

Chapter 7

Advanced Applications

Content

- Web-mining: Web Content Mining, Web Usage Mining
- Time-series data mining

- 3 LH

Web-mining:

Web Content Mining, Web Usage Mining

7.1 Web Mining

- Web mining is the application of data mining techniques to extract knowledge from Web data, i.e. Web Content, Web Structure and Web Usage data.
- **Web Mining** is the process of Data Mining techniques to automatically discover and extract information from Web documents and services.
- The main purpose of web mining is discovering useful information from the World-Wide Web and its usage patterns.

7.1 Web Mining

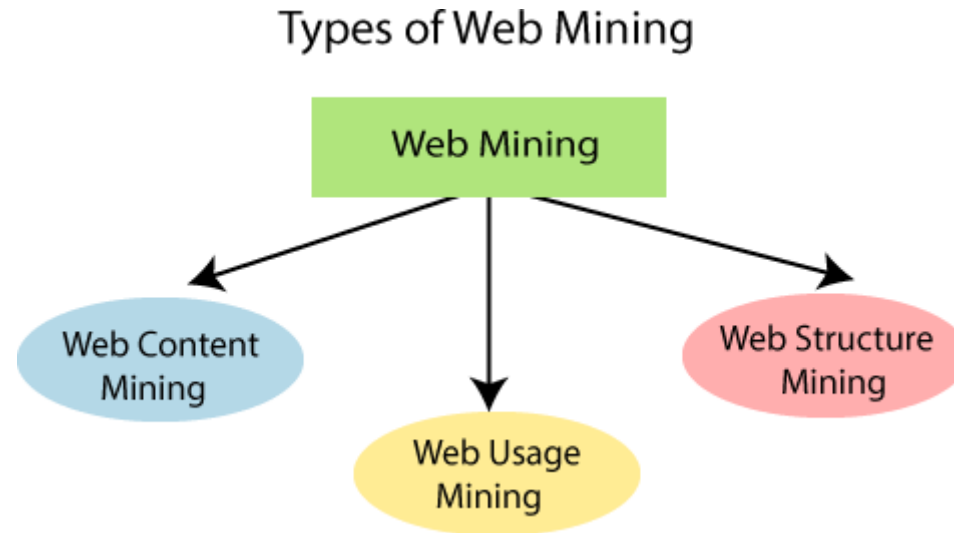
- Web mining is the art and science of discovering patterns and insights from the WWW so as to improve it.
- The WWW is at the heart of the digital revolution. More data is posted on the Web everyday. Billion of users are using it everyday for a variety of purposes, like as ecommerce, business communication and many other application.
- Web mining analyzes data from the web and helps find insights that could optimize the web content and improve user experience.
- Data for web mining is collected via web crawlers, web logs and other means.

Applications of web mining

- Web mining helps to improve the power of web search engines by classifying web documents and identifying web pages.
- It is used for Web Searching e.g., Google, Yahoo, etc,
- Web mining is used to predict user behavior.

Types of web mining

- Web mining can be broadly divided into three different types of techniques of mining: Web Content Mining, Web Structure Mining, and Web Usage Mining.



i) Web Content Mining

- Web content mining is a technique used to extract useful information and knowledge from web pages and other online content.
- The goal of web content mining is to automatically analyze and extract information from the vast amount of unstructured data available on the web, such as text, images, audio, and video.
- Web content mining involves several steps, including data collection, data preprocessing, feature extraction, and data analysis.
 - The first step is to collect data from the web, either by web scraping or web crawling. Web scraping involves extracting data from specific web pages, while web crawling involves automatically following links to gather data from multiple pages.
 - Once the data is collected, it is preprocessed to clean and transform the data into a suitable format for analysis. This may involve removing noise, correcting spelling errors, and converting the data into a structured format.
 - The next step is to extract features from the data, such as keywords, topics, or sentiment. Feature extraction may involve natural language processing techniques, such as text classification, entity extraction, or sentiment analysis.
 - Finally, the extracted features are analyzed using statistical and machine learning techniques to uncover patterns and relationships in the data. Web content mining can be used for a variety of applications, such as web search, recommendation systems, and online advertising.

ii) Web Structure mining

- Web structure mining is a technique used to analyze the structure of the web and extract useful information from it.
- The goal of web structure mining is to automatically discover patterns and relationships in the way web pages are linked together.
- Web structure mining involves analyzing the links between web pages, including hyperlinks, anchor texts, and other structural features. The analysis can be done at different levels, including individual web pages, groups of pages, and entire websites.
- Web structure mining can be used for a variety of applications, such as search engine optimization, web navigation, and web analytics.
- For example, web structure mining techniques can be used to identify the most important pages on a website, optimize the structure of a website for search engines, or identify patterns in user navigation behavior.

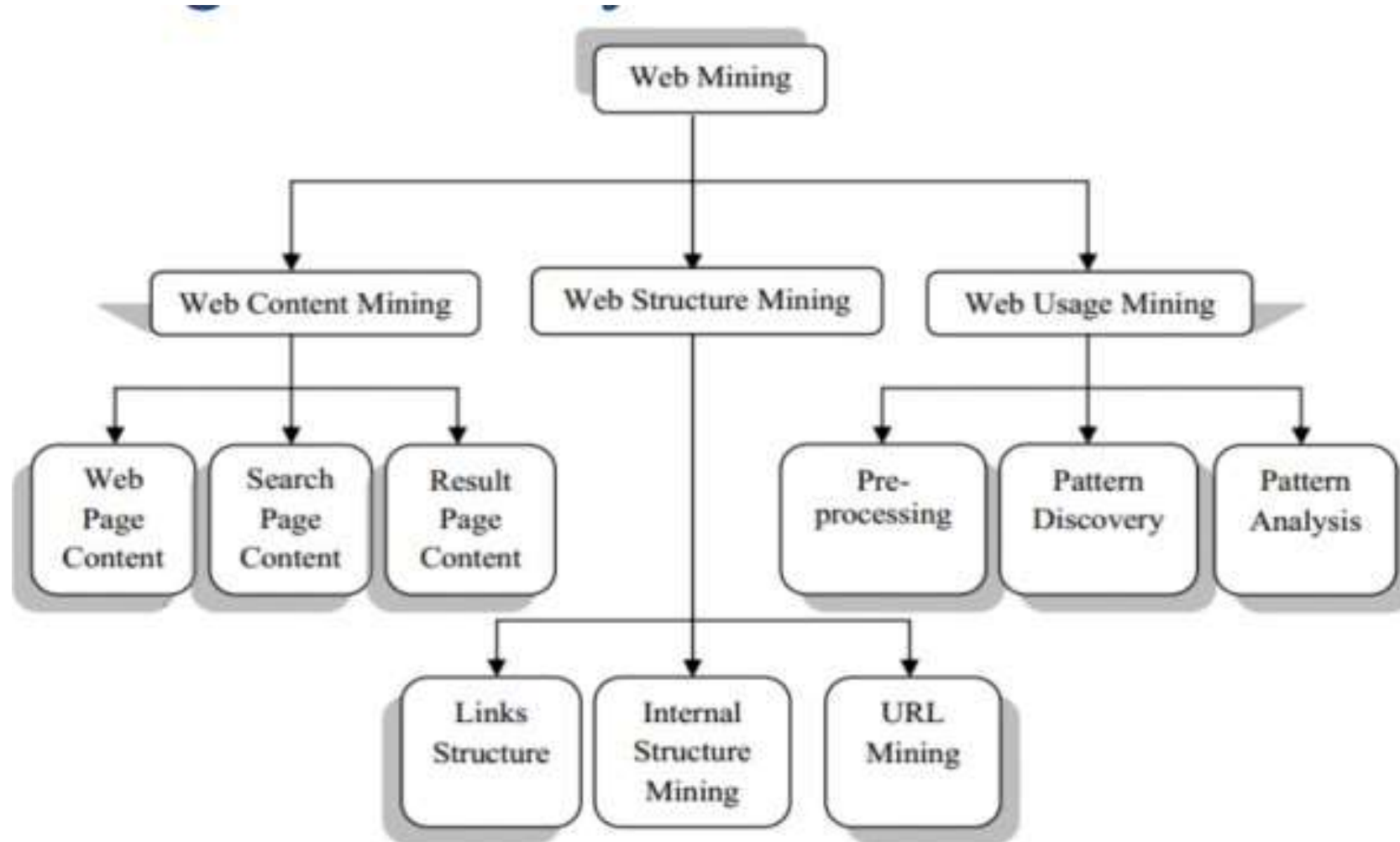
ii) Web Structure mining

- There are two main approaches to web structure mining: graph-based and model-based.
 - Graph-based approaches involve representing the web as a graph, where web pages are nodes and links between pages are edges. The graph can then be analyzed using graph theory techniques to identify patterns and relationships in the web structure.
 - Model-based approaches involve creating models of the web structure and using them to make predictions or classify web pages. For example, a model-based approach might involve using machine learning techniques to classify web pages based on their link structure or other structural features.

iii) Web Usage Mining:

- Web usage mining is a technique used to analyze and extract useful information from user interactions with web-based applications, such as websites or mobile apps.
- The goal of web usage mining is to understand user behavior, preferences, and interests based on their interactions with these applications.
- Web usage mining involves analyzing web server log files or user tracking data to extract patterns and relationships in user behavior. This may include analyzing the pages visited by users, the time spent on each page, the search queries used, and other user interactions, such as clicks or form submissions.
- The analysis of user behavior may involve various techniques, such as clustering, classification, and association rule mining.
- For example, clustering techniques may be used to group users with similar behavior, while classification techniques may be used to predict user preferences or interests based on their behavior.
- Web usage mining can be used for a variety of applications, such as improving website design, personalizing content for users, and predicting user behavior.
- For example, web usage mining techniques can be used to identify the most popular pages on a website, optimize the layout of a website for better user engagement, or recommend relevant content to users based on their past behavior.

Web mining Taxonomy:



Challenges in Web mining:

Web mining is the process of extracting useful information and knowledge from web data sources, which can present several challenges. Some of the major challenges in web mining include:

1. **Heterogeneous data:** Web data can be in different formats, such as text, images, audio, and video, which makes it difficult to process and analyze. Techniques such as feature extraction and fusion may be required to handle the heterogeneity of web data.
2. **Data volume:** Web data sources can be very large and complex, which makes it difficult to process the data in a timely manner. Scalable algorithms and techniques are required to handle large amounts of web data.
3. **Dynamic nature of the web:** Web data is constantly changing and evolving, which makes it difficult to maintain up-to-date data sets. Techniques such as web crawling and indexing may be required to keep up with the dynamic nature of the web.
4. **Data quality:** Web data can be noisy, incomplete, or inconsistent, which can make it difficult to extract useful information. Data cleaning and preprocessing techniques may be required to ensure the quality of the data.
5. **Privacy concerns:** Web mining can involve collecting personal information about users, which raises privacy concerns. Appropriate measures must be taken to ensure that user privacy is protected.
6. **Interpretability:** Web mining results can be difficult to interpret, especially if the algorithms used are complex or black-box. Techniques such as visualization and explanation can be used to improve the interpretability of the results.
7. **Bias:** Web mining can be biased if the data is not representative of the target population. Sampling techniques may be required to ensure that the data is representative.

Issues in Web mining:

1. **Data quality:** Web data can be noisy, incomplete, or inconsistent, which can make it difficult to extract useful information. Data cleaning and preprocessing techniques may be required to ensure the quality of the data.
2. **Privacy concerns:** Web mining can involve collecting personal information about users, which raises privacy concerns. Appropriate measures must be taken to ensure that user privacy is protected.
3. **Bias:** Web mining can be biased if the data is not representative of the target population. Sampling techniques may be required to ensure that the data is representative.
4. **Scalability:** Web data sources can be very large and complex, which can make it difficult to process the data in a timely manner. Scalable algorithms and techniques are required to handle large amounts of web data.
5. **Interpretability:** Web mining results can be difficult to interpret, especially if the algorithms used are complex or black-box. Techniques such as visualization and explanation can be used to improve the interpretability of the results.
6. **Ethical issues:** Web mining can raise ethical issues, such as the use of user data for marketing or other purposes without their consent. Appropriate ethical guidelines and regulations must be followed to ensure that web mining is conducted ethically.

Data Mining vs Web mining

Data Mining:

1. The process of discovering hidden patterns, trends, and relationships in large data sets.
2. Can be applied to a wide range of data sources, including databases, transactional data, and sensor data.
3. Focuses on structured data that can be represented in tables or databases.
4. Uses statistical and machine learning techniques to analyze the data.
5. Can be used for a variety of applications, such as fraud detection, customer segmentation, and predictive modeling.

Web Mining:

1. A specific type of data mining that focuses on discovering patterns and trends from web data sources.
2. Focuses on unstructured and semi-structured data from web pages, social media, and online user behavior.
3. Requires specialized techniques to deal with the unique characteristics of web data, such as its unstructured nature, heterogeneity, and dynamic nature.
4. Involves the use of web scraping and web crawling techniques to gather data from the web.
5. Can be used for applications such as web personalization, recommendation systems, and online advertising.

Time-series data mining

7.2 Time series data

- Time series data are a sequence of data recorded over a period of time.
- They are generated by an economic process like Stock Market analysis, Medical Observations.
- They are useful for studying natural phenomena.
- Time Series Analysis comprises methods for analyzing time-series data in order to extract meaningful *statistics*, *rules* and *patterns*.
- These rules and patterns might be used to build forecasting models that are able to predict future developments.



Time series data mining

- Time series analysis is a **specific way of analyzing a sequence of data points collected over an interval of time.**
- In time series analysis, analysts record data points at consistent intervals over a set period of time rather than just recording the data points intermittently or randomly
- In general, there are two goals in time series analysis.
 - *Modeling Time Series:* Generating the time series with underlying mechanism.
 - *Forecasting Time Series:* Predict the future values of the time series variables.

Time series data mining

- Time series data mining is the process of extracting patterns, trends, and relationships from time series data.
- It involves visualizing the data, decomposing the different components, and using statistical and machine learning techniques to analyze the data and make predictions about future trends.
- An example of time series data mining could be analyzing hourly electricity consumption data to identify daily and seasonal patterns, decomposing the data to separate out the trend and seasonality, and using statistical or machine learning techniques to make predictions about future consumption values.

Application of Time Series Mining:

1. Financial:

- Used for stock price evaluation
- For the measurement of Inflation

2. Industry:

- Determine the power consumption

3. Scientific:

- Used for experiment results

4. Meteorological:

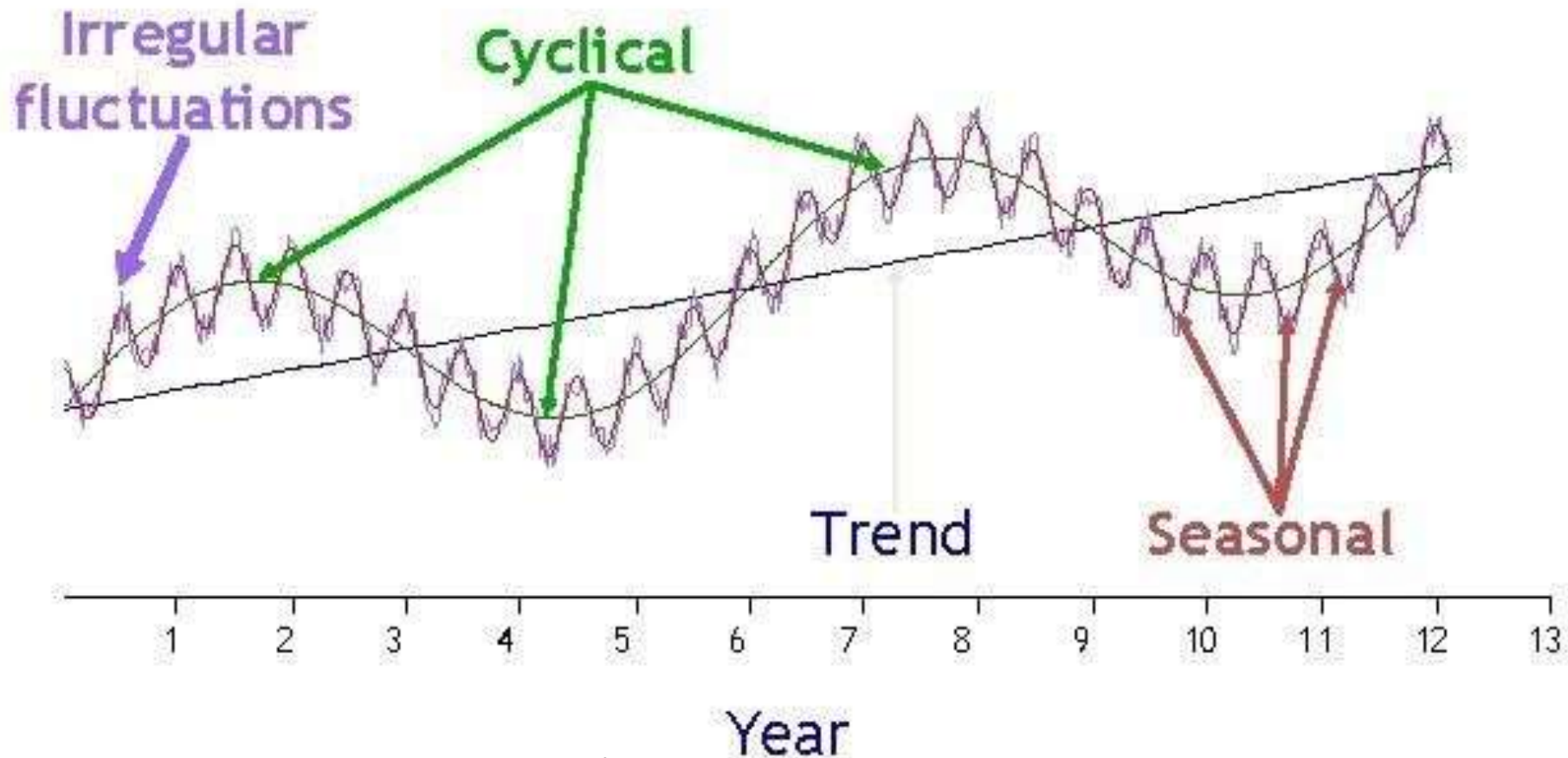
- Concerned with the processes and phenomena of the atmosphere, basically for forecasting weather

Importance of Time Series Analysis in Business

- Business owners use time series analysis to see seasonal trends and understand the underlying reasons for their occurrence.
- Businesses can use time series forecasting to predict the probability of upcoming events.
- Time series forecasting can highlight possible modifications in the data such as cyclic or seasonal behavior, which helps in better forecasting by offering a clear idea of data variables.

Category of Time-series Movements (Trend analysis)

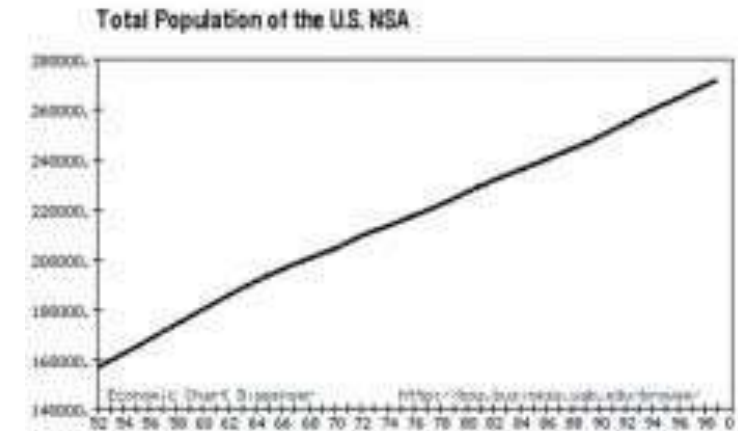
Components of Time-Series Data



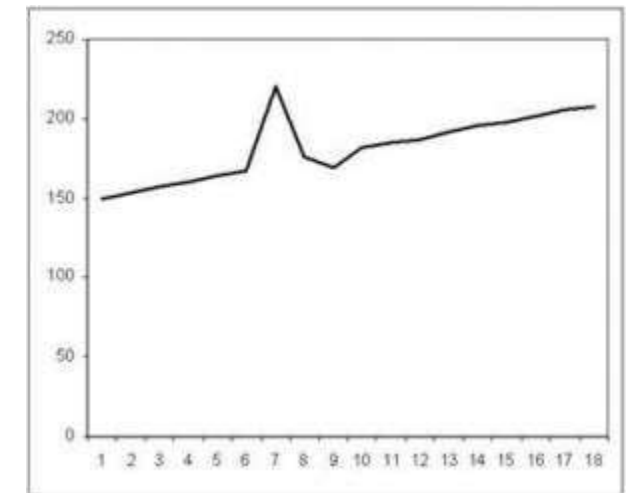
Category of Time-series Movements (Trend analysis)

There are three main categories of time series movement:

1. **Trend:** The trend component represents the long-term movement of the time series data in a particular direction. Trends can be upward (increasing over time), downward (decreasing over time), or stationary (remaining constant over time).
2. **Seasonality:** The seasonality component represents the pattern of the time series that repeats itself at regular intervals, such as daily, weekly, or monthly. Seasonality is typically associated with calendar events, holidays, or weather patterns.



3. **Cyclical:** The cyclical component represents the fluctuations in the time series that do not have a fixed periodicity. These fluctuations are often influenced by economic, political, or other external factors that affect the overall movement of the time series. For example, business cycles
4. **Irregular or random movements:** These fluctuations are unforeseen, uncontrollable and unpredictable. They are not regular variations and are purely random or irregular.



Approaches for Time series Data Analysis

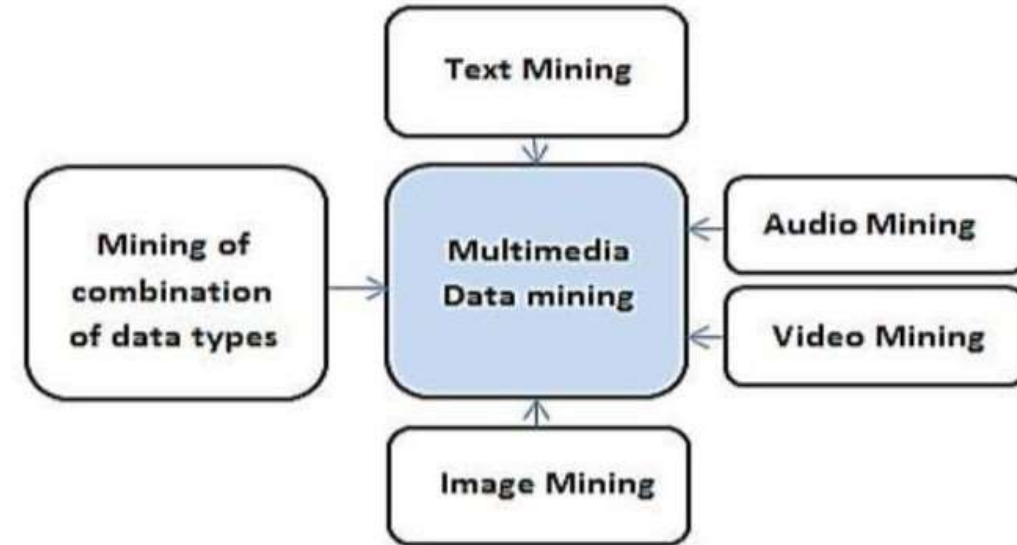
1. **Descriptive analysis:** This approach involves examining the properties of the time series data, such as trend, seasonality, and cyclical patterns. Descriptive analysis can be done using techniques such as time plots, autocorrelation plots, and periodograms.
2. **Time series decomposition:** Time series decomposition is the process of separating the different components of the time series, such as trend, seasonality, and cyclical patterns. This can be done using techniques such as moving averages, exponential smoothing, and Fourier analysis.
3. **Regression analysis:** Regression analysis involves examining the relationship between the time series data and one or more predictor variables. This can help identify factors that influence the movement of the time series.

4. **Machine learning:** Machine learning techniques, such as neural networks, support vector regression, and decision trees, can be used to analyze time series data. These techniques can help identify complex relationships between variables and make predictions about future trends in the data.
5. **Signal processing:** Signal processing techniques, such as wavelet analysis, can be used to analyze the frequency content of time series data. This can help identify patterns and anomalies that may not be visible in the time domain.
6. **Bayesian analysis:** Bayesian analysis involves using probability theory to model the uncertainty in time series data. This approach can help make predictions about future trends and estimate the likelihood of different outcomes.

Object/Image/Multimedia mining

7.3 Object/Image/Multimedia mining

- Multimedia data mining is the discovery of interesting patterns from multimedia databases.
- Multimedia database system stores and manages a large collection of multimedia data such as audio, video, images, graphics, speech, text etc.



7.3 Multimedia mining Process

The multimedia data mining process involves several steps, which are as follows:

1. **Data collection:** The first step in multimedia data mining is to collect the multimedia data from various sources. The data can be in different formats such as images, audio, video, and text.
2. **Preprocessing:** Preprocessing involves cleaning and transforming the raw data into a suitable format for analysis. This includes removing noise, normalizing, and reducing the dimensionality of the data.
3. **Feature extraction:** Feature extraction is the process of identifying and selecting relevant features from the preprocessed data. These features could be color, texture, shape, audio features, or any other relevant attributes.
4. **Data mining:** In this step, data mining techniques such as clustering, classification, association rule mining, and anomaly detection are applied to the extracted features. The objective of data mining is to identify patterns, trends, and relationships within the data.
5. **Evaluation:** In this step, the results of the data mining process are evaluated to determine their usefulness and effectiveness. The evaluation is done using metrics such as accuracy, precision, and recall.
6. **Interpretation and Visualization:** The final step involves interpreting and visualizing the results of the data mining process. The insights gained from this step can be used to make informed decisions or to further refine the mining process.

7.3 Complexities in Multimedia mining

1. **Large volumes of data:** Multimedia data can be massive, making it challenging to process and analyze efficiently. The sheer size of multimedia data can cause memory and processing issues.
2. **Heterogeneity:** Multimedia data comes in different formats and types, making it difficult to analyze using traditional data mining techniques. There are also variations in quality, resolution, and encoding formats, which adds to the complexity.
3. **Semantic Gap:** The semantic gap refers to the difference between the low-level features extracted from multimedia data and the high-level concepts that humans associate with the data. Bridging this gap is a significant challenge in multimedia data mining.
4. **Subjectivity:** Multimedia data can have a subjective nature and can be interpreted differently by different people. This subjectivity makes it challenging to develop algorithms that can accurately analyze multimedia data.
5. **Privacy and Ethical Concerns:** With multimedia data, there are issues of privacy, data ownership, and ethical concerns that need to be taken into account when developing algorithms for multimedia data mining.
6. **Contextual Information:** In multimedia data mining, it is essential to consider the context in which the data was created or used, such as the user's intent, background, and environment. The context can affect the interpretation of the data and the mining results.

Application of Multimedia mining

- Digital Library
- Traffic Video Sequences
- Medical Analysis
- Customer Perception
- Media Making and Broadcasting
- Surveillance system



Reference: <https://www.javatpoint.com/what-is-multimedia-data-mining>

Exercise:

1. What is web mining? Mention its types.
2. Mention the applications of web-mining.
3. Explain the different types of web-mining.
4. What are the challenges in web mining?
5. Differentiate between data mining and web mining.
6. Define Time-series data. Give few examples of such.
7. Define Time-series data mining.
8. Mention the aim of Time series data mining.
9. Mention the application of Time-series data mining.
10. Discuss on the importance of Times series analysis in business.
11. Explain on the categories of Time-series movements.
12. How multimedia mining is performed?
13. What are the complexities in multimedia data mining?

Exam questions:

1. Give any two application of web mining. [BIM 2022, Group A]
2. Give any two applications of time series data mining. [BIM 2021, Group A]
3. What is time series data mining? [BIM 2018, Group A]
4. What is web structure mining? [BIM 2018, Group A]
5. What are the complexities in multimedia data mining? [BIM 2021, Group A]
6. Describe about web usage mining and web content mining. [BIM 2021, Group C]

End of chapter