

Challenges in BA

- * Data quality & integration issues
- * Lack of skilled personnel & analytical expertise
- * Privacy & security concerns
- * Resistance to adopting data-driven decision making

Future Trends in BA

- * AI & ML integration
 - * Real-time Analytics
 - * Cloud-Based Analytics
 - * Natural Language Processing (NLP)
 - * Ethical AI & Analytics
- Brief Introduction To Big Data Analytics
- * Big data is data that cannot be stored in a single storage unit.
 - * Big data refers to data that is arriving in many different forms, be they structured, unstructured or in a stream.
 - * Two aspects to big managing data on this scale: storing & processing.
 - * 3V's (Volume, Velocity, Variety)

Machine versus Men On Jeopardy!

- * Watson Overview
 - * Watson is an advanced Computer System
 - * designed to answer natural human language questions
 - * Developed in 2010 by IBM Research for the Deep Blue project, named after IBM's first president, Thomas J. Watson
- Background
- * IBM Research pursued a challenge to rival Deep Blue & advanced computer science, benefiting science, business & society.
 - * The challenge: build a real-time Jeopardy! contestants capable of listening, understanding & responding
 - * Competing against the Best

- * In 2011, Watson completed its first human versus-machine match.
- * Watson won a two-game match, defeating top players Brad Rutter & Ken Jennings.
 - * Watson excelled in signaling but had trouble with short few word clues.

- * It processed 800 million pages of data using four terabytes of storage without Internet access.
- * Watson employed multiple QA technologies like text mining, NLP & question-answering techniques.

* It computed confidence levels to decide whether to buzz in within 1.6 seconds, averaging 3 seconds.

How does Watson do it?

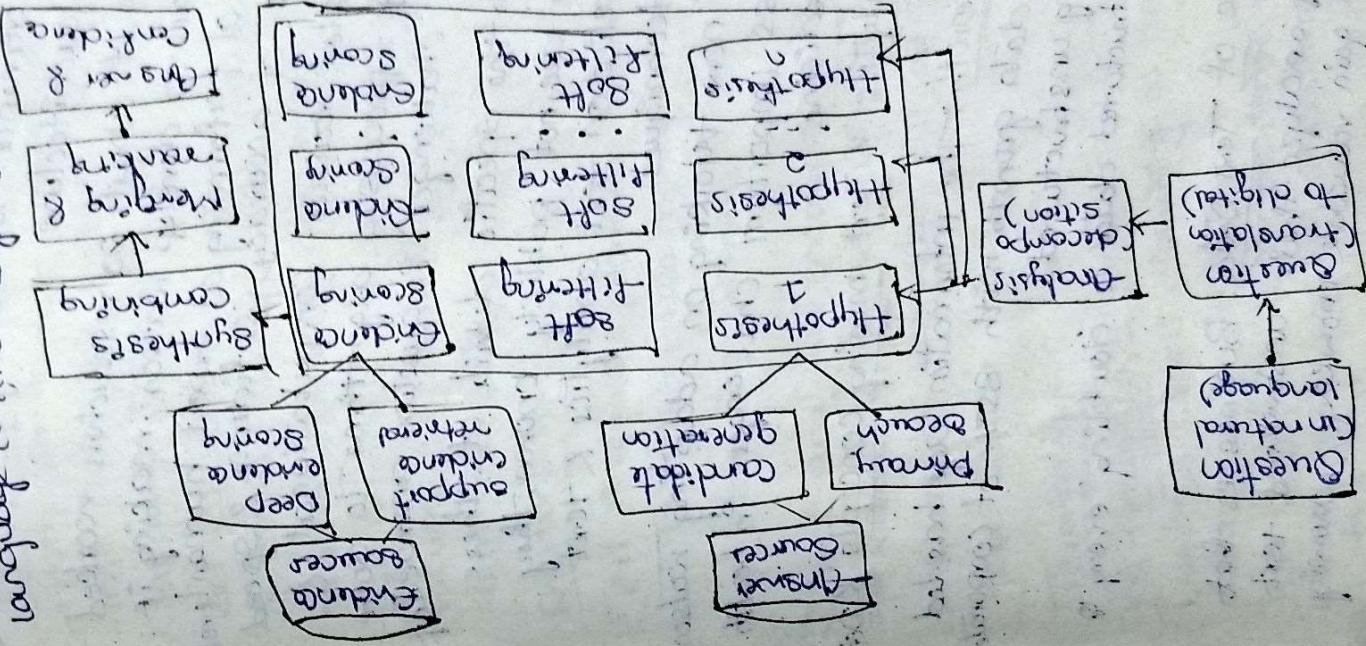
* Watson's system called DeepQA is a powerful, parallel-processing architecture that uses over 100 techniques to analyze language, generate answers & evaluate them.

* The key to its success lies in how these techniques work together to improve accuracy, speed & confidence.

Key Principles Of DeepQA

- * Massive parallelism: It examines multiple interpretations & answers at once.
- * Many experts: uses a range of probabilistic methods for analyzing questions & content.
- * Pervasive, confidence estimation: Each part of the system scores potential answers with confidence levels, which are combined to find the best answer.

* Shallow & deep knowledge: it balances detailed & broad understanding of language & knowledge structures.



Conclusion

- * the Jeopardy challenge helped IBM define the requirements for creating Watson & the DeepQA architecture
- * A team of about 20 researchers worked for 9 years to develop Watson, making it capable of performing at expert human level, in terms of accuracy, confidence & speed
- * IBM created various language processing algorithms to solve different challenges, in question answering
- * Although the specific details of these algorithms are not publicly known, they heavily relied on text analytics & text mining techniques
- * IBM is now talking on adapting Watson to address important challenges in health care & medicine
- * Text analytics & Text mining concepts & Text mining growth
- * Data growth: the information age has led to rapid data growth with 85% of corporate data being unstructured
- * this unstructured data is doubling every 18 months

Importance of Text Data: Businesses that effectively analyze their unstructured text data will gain valuable knowledge, allowing

from to make better decisions & achieve a competitive edge.

Text Analytics & Text Mining
Both aim to turn unstructured text data into actionable insights using natural language processing (NLP) & analytics

Text Analytics broader concept that includes information retrieval, extraction, data mining & web mining

Text mining focuses on discovering new, useful knowledge from text data

Relationship

Text Analytics = Information Retrieval + Information Extraction + Data mining + Web mining = Information Retrieval + Text mining

* Text analytics encompasses multiple processes, while text mining specifically focuses on knowledge discovery from text sources

Terminology

Text Analytics: A newer term, often used in business context, focusing on analysing text data to gain insights

Text Mining: used in academic research, involving the extraction of patterns & knowledge from text

Text Mining Definition

Text semi-automated process to extract useful patterns & knowledge from large amounts of structured data

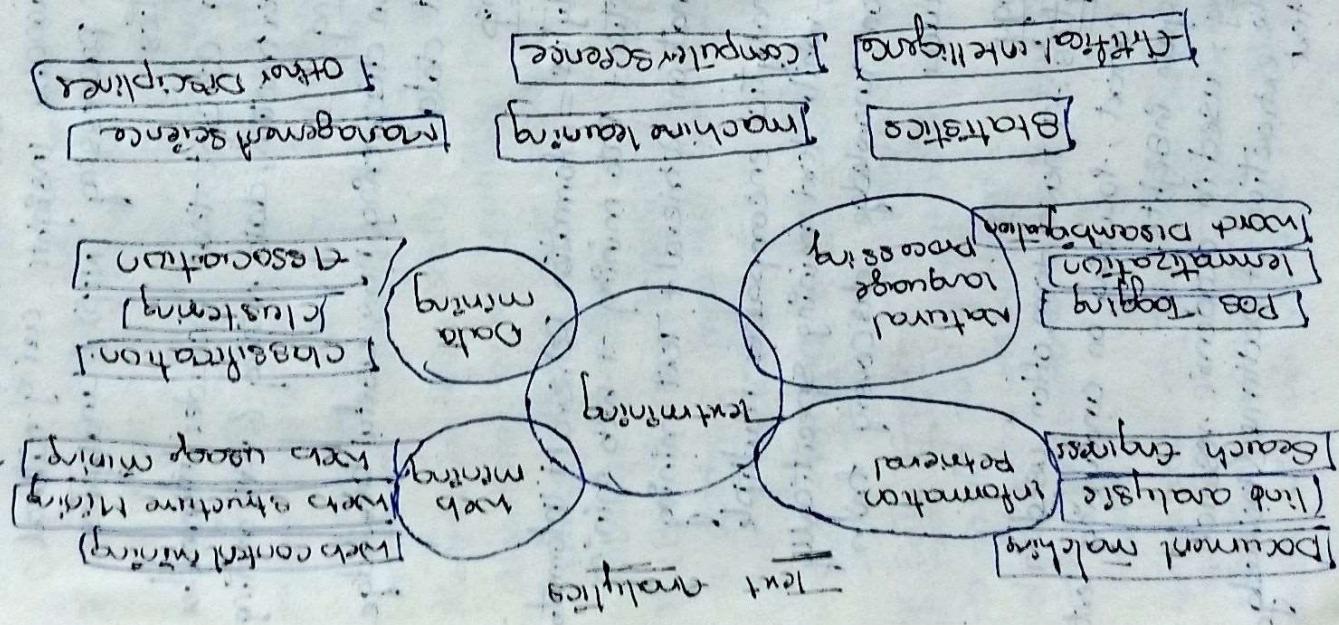
comparison with Data mining identifies patterns in structured data DM uses similar techniques but focuses on unstructured data.

Process of Text mining :
impose structure on unstructured text data
use data mining techniques to extract & analyze information from this structured data

1. Information extraction
2. Topic tracking
3. Summarization
4. Categorization
5. Clustering
6. Concept linking
7. Question answering

① Information extraction : It identifies key phrases & relationships in text using pattern matching. The most common method is named entity extraction, which includes:

Purpose: Both teams aim to convert unstructured text data into actionable insights using methods like NLP & analytics.



- * Named Entity Recognition: Detects entities, defined set of categories based on those themes
e.g. A blog post discussing healthy eating habits is categorized as "Health" & "Nutrition" based on its main theme
- * co-reference resolution: links entities that refer to the same thing
- * relation extraction: identifies relationships between entities
 - e.g. "John works at Google. In California, it enacts "John" (person), "Google" (organization), "California" (place), plus the relationship "works at"
- ⑤ Topic tracking
 - * Based on a user profile & documents that a user views, text mining can predict other documents of interest to the user
 - * It is used in search engines to recommend content based on user searches
 - e.g. If a user frequently searches for deep learning & AI, the engine will provide more related suggestions to enhance their experience
- ⑥ Summarization
 - * Summarizing a document to "Save time on the part of the reader, or to make it easier to store"
 - e.g. Paragraph into sentences
- ④ Classification
 - * Identifying the main themes of a document & then placing the documents into a pre-existing categories

- 5. clustering: grouping similar documents
- 6. concept linking: connects related documents by identifying their shared concepts & by doing so, helps user find information that they would not have found using traditional search methods
- 7. It user reads an article about climate change
 - * the system shows a link to another article about renewable energy.
 - * This helps the user find more related information they might not have seen before
- 8. Question answering
 - * It finds the best answer to a given question through knowledge-driven pattern matching
 - e.g. "What causes rain?"
↳ caused by condensation of water vapour in the atmosphere.

- Natural Language Processing (NL)
- Def: A branch of artificial intelligence & computational linguistics that focuses on the interaction between computers & human language.

Goal: convert human language into structures data that computers can process, aiming for a deeper understanding.

Objectives:

- * Moves past simple word counting to include understanding of grammar, semantics & context.
- * True understanding of natural language is difficult because it requires extensive knowledge that goes beyond the text itself.
- * NLP has evolved from basic text processing methods to more advanced techniques that attempt to understand language more deeply.

Some challenges commonly associated with the implementation of NLP

1. part-of-speech tagging
2. Text Segmentation
3. word Sense disambiguation
4. Syntactic ambiguity
5. imperfection of irregular input.
6. Speech acts

1. parts of speech tagging

It is difficult to mark up terms in a text to a particular part of speech because the part of speech depends not only on the text of the term but also on the context within which it is used.

* As a verb:

- * "please light the candle."
- * in this case, "light" is an action

* As an adjective:

- * "she carried a light backpack."
- * "light" describes the weight of the backpack

2. Text Segmentation

* Text-parsing task requires the identification of word boundaries, which is often a difficult task.

* In speech segmentation emerge when analyzing spoken language, because sounds representing successive letters & words blend into each other.

Ex: "I scream," it can sound like "ice cream" when spoken quickly. The blending of sounds makes it challenging to tell where one word ends & the other begins, causing confusion in understanding.

3. Word sense disambiguation

- * Many words have more than one meaning
- * Selecting the meaning that makes the most sense can only be accomplished by taking into account the content within which the word is used
- * Financial institution: "I need to deposit money at the bank"
- * "River bank: "we had a picnic by the river bank."

4. Syntactic ambiguity

- * The grammar for NLP is ambiguous. That is, the grammar usually requires a fusion of Semantic & contextual information
- * need to be considered. choosing the most appropriate structure
- * usually requires a fusion of Semantic & contextual information

- * "I looked at the man using the telescope."

possible structures:

- * "I looked at the man using the telescope"
- * "I looked at the man who was using the telescope"
- * "I looked at [the man]. Using the telescope observed the man through the telescope"

5. Imperfect or irregular input.

- * Foreign or regional accents & vocal imitations in speech & typographical (or grammatical) errors in texts make the processing of the language an even more difficult task.

- * speech with accents: "I need to book a flight" might be pronounced differently by some one with a strong regional accent, making it harder for speech recognition systems to accurately transcribe it

- * typographical errors: "I need a research paper" instead of "I need a research paper" can confuse text processing systems & lead to incorrect or incomplete interpretation

6. Speech acts

- * It is something expressed by an individual that not only presents information but performs an action as well
- * the sentence structure alone may not contain enough information to define this action.
- * This might be requests, warnings, promises, apologies, physical actions, greetings.

Q: "Can you pass the class?"

These are tough for computers because understanding language involves more than just recognizing words.

11. optical character recognition
the automatic translation of images of handwritten, typewritten or printed text into machine editable textual documents

Text Mining Applications

1. Marketing Applications
2. Security "
3. Biomedical "
4. Academic "

Marketing Applications

1. Cross-Selling & Up-Selling

- * Text mining examines unstructured data from call centre notes & voice transcripts to understand customer perceptions
- * This analysis helps identify opportunities for cross-selling & up-selling

2. Customer Sentiment Analysis

- * Text mining processes blogs, product review & discussion boards posted to capture customer sentiments
- understanding this rich data helps enhance customer satisfaction & increase their overall lifetime value with the company

3. Customer Relationship Management (CRM)
* Companies combine unstructured text data with structured data from their databases to gain insights into customer behavior

- * Text mining improves the ability to predict customer churn, enabling companies to identify at-risk customers for targeted retention efforts.

4. Product Attribute Analysis

- * Text mining systems can identify both explicit & implicit product attributes allowing retailers to analyze product databases more effectively.

- * Explicit attributes are clearly defined features such as color, size or brand, while implicit attributes are inferred from data, like customer sentiment (or) usage context.

- * Treating products as sets of attribute-value pairs enhances effectiveness in demand forecasting, product recommendations & supplier selection.

Product: Smartphone

Implicit Attribute

These are clearly defined features that can directly be observed or stated about the product.

Brand: Samsung

color: Black

storage: 128GB

Price: \$699

Implicit Attribute

They require interpretation often using methods like text mining, sentiment analysis.

Customer Sentiment

Positive reviews about the smartphone's battery life or design usage context. The smartphone is often used for gaming or photography.

Security Applications

a) Echelon Surveillance System

- * Echelon is major classified text mining application in the security field
- * It is classified system rumored to analyze content from phone calls, faxes, emails & other communications
- * It can intercept information via various channels including satellite & public networks

b) EUROPOL's OASIS System

- * It was developed in 2007 to track trans-national organized crime by analyzing large volumes of structured & unstructured data

* It integrated advanced data & text mining technologies enhancing law enforcement efforts internationally.

(c) FBI and CIA Supercomputer System

The FBI & CIA are working together to create a comprehensive data warehouse for law enforcement.

(d) Supercomputer System goal is to improve knowledge discovery by connecting previously separate databases enhancing data accessibility for federal, state & local agencies.

(e) Deception Detection

* Text mining is used to analyze statements from persons of interest in criminal investigations

* It developed model that achieved 40% accuracy in distinguishing between deceptive & truthful statements relying solely on textual cues

* This method can be applied to both text and transcriptions of voice recordings, offering an alternative to traditional techniques like polygraphs

Biomedical Applications

Medical literature

Developing medical literature is the field of biomedical literature, with the expanding rapidly, particularly, with the use of open-source journals. Medical literature is well-organized & standardized, making it easier to mine for valuable insights.

Data from experimental Techniques

Techniques like DNA microarray analysis, mass Spectrometry, produce vast amounts of data. Text mining tools are essential for analyzing & interpreting this data, by cross-referencing it with existing literature. Text mining techniques are enhancing understanding of biological processes by predicting protein locations, extracting disease-gene relationships & discovering gene-protein interactions.

Protein location prediction

Knowing a protein's location in a cell is crucial for understanding its biological role & potential as a drug target. Shatkay (2007) developed a system that integrates sequence-based & text-based features to predict protein locations. Their system outperformed previous models.

Disease-gene Relationship Extraction

Chen, et. al (2006) created a system to extract relationships between diseases & genes from biomedical literature using a dictionary for disease/gene names.

To reduce false positives, they used machine learning-based named entity recognition filtering, which improved precision by 26.7%

Gene-protein Relationship Discovery

Nakov (2005) demonstrated a 'middle-level' text analysis process to discover gene-protein & protein-protein relationships from biomedical texts. This process involves tokenizing text, part-of-speech tagging & shallow parsing, then matching terms against domain ontologies to interpret relationships.

Academic Applications

To analyse large volume of scholarly literature, identifying research trends, key topics, emerging areas of study, research gaps & patterns within a field, allowing researchers to quickly summarize & extract relevant information from published papers, ultimately accelerating the research process.

Text Mining Tools

Commercial Software Tools

(i) **ClearForest**: offers text analysis & visualization, provides Matador & data & text analysis toolsuite.

(ii) **IBM** offers SPSS Miner provides a rich suite of text processing & analysis tools.

(iii) **rxEN Text Coder (KTC)** offers a text analysis solution for automatically preparing & transforming unstructured text attributes into a structured representation for use in KXEN analytic framework.

(iv) **Statistica Text mining engine** provides easy-to-use-text mining functionality with exceptional visualization capabilities.

(v) **Vantage point** provides a variety of interactive graphical views & analysis tools with powerful capabilities to discover knowledge from text databases.

(vi) **WordStat analysis module** from provaxis Research analyzes textual information such as responses to open-ended questions, interviews etc.

Free Software Tools

1. **Apache open source tools** which are open source, are available from a no of non profit organization.

1. **RapideMiner**: one of the most popular free open source software tools for data mining & text mining is tailored with a graphical appealing, drag & drop user interface.

2. **open calais** is an open source toolkit for including semantic functionality within your blog, content management system, website etc.

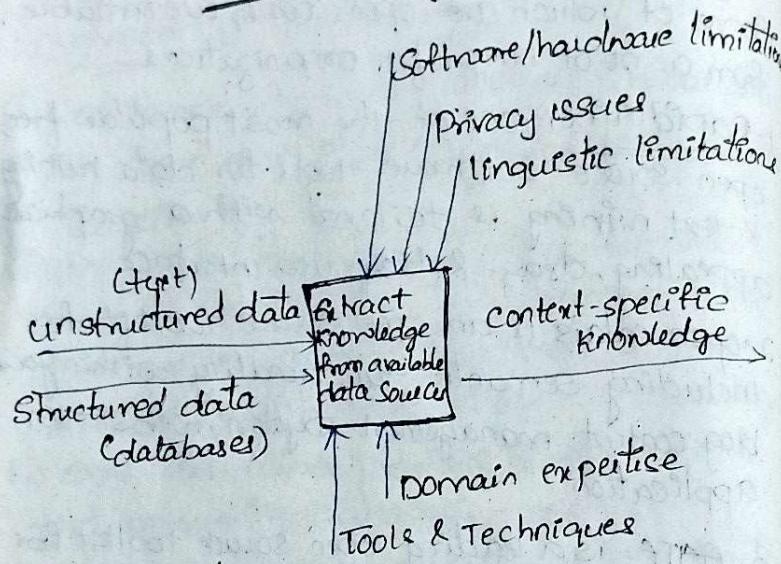
3. **SpaTe** is a leading open source toolkit for text mining. It has a free open source framework & graphical development environment.

4. **lingpipe** is a suite of Java libraries for the linguistic analysis of human language.

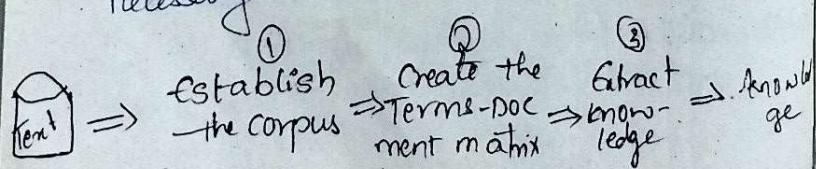
5. **stem (spstem)** is a text classification system that learns from positive & unlabelled examples.

6. **vinersmoldust**: is a web search and text clustering engine.

Text Mining Process



- * The text-based knowledge discovery process involves analyzing both unstructured & structured data to generate content-specific knowledge for decision-making.
- * Inputs include collected data, while outputs are actionable insights.
- * The process is governed by constraints like software limitations & privacy issues, & relies on techniques & expertise.
- * It consists of three tasks, with feedback loops for adjusting outputs if necessary.



Three step text mining process

1. Establish the corpus: collect & organize domain-specific unstructured data (Text, XML, HTML)

2. Create the Term-Document Matrix
Structure the data into a matrix where each cell represents term frequency in documents.

3. Extract knowledge:
use the matrix to find patterns using classification, clustering (or) association techniques.

1. Establish the corpus.

The main purpose of the first task actually is to collect all of the documents related to the context being studied. This collection may include textual documents, XML files, emails, web pages & short notes, the readily available textual data voice recordings may also be transcribed using speech recognition algorithms & made a part of the text collection.

Once collected, the text documents are transformed & organized in a manner such that they are all in the same representation form for computer processing. The organization of the documents can be simple as a

collected of digitized text except stored in a file folder or it can be a list of links to a collection of web pages in a specific domain. Many commercially available text mining software tools could accept them as input & converts them into a flat file for processing. Alternatively, the flat file can be prepared outside the text mining software & then presented as the input to the text mining application.

2. Create the Term-Document

The digitized & organized documents are used to create the term-document matrix. In the TDM rows represent the documents & columns represent the terms.

The relationships between the term & document are characterized by indices

Document	Term 1	Term 2	Term 3	Term 4	Term 5	Term 6
DOC 1	1				1	
DOC 2		1				
DOC 3			3	1		
DOC 4		1				
DOC 5			2	1		
DOC 6	1				1	
..						

Steps

1. Collect the Text Data: Gather the text documents that need to be analyzed. They could be in any format.
2. Preprocessing the Text: Clean & prepare the text data by performing
 - **Tokenization**: Split the text into words
 - **lowercasing**: Convert all text to lowercase to avoid treating the same word in different cases as distinct
 - **Removing stop words**: Remove common words that do not add much meaning
 - **stemming/lemmatization**: Reduce words to their base or root form
 - **Removing punctuation & special characters**: Clean up the text to focus only on words
3. Build a vocabulary: Identify all the unique terms in the corpus.
4. Create the Matrix structured
 - Rows: Each row corresponds to a document
 - Columns: Each column corresponds to a unique term from the vocabulary
 - Cells: The value in each cell represents the frequency or weight of the term in the respective document

Several ways to fill the cells

Binary encoding: Assign 0 or 1

Term Frequency (TF): count how many times a term appears in a document

Term Frequency - inverse document frequency (TF-IDF): A weighted term frequency, where the TF is adjusted by how commonly the term appears across all documents

Example

DOC 1 : "Text Mining is fun!"

DOC 2 : "Text Mining involves extracting useful information."

DOC 3 : "Mining data can be fun & useful!"

Step 0: preprocess the text

* after preprocessing, the tokenized & cleaned text for each document

DOC 1: ["text", "mining", "fun"]

DOC 2: ["text", "mining", "involves", "extracting", "useful", "information"]

DOC 3: ["mining", "data", "can", "be", "fun", "useful"]

④ vocabulary

["text", "mining", "fun", "involves", "extracting", "useful", "information", "data", "can", "be"]

③ create the term document matrix.

	DOC 1	DOC 2	DOC 3
Term	1	1	0
Text	1	1	1
mining	1	0	1
fun	1	1	0
involves	0	1	0
extracting	0	1	0
useful	0	1	1
information	0	0	1
data	0	0	1
can	0	0	0
be	0	0	1

③ Extract knowledge

Goal: Discover meaningful patterns or insights from the term-document matrix

Action: use text mining algorithms to uncover patterns, such as key topics or word relationships

Output: Extracted knowledge or insights such as clusters of similar documents or terms, which can be used to solve specific problems

Apply text mining techniques to analyze the structured term-document matrix. common algorithm used included.