MASTER MEIM 2021-2022

# BIG DATA ANALYTICS
# Master 2021-2022

**TEXT MINING**

Giorgia Rivieccio

Professor of Economic Statistics at Parthenope University

# BIG DATA ANALYTICS

**BIG DATA ERA**

*In the data-driven era, providing value means being able to translate data rapidly into value-add information.*

In today's fast growly and highly complex world, the main challenge is to find meaning in data, so that the derived knowledge can be used to make informed decisions.

The Big Data era is characterized by:

➢ A large and complex quantity of data

➢ The development of IT

Examples: financial services, fraud detection, implementation of algorithms for trading, risk analysis, retail, CRM.

# BIG DATA ANALYTICS

BIG DATA ANALYTICS

Big data analytics includes all the activities to help managers in coping with strategic decisions, in achieving major goals and in solving complex problems, by collecting, analyzing and reporting the most relevant information.

Information could be about the causes of the current situation, the most likely trends to occur, and what should be done as a result.

The purpose of the collection and processing of large volumes of complex data is to understand the trends of the phenomena of interest, uncover hidden trends, detect anomalies, in order to take data-driven decisions.

# BIG DATA ANALYTICS

ARTIFICIAL INTELLIGENCE

Inundated with information and cognitive stimuli in a world where about 80% of data consists into free text, written in natural language, not classified or structured, treating and analyzing such a large amount of data can simultaneously represent a opportunity and a challenge.

Artificial intelligence techniques such as text mining, a technique aimed at extracting useful information hidden in a text, respond to this challenge.

This technique uses natural language processing (NLP) to transform free text, also known as "unstructured", into structured and normalized data.

A brief introduction is therefore required to distinguish structured data from unstructured ones.

# BIG DATA ANALYTICS

STRUCTURED DATA VS UNSTRUCTURED DATA

Structured data are classical data, which are organized and well formatted, and conform to the formal system of relational databases and traditional spreadsheets. In summary, they are all those quantitative data that can be stored in a SQL (Structured Query Language) database.
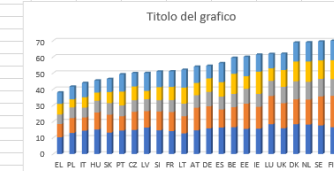
Structured data are analyzed by means of data mining.

Unstructured data, not being organized or correctly formatted, implies many difficulties in their collection, processing and analysis.

They are qualitative data processed by text mining.

Examples: Newspaper articles, market place reviews, web pages, e-mails, messaging are just a few examples.

# BIG DATA ANALYTICS

DATA  MINING VS TEXT MINING

The challenge, therefore, with text mining consists precisely in the creation of algorithms (image), and computer and automated procedures capable of preparing this data source for an analytical use, accessible, in QUANTITATIVE form, to the various data mining algorithms.

Text mining and data mining are considered complementary techniques.

Specifically, text mining extracts, starting from NON-STRUCTURED textual information, STRUCTURED numerical data to be processed later with DATA MINING techniques.

To transform unstructured textual resources into structured information text mining adopts linguistic and statistical techniques capable of both analyzing free text formats and combining each document with metadata.

# BIG DATA ANALYTICS

## NATURAL LANGUAGE PROCESSING (NLP)

Therefore, Text Mining is thus an Artificial Intelligence (AI) technique that uses the Natural Language Processing (NLP) to transform the free, unstructured text of documents / databases such as web pages, newspaper articles, e-mails, press, post / comment on social media etc. in structured and normalized data.

With the expression Natural language processing (NLP), we mean that part of computer science and artificial intelligence that deals with the processing of human languages, capable of performing a linguistic analysis, which, essentially, helps a machine to "read" the text.

# BIG DATA ANALYTICS

TEXT MINING

Typically, text mining is used for two main functions:

classifying data by segment and predicting its behavior.

The most commonly used features are:

• Cleaning, literally the process of cleaning the text aimed at organizing the material for analysis.

• Extraction, the extraction of keywords that involves the choice of words and phrases that best identify the nature of the text.

• Categorization, the classification of the text into one or more categories.

• Modeling, which combines the essential parts of the text with structured data attributes and is useful for optimizing the forecast.

# BIG DATA ANALYTICS

OBJECTIVES OF TEXT MINING

o   Identify thematic groups

o   Classify documents into predefined categories

o   Discover hidden associations (links between topics, or between authors, temporal trends, ...)

o   Extract specific information (ex: names of geniuses, names of companies, ...)

o   Train search engines

o   Extract concepts for the ontology learning.

o   Identify attitudes and opinions on opposite emotional states (sentiment analysis)

# BIG DATA ANALYTICS

**Text mining applications**

Among all the possible applications, the most frequent are:

• Digital marketing for contextual retargeting. Compared to the traditional cookie-based approach, contextual advertising offers greater accuracy, fully preserving the user's privacy.

• Customer care, as it allows to monitor and improve the customer experience using various sources of information such as surveys, trouble tickets and customer call notes, to optimize problem resolution

• Engagement analysis and brand reputation: to analyze or predict the needs of customers and understand the perception of their brand, it is possible to analyze both large volumes of unstructured data, and by extrapolating opinions, emotions and relationships with brands and products and feelings.

# BIG DATA ANALYTICS

STEPS OF TEXT MINING

**1. Pre – Processing:**

- Corpus

- Tokenization

- Stemming

**2. Document term matrix**

**3. Exploratory Analysis**

- Word frequency e word cloud

- Word association

- Word clustering

# BIG DATA ANALYTICS

PRE-PROCESSING STEPS:

- Data selection (creating corpus)

- Data manipolation (tokenization)

- Reduction of the variants associated with each word (stemming)

- Generation of characteristics and vectors

- Classification

- It requires solving numerous difficulties in data processing

- Ambiguity of language: the same word can take on different meanings depending on the context  and different words can mean the same thing (synonyms)

- Language sensitivity (sensitive topics)

- Numerous dimensions involved in the extraction of concepts / words

- Difficulties due to: spelling errors, abbreviations, language variants, etc…

# BIG DATA ANALYTICS

## PRE-PROCESSING - CORPUS

It transforms a group of separate text documents into a single text that merges them all.

❑ Client 1: Good morning, I would like to have information on product X.

❑ Client 2: I would like to know where product X can be purchased.

↩

Good morning, I would like to have information on product X. I would like to know where product X can be purchased.

# BIG DATA ANALYTICS

It is the manipulation of data to extract the relevant elements, namely words, phrases or even letters.

The text is divided into tokens, which are blocks of atomic text, made up of indivisible characters, as a sequence of characters surrounded by delimiters

There is great demand on the market for this product

There | is | great |demand |on | the | market| for | this | product

# BIG DATA ANALYTICS

PRE-PROCESSING - TOKENIZATION

**DELETE STOPWORDS**

Stopwords are non-informative words, like articles, prepositions.

In the text processing, it is therefore necessary to load, for each language, a list of "stopwords" that the system will eliminate before proceeding with the processing of the text. Example on a facebook comment.

| Description | Before | After |
|---|---|---|
| Remove capital letter | The product is fantastic,      the NUMBER 1! | the product is fantastic,     the number 1! |
| Remove punctuation | the product is fantastic,      the number 1! | the product is fantastic     the number 1 |
| Remove numbers | the product is fantastic     the number 1 | the product is fantastic     the number |
| Remove spaces | the product is fantastic     the number | the product is fantastic the number |
| Remove specific terms | the product is fantastic the number | product fantastic number |

# BIG DATA ANALYTICS

PRE-PROCESSING - STEMMING

Process of reducing the inflected form (any morphological variation) of words to the basic form, called root or theme.

Extraction of the root of a word, removing affixes and endings (e.g. playing, playing games, players, play…).

# BIG DATA ANALYTICS

Identification of the NLP vocabulary (lemma) starting from a word with ending. Determining the part of speech of a word, and then applying the different normalization rules.

Given a wordform, stemming is a simpler way to get to its root form. Stemming simply removes prefixes and suffixes.

Lemmatization on the other hand does morphological analysis, uses dictionaries and often requires part of speech information.

Thus, lemmatization is a more complex process

**Stemming**

adjustable → adjust
formality → formaliti
formaliti → formal
airliner → airlin ⚠

**Lemmatization**

was → (to) be
better → good
meeting → meeting

# BIG DATA ANALYTICS

## 2. DOCUMENT-TERM MATRIX (DTM)

It consists of one of the most common formats for representing a corpus of text in a bag-of-words format.

The Document matrix is a table containing the frequency of each word occurring in the text.

|  | word 1 | word 2 | word 3 | ... | word k |
|---|---|---|---|---|---|
| document 1 |  |  |  |  |  |
| document 2 |  |  | $n_{23}$ |  |  |
| document 3 |  |  |  |  |  |
| ... |  |  |  |  |  |
| document n |  |  |  |  |  |

n23= Number of times which the word 3 occurs in the document 2

The resulting data frame is

| Term | Frequency |
|---|---|
| Bad | 77 |
| Beautiful | 58 |
| Need | 30 |
| Assistance | 26 |
| break | 24 |
| small | 20 |
| negative | 18 |
| failure | 10 |
| problem | 7 |
| price | 5 |

# BIG DATA ANALYTICS

**EXPLORATORY ANALYSIS -**

WORD FREQUENCY

The bar-plot allows you to compare the frequency of the most used words in a text

**Most frequent words**

# BIG DATA ANALYTICS

EXPLORATORY ANALYSIS -

WORD CLOUD

Once the frequency associated with each word detected in the text has been identified, the so-called word cloud can be presented, showing the most frequent terms of the text.

Word clouds are a powerful communication tool. They're easy to understand, easy to share, and they're impactful.

Word clouds add simplicity and clarity. The most used keywords stand out best in a word cloud. By size and color.

Word clouds are more visually appealing than table data

# BIG DATA ANALYTICS
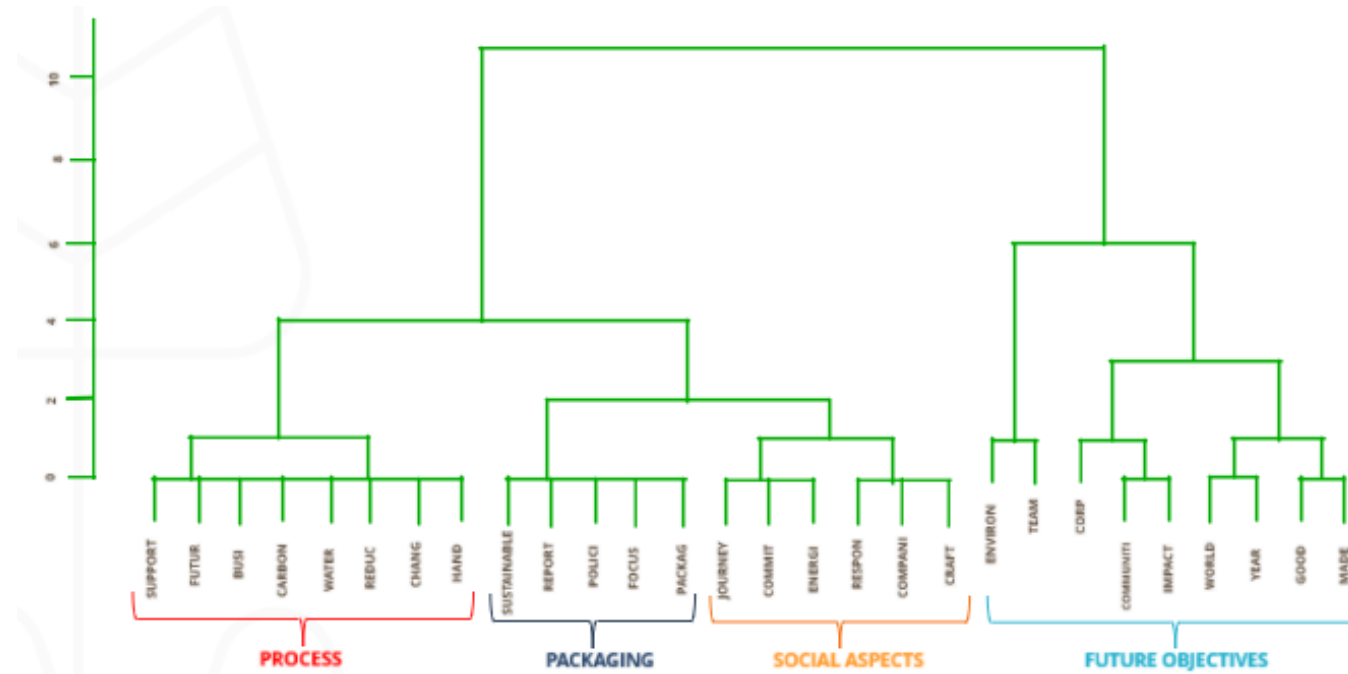
**EXPLORATORY ANALYSIS - WORD ASSOCIATION**

The link between the words can be summarized through a plot of the network that can be created between them.

# BIG DATA ANALYTICS

EXPLORATORY ANALYSIS - WORD CLUSTERING

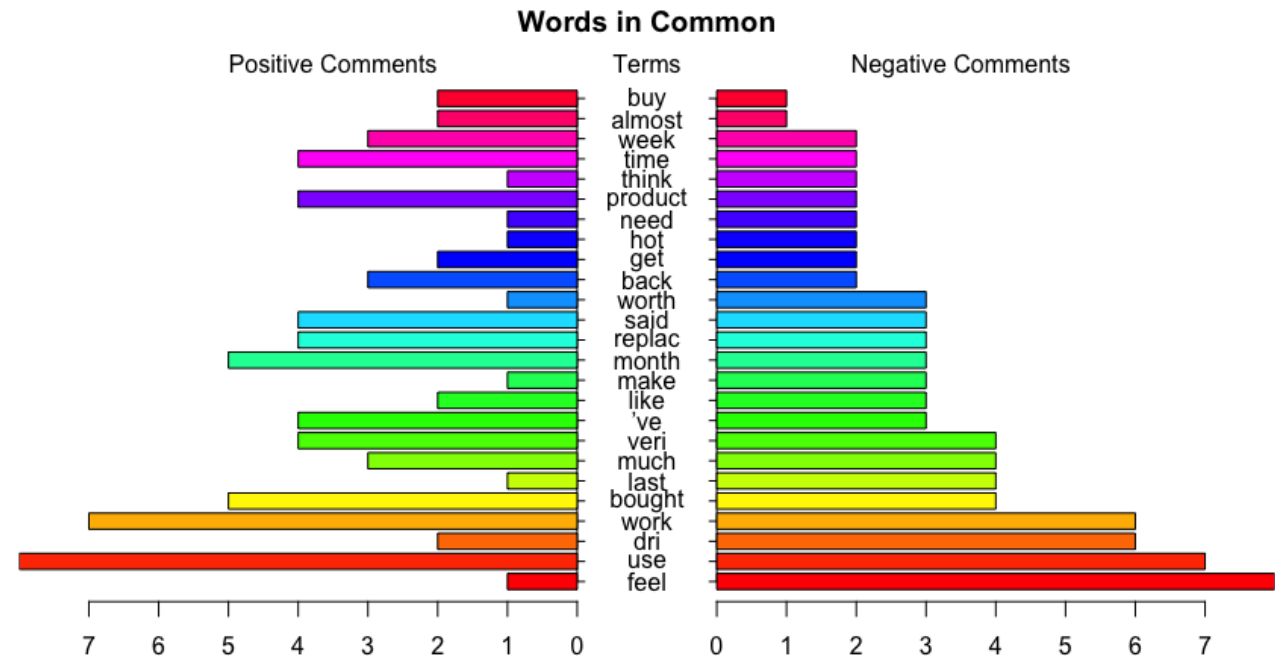THE HIERARCHICAL CLUSTER OF WORDS IS REPRESENTED
BY THE DENDROGRAM

# BIG DATA ANALYTICS

EXPLORATORY ANALYSIS - POLARYZED TAG PLOT

The polarized tag-plot (pyramid plot) allows to identify the frequency of a term used into two different documents

It is created starting from a data matrix with all common words occurring in both corpora.

Adding another matrix for the absolute difference between both corpus for each word and the graph is created.



Words in Common

# BIG DATA ANALYTICS

EXPLORATORY ANALYSIS – SENTIMENT ANALYSIS

Sentiment Analysis is a technique aimed at identifying the opinions expressed in online texts on a product or service, on a company, on a brand or on an event. This type of analysis allows to understand the nature of the interactions carried out between users, in a precise context and in a given period of time.

It is a multifunctional tool, it can be used in a different but functional way in various fields, even very different from each other.

EXAMPLES: Companies that want to know directly the opinion of their users, political parties, sociologists, museums, research bodies, computer scientists, even seismologists, epidemiologists, …
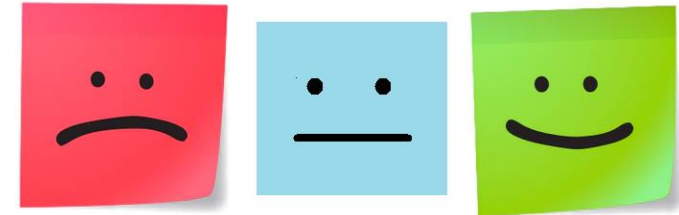
# BIG DATA ANALYTICS

EXPLORATORY ANALYSIS – SENTIMENT ANALYSIS

Sentiment Analysis is a useful tool for market research, such as understanding the opinion of the chosen target on a specific product or topic of interest or it can facilitate the work of market segmentation in order to get to know your customers or main stakeholders better.

Furthermore, it can be a useful tool to analyze Brand Reputation, through Social Networks, to understand the general opinion of stakeholders and in particular of customers regarding their brand.

It still allows you to monitor marketing campaigns by evaluating the effectiveness of a specific marketing activity.

Evaluating digital *word of mouth* thus becomes an important aspect in evaluating the brand reputation of a brand or customer satisfaction, in relation to a service or product.
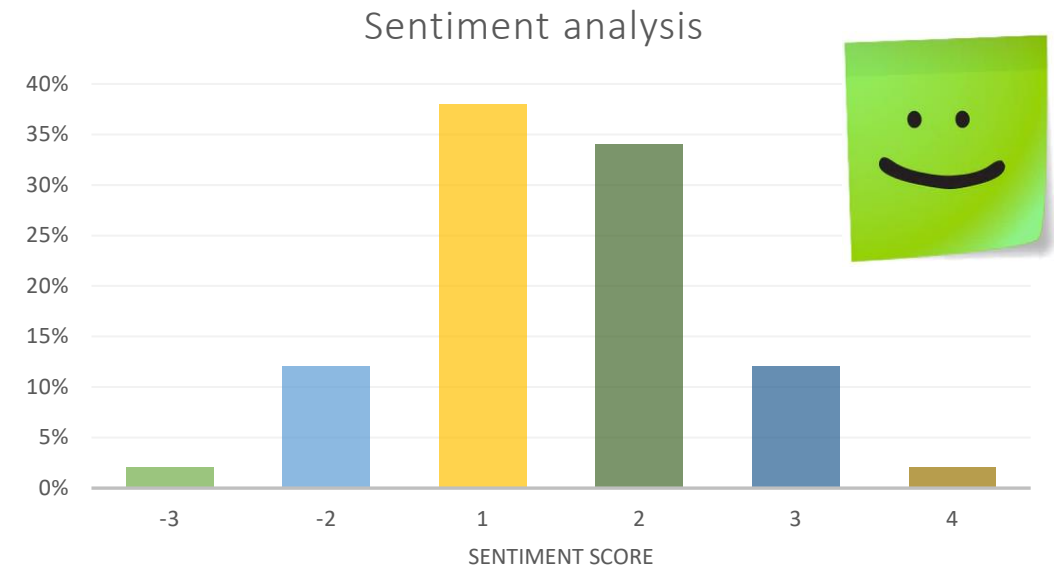
# BIG DATA ANALYTICS

EXPLORATORY ANALYSIS – SENTIMENT ANALYSIS

Furthermore, given that the discussion on social media can take place at the same time as the choice process is matured or even coinciding with the purchase process, the monitoring of social channels has an important impact on the success of new products, and on the effectiveness of communication or marketing campaigns (Jansen et al., 2009).
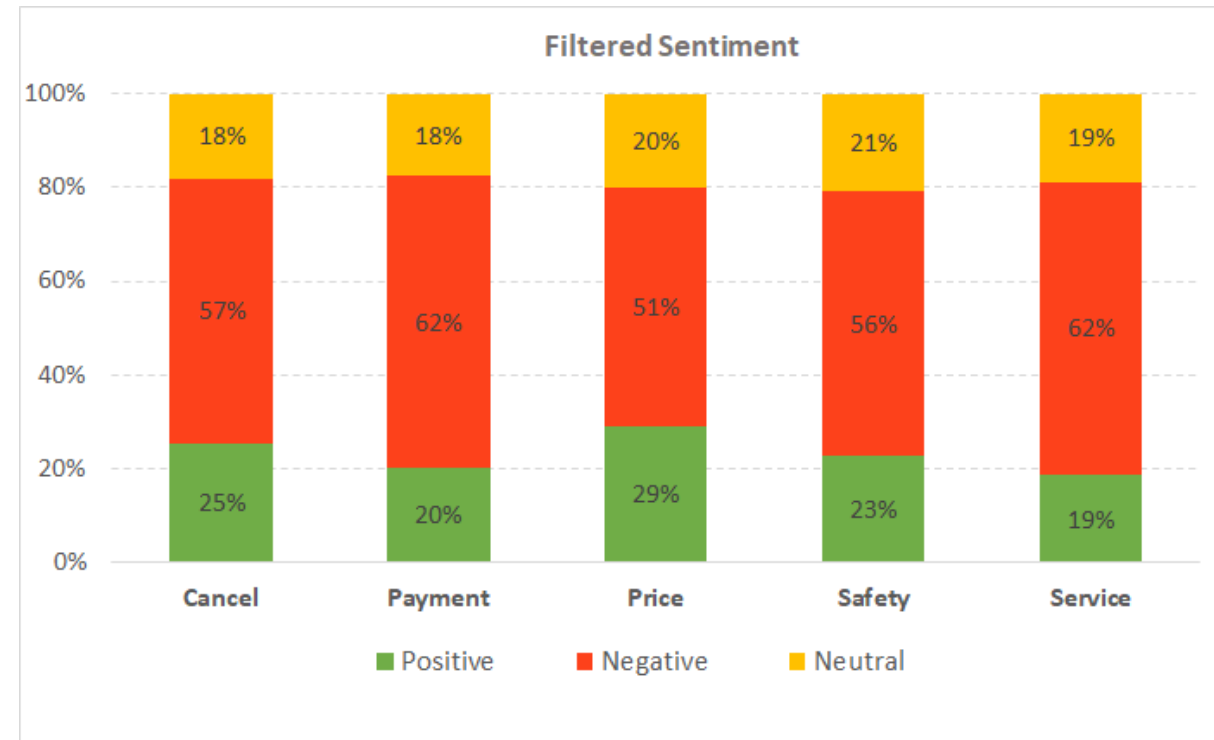
However, the judgment can also be expressed subsequently with respect to the consumption of the same, thus allowing the satisfaction of the buyer to be assessed. Social media becomes an important brand management tool. With this approach, therefore, consumer sentiments and opinions become the fulcrum of analysis. There are two approaches used by companies to monitor the mood of their market: Top Down and Bottom Up.



Sentiment analysis

# BIG DATA ANALYTICS

EXPLORATORY ANALYSIS – SENTIMENT ANALYSIS

Identify a sort of polarity that people show towards a topic by creating an index that associates numerical values from "completely positive" to "completely negative", passing through the neutral position.

MASTER MEIM 2021-2022

# THANKS FOR YOUR ATTENTION AND LET'S START TO APPLY TEXT MINING!