


# Practical Text Analytics: Introduction

Fan Dai, Subrata Pal

Iowa State University

September 24, 2019

- Github repo for tutorials and coding materials  
<https://github.com/fanstats/text-analytics>
- References
  -  Murugan Anandarajan, Chelsey Hill, Thomas Nolan  
Practical text analytics: maximizing the value of text data  
*Springer International Publishing*, 2019  
<https://link-springer-com.proxy.lib.iastate.edu/book/10.1007%2F978-3-319-95663-3>
- Practical skills with R  
Prepare a corpus of documents, convert them to a data matrix and illustrate statistical methods available for different analysis.  
Packages: `tm`, `tidytext`, `tidyr`, `tidyverse`, `hunspell`, etc.

# Introduction: Text data

- Text data
  - Product reviews, social media posts, customer feedback, etc.



- \* Concepts & topics
- \* Sentiment & Emotions
- \* Behaviors

Source: <https://www.wordstream.com/social-media-marketing>

# Introduction: Text data

- Text data
  - Product reviews, social media posts, customer feedback, etc.
  - **Unstructured data**
    - \* Structured data: stored and organized in typical databases like Excel, Google Sheets, and SQL.
    - \* Unstructured data: everything else..

# Introduction: Text data

- Text data
  - Product reviews, social media posts, customer feedback, etc.
  - **Unstructured data**
    - \* Structured data: stored and organized in typical databases like Excel, Google Sheets, and SQL.
    - \* Unstructured data: everything else..

Traveller Id	Traveller star	location	cleanliness	Service	Sleep Quality	Value	Rooms	Trip Type	Text reviews
Filippo S	5		5	5	5			business	I had the most incredible experience in this hotel. I
Explorer078	3	1	3	2	2	3	3	business	Getting in and out of hotel is extremely hazardous I
hazimismail	4	3	3	4	4	3	4	business	Wverything was fine with the hotel. Rooms were a
Barranquita	5	5	5	5	5	5	5	family	Free airport shuttle on time both at arrival and dep
									HGI Miami APW was a very nice place to unwind af Cruise. They were very accommodating (able to che free water bottles). It is not too far from the airport free shuttle) and has many fast-food restaurants (M Starbucks, Dunkin Donuts) within walking distance. Also appreciate the onsite washing machines and ci a small fee) and morning breakfast (another fee). Highly recommend.
Josephine B	5								
Isabelle W	5	5	5	5	5	5	5	couple	Just check-out of my favorite Miami hotel - again, th
mikeslady	4								Great breakfast buffet, nice seating area in lobby fo
Frankie_Sur	4	3	4	4	5	4	5	family	Very Nice Hotel , Clean rooms, free shuttle to Airpo

Figure: Customer reviews for Hilton Garden Inn-Miami Airport (Dec-Jan, 2019)

# Introduction: Text data

- Text data
  - Product reviews, social media posts, customer feedback, etc.
  - **Unstructured data**
    - \* Structured data: stored and organized in typical databases like Excel, Google Sheets, and SQL.
    - \* Unstructured data: everything else..
  - **Big data**
    - \* Massive volumes
    - \* Many different file types
    - \* High-speed data creation

# Introduction: Text analytics

- Text analytics
  - Uses: including but not limited to email filtering, product suggestions, fraud detection, opinion mining, trend analysis, search engines, etc.
  - Goals and objectives: understanding semantic information, text summarization, classification, clustering, etc.
  - Information retrieval and extraction
  - Data mining
    - \* Quantitative results
    - \* Graphs, tables and other sorts of visual reports.





# Introduction: Text analytics



Figure: Word cloud of text analytics job titles

Source: Practical text analytics: maximizing the value of text data

## Introduction: Text analytics



**Figure:** Word cloud of the skills required for text analytics jobs

Source: **Practical text analytics: maximizing the value of text data**

# Introduction: Text analytics process

## Text analytics process

- Planning the text analytics projects
- Preparing and preprocessing text
- Analyzing data
- Interpreting results

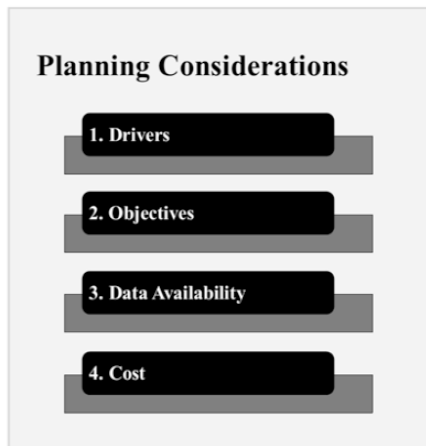
# Introduction: Planning the text analytics

## Planning the text analytics projects

- Initial consideration
- Problem framing
- Data generation
- Method and implementation selection
- Inference and decision-making

# Introduction: Planning the text analytics

## Initial considerations



**Figure:** Initial planning considerations

Source: **Practical text analytics: maximizing the value of text data**

# Introduction: Planning the text analytics

## Problem framing

- Identify problems of analysis
- Decide types of inference
  - **Deductive inference**
    - Research questions  $\Rightarrow$  Theory  $\Rightarrow$  Hypotheses  $\Rightarrow$  Conceptualization
    - Theory-driven coding schemes

# Introduction: Planning the text analytics

## Problem framing

- Identify problems of analysis
- Decide types of inference
  - **Deductive inference**
    - Research questions  $\Rightarrow$  Theory  $\Rightarrow$  Hypotheses  $\Rightarrow$  Conceptualization
    - Theory-driven coding schemes
  - **Inductive inference**
    - Data-driven coding schemes
    - Create a generalizable theory

# Introduction: Planning the text analytics

## Data generation

- Define research scope and purpose
  - Unitization
    - Sampling unit: each row of data or text document
    - Recording unit: piece of text that is categorized in the coding step
    - Context unit: define the amount of text that can be used



# Introduction: Planning the text analytics

## Data generation

- Define research scope and purpose
  - Unitization
    - Sampling unit: each row of data or text document
    - Recording unit: piece of text that is categorized in the coding step
    - Context unit: define the amount of text that can be used
- Collect text data
  - Sampling techniques
    - Non-probability sampling: convenience sampling, relevance sampling
    - Probability sampling: systematic sampling, stratified sampling

# Introduction: Data generalization

- Data quality



**Figure:** Quality dimensions of text data

## Example: Business policies on human rights (Preuss and Brown, 2012)

- Inductive inference
- Research questions: frequency and content of human rights policies in large corporations
- Relevance sampling: all companies on the Financial Times Stock Exchange 100 (FTSE 100) list
- Sampling unit: each firm's website
- Context unit: corporate codes of conduct and supporting documents
- Recording unit:  
For example, to answer whether a company had a human rights policy, the recording unit is each firm's website

# Introduction: Text analytics process

Different methods based on our objective(s)

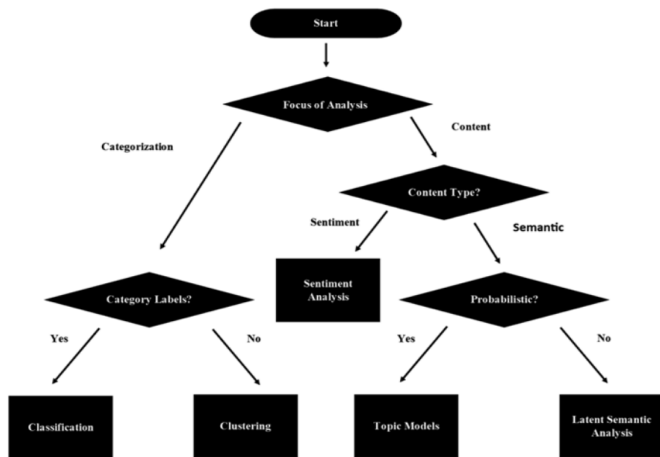


Figure: Flow chart of what-to-do based on our objective

## Different methods based on our objective(s):(contd.)

- If our objective is to categorize or classify the documents, go for **Classification** or **Clustering** depending on levels exists or not.
- If our objective is not to categorize or classify the documents, rather we're concerned with the content of the documents, we should go for the other type:

## Different methods based on our objective(s):(contd.)

- If our objective is to categorize or classify the documents, go for **Classification** or **Clustering** depending on levels exists or not.
- If our objective is not to categorize or classify the documents, rather we're concerned with the content of the documents, we should go for the other type:
  - Identify the sentiment of the documents in the document collection → sentiment analysis.  
If the data have sentiment labels, classification methods can be used also.

## Different methods based on our objective(s):(contd.)

- If our objective is to categorize or classify the documents, go for **Classification** or **Clustering** depending on levels exists or not.
- If our objective is not to categorize or classify the documents, rather we're concerned with the content of the documents, we should go for the other type:
  - Identify the sentiment of the documents in the document collection → sentiment analysis.  
If the data have sentiment labels, classification methods can be used also.
  - To identify semantic ("*linguistic and philosophical study of meaning in language...*" - Wikipedia) information in the documents, there are two classical analysis methods available: **latent semantic analysis (LSA)** and **topic models**. Topic models are probabilistic, generative models; while LSA deals with dimension reduction methods.

## Different methods based on our objective(s):(contd.)

- If our objective is to categorize or classify the documents, go for **Classification** or **Clustering** depending on levels exists or not.
- If our objective is not to categorize or classify the documents, rather we're concerned with the content of the documents, we should go for the other type:
  - Identify the sentiment of the documents in the document collection → sentiment analysis.  
If the data have sentiment labels, classification methods can be used also.
  - To identify semantic ("*linguistic and philosophical study of meaning in language...*" - Wikipedia) information in the documents, there are two classical analysis methods available: **latent semantic analysis (LSA)** and **topic models**. Topic models are probabilistic, generative models; while LSA deals with dimension reduction methods.
- Other methods, like text summarization also exists.