the consequence. Since the hub nodes often link to several distinct areas of the network, they might confuse many network clustering methods. Because it is not particular to a query or subject, the concept of hubs as presented here differs somewhat from the concept of hubs as mentioned in the HITS algorithm. However, in both instances, the common understanding that nodes serve as the primary hubs of connection in a network is maintained.

**Community Detection**

In the domain of social network analysis, "clustering" and "community detection" are roughly equivalent terms. In the classical work on network analysis, the clustering of networks and graphs is also frequently referred to as "graph partitioning." As a result, there is a wealth of literature in this topic that draws from many other disciplines. Prior to the official study of social network analysis, much work has been done in the area of graph partitioning. But it still has application to the world of social networks. One of the most essential issues in social network analysis is community discovery. After example, one of the clearest and most accessible methods to describe social systems is to summarise closely connected social groupings. Due to several inherent characteristics of typical social networks, network clustering algorithms often struggle to clearly distinguish between various clusters in the social network space.

It is difficult to extend multidimensional clustering techniques to networks, such as the distance-based k-means algorithm. Small-world networks do not allow for a fine-grained indication of similarity because the distances between distinct pairs of nodes are too close together. It is more crucial to overtly or implicitly exploit the triadic closure features of actual networks throughout the clustering process.

Unlike social networks, which often have unique community structures, high-degree hub nodes link several communities together. The edge densities of various social network components vary. In other words, the local clustering coefficients in various areas of the social network are often quite different from one another. Because a single global parameter option is irrelevant in many network locales, using particular selections of parameters to quantify the clusters worldwide results in imbalanced clusters. Real social networks often include a large, intricately interwoven component. This makes it more likely for community discovery methods to produce unbalanced clusters, where one cluster benefits from the majority of the network's nodes.

Such problems are handled by built-in features in several network clustering techniques. The most well-known network clustering methods will be covered in the discussion that follows.

Assume that G = represents the undirected network (N, A). $W_{ij} = w_{ji}$ stands for the weight of the edge I j) between nodes I and j. In certain instances, edge costs (or lengths) are expressed as the inverse notion rather as weights. In these circumstances, we suppose that $c_{ij}$ represents the edge cost. $W_{ij} = 1/c_{ij}$ or another suitable kernel function may be used to translate these values into one another. The challenge in network clustering, also known as community discovery, is to divide the network into k sets of nodes while minimising the weighted sum of the edges with end points in each division. This fundamental objective function is utilised in many different ways in reality to accomplish various application-specific objectives, such as partition balancing, where different clusters have about equal numbers of nodes.

The 2-way cut problem can be solved polynomially in the particular situation when wij = 1 and there are no balancing requirements on partitions. It is suggested that the reader look into Karger'srandomised minimum cut method in the bibliographic notes. For a network of n nodes, this approach may calculate the lowest cut in $O(n2logr (n))$ time, where r is a parameter controlling the desired degree of probabilistic accuracy. However, the cut that results is often unbalanced. The issue becomes NP-hard when arbitrary edge weights or balance constraints are added. The balanced 2-way splitting of graphs is the subject of several network clustering techniques. Recursively creating a k-way partitioning using a 2-way partitioning is possible.

### Increasing speed Kernighan–Lin

The improvements made by Fiduccia and Mattheyses serve as the foundation for a quick version of Kernighan-Lin.

Additionally, this version can manage weights related to both nodes and edges. The method also enables the definition of the ratio representing the degree of balance between the two divisions. To make the benefit Di of Eq. 19.14 as great as feasible, one may just shift a single node I from one partition to the other rather than pairing nodes in an epoch to swap them. At each stage, only nodes that can move without breaking the balancing constraint2 are taken into consideration. Node I is moved and then tagged to ensure that it won't be taken into account again during the current era. To account for this modification, the values of Dj on the remaining vertices j N are adjusted. Until the balance criteria forbids additional movements or until all nodes have been considered for a move in an epoch, this procedure is repeated. The latter is feasible if the nodes lack unit weights or the required partition ratios are not balanced.

Keep in mind that many possible movements in an era could have a loss. As a result, the remaining movements are undone, just as in the original Kernighan-Lin method, and only the best partition generated during each epoch is deemed permanent. Fiduccia and Mattheyses additionally developed a unique data structure to implement each epoch in $O(m)$ time, where m is the quantity of edges. In most real-world networks, a very limited number of epochs are often needed for convergence, however there is no assurance of this.

Karypis and Kumar noted that it is not required to shift as many vertices as feasible in an epoch, contrary to Fiduccia and Mattheyses' original innovation.Instead, an epoch may be terminated if the partitioning objective function does not improve after np moves. The epoch ends once these np moves are undone. The usual choice for np is 50. Furthermore, as long as the gain is positive, it is not necessarily required to relocate the vertex with the biggest benefit. The per-move cost is greatly enhanced by removing the need of selecting the vertex with the biggest benefit.

In many situations, the gains brought about by these simple changes are considerable.

### Girvan-Newman Algorithm

Instead of using edge weights wij, this technique makes use of edge lengths cij. You may think of the edge lengths as being the inverse of the edge weights. Edge weights that are supplied may be heuristically converted to edge lengths using cij = 1/wij or another appropriate application-specific function.Based on the hunch that edges with strong betweenness tend to link various clusters, the Girvan-Newman method was developed. As an example, the edges in Fig. 19.2 that are incident on the hub nodes have a high betweenness. Because so many pairwise shortest routes between nodes in various communities cross via

these edges, they have a high betweenness. The natural clusters in the original network will be matched by a collection of linked components as a consequence of the disconnection of these edges. The Girvan-Newman algorithm is built on this disconnection strategy.

The Girvan-Newman method is a top-down hierarchical clustering technique that divides the network into the necessary number of linked components by gradually deleting the edges with the greatest betweenness. The betweenness values of these edges must be recalculated after each edge removal since each removal affects the betweenness values of some of the other edges.

The calculation of the edge betweenness values is the biggest difficulty in the Girvan-Newman method. An intermediate stage in the calculation of edge-betweenness is the computation of node betweenness values. Remember that the comprehensive set of shortest routes between all source-sink pairs defines all node and edge-betweenness centrality values. Therefore, these betweenness centrality values may be divided into a number of additive components, each of which is specified by the subset of shortest routes leading to a certain source node. For each potential source node s, a two-step procedure is utilised to calculate these betweenness components:

It is calculated how many shortest routes there are from source node s to each other node.

Using the calculations from step one, the component $Bs(i)$ of the node betweenness centrality of node I and the component $bs(i, j)$ of the edge betweenness centrality of edge I j) are calculated. These components correspond to the subset of shortest routes coming from a certain source node s.

**Classification of All**

There are several social networking programmes where labels and nodes may be joined. Consider the scenario of a social networking programme where it would be useful to identify everyone who is interested in golf. Some actors' names may already be known under their own labels. It is preferable to classify nodes for which the label is unknown using the labels that are already accessible.The idea of homophily is essential to the solution to this model. It seems sense to presume that node labels are related since nodes with comparable attributes are often connected. This issue can be easily solved by looking at the k nearby labelled nodes and reporting the majority label. In actuality, this method is a closest neighbour classifier on the network. However, because to the lack of node labels in collective categorization, such a method is often not feasible.

It is thus clear that in addition to using indirect connections via unlabeled nodes, one must also employ direct connections to labelled nodes. In a transductivesemisupervised scenario, where test cases and training instances are categorised simultaneously, collective classification in networks is therefore always carried out. In reality, by converting the data into a similarity graph, collective classification techniques may be utilised for semisupervised classification of any kind of data.

To make it easier to apply a multidimensional classifier like a k-nearest neighbour classifier, the first technique directly converts the graph to multidimensional data. Except for the inclusion of the class information, the embedding method is the same as that used in spectral clustering. The second technique employs a spectral clustering-related optimization formulation to directly train an n k class probability matrix Z. This class probability matrix Z was obtained by label propagation, same as that matrix. It's interesting how tightly tied to label propagation the second technique is.

**Link Forecast**

Predicting future connections between network node pairs is useful in many social networks. For instance, for-profit social networks like Facebook often suggest members as possible buddies. In general, relationships between pairs of nodes may be predicted using both structure and content similarity. These standards are covered below:

Structural actions: The triadic closure concept is often used by structural measurements to make predictions. The theory states that two nodes are more likely to link in the future, if they are not currently connected, if they share comparable nodes in their surroundings.

Measures depending on content: In these situations, predictions are based on the homophily concept. According to the theory, nodes with comparable content are more likely to join together. For instance, a node holding the term "data mining" is more likely to be related to another node containing the keyword "machine learning" in a bibliographic network showing relationships between scientific co-authors.

Although it has been shown that content-based measurements may improve link prediction, the outcomes are fairly network-specific. For instance, content-based approaches are ineffective in a network like Twitter where the content takes the form of brief, obnoxious messages with a lot of unusual acronyms. Furthermore, whereas structural connection often suggests homophily based on content, this is not always the case.

As a consequence, the application of content similarity has varied outcomes across various network domains. However, in various kinds of networks, structural measures are usually always successful. This is due to the fact that triadic closure is commonplace across several network domains and applies more directly to link prediction.

**Neighbor-Based Interventions**

Neighborhood-based metrics, in various ways, estimate the probability of a connection between a pair of nodes I and j by counting the shared neighbors between them. For instance, Alice and Bob have four neighbors in common. Consequently, it is conceivable that a connection between them may one day develop. They have separate groups of neighbors in addition to their usual neighbors. Different approaches of normalise neighborhood-based measurements exist to take into consideration the quantity and relative significance of various neighbours.

**Katz Scale**

Although neighborhood-based metrics provide a reliable assessment of the possibility of a connection establishing between two nodes, they are not as useful when there are few common neighbours between the nodes. Jim and Alice also have a mutual neighbour. As a result, in these scenarios, neighborhood-based metrics find it challenging to discern between various pairwise predictions strengths. However, it seems that longer pathways also contribute significantly to the indirect connectedness in these situations. In these circumstances, walk-based interventions are more suitable. The Katz measure is a specific walk-based metric that is often used to assess the link-prediction potency.

**Walk-Based Random Measures**

An alternative method of determining connectedness between pairs of nodes is using random walk-based metrics. PageRank and SimRank are examples of such metrics. These techniques won't be covered in depth here. The customised PageRank of node j is used in the first method of calculating the similarity between nodes I and j, where node I serves as the restart node.

According to the theory, when the restart is carried out at node I j will have a very high customised PageRank measure if it is structurally close to i. This suggests that nodes I and j have stronger predicted links. A measure that is asymmetric between nodes I and j is the customised PageRank. One may take the average of the values of P ersonalizedP ageRankand PersonalizedP ageRank since the discussion in this section is for the situation of undirected graphs (j, i).The SimRank measure, which is already symmetric, is another option. This measurement determines the inverse function of the distance that two randomly travelling surfers would have to travel to cross paths. The link prediction measure is presented as the equivalent value.

**Link Prediction as a Problem of Classification**

Unsupervised heuristics make up the aforementioned measurements. One of these methods may be more successful for a certain network, but another one may be more beneficial for a different network. How can this conundrum be solved, and what are the best solutions for a certain network?

By interpreting the existence or absence of a link between a pair of nodes as a binary class indication, it is possible to approach the link prediction issue as a classification problem. As a result, for each pair of nodes, a multidimensional data record may be retrieved. All of the distinct neighborhood-based, Katz-based, or walk-based commonalities between nodes are among the characteristics of this multidimensional record. The node-degrees of each node in the pair as well as many additional preferential-attachment properties are employed. Consequently, a multidimensional data record is built for each pair of nodes. The outcome is a classification issue with positive and unlabeled instances, where node pairs with edges are the positive examples and the other pairs are unlabeled examples. For training purposes, the unlabeled instances may be roughly regarded as negative examples. Only a sample of the negative examples is utilised since there are too many negative example pairings in big and sparse networks. Consequently, the following is how the supervised link prediction algorithm operates:

Create a multidimensional data set for the training phase that contains a sample of data records from pairs of nodes with and without edges in between them, as well as one data record for each pair of nodes with an edge between them. The retrieved structure and similarity characteristics between node pair are the features. The existence or absence of an edge between the pair serves as the class label. Create a training model using the information.

The multidimensional record for each test node pair is created. To predict labels, use any standard multidimensional classifier. The underlying classification issue is unbalanced, cost-sensitive variants of many classifiers are often utilized. This strategy has the benefit of allowing for seamless usage of content features.

A pair of nodes' content similarity, for instance, may be employed. During the training phase, the classifier will automatically determine if these properties are relevant. The technique can also handle directed networks by extracting characteristics in an asymmetric fashion, unlike many other link prediction algorithms.

For instance, indegrees and outdegrees might be used as features in place of node degrees. On directed networks, random walk characteristics may also be specified asymmetrically. For example, computing the PageRank of node j with a restart at node I and vice versa. The supervised model is more adaptable in general because it can learn correlations between linkages and properties of different types.

**Link Prediction as a Missing-Value Estimation Problem**

In general, it is possible to consider the link prediction issue and the recommendation problem as examples of missing value estimate on various kinds of matrices. On user-item utility matrices, recommendation techniques are used, but on incomplete adjacency matrices, link prediction methods are used. The edges of the matrix are all represented by 1s. The remaining entries are only somewhat randomly sampled, and the remainder are presumed to be nonspecific. The values of the missing entries may be estimated using any of the missing-value estimating techniques. Matrix factorization techniques are among the most often utilised techniques in this class. These techniques have the benefit of not requiring symmetry in the given matrix. In other words, directed graphs may also be employed with the method.

**Conversation**

Over various data sets, it has been shown that the various measures are beneficial to variable degrees. Neighborhood-based metrics have the benefit of being efficiently calculated for extremely large data sets. Additionally, their performance is practically on par with that of the other unsupervised metrics. However, random walk-based and Katz-based measurements are particularly helpful in highly sparse networks where it is difficult to determine the number of common neighbours. Although supervision offers more accuracy, it requires costly computing resources. However, in terms of flexibility across different social network domains and readily accessible ancillary data, such as content characteristics, supervision offers the best results.In recent years, link prediction has also benefited from the utilisation of content. Although content may considerably enhance link prediction, structural approaches are still far more effective. Due to the fact that structural metrics directly depend on the triadic characteristics of actual networks. Networks have the triadic characteristic, which holds true for almost all data domains.Contrarily, content-based metrics rely on "reverse homophily," which uses comparable or link-correlated material to forecast links. This has very network domain-specific efficacy. As a result, content-based measurements are often utilised to assist in link prediction and are seldom used independently of link prediction.

**Analyzing social influence**

Individuals' degrees of influence on one another change as a consequence of every social encounter. This is frequently referred to as "word of mouth" effect in more formal social interactions. Online social networks are likewise subject to this fundamental idea. For instance, when an actor tweets a message on Twitter, the message is seen by the actor's followers. The message on the network may often be retweeted by the followers. As a consequence, information, concepts, and viewpoints propagate across the social network. This sort of information dissemination is valued as a route for advertising by many businesses. If a popular message is tweeted to the appropriate audience and spreads like a cascade via the social network, millions of dollars' worth of advertising may be produced. During the Super Bowl game between the San Francisco 49ers and the Baltimore Ravens, the power went out. No issue. Even in the dark, you can still dunk. Thousands of times were retweeted by viewers who adored Oreo's message. As a result, Oreo was able to produce millions of dollars in free advertising that seemed to have more of an effect than expensive television ads during the Super Bowl.

The capacity of various actors to affect their social network peers varies. The following are the two most frequent elements that control an actor's influence:

The degree of their effect is significantly influenced by their centrality inside the social network structure. As an example, important actors are more likely to have high centrality

values. Actors with a strong reputation are more likely to be influential in directed networks. The possibility that the respective pair of actors may affect one another is often used to determine the weights attached to the edges in a network. These weights may sometimes be explicitly interpreted as influence propagation probabilities, depending on the diffusion model in use. These probabilities might be affected by a number of variables. A well-known person, for instance, could have more influence than someone who is less well-known. Similar to this, long-term buddies are more likely to have an impact on one another. Although some new techniques demonstrate how to estimate these probability using data-driven approaches, it is often believed that the impact propagation probabilities are already accessible for analytical reasons. An influence propagation model is used to calculate the exact effect of the aforementioned components. They are also known as diffusion models. Such models' major objective is to identify a set of seed nodes in the network where information propagation optimizes impact.

**Model of Independent Cascade**

Once a node becomes active in the aforementioned linear threshold model, it has several opportunities to affect its neighbors. The threshold that the random variable I was linked with was a node. The independent cascade model, on the other hand, gives each active node a single opportunity to activate its neighbors, with propagation probability linked to the edges. Pij represents the propagation probability connected to an edge. Only the newly active nodes that have not previously been activated are permitted to impact their neighbors throughout each iteration. Each edge I j) connecting a particular node to its recently active neighbors I separately flips a coin with a success probability of pij. The node j is triggered if the coin toss for edge I j) is successful. In the next cycle, if node j is engaged, it will only have one opportunity to affect its neighbors. When no new nodes are activated during an iteration, the algorithm ends.

The impact function value at termination is the same as the quantity of active nodes. A coin is thrown for each edge no more than once throughout the algorithm's runtime since nodes are only permitted to affect their neighbors once.

**Evaluation of Influence Function**

The influence function f(S) may be calculated using a model using either the linear threshold model or the independent cascade model. In most cases, simulation is used to estimate f(S). Consider the linear threshold model, where the thresholds at the nodes may be determined using a random number generator for a given seed node set S. Any deterministic graph search method may be used to identify the active nodes once the thresholds have been defined, beginning with the seed nodes in S and gradually activating nodes when the threshold condition is met.

To provide more reliable estimates, the process might be performed for other sets of thresholds that were chosen at random. Another simulation may be employed in the independent cascade model. For each edge, a coin with probability pij may be flipped. If the coin toss was successful, the edge is shown as live. When a live edge route exists from at least one node in S to a node, it can be shown that the independent cascade model will ultimately activate that node.

By simulating, the size of the (final) active set may be estimated using this. The calculation is carried out several times, and the outcomes are averaged. Pointers in the bibliographic notes provide evidence that the linear threshold model and the independent cascade model are sub modular optimization issues. These models are not the only ones that have this trait, however.

When the rules of diminishing returns are applied to the incremental effect of individual influence in bigger organizations, sub modularity is a fairly natural result. As a consequence, sub modularity will be satisfied by the majority of viable models for impact analysis.

**----------------------------**

<center>**CHAPTER 13**</center>

<center>**APPLICATIONS AND TRENDS IN DATA MINING**</center>

<center>Sachin Jain</center>
<center>Assistant Professor, School of Computer & Systems Sciences, Jaipur National University, Jaipur, India,</center>
<center>Email Id-sachin.jain@jnujaipur.ac.in</center>

The principles and techniques for mining relational data, data warehouses, and complex types of data (such as stream data, time-series and sequence data, complex structured data, spatiotemporal data, multimedia data, heterogeneous multi-database data, text data, and Web data) have been covered in earlier chapters of this book. Since data mining is a relatively new field with many applications, there is currently a substantial gap between generic data mining concepts and efficient data mining tools. In this part, we look at a few application areas and talk about the need for developing tailored data mining tools.

**Analyzing Financial Data using Data Mining**

The majority of banks and financial organisations provide a broad range of banking services, including business and individual checking and savings accounts, credit services including business, mortgage, and vehicle loans, and investment services like mutual funds. Some also provide stock investing and insurance services.The banking and financial sector often collects financial data that is comparatively full, trustworthy, and of high quality, which makes systematic data analysis and data mining easier. Here, we outline a few common situations:

Data warehouse design and development for multidimensional data processing and data mining For banking and financial data, data warehouses must be built, much as for many other applications. The general characteristics of such data should be examined using multidimensional data analysis techniques. One may find it interesting to see the changes in debt and income by month, by area, by industry, and by other parameters, as well as the maximum, minimum, total, average, trend, and other statistical data. The analysis and mining of financial data makes extensive use of data warehouses, data cubes, multi-featured and discovery-driven data cubes, characterization and class comparisons, and outlier analysis. Analysis of client credit policies and loan payment forecasting: A bank's operations depend heavily on credit research and loan payment forecasting. Both loan payment performance and client credit rating may be substantially or weakly influenced by a variety of circumstances. Finding significant elements and removing unimportant ones may be aided by data mining techniques like attribute selection and relevance ranking. For instance, loan-to-value ratio, loan length, debt ratio (total monthly debt compared to total monthly income), payment-to-income ratio, customer income level, education level, residency location, and credit history are all elements that affect the risk of loan payments. If the payment history of a client is examined, it may be discovered that, for example, the payment-to-income ratio dominates the analysis while the debt-to-income ratio and degree of education do not. The bank may then choose to modify its loan-granting policy to provide loans to clients whose applications had previously been turned down but whose profiles, as determined by the key factor analysis, showed comparatively low risks.

**Customer classification and grouping for focused marketing**: For client group identification and focused marketing, classification and clustering algorithms may be applied. For instance, categorization may be used to determine the most important elements that may

affect a customer's choice of bank. Multidimensional clustering algorithms may be used to find customers that have similar loan payment habits. These may aid in customer group identification, the matching of a new client with the proper customer group, and targeted marketing.Integrated information from several databases, such as bank transaction databases and federal or state criminal history databases, is crucial for the detection of money laundering and other financial crimes, provided that they may be relevant to the research. Then, other data analysis methods may be utilised to find odd patterns, including big cash flows by certain consumer groups at particular times. Tools that are helpful include data visualisation tools (used to display transaction activities using graphs by time and by groups of customers), linkage analysis tools (used to identify links between various customers), clustering tools (used to group various cases), outlier analysis tools (used to find outliers in unusually high or low amounts of fund transfers or other activities), and sequential pattern analysis tools (to characterise unusual access sequences). These techniques may assist investigators concentrate on suspect cases for further in-depth investigation by highlighting significant links and patterns of behaviour.

## Retail Industry Data Mining

Since it gathers enormous volumes of data on sales, customer shopping histories, the transportation of products, consumer consumption, and service, the retail sector is a key application area for data mining. The amount of data being gathered is growing quickly, particularly in light of how simple, accessible, and well-liked online shopping has become. Today, a lot of retailers provide online shopping via their websites. Some companies, like Amazon.com (www.amazon.com), are entirely online-only enterprises with no physical stores. A valuable source for data mining is retail data.Retail data mining can assist in identifying customer purchasing patterns, uncovering customer shopping trends and patterns, enhancing customer retention and satisfaction, increasing product consumption ratios, creating more efficient product transportation and distribution plans, and lowering business costs.

Following are a few instances of data mining in the retail sector.

Based on the advantages of data mining, design and construction of data warehouses: There are several methods to create a data warehouse for the retail sector since the data it generates is so diverse (covering sales, customers, staff, products transit, consumption, and services, among other things). The amount of information that should be included might likewise differ greatly. Data warehouse architectures may be designed and developed using the results of basic data mining operations as a guide. Multidimensional analysis of sales, customers, products, time, and region: The retail industry requires timely information regarding customer needs, product sales, trends, and fashions, as well as the quality, cost, profit, and service of commodities. This involves choosing which dimensions and levels to include and what preprocessing to perform in order to facilitate effective data mining. Therefore, it is crucial to provide robust tools for multidimensional analysis and visualization, including those for building complex data cubes in accordance with the requirements of data analysis. In retail data analysis, the multi-feature data cube is a helpful data structure since it makes it easier to analyze aggregates with complicated circumstances. Analysis of sales campaign effectiveness: To promote items and draw consumers, the retail sector runs sales campaigns employing adverts, coupons, and different types of discounts and incentives. The profitability of a firm may be increased by carefully analyzing the performance of sales initiatives. This may be accomplished using multidimensional analysis by contrasting the volume of sales and the number of transactions involving the sale products during the sales period with those involving the same items prior to or after the sales campaign. Additionally, association

research may reveal which products are most likely to be bought alongside the discounted goods, particularly when compared to purchases made before or after the campaign.

**Analysis of customer loyalty and customer retention**: Using information from customer loyalty cards, one may track specific consumers' purchase sequences. Trends in customer loyalty and purchases may be carefully examined. Sequences may be created from products that the same consumers bought at several times. The use of sequential pattern mining may then be utilized to analyze shifts in consumer behavior or brand loyalty and make recommendations for modifications to product price and variety that will both help maintain current consumers and draw in new ones.

**Cross-referencing of goods and product recommendations:** One may learn that a consumer who purchases a digital camera is likely to purchase another set of goods by mining associations from sales data. Product suggestions may be made using this data. In real-time consumer interactions, collaborative recommender systems employ data mining methods to provide tailored product suggestions based on feedback from other customers. To enhance customer service, assist consumers in making purchases, and boost sales, product recommendations may also be published online, in weekly flyers, or on sales receipts. To encourage purchases, information like "popular goods this week" or alluring discounts might be provided with the associated information.

### The Telecommunications Industry and Data Mining

The telecommunications sector has quickly developed from providing local and long-distance phone services to offering a wide range of other comprehensive communication services, such as fax, pager, cell phone, Internet messenger, images, e-mail, computer and Web data transmission, and other data traffic. Additionally, efforts are being made to integrate telecommunication, computer networks, the Internet, and various other forms of communication and computation. In addition, the telecommunications market is increasing quickly and is very competitive due to the deregulation of the sector in many nations and the development of new computer and communication technology. This increases the need for data mining to better understand the company at hand, discover communications trends, uncover fraudulent activity, optimize resource use, and enhance service quality. Here are a few situations when data mining might enhance telecommunications services:

**Data from communications are analyzed on several dimensions:** The location of the caller, the location of the caller, the kind of call, and the calling time are just a few of the aspects that make telecommunication data inherently multidimensional. It is possible to detect and compare the data flow, system burden, resource utilization, user group behavior, and profit using multidimensional analysis of such data. For instance, industry analysts may find it useful to often review graphs and charts that show use trends for the calling source, destination, volume, and time of day. Consolidating telecommunication data into sizable data warehouses and consistently doing multidimensional analysis utilizing OLAP and visualization tools are therefore often beneficial.

**Detecting strange patterns and analyzing fraudulent patterns**: The telecommunications sector suffers annual losses of millions of dollars due to fraud. It is crucial to

(1) Identify potentially fraudulent users and their unusual usage patterns,

(2) Spot attempts to access customer accounts fraudulently,

(3) Find out-of-the-ordinary patterns that might require special attention, such as busy-hour frustrated call attempts, switch and route congestion patterns, and periodic calls from

incorrectly programmed automatic dial-out equipment (like fax machines). Multidimensional analysis, cluster analysis, and outlier analysis may all be used to identify many of these patterns.

Analyzing sequential patterns while using several dimensions: The promotion of telecommunication services may benefit from the identification of association and sequential patterns in multidimensional analysis. Consider the scenario where you want to discover utilization trends for a collection of communication services by client segment, month, and hour of the day.

**Services for mobile communications:** Mobile communications, Web and information services, and mobile computers are all becoming more integrated and widespread in our daily lives and work. The connection of mobile telecommunication data with spatiotemporal data is a crucial aspect. Finding specific patterns may require the use of spatial and temporal data mining. For instance, extremely heavy mobile phone traffic at certain regions might be a sign of aberrant activity there. Additionally, making new mobile services simple to use is essential for luring customers to utilise them. The development of adaptable solutions that let consumers get valuable information with just a few keystrokes will probably heavily rely on data mining.

Using visualisation tools for study of telecom data for the study of telecommunication data, it has been shown that tools for OLAP visualisation, linkage visualisation, association visualisation, clustering, and outlier visualisation are particularly helpful.

**Analysis of Biological Data Using Data Mining**

Research in genomes, proteomics, functional genomics, and biomedicine has grown rapidly during the last ten years. Examples include the identification and comparative analysis of the human and other species' genomes (by identifying sequencing patterns, gene functions, and evolutionary pathways), the study of genetic networks and protein pathways, the creation of novel medications, and advancements in the treatment of cancer. Bioinformatics, a recent branch of science, now heavily relies on biological data mining. It is difficult to fully address such a significant and developing issue in a single subsection given the breadth, depth, and dynamic nature of the area of biological data mining. Only a few intriguing subjects in this area are discussed here, with a focus on the analysis of genomic and proteomic data. A thorough overview of biological data mining might take up many volumes. Numerous books on bioinformatics and the analysis of biological data have already been released, and more are anticipated. We list citations in our bibliographic notes.

The genetic instructions of all living things are built on DNA sequences. The four fundamental building units of DNA, known as nucleotides, are adenine (A), cytosine (C), guanine (G), and thymine (T). These four bases (or nucleotides) are joined to create lengthy chains that resemble twisted ladders. The information and metabolic machinery that can be passed down from generation to generation are carried by the DNA. Insertions, deletions, or mutations (also known as substitutions) of nucleotides are introduced into the DNA sequence throughout the "copying" processes, resulting in various evolutionary pathways. A gene typically consists of hundreds of distinct nucleotides organised in a certain order. To create unique genes, the nucleotides may be arranged and sequenced in an almost infinite number of different ways. The whole of an organism's genes make up its genome. An estimated 20,000–25,000 genes make up the human genome. The study of genome sequences is known as genomics.

Any creature needs proteins as necessary chemicals. They are the bulk of cellular structures and carry out life processes. Through a variety of translational changes and gene splicing processes, the about 25,000 human genes produce nearly 1 million proteins. The constituent parts of proteins are known as amino acids (or residues). There are 20 amino acids, each represented by a different alphabet letter. One or more triplets of DNA's nucleotides are responsible for encoding each amino acid. Another set of triplet's codes for the chain's finish.

It is difficult to pinpoint DNA or amino acid sequence patterns that contribute to numerous biological processes, genetic disorders, and evolution. To provide efficient genomic and proteomic data processing tools, a lot of research in computing algorithms, statistics, mathematical programming, data mining, machine learning, information retrieval, and other fields is needed.

**In the following ways, data mining might support the study of biological data:**

Diverse, dispersed genomic and proteomic databases semantically integrated: Data sets for genomic and proteomic analyses are often produced using various techniques at various laboratories. They are dispersed, heterogeneous, and come in many different forms. The cross-site analysis of biological data requires the semantic integration of such data. Finding accurate connections between research literature and the biological entities that go along with it is also crucial. The systematic and coordinated study of genetic and biological data would be made easier by such integration and linkage analysis.

To store and manage the original and derived biological data, this has encouraged the creation of integrated data warehouses and distributed federated databases.Methods for data cleansing, data integration, reference reconciliation, classification, and clustering will make it easier to integrate biological data and build data warehouses for the study of biological data.Multiple nucleotide/protein sequence alignment, indexing, similarity searching, and comparative analysis techniques have been developed during the last two decades. Tools for the systematic study of genomic and proteomic data include BLAST and FASTA, in particular. The many sequential pattern analysis methods suggested in data mining research are different from biological sequence analysis techniques. To handle with insertions, deletions, and mutations, they should allow for gaps and mismatches between a query sequence and the sequence data to be searched. Two amino acids should also be regarded as "matching" in protein sequences if they may be formed from one another via substitutions that are likely to happen in nature. When creating alignment algorithms, sophisticated statistical analysis and dynamic programming techniques often play a crucial role. Such data sets may be used to create indices, which will facilitate efficient similarity and precise searches.

There are an infinite number of possible combinations for approximation alignment of several sequences. Multiple sequence alignment is thus thought to be a more difficult process. Using Hidden Markov Models, often known as HMMs, and breaking down a multiple alignment into a series of pairwise alignments are two techniques that might be helpful. The effective and methodical alignment of several biological sequences is still a subject of continuous study. Using phylogenetic trees to infer evolutionary links between species, multiple sequence alignments may be utilised to find highly conserved residues across genomes. Furthermore, it could assist in revealing the genetic mysteries of evolution.

Medical scientists may analyse genomic and proteomic sequences obtained from sick and healthy tissues to pinpoint significant differences between them. Sequences that appear more often in the afflicted samples may point to the disease's genetic causes. Those showing up more often solely in the healthy samples may point to defences the body has against the

illness. Although similarity search is necessary for genetic analysis, the method used here is different from that for time-series data. Since genetic data are nonnumeric, data transformation techniques including scaling, normalisation, and window stitching are often used in time series analysis but are useless for this kind of data. These approaches ignore the interactions between nucleotides, which are crucial for biological activity. It is crucial to continue creating effective sequential pattern analysis techniques for comparing biological sequences.

Identification of structural patterns and examination of protein pathways and genetic networks Protein sequences are folded into three-dimensional structures in biology, and these structures communicate with one another depending on their locations in relation to one another and how far apart they are from one another. Intricate genomic networks and protein pathways are built on the foundation of such complex connections. Finding structural patterns and regularities in such vast yet complicated biological networks is vital. In order to find approximate and common structural patterns and to research the regularities and irregularities among such linked biological networks, it is crucial to build effective and scalable data mining tools.

Identifying co-occurring gene sequences and connecting genes to various phases of disease development using association and path analysis: Comparing one gene to another has been the topic of several recent research. But rather than being caused by a single gene, most illnesses are caused by many genes working together. It is possible to identify the kind of genes that are most likely to co-occur in target samples by using association analysis approaches. The identification of gene families and the investigation of their interactions and linkages would be made easier by such analysis.

Various genes may become active at different phases of the illness, even if a set of genes may contribute to the development of a disease. It may be feasible to design pharmacological treatments that target the various phases of disease formation independently, leading to more effective disease therapy, if the sequence of genetic activity throughout the many stages of disease development can be discovered. Such route analysis is anticipated to be crucial in genetic research.

**Tools for visualizing genetic data analysis:** The best way to convey alignments between genomic or proteomic sequences and the relationships between intricate biological structures is graphically, using different types of simple visual presentations. Such aesthetically pleasing patterns and structures aid in the comprehension of patterns, information discovery, and interactive data exploration. Therefore, visual data mining and visualization are crucial in the study of biological data.

**Applications of Data Mining in Other Scientific Fields**

Prior until now, the majority of scientific data processing activities tended to deal with homogeneous, tiny, and compact data sets. The "formulate hypothesis, develop model, and assess findings" approach was often used to examine such data. For their analysis in these situations, statistical approaches were suitable and often used. Recent advancements in data gathering and storage technology allow for the faster and more affordable acquisition of scientific data. Huge amounts of high-dimensional data, stream data, and heterogeneous data, comprising rich geographical and temporal information, have accumulated as a consequence of this. As a result, scientific applications are moving away from the "hypothesize-and-test" paradigm and toward a process where data is collected and stored, mined for new hypotheses, and then confirmed by data or testing. Data mining now faces additional difficulties as a result of this change. Using powerful telescopes, multispectral high-resolution remote

satellite sensors, and global positioning systems, enormous volumes of data have been gathered from scientific fields (including geosciences, astronomy, and meteorology). Fast numerical simulations in a variety of disciplines, including chemical engineering, fluid dynamics, structural mechanics, and climate and ecosystem modelling, are generating large data sets. In addition, the fields of biomedical engineering and telecommunications need the processing of significant volumes of complicated data.

In this part, we examine some of the difficulties caused by new scientific data mining applications, including the following:

**Data preprocessing and data warehouses:** Data warehouses are essential for data mining and information sharing. However, there is currently no actual geographic data warehouse available. Finding ways to resolve spatial and temporal data incompatibilities, such as balancing semantics, reference systems, geometry, accuracy, and precision, is necessary for building such a warehouse. Methods are required for recognising events and integrating data from diverse sources (such as data spanning several time periods) for scientific applications in general. The issue arises when there are too many occurrences in the geographical domain and not enough in the temporal domain, as in the case of climate and ecosystem data, which are both spatial and temporal. (For instance, El Nino episodes only happen every four to seven years, and earlier data may not have been gathered as methodically as it is now.) Methods are required for the management of spatially connected data streams as well as the efficient computing of complex spatial aggregates.

Complex data types mining scientific data sets are often diverse and include semi-structured and unstructured data, including multimedia and georeferenced stream data. Spatiotemporal data processing, linked idea hierarchies, and complicated geographic linkages (such non-Euclidian distances), all need robust approaches.

Mining using graphs Due to the shortcomings of current modelling methodologies, it is sometimes difficult or impossible to simulate a variety of physical events and processes. As an alternative, labelled graphs may be utilised to represent many of the topological, geometrical, topographical, and other relationship aspects of scientific data sets. Each item that has to be mined in graphmodeling is represented by a vertex in a graph, and the edges between vertices signify the connections between the objects. Graphs may be used, for instance, to represent the data produced by numerical simulations of fluid flow and chemical structures. However, for graph modelling to be successful, several traditional data mining activities, including classification, frequent pattern mining, and clustering, must be made more scalable and effective.

Tools for visualisation and expertise in a certain domain: For scientific data mining systems, high-level graphical user interfaces and visualisation tools are necessary. These should be combined with already-existing database and domain-specific information systems to aid researchers and everyday users in finding patterns, understanding and displaying patterns found, and using knowledge gained in decision-making.

### Data Mining to Find Intruders

Our computer systems' and our data's security is always under danger. Intrusion detection has evolved into a crucial part of network management as a result of the Internet's rapid expansion and the proliferation of tools and techniques for hacking and attacking networks. Any series of acts that jeopardise the availability, integrity, or secrecy of a network resource (such as user accounts, file systems, system kernels, and so on) is referred to as an intrusion.

The majority of commercial intrusion detection systems have limitations and don't provide a full answer. Typically, abuse detection is used by such systems. Misuse detection looks for user or programme activity patterns that match recorded signatures for recognised intrusion situations. Based on their in-depth understanding of infiltration strategies, human specialists laboriously produce these hand-coded signatures.

An event for which an alert is triggered is signalled if a pattern match is discovered. The actions that should be taken, like as shutting down a portion of the system, informing the appropriate Internet service provider of suspicious traffic, or just noting odd activity for later reference, are determined by human security experts after they examine the alerts. A highly sophisticated network's intrusion detection system often generates hundreds or millions of alerts each day, which is a burdensome chore for the security analysts. Systems are dynamic, thus the signatures must be updated anytime new software releases or modifications to network setup take place. The fact that abuse detection can only find situations that match the signatures is another significant limitation. That is to say, it cannot recognise novel or previously undiscovered intrusion strategies. Anomaly detection techniques may discover new invasions. In order to find novel patterns that substantially differ from the profiles, anomaly detection creates models of typical network activity (referred to as profiles). These variations can constitute real invasions or they might just be brand-new behaviours that need to be incorporated to the profiles. The primary benefit of anomaly detection is that it has the potential to identify brand-new invasions that have not previously been noticed. The variations are often sorted by a human analyst to determine which ones are actual incursions. The high rate of false positives is a limiting issue in anomaly detection. The collection of signatures for misuse detection may be expanded to include new patterns of infiltration.

This debate has shown that standard intrusion detection systems currently in use have significant drawbacks. Data mining for intrusion detection is becoming more popular as a result of this. Data mining technology may be used or improved in the following areas for intrusion detection:

**Data mining methods for intrusion detection development:** Algorithms for data mining may be used to find anomalies and identify abuse. Training data are classified as "normal" or "intrusion" in misuse detection. Then, a classifier to find known invasions may be developed. Application of classification algorithms, association rule mining, and cost-sensitive modelling have all been studied in this field. Anomaly detection creates models of typical behaviour and swiftly finds notable departures from it. You may utilise supervised or unsupervised learning. In a supervised method, training data that are considered to be "normal" are used to build the model. An unsupervised technique provides no details on the training data. Application of classification algorithms, statistical methods, clustering, and outlier analysis have all been used in anomaly detection research. The approaches must be effective, scalable, and able to handle heterogeneous, large volume, and dimensional network data. To choose and create discriminating qualities, use association and correlation analysis, as well as aggregation: In order to discover connections between the system properties defining the network data, association and correlation mining may be used. Such details may provide light on the choice of practical qualities for intrusion detection. It may also be useful to create new properties from aggregated data, such as summary counts of traffic that fits a certain pattern.

**Data analysis on streams:** It is essential to carry out intrusion detection in the data stream environment because malicious assaults and intrusions have a transitory and dynamic character. Additionally, a single event may be regarded innocent if it occurs randomly, but malicious if it occurs as part of a chain of occurrences. As a result, it's important to research

the common event sequences that occur together, look for sequential patterns, and spot outliers. For real-time intrusion detection, further data mining techniques are required for identifying changing clusters and creating dynamic classification models in data streams.

**Distributed data mining:** Attacks may be conducted from several places and directed to numerous targets. To identify these spread assaults, network data from several network locations may be analysed using distributed data mining techniques.

**Tools for visualisation and querying:** Any observed aberrant patterns should be viewable using visualisation tools. These tools could provide options for associations, clusters, and outliers. Additionally, intrusion detection systems must to feature a graphical user interface that enables security analysts to ask questions about network information or the outcomes of intrusion detection.Intrusion detection systems based on data mining are often more accurate and need far less manual processing and input from human specialists than standard intrusion detection systems.

**Products for Data Mining Systems and Research Prototypes**

There are several off-the-shelf data mining system solutions and domain-specific data mining application softwares available, despite the fact that data mining is still a relatively new topic with many challenges that need more investigation. Data mining is a discipline with a relatively recent history that is constantly changing. Every year, new data mining systems hit the market, and existing systems are routinely updated with new features, functions, and visualisation tools. There are also ongoing efforts to standardise data mining language. As a result, we have no intention of describing commercial data mining techniques in great depth in this book. Instead, we provide a brief overview of a few common data mining systems and discuss the features to take into account when choosing a data mining solution. The bibliographic Notes provide a collection of citations, webpages, and contemporary studies of data mining systems.

**Selecting a Data Mining System**

What kind of system should I pick given the wide range of data mining system goods on the market? Some individuals may believe that data mining systems operate similarly too many commercial relational database systems when it comes to typical features and have the same well-defined processes. If so, the decision would be made more in light of the hardware platform, compatibility, robustness, scalability, affordability, and service of the systems. Unfortunately, reality is far different from this. Many commercial data mining solutions deal with entirely distinct types of data sets and share nothing in terms of data mining capabilities or technique.

A multidimensional perspective of data mining systems is crucial if you want to choose one that is suitable for your work. In general, the several factors listed below should be used to evaluate data mining systems data kinds the majority of commercially available data mining solutions work with formatted, record-based, relational-like data that has numerical, category, and symbol properties. The information may take the form of ASCII text, relational database information, or data from a warehouse. It's critical to confirm the precise format(s) that any system you're thinking about can support. Certain types of data or applications can need specialised algorithms to look for patterns; as a result, their needs might not be met by commercially available, generic data mining solutions. Instead, specialised data mining systems that mine text documents, geographical data, multimedia data, stream data, time-series data, biological data, or Web data may be employed. These systems are either devoted to particular applications or mine a variety of data types (such as finance, the retail industry,

or telecommunications).Additionally, a lot of data mining businesses provide unique data mining solutions that include crucial data mining features or approaches.

**System problems:** One operating system or a number of them may be supported by a certain data mining system. Microsoft Windows and UNIX/Linux are the two most widely used operating systems for hosting data mining applications. Data mining programmes are also available for OS/2, Macintosh, and other platforms. A client/server design is often used by large, industry-focused data mining systems, where the client is typically a personal computer and the server is typically a group of powerful parallel computers. Data mining solutions now often provide Web-based user interfaces and support XML data for input and/or output.data resources This is a reference to the particular data types that the data mining system will use. While many systems operate with relational data, or data warehouse data, accessing several relational data sources, many others merely deal with ASCII text files.A data mining system must support ODBC or OLE DB connectors for ODBC connections. All relational data, including that found in IBM/DB2, Microsoft SQL Server, Microsoft Access, Oracle, Sybase, and other databases, as well as structured ASCII text data, are accessible thanks to these open database connections.

**Data mining techniques and purposes:** A data mining system's core consists of data mining operations. Some systems for data mining only provide a single data mining function, like classification. Others may provide support for a variety of data mining operations, including concept description, discovery-driven OLAP analysis, association mining, linkage analysis, statistical analysis, classification, prediction, clustering, outlier analysis, similarity search, sequential pattern analysis, and visual data mining. Some systems might only support one method for a specific data mining function (like classification), whereas others might support a wide range of methods (like decision tree analysis, Bayesian networks, neural networks, support vector machines, rule-based classification, k-nearest-neighbor methods, genetic algorithms, and case-based reasoning). Different data mining functions and multiple methodologies for each function are supported by data mining systems, giving users more flexibility and analytical capacity. Users may need to combine many mining functions to solve many issues, and some approaches may work better with certain types of data than others. Users may need more education and expertise, nevertheless, in order to benefit from the increased flexibility. Therefore, such systems should also make it simple for new users to access the most used function and method or default settings.

**Integrating database and/or data warehousing systems with data mining:** A database and/or data warehouse system should be connected with a data mining system, and the related components should be easily integrated into a standardised information processing environment. These couplings may generally be divided into four categories: no coupling, loose coupling, semitight coupling, and tight coupling. Some data mining systems are completely uncoupled from database or data warehouse systems and just deal with ASCII data files. These systems struggle to effectively handle huge data volumes and use the data held in database systems. Data are retrieved into a buffer or main memory by database or warehouse operations in data mining systems that are loosely coupled with database and data warehouse systems, after which mining functions are used to evaluate the returned data. When processing data mining queries, these systems may not have scalable methods to handle big data volumes.

A data mining system may be semitightly coupled to a database or data warehouse system, enabling the effective implementation of a few key data mining primitives (such as sorting, indexing, aggregation, histogram analysis, multiway join, and the precomputation of specific statistical measures). An ideal relationship between a data mining system and a database

system would see the integration of the data mining and data retrieval processes by improving data mining queries throughout the iterative mining and retrieval process. Data mining and OLAP operations should be tightly coupled in order to offer OLAP-mining characteristics. This may be done by integrating data mining and OLAP operations.

The two types of scalability problems in data mining are row (or database size) scalability and column (or dimension) scalability. If a data mining system can perform the identical data mining queries no more than 10 times when the number of rows is increased by 10, it is said to be row scalable. If the amount of time it takes to execute a mining query rises linearly with the number of columns, a data mining system is said to be column scalable (or attributes or dimensions). Making a system column scalable is significantly harder than making it row scalable because of the curse of dimensionality. Visualization tools: In data mining, the adage "A picture is worth a thousand words" is quite accurate. Data visualization in data mining may be divided into four categories: data visualization, mining result visualization, mining process visualization, and visual data mining. The usefulness, interpretability, and attractiveness of a data mining system may be significantly impacted by the range, caliber, and flexibility of visualization tools.

Data mining's graphical user interface and query language: Data mining is an exploratory activity. To encourage user-guided, highly interactive data mining, a high-quality graphical user interface is necessary. The majority of data mining technologies provide simple user interfaces for mining. However, unlike relational database systems, most data mining systems do not share any underlying data mining query language. In relational database systems, most graphical user interfaces are built on top of SQL (which acts as a standard, well-designed database query language). Standardizing data mining goods and processes is challenging in the absence of a common data mining language.

### Additional Data Mining Themes

Not all of the subjects related to data mining can be fully addressed in this book due to the extensive scope of data mining and the wide diversity of data mining approaches. We quickly touch on a number of intriguing topics in this part that were not completely covered in the book's earlier chapters.

### Foundations of Data Mining Theory

Theoretical studies behind data mining are still in their infancy. Because it may contribute to the creation of a cogent framework for the advancement, assessment, and use of data mining technologies, a sound and methodical theoretical basis is crucial.

The following are a few possibilities about data mining's origins:

Data compression According to this approach, data mining starts by reducing the data representation. In order to provide rapid approximations to queries on extremely large databases, data reduction compromises accuracy for speed. Singular value decomposition, which serves as the basis for principal components analysis, wavelets, regression, log-linear models, histograms, clustering, sampling, and the creation of index trees are all examples of data reduction techniques.

According to this theory, the fundamental step in data mining is to compress the input data using techniques like bit encoding, association rules, decision trees, clusters, and so on. The "optimal" theory to infer from a collection of data is the one that minimises the length of the theory and the length of the data when encoded, using the theory as a predictor for the data.

This is known as encoding based on the minimal description length concept. Usually, this encoding is done in bits.

The goal of data mining, according to this idea, is to find patterns in a database, such as relationships, classification models, sequential patterns, and so on. This theory draws on a number of different subfields, including machine learning, neural networks, association mining, sequential pattern mining, clustering, and others.Probabilistic reasoning Based on statistical theory, this. According to this theory, the goal of data mining is to identify joint probability distributions of random variables using techniques like hierarchical Bayesian models or Bayesian belief networks.Microeconomic perspective The work of uncovering patterns that are interesting solely to the degree that they may be employed in the decision-making process of some organisation (for example, about marketing strategies and production plans) is referred to as data mining in the microeconomic perspective. According to this utilitarian viewpoint, patterns are only deemed intriguing if they can be put to use. The goal of optimization challenges, which businesses are said to face, is to maximise the utility or value of a choice. This theory states that data mining turns into a nonlinear optimization issue.

According to the principle of inductive databases, a database schema is made up of the data and patterns that are kept there. Since the objective is to query the data and the theory (i.e., patterns) of the database, data mining is therefore the issue of conducting induction on databases. Many database system researchers share this viewpoint.These ideas do not conflict with one another. For instance, pattern discovery may be thought of as a kind of data compression or data reduction. An ideal theoretical framework would be probabilistic in nature, able to handle various types of data, and take into account the iterative and interactive nature of data mining. It would also be able to simulate common data mining tasks (such as association, classification, and grouping). The development of a clear framework for data mining that complies with these demands calls for further effort.

**Analyzing Statistical Data**

The methods for data mining presented in this book are essentially database-oriented, that is, they are created for the effective management of enormous volumes of data, most of which are multidimensional and may be of different complicated sorts. However, there are several tried-and-true statistical methods for analysing data, especially for numerical data.Several forms of scientific data, including those from experiments in physics, engineering, manufacturing, psychology, and medicine, as well as data from economics and the social sciences, have been subjected to considerable use of these methodologies. This book has previously covered several of these methods, including principal component analysis, regression, and clustering. Although it is beyond the scope of this book to examine all relevant statistical techniques for data analysis, a few are presented here for completeness. There are references to these methods in the bibliographic notes.

**Regression:** When the variables are numerical, these approaches are often used to predict the value of a response (dependent) variable from one or more predictor (independent) variables. Regression may take many different forms, including linear, multiple, weighted, polynomial, nonparametric, and robust (robust techniques are helpful when errors don't conform to normality assumptions or when the data include large outliers).

**Generalized linear models:** These models and its generalisation (generalised additive models) enable a categorical response variable to be associated to a collection of predictor variables in a way comparable to the modelling of a numeric response variable using linear regression. Poisson regression and logistic regression are examples of generalised linear

models.Analyzing experimental data for two or more populations using a numeric response variable and one or more category variables using an analysis of variance (factors). ANOVA (single-factor analysis of variance) problems often require comparing the means of k populations or treatments to see whether at least two of the means vary from one another. There are other more difficult ANOVA issues.

Mixed-effect models are used to analyze data that may be categorized based on one or more grouping factors, or grouped data. In data organized according to one or more variables, they often indicate relationships between a response variable and certain covariates. Multilevel data, repeated measurements data, block designs, and longitudinal data are examples of common application areas. Determine which variables are combined to produce a certain factor using factor analysis. For instance, it is often feasible to measure other values that represent the factor of interest, such as student test scores, for many psychiatric data even while it is difficult to assess a particular component of interest directly (such as intellect). None of the variables in this situation are marked as dependent.

**Discriminant analysis:** This method is used to forecast a response variable that is categorical. It makes the assumption that the independent variables have a multivariate normal distribution, unlike generalized linear models. The process looks for various discriminant functions (linear combinations of independent variables) that may distinguish between the groups the response variable defines. In social sciences, discriminant analysis is often used.

Time series analysis: A variety of statistical approaches, including long-memory time-series modelling, unilabiate ARIMA (autoregressive integrated moving average) modelling, and auto regression methods, may be used to analyze time-series data.

There are several accepted statistical methods for conducting a survival study. These methods were first intended to estimate the likelihood that a patient receiving medical care would live at least until time t. However, techniques for survival analysis are often used in manufacturing environments to predict the lifespan of industrial machinery. Popular techniques include Cox proportional hazards regression models and their expansions as well as Kaplan-Meier estimates of survival.Quality assurance Charts for quality control may be created using a variety of statistics, including Shewhart charts and cusum charts (both of which provide group summary statistics). The mean, standard deviation, range, count, moving average, moving standard deviation, and moving range are among these statistics.

**Audiovisual Data Mining**

Using data and/or knowledge visualisation approaches, visual data mining extracts implicit and practical information from enormous data sets. The brain, which can be thought of as a sophisticated, highly parallel processing and reasoning engine with a substantial knowledge base, controls the visual system in humans together with the eyes. In essence, visual data mining combines the strength of these elements, making it a very appealing and useful tool for understanding data distributions, patterns, clusters, and outliers.

Data visualisation and data mining may be seen as two disciplines that are combined in visual data mining. It is also strongly connected to high-performance computing, multimedia systems, pattern recognition, human computer interaction, and computer graphics.

Data mining and data visualisation may generally be combined in the following ways:

**Data visualisation:** Data in a database or data warehouse may be displayed as various combinations of characteristics or dimensions, or at various degrees of granularity or

abstraction. Many other visual representations of data are possible, including boxplots, 3-D cubes, data distribution charts, curves, surfaces, link graphs, and more. The StatSoft figures 11.2 and 11.3 show data distributions in multidimensional space. Users may get a clear idea and overview of the data properties in a database with the use of a visual presentation.

Visualizing data mining results is the process of presenting the findings or information discovered by data mining in a visual manner. These formats might include decision trees, association rules, clusters, outliers, generalised rules, scatter plots, boxplots (obtained through descriptive data mining), and more.

## Impacts of Data Mining on Society

Although we may not always be aware of it, data mining is a part of most people's everyday life. There are various instances of "ubiquitous and invisible" data mining that have an impact on daily activities, such as the items available at our neighborhood supermarket, the advertisements we encounter when browsing the Internet, and crime prevention. By enhancing customer service, contentment, and lifestyle in general, data mining may provide the person with a host of advantages. However, it also has important consequences for data security and one's right to privacy.

### Data mining that is pervasive and invisible

Whether we recognise it or not, data mining is prevalent in many facets of our everyday life. It has an impact on the way we work, shop, and look for information. It may also have an impact on our leisure time, health, and general well-being. Several instances of this pervasive (or constant) data mining. The use of data mining in "smart" software such as Web search engines, customer-adaptive Web services (using recommender algorithms, for example), "intelligent" database systems, e-mail managers, ticketing systems, and other similar applications frequently occurs without the user being aware of it. Several of these examples also represent invisible data mining.

Data mining has creatively changed what we purchase, how we shop, and how we experience buying, from supermarket shops that print tailored discounts on customer receipts to online businesses that suggest extra things based on consumer preferences. Wal-Mart is one such; each week, more than 100 million people visit its more than 3,600 locations in the United States. On NCR Teradata mainframes, Wal-Mart has 460 terabytes of point-of-sale data stored. Experts believe that the Internet only has less than half of this data, so that puts things into perspective.

Suppliers are given access to data about their items by Wal-Mart so they may conduct analysis utilizing data mining technologies. This enables suppliers to track inventory and product placement, determine new merchandising possibilities, and discover consumer purchasing habits. Consider all of these the next time you browse the aisles at Wal-Mart since they all have an impact on the products that get placed on the shelves of the shops (as well as how many of them).

The online purchasing experience has been altered by data mining. Online retailers are often used by consumers to buy toys, music, movies, and books. Collaborative recommender systems, which provide individualized product suggestions based on the feedback of other consumers. The use of such a customised, data-driven approach as a marketing tactic was pioneered by Amazon.com. The toughest aspect of running a conventional brick-and-mortar business, according to CEO and founder Jeff Bezos, is getting customers inside. Due to the

high expense of visiting another business, once the client is there, she is more likely to make a purchase.

As a result, rather than focusing on the actual in-store customer experience, marketing for physical and mortar shops often focuses on attracting customers. In contrast, clients at online retailers may "walk out" and visit another online business with the simple click of a mouse. Amazon.com benefited from this distinction by providing a "customised shop for every consumer," using a variety of data mining tools to determine their preferences and provide accurate suggestions.

While we're talking about shopping, let's say you've been using your credit cards a lot. Nowadays, receiving a call from one's credit card provider questioning suspicious or odd spending habits is not uncommon. Data mining is used by credit card firms (and long-distance telephone service providers, for that matter) to identify fraudulent activity, which results in annual savings of billions of dollars.

In place of mass marketing, many businesses are turning to data mining for customer relationship management (CRM), which enables them to provide more individualised, personalised service that addresses the demands of each individual consumer. Companies may better target their marketing and promotions to individual consumers by researching their online shopping and browsing habits. This reduces the likelihood that customers will get irritated by unsolicited bulk mailings or junk mail. Companies may save a lot of money by taking these steps. Customers also gain since they are more likely to hear about offers that are truly interesting, which leads to less personal time wasted and higher satisfaction. As we'll see in a moment, this reoccurring topic might appear throughout our day multiple times.

Data mining has had a significant impact on how individuals use computers, do information searches, and go about their daily lives. Consider that you have just connected on to the Internet while using your computer. It's likely that you have a tailored gateway, which means that the first Web page that your Internet service provider displays is made to see and feel like it caters to your individual interests. This idea was initially introduced by Yahoo (www.yahoo.com). MyYahoo use records are mined to provide Yahoo useful knowledge about a user's Web usage patterns, allowing Yahoo to deliver personalised content. According to the Media Power 50 list published by Advertising Age's BtoB magazine (www.btobonline.com), which honours the 50 most effective and focused business-to-business advertising sources each year, Yahoo has been consistently ranked as one of the top Web search providers for years.

You choose to check your email after connecting to the Internet. Unbeknownst to you, a spam filter that employs classification algorithms to identify spam has already eliminated a number of irritating emails. Once your email has been processed, you visit Google (www.google.com), which gives you access to data from over 2 billion Web sites that have been indexed on its server. One of the most well-known and often used Internet search engines is Google. For many individuals, searching for information on Google has become second nature. A new verb in English, meaning "to search for (anything) on the Internet using the Google search engine or, by extension, any comprehensive search engine," was created as a result of Google's immense popularity. 1 You choose to provide a few keywords for an appealing subject. Google provides a list of websites related to your search query, mined and sorted by PageRank. Instead than relying exclusively on Web content to determine which sites were relevant to a search, older search engines used structural link data from the Web graph to determine a page's relevance. It serves as Google's Web mining technology's brain.

Various advertising that are related to your search appear while you are browsing the results of your Google inquiry. The effectiveness of Google's approach of customising advertising to users' interests may be seen in the four- to five-fold increase in clicks for the firms concerned.

A tool called "web-wide tracking" follows a person to all the websites she visits. Therefore, when browsing the Internet, information about each website you visit may be logged, giving advertisers information about your interests, way of life, and habits. Using Web-wide monitoring, DoubleClick Inc.'s DART ad management system targets advertising based on demographic or behavioural characteristics. On their websites, businesses pay to utilise DoubleClick's service. The clickstream data from all of the DoubleClick-enabled websites is combined and analysed to discover user profiles for those who visit these websites.

Then, on behalf of its customers, DoubleClick may modify adverts to appeal to end users. In general, customer-tailored advertising is not only confined to mail-outs from businesses or ads posted on websites with storefronts. Future digital television, online books, and online newspapers may potentially provide adverts that are created and chosen precisely for a certain viewer or viewing group based on demographic and consumer profile data. You are delighted with the offers that have been submitted so far, supposing that they are genuine. Fortunately, eBay now employs data mining to differentiate between genuine and fake bids.

Data mining and OLAP technologies, as we have seen throughout this book, may be quite useful to us in our job. Governments, scientists, and business analysts may all utilise data mining to examine and understand their data. They may utilise OLAP and data mining technologies without having to be familiar with the specifics of any underlying techniques.

The user only cares about the final output that these systems deliver, which they may then analyse or utiliseto guide their decision-making.Data mining may also affect how we use our free time, including when we eat and watch movies.

Let's say you decide to stop for fast food on your way home from work. A well-known fast food chain uses data mining to analyse time series data and market baskets to identify consumer behaviour. As a result, a push to turn "drinkers" into "eaters" was started by giving hamburger-drink combos for not much more than the cost of the drink alone. The next time you order a meal combination, keep that in mind. The restaurant may even know what you want to order before you get to the counter with a little data mining assistance. Bob, an automated management system for fast-food restaurants created by HyperActive Technologies (www.hyperactivetechnologies.com), makes predictions about what customers will likely order based on their height and the sort of automobile they drive to the business. For instance, the consumer is likely to purchase a quarter pounder if a pick-up truck drives up. Children are likely to be in a family automobile, which calls for fries and chicken nuggets. The goal is to offer the cooks with guidance on the best cuisine to prepare for arriving clients in order to speed up service, provide better meals, and minimise food waste.

### Data Security, Privacy, and Mining

There are growing worries that data mining might be a danger to our privacy and data security as more and more information is available online and in electronic form, and as more potent data mining techniques are created and used. It is crucial to remember that the majority of the most popular data mining tools do not even touch personal information. The use of natural resources, the forecasting of floods and droughts, meteorology, astronomy, geography, geology, biology, and other scientific and engineering data are a few prominent examples. Additionally, the majority of data mining projects do not use personal data and instead concentrate on creating scalable algorithms. Data mining technology focuses on

finding broad trends rather than detailed information about specific people. In this regard, we think that the actual privacy issues are with unrestricted access to individual data, such as credit card and banking apps, which must access information that is sensitive to privacy, as an example. In many instances, simple techniques like deleting sensitive IDs from data may be sufficient to safeguard the privacy of the majority of people for those data mining applications that do use personal data. Recently, several methods for strengthening data security have been created. Additionally, a lot of work has been put into recent months to create data mining techniques that protect privacy. We examine some of the developments in data mining's privacy and data security in this part.Fair information practises are a collection of international regulations that were developed in 1980 by the Organization for Economic Co-operation and Development (OECD).These rules are intended to safeguard data accuracy and privacy. Aspects of data collection, utilisation, transparency, security, quality, and accountability are covered. They include the following guidelines:

**Specification of the purpose of collection and restriction of its use:** Personal data shall only be gathered for the reasons that are disclosed at the time of collection. Data collecting generally serves a secondary goal, data mining. According to one argument, a caution that the data may potentially be used for mining is not commonly seen as being a sufficient indication of purpose. The exploratory nature of data mining makes it hard to predict what patterns could be found, and as a result, there is no assurance about their potential applications.

**Openness:** There has to be a general policy of transparency about personal data developments, practises, and policies. People have a right to information about the kind of information gathered on them, who is in charge of upholding the rules as the data controller, and how the information is being used.Security Measures: Reasonable security measures should be used to safeguard personal data from risks including loss or unauthorised access, as well as against use, modification, and disclosure.

**Individual Participation:** Every person has the right to know if the controller of their data possesses any information on them, and if so, what exactly that information is. The person may also contest such information. The person has the right to request that the information be changed, finished, or removed if the challenge is successful. Inaccurate facts are often only discovered after they have some negative effects on a person, such being denied credit or having money withheld from them. Usually, the implicated organisation lacks the requisite contextual information to identify such errors.

**Data mining has a new application field that is gaining popularity: counterterrorism.**

To identify anomalous trends, terrorist activity (including bioterrorism), and dishonest conduct, data mining for counterterrorism may be employed. This application field is still developing since it has so many obstacles to overcome. These include creating algorithms for multimedia data mining (which includes mining audio, video, and image data in addition to text data), real-time mining (building models in real time to detect real-time threats such as that a building is scheduled to be bombed by 10 a.m. the following morning), and finding unclassified data to test such applications. Although this new method of data mining raises concerns about the privacy of specific individuals, it is crucial to remember that the goal of the research is to create a tool for the detection of abnormal patterns or activities, and that only authorised security agents are allowed to access specific data in order to find terrorist patterns or activities.

To assist safeguard data, several data security-enhancing methods have been created. Databases may utilise a multilayer security model to categorise and limit data according to different degrees of protection, with users only being allowed access to the level to which

they are authorised. However, it has been shown that data mining offers a comparable opportunity and that people running specialised queries at their permitted security level may still deduce more sensitive information. Another method for encoding specific pieces of data is encryption. This may involve biometric encryption (where a person's iris or fingerprint image is used to encode his or her personal information), blind signatures (which build on public key encryption), and anonymous databases (which allow the consolidation of various databases but restrict access to personal information to only those who need to know; personal information is encrypted and stored at different locations). Another prominent field of study that contributes to the privacy of personal data is intrusion detection.In response to privacy protection during mining, a new field of data mining study called privacy-preserving data mining is developing. Additionally known as privacy-sensitive data mining or privacy-enhanced data mining. It deals with getting reliable data mining findings without knowing the values of the underlying data. Secure multiparty computation and data obscuration are the two popular methods. Data values are encoded in safe multiparty computation utilising simulation and cryptographic methods to prevent parties from learning one another's data values. When mining big datasets, this method may not be feasible.In data obscuration, the real data are falsified by aggregation (such as utilising the neighborhood's average income rather than the individual inhabitants' incomes) or by the addition of random noise. A reconstruction method may be used to determine a close approximation of the original distribution of a set of distorted data values. Instead of utilising the real numbers, mining may be done using these approximations. Although there still needs to be a standard framework for defining, measuring, and assessing privacy, significant progress has been done.The sector is anticipated to prosper.

Data mining may be abused, just like any other technology. We must not lose sight of the many advantages that data mining research may provide, from understandings acquired from medical and scientific applications to higher consumer satisfaction through assisting businesses in better meeting the demands of their customers. We anticipate that computer scientists, policy experts, and counterterrorism specialists will continue to collaborate with social scientists, attorneys, businesses, and customers to develop ways to maintain data security and privacy. In this approach, we may keep taking advantage of data mining's time and money savings and knowledge finding advantages.

**Trends in Data Mining**

Numerous difficult research problems in data mining are caused by the variety of data, data mining tasks, and data mining methodologies. Important tasks for data mining researchers and developers of data mining systems and applications include the creation of efficient and effective data mining methods and systems, the development of interactive and integrated data mining environments, the design of data mining languages, and the use of data mining techniques to address complex application problems. In this section, certain data mining trends that indicate the resolution of these problems are discussed. Application investigation early uses of data mining generally aimed to provide corporations a competitive advantage. As e-commerce and e-marketing have become staples of the retail sector, research into data mining for companies keeps growing. Data mining is rapidly being utilized to investigate potential applications in fields outside of research, biology, telecommunications, and financial analysis. Data mining for counterterrorism (including and beyond intrusion detection) and mobile (wireless) data mining are two newer application fields. We may see a trend toward the creation of increasingly application-specific data mining systems since generic data mining systems may have limits in addressing application-specific issues.

**Interactive and scalable data mining techniques**: Data mining must be able to process enormous volumes of data effectively and, if at all feasible, interactively, in contrast to conventional data analysis techniques. Scalable algorithms are needed for both standalone and integrated data mining operations since the quantity of data being gathered is growing fast. Constraint-based mining is a crucial step in boosting user participation while enhancing the overall efficiency of the mining process. By enabling the design and usage of limitations to direct data mining tools in their search for intriguing patterns, this gives users more control over the process.

**Data mining integration with database, data warehouse, and web database systems:** The Web, database systems, and data warehouse systems are becoming commonplace information processing tools. It is crucial to make sure data mining functions as a crucial part of data analysis that can be easily included into such a processing environment for information. A data mining system should be closely connected to database and data warehouse systems, as was previously stated. It is important to combine transaction management, query processing, online analytical processing, and online analytical mining into a single, cohesive architecture. In order to support multidimensional data analysis and exploration, this will provide data accessibility, data mining portability, scalability, high performance, and an integrated information processing environment. Language for data mining must be standardised: The methodical creation of data mining solutions will be made easier, the interoperability of various data mining systems and functions will be improved, and the teaching and use of data mining systems in business and society will be encouraged. Microsoft's OLE DB for Data Mining (an introduction is provided in the appendix of this book), PMML, and CRISP-DM are recent initiatives in this area.

**Visual data mining:** This method of extracting information from vast volumes of data is effective. The promotion and use of data mining as a tool for data analysis will be made easier by the methodical research and development of visual data mining methods.

Innovative techniques for mining complicated forms of data: Mining complicated forms of data is a significant research horizon in data mining.

Despite advances in mining stream, time-series, sequence, graph, spatiotemporal, multimedia, and text data, there is still a substantial technological gap between what is required for these applications and what is now possible. More study is needed, particularly on how to combine data mining methods with already used data analysis methods for various kinds of data.

Although mining complicated forms of data or application exploration are appropriate categories for biological data mining, the combination of complexity, richness, scale, and relevance of biological data calls for particular consideration in data mining. Intriguing areas for biological data mining research include mining DNA and protein sequences, mining high-dimensional microarray data, biological pathway and network analysis, link analysis across heterogeneous biological data, and information integration of biological data through data mining. Software engineering and data mining: The work of ensuring software resilience and dependability is becoming more difficult as programs grow in size, complexity, and tendency to be the result of the fusion of several components created by various software teams. The examination of a defective software program's executions is basically a data mining operation; tracking the data produced during program executions may reveal significant trends and outliers that may eventually enable automated software bug identification. We anticipate that as data mining approaches for software debugging continue to advance, software engineering will become more robust.

**Web mining:** Additionally covered topics pertaining to Web mining. Web content mining, Weblog mining, and data mining services on the Internet will become one of the most significant and prospering subfields in data mining given the vast quantity of information accessible on the Web and the growing importance that the Web plays in modern life. The Internet, intranets, local area networks, high-speed wireless networks, and sensor networks are just a few examples of the distributed computing environments that are currently in use. Traditional data mining techniques, which were created to operate at a centralized location, do not perform well in these environments. It is anticipated that distributed data mining techniques will advance.

Time-sensitive or real-time data mining Dynamic data mining models must be constructed in real-time for many applications that use stream data, including e-commerce, Web mining, stock analysis, intrusion detection, mobile data mining, and data mining for counterterrorism. In this field, further advancement is required. Graph mining, link analysis, and social network analysis are helpful for capturing the sequential, topological, geometric, and other relational characteristics of many scientific data sets, including those for biological networks and chemical compounds, as well as social data sets, like those used to analyses covert criminal networks. The analysis of connections in Web structure mining may also benefit from such models. A problem for data mining is the creation of effective graph and connection models.

**Data mining that uses several relational and database systems**: The majority of data mining techniques look for patterns in only one relational table or one database. However, the majority of information and data in the actual world are dispersed over several databases and tables. Multiple tables (relations) from a relational database are searched for patterns using multi-relational data mining techniques. Multiple datasets are mined for patterns using multi-database mining. It is anticipated that more study will be done on effective and efficient data mining across several relations and datasets.

**Data mining and information security and privacy protection:** Our privacy and data security are in danger because there is a surplus of recorded personal information that is readily accessible in electronic form and on the Internet and because data mining techniques are becoming more potent. The vulnerability is further increased by the growing interest in data mining for counterterrorism. It is anticipated that privacy-preserving data mining techniques will continue to advance. A rigorous definition of privacy and a formalization to demonstrate privacy-preservation in data mining need the cooperation of engineers, social scientists, legal experts, and businesses.

**----------------------**

# Questions for practices

1. How is a data warehouse different from a database?

2. What is the procedure of cleansing data?

3. How about mining geographical data using statistical methods?

4. What is data mining architecture of data mining?

5. What are the data mining techniques?

6. What types of relationships are there in multimedia data that can be mined?

7. A multimedia database is what, exactly?

8. What information retrieval techniques are there?

9. What data is needed for social network analysis?

10. What are the applications and issues in data mining?

---------------------

# References book for further Reading

1.  "Data Mining: Concepts and Techniques" by Han

2.  "Data Warehousing" by ReemaThareja

3.  "Data Warehousing and Data Mining" by Singh M

4.  "Data Mining and Warehousing" by S Prabhu

5.  "Data Mining and Warehousing" by Khushboo and Sandeep

6.  "Data Warehousing: OLAP and Data Mining" by Nagabhushana S

7.  "Data Mining and Warehousing" by M SudheepElayidom

8.  "Python for Data Science for Dummies" by John Paul Mueller and Luca Massaron

9.  "The Encyclopedia of Data Warehousing and Mining" by John Wang

--------------------