## Tipo actividad:  Reading comprehension: "Text mining: Definition, techniques, use cases" and related activities

**Socialization key vocabulary: "Text mining: Definition, techniques, use cases"**

1. **Text Mining:** Text mining, or text analysis, is the process of using Machine Learning for text analysis. It involves transforming unstructured text into structured data to identify meaningful patterns and gain insights.

2. **Natural Language Processing (NLP):** Natural Language Processing is a technology that enables machines to understand and process human language automatically. It forms the basis of text mining, allowing for the analysis of texts by sentiment, subject, or intent.

3. **Information Retrieval (IR):** Information retrieval is the process of finding relevant information from a predefined set of queries or phrases. It is commonly used in library catalog systems and web search engines, utilizing algorithms to track user behavior and identify relevant data.

4. **Text Classification:** Text classification involves assigning labels to unstructured text data. It is crucial for organizing and structuring complex text, enabling the extraction of relevant data. It includes sub-tasks like Topic Analysis, Sentiment Analysis, Language Detection, and Intention Detection.

5. **Text Analytics:** Text analytics is the process of discovering patterns across large datasets through the analysis of text. It provides quantitative insights and is used to create tables, charts, graphs, or visual reports. Text analytics complements text mining by visualizing results from text mining analyses.

**7) Reading comprehension activity #2: "Text mining: Definition, techniques, use cases"**

**Text mining: Definition, techniques, use cases**

Text mining consists in using Machine Learning for text analysis. Discover all you need to know: definition, functioning, techniques, advantages, use cases... Modern companies have a lot of data on their customers or their business sector. New digital technologies such as social networks, e-commerce, or mobile applications for smartphones give access to a vast amount of information.

By analyzing this data, it is possible to discover untapped opportunities or alarming problems that need to be addressed urgently. However, some types of data are more difficult to exploit than others.

Data from social networks or other websites are mainly texts: comments on posts, product reviews, and complaints on community forums…

However, texts are part of the so-called "unstructured" data. This information cannot be properly processed by traditional data analysis software and tools. It is, therefore, necessary to rely on "Text Mining".

Text mining, or text analysis, consists of transforming unstructured text into structured data and then proceeding with the analysis. This practice is based on the technology of "Natural Language Processing", which allows machines to understand and process human language automatically.

Artificial intelligence is now able to automatically classify texts by sentiment, subject, or intent. For example, a text mining algorithm can review product reviews to determine whether they are mostly positive, neutral, or negative. It is also possible to identify the most frequently used keywords.

In this way, companies can analyze large and complex data sets in a simple, fast, and efficient way. This discipline also reduces time wasted on manual and repetitive tasks.

Teams save time and can focus on more important tasks that require human intervention. And business leaders can leverage data to make better decisions.

**How does Text Mining work?**

Text mining is based on Machine Learning, a subcategory of artificial intelligence, which encompasses many techniques and tools that enable computers to learn to perform tasks autonomously.

Machine Learning models are trained on data to be able to make accurate predictions. Text mining is the automation of text analysis using Machine Learning. To achieve this, the algorithms are trained using text as example data.

The first step is to assemble data. This data can come from internal sources, such as chat interactions, emails, surveys, or company databases. It can also come from external sources such as social networks, review sites, or news articles.

The data must then be prepared using various Natural Language Processing techniques. This "data pre-processing" aims to clean and transform the data into a usable format.

This is an essential aspect of Natural Language Processing, involving the use of different techniques such as language identification, tokenization, part-of-speech labeling, chunking, and syntax analysis. The objective of these different methods is to format the data for analysis.

After completing this "pre-processing" of the text, it is time for data analysis. Various text-mining algorithms are used to extract information from the data.

### Text mining methods and techniques

There is a wide variety of text mining techniques and methods. Here are the most commonly used.

### Analysis techniques

The "word frequency" technique consists of identifying the most recurrent terms or concepts in a data set. This can be very useful, especially when analyzing customer reviews or conversations on social networks.

For example, if terms such as "too expensive" or "overpriced" recur frequently, the analysis may suggest that the product is too expensive. It is, therefore, necessary to adjust the price if possible.

The collocation method, on the other hand, consists of identifying sequences of words that frequently appear close to each other. Some words appear together very often. These may be bigrams or trigrams, combinations of two or three words. By identifying these collocations, it is possible to better understand the semantic structure of a text and to obtain more reliable Text Mining results.

### Information retrieval

Information retrieval is the process of finding relevant information from a predefined set of queries or phrases. This approach is often used in library catalog systems or web search engines.

IR (information retrieval) systems use different algorithms to track user behavior and identify relevant data. Tokenization" consists of breaking down a long text into sentences or words called "tokens". These tokens are then used in models for text clustering or document association tasks.

Stemming, on the other hand, consists of separating the prefixes and suffixes of words to derive the root word and its meaning. This technique reduces the size of index files.

### Text classification

There are also more advanced methods of text mining. Text classification consists of assigning labels to unstructured text data. This is an essential and indispensable step for Natural Language Processing.

It allows the organization and structure of a complex text to extract relevant data. It is thanks to this technique that companies can analyze all kinds of textual information to extract valuable information.

There are different forms of text classification. Topic Analysis is used to understand the main themes or topics of a text. This is one of the main ways to organize text data.

Sentiment Analysis is the analysis of the emotions in a text. This allows for a better understanding of customer opinions, for example, by reviewing comments about a product. Text can be classified as positive, negative, or neutral.

Language detection consists of classifying a text according to its language. For example, it will be possible to sort customer service requests and redirect them to an advisor or agent who masters the appropriate language. This saves precious time.

Finally, intention detection allows for the automatic recognition of the intentions of a text. For example, the analysis of different responses to an advertising email can determine which interlocutors are interested in a product.

### Information extraction

Another text mining technique is text extraction. It aims at extracting specific data from a text, such as keywords, proper names, addresses, or emails. This avoids having to sort the data manually and therefore saves time.

One can select the features that contribute most to the results of a predictive analysis model, extract features to improve the accuracy of a classification task or detect and categorize specific entities in a text.

It is of course possible to combine text mining and text classification, or other text mining methods in the same analysis.

### Text Mining vs. Text Analytics: what is the difference?

Text mining is often confused with text analytics. In reality, they are two slightly different concepts.

Both aim at automatically analyzing texts but are based on different techniques. Text mining identifies relevant information in text, while text analytics aims to discover patterns across large datasets.

One provides qualitative analysis and the other quantitative analysis. In general, Text Analytics is used to create tables, charts, graphs, or other visual reports.

Text mining combines statistics, linguistics, and machine learning to automatically predict outcomes from past experiences. Text Analytics, on the other hand, is about creating data visualizations from the results of Text Mining analyses. It is of course possible to combine these two approaches.

### The Benefits of Text Mining

Text mining has many advantages, at a time when companies and individuals generate huge volumes of data every day. Indeed, nearly 80% of text data is unstructured. It is therefore impossible to analyze it without using text mining.

For example, emails, social media posts, messenger discussions, customer service requests, surveys… It is very difficult to sort out this information manually.

Text analytics allows you to analyze large volumes of data in just a few seconds, thus increasing productivity. These analyses can be performed in real-time, and it is, therefore, possible to intervene immediately if a problem is detected.

### How can Text Mining be used?

Text mining can be used in many ways by companies. The applications of this technology are limitless and extend to all industries.

It can be used to automate text analysis for marketing, product development, sales, and customer service. Teams can become more efficient and productive by focusing on more important tasks.

### Customer service

In the field of customer service, it is for example possible to automatically sort requests. Text mining automatically identifies the topics, intent, complexity, and language of the requests to organize them. This allows agents to focus on helping customers.

If a request is more important or urgent than another, it can be automatically prioritized and processed before others. In addition, text analytics can also be used to measure customer service efficiency and user satisfaction.

Text mining is also very useful for analyzing customer feedback and opinions about the brand and its products. This allows you to understand their opinions, but also their expectations and the quality of their experience with your company.

Product reviews, comments on social networks, and survey responses can be scrutinized. In this way, it is possible to use the data to make the right decisions and improve weak points.

### Risk management

Text mining is used in the field of risk management. It can be used to extract information about industry trends or financial markets by monitoring changes in sentiment or extracting information from analytical reports and white papers.

This can be very useful within banking institutions. This is because the data allows them to approach investments in different sectors with more confidence. Many banks are now taking this approach.

### Maintenance

Text mining offers a complete overview of the activity and operation of industrial equipment and machinery. It allows for the automation of maintenance decisions.

For example, it is possible to highlight patterns and trends suggesting the occurrence of a problem. In this way, it is possible to implement predictive maintenance measures to intervene before it is too late. Maintenance operations can then be carried out proactively.

### Healthcare

In the field of health, Text Mining techniques are increasingly used by researchers. For example, information clustering allows us to extract information from medical books in an automated way.

This saves time and money. Thus, this approach is proving to be of great help to the world of medicine and health.

### Cybersecurity

Text analysis can also be particularly useful for cybersecurity. For example, it is possible to detect and filter spam automatically in email boxes.

This way, hackers can no longer use the spam method to hack into computer systems. The risk of cyber attacks is drastically reduced, and the user experience is also improved.

**8) Multiple choice activity.**

**1. What is the primary goal of Text Mining?**

- A. Creating data visualizations
- B. Transforming unstructured text into structured data
- C. Analyzing numerical datasets
- D. Developing machine learning models

**2. Which technology forms the basis of Text Mining for language understanding?**

- A. Machine Learning
- B. Data Analysis
- C. Natural Language Processing (NLP)
- D. Artificial Intelligence

**3. What is the main purpose of Information Retrieval (IR) in Text Mining?**

- A. Classifying texts
- B. Extracting information from documents
- C. Finding relevant information from predefined queries
- D. Analyzing sentiment in text

**4. What does Text Classification involve?**

- A. Identifying word frequency
- B. Assigning labels to unstructured text data
- C. Extracting specific data from a text
- D. Creating data visualizations

**5. How does Text Analytics differ from Text Mining?**

- A. Text Analytics focuses on creating data visualizations, while Text Mining transforms unstructured text.
- B. Text Mining is quantitative, while Text Analytics is qualitative.
- C. Both are synonymous and have no differences.
- D. Text Analytics is based on Natural Language Processing, while Text Mining uses machine learning.