

# Text mining / Text Analytics

Data in the form of Numbers, characters, special characters, white space etc

→ Text mining is a process of deriving high quality information from text.  
Used in Analytics.

→ Converting data into Tidy data

→ Each variable is a column

→ Each observation is a Row

→ Each type of observation unit is a table

India is my Country!!! → Unstructured

↓ transform  
india / is / my / country / !!!

→ india / my / country →

and, or, not, its, the, a → filtered

Chapt	page	chr
1	1	india
1	1	my
1	2	country
1	2	

Structured

Tidy text is a table with One token per row

Tokenization is a process of splitting text into tokens.

Tokens can be → Words → india / is / my / country / india / is / great /

→ Phrases / sentence → india is my country / india is great /

Text is often stored in below 3 types

1) Strings → Text can be a character vector. → Read in to Memory  
with Metadata (Description about data)

- 2) Corpus → Raw strings with Metadata (Description about data)
- 3) DTM → Document Text matrix → sparse matrix → Rows × Columns.

Steps → Unnest - tokens.

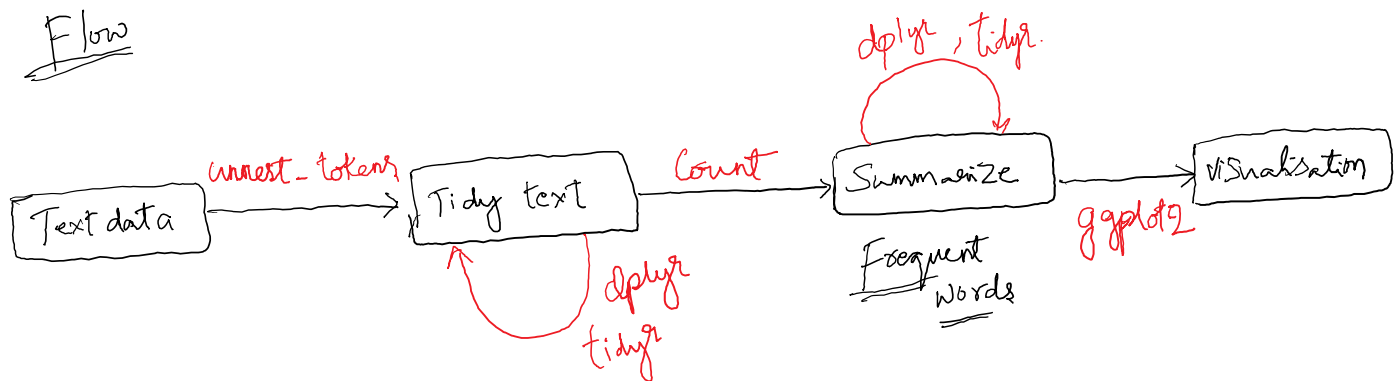
→ Convert the data into a data-frame → tibble dataframe

→ Split each sentence into Tokens → Assign each to a row.

I am Good! → i/am/good → i  
→ am  
→ good

→ Punctuations are removed  
→ Converted to lower case

Flow



Tokens → It is a meaningful unit of text, most often a word, that we are interested in using for further analysis.

Word frequency → Most Common task in Text Analysis

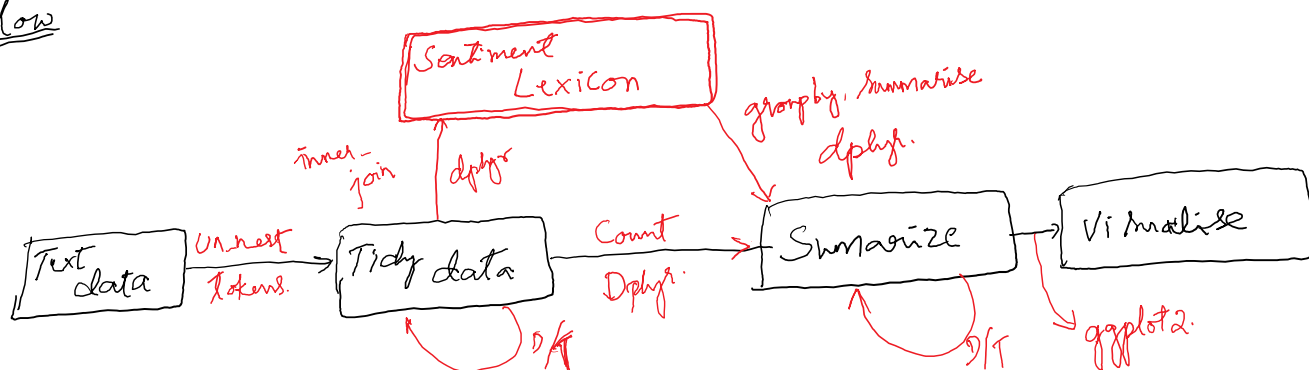
→ Which word occurs the most in a given text.

Sentimental Analysis → Process of Computationally identifying & Categorizing opinions expressed in a piece of text, to determine the authors attitude towards particular topic, product etc is Positive/Negative.

Flow



## Flow



Lexicon → like dictionaries. → Words are associated with subjects

In R  
 get.sentiment()  
 → sentiment of Lexicon  
 inner join  
 in sentiment analysis  
 Explained in code

	<u>Dictionary</u>	<u>Lexicon</u>
Happy	Meaning	positive
Sad	_____	Negative
Angry	_____	Negative
Cherish	_____	Positive
Good	_____	positive
Word	English meaning	Sentiment Associated.

Lexicons → AFINN from Arup Nielsen  
Bing from Bing Lin  
NRC from Saif Mohammed & Peter Turney

Compared in R code

Visualisation → Frequent words. → Most frequently occurring words  
 → Sentiment score → Word variation along the trajectory of the story  
 → Word cloud → Frequent words appear big in the cloud  
 → Sentiment word cloud → Combined sentiment scale & word frequency.  
 We can easily identify the most important +ve/-ve words & use

them in any of the further Analysis

---