

Text Mining with fi:-

Tokens and Tokenization:-

- Tokens are building blocks of text in text analysis.
- They represent individual units of text that have a meaning and can be analyzed.
- Often a token is a word, but it could also be a sentence, paragraph or even a group of words like n-gram.
- n-gram: a sequence of n words, such as "New york" for 2 words or "once upon a time" for 4 words.
- The act of breaking down a large piece of text into these tokens, is called Tokenization.
- It the act of taking a sentence, paragraph or document and splitting it into these smaller pieces.

why tokenization is important:-

- Tokenization allows us to examine the content and structure of text data in a manageable and analyzable format.

For instance, we can:

- Count the occurrence of each word in a text to identify frequently used terms.
- Filter out common but less ~~meaningful~~ meaningful words (like "the" or "and") to focus on keywords.
- Analyze patterns in text, like word frequency or sequence of words.

Corpus:- (plural- corpora)

A corpus is a structured collection of text documents. It's often created for text analysis or natural language ~~task~~ processing tasks.

- Each document in a corpus is typically annotated with metadata, which provides additional info like author, publication date, language or genre.
-

Document-Term Matrix (DTM)

A document-Term matrix is a type of matrix representation for a corpus, where

- Rows represent documents in corpus.
- Columns represent terms found in corpus.
- Values in the matrix are frequencies or important measures, such as term frequency (count of a term in a document) or tf-idf (term frequency-inverse document frequency, which reflects importance).

StopWords:-

Stopwords are common words, like "the", "of", "and", "is", which frequently appear but don't add much meaningful information to analysis.

Sentiment Analysis:- (Opinion Mining)

- A method to evaluate the emotional tone of a text by identifying and measuring positive, negative or nuanced emotions (joy, anger etc).

Sentiment lexicons:-

- Sentiment lexicons are predefined lists of words that are associated with specific emotional or opinion-based meanings:-

(1) AFINN lexicon:- Score words on a scale from -5 (very negative) to +5 (very positive)

(2) Bing Lexicon Classifies words into binary categories of "positive" or "negative".

(3). NRC Lexicon:- Categorizes words into various emotion (like joy, anger, trust) as well as binary "positive" and "negative" categories.

Sentiment Scores and Calculations

- Unigram:- A single word, most lexicons use unigram-based sentiment scores.
- Sentiment Score:- The numeric or categorical value assigned to each word based on its emotional content (eg. -5 to +5 in affin).
- Net sentiment:- The overall sentiment for a section of text, often calculated as the difference b/w positive and negative sentiment score.

Chunking text for Sentiment Analysis:-

Chunking: Dividing text into manageable segments (e.g. by paragraph or sentences) for sentiment scoring, as overall sentiment can vary across the text.

Integer division: A technique used to define sections in the text by dividing ~~the~~ line numbers to create groups ~~sections of~~ ~~sentences~~.

One-Row-per-Token-Format:-

The structure for tidy text data where each token (e.g. word) occupies its own row in a table which makes the data compatible with tidy tools.

Chapter 3 :- Analyzing word and document frequency.

(1) Term Frequency:-

- Term frequency measures how often a word or "term", appears within a single document.
- It's calculated as the count of a word in a document divided by total number of words in that document.
- TF is useful for identifying the most frequent terms within a document but does not distinguish b/w common, less meaningful words and content specific words.

(2) Inverse Document Frequency:-

- Inverse document frequency (IDF) adjusts for words that are common across a collection of documents, giving a higher weight to

terms that are less frequent across the entire collection.

- IDF helps downweight common words and upweight unique words that might be more meaningful in certain documents.

- $\text{idf}(\text{term}) = \ln \left(\frac{\text{no. of documents}}{\text{no. of documents containing term}} \right)$

• 3. TF-IDF :

- TF-IDF is a metric that combines Term frequency and Inverse document frequency to highlight words that are important within individual documents but not too common across the document collection.

- It calculated as

$$\text{tf.idf}(\text{term}, \text{document}) = \text{tf}(\text{term}, \text{document}) * \text{idf}(\text{term})$$

H. Zipf's Law

- Zipf's Law describes a pattern where frequency of a word is inversely proportional to its rank, in terms of frequency.
- This law results in a "long-tailed" distribution where a few words are very common, while many words are used infrequently.
- In a graph of log-transformed term frequencies against ranks, Zipf's Law is observed as a line with negative slope, often around -1.

Chapter 4:- Relationship between words:-

n-grams and correlations:-

(i) N-grams:-

- N-grams are sequences of n consecutive words in a text, often used to capture the context in which words occur:-
 - Bigram - An n-gram of two words.
(eg, "happy day").
 - Trigram - An n-gram of three words.
(eg, "beautiful sunny day").
- N-grams allow us to understand common word pairs or sequences in a document, providing more context than analyzing individual words alone.

Important functions:-

(1) unnest_tokens() :-

A function from tidytext that converts a text column into tokens, creating a row for each token.

(2) anti_join() :-

A function in dplyr used to filter out stop words from text data by joining it with a list of stopwords and keeping only the non-matching terms.

(3). wordcloud(): (Wordcloud):-

- Visual display of word frequencies where the size of each word represents its frequency, often used to show prominent words by sentiment.

(previous chapter functions) ↑, from now functions used in chapter H.

(4) `separate()` and `unite()` from tidyverse

- `separate()` - splits a single column into multiple columns based on a delimiter.
- `unite()` - combines multiple columns into a single column.

(5) ~~gg~~ ggraph package:- (`ggraph()`)

- ggraph is an extension of ggplot2 designed for creating network plots, allowing us to visualize relationships b/w words (e.g. bigrams) as graphs.
- Network visualizations can show how words connect, helping us see clusters, common phrases or key phrases that structure a text.
- `geom_edge_link()`, `geom_node_point()`, `geom_node_text()`, `ggraph()`.

6. igraph package :- (graph_from_data_frame())

- igraph is a package used to create and manipulate graph structures.
- It helps prepare network data before visualizing it with ggplot.
- Graphs in igraph are connected with nodes (individual words) and edges (connections between words) with optional weights (such as bigram frequency).

7. pairwise_count() from wigr :-

- pairwise_count() counts the no. of times each pair of items (e.g. words) appear together within the same unit of text (e.g. a document or section). This is useful for exploring co-occurrences of words that appear within the same text sections, but not necessarily as adjacent words.

8. pairwise_cor() from widyr:

- pairwise.cor() calculates the phi coefficient or binary correlation b/w two items (e.g. words) to measure their likelihood of appearing together (~~relative~~) relative to how often they appear separately, this is helpful for identifying strong for identifying strong relationship/association within text, even if they are adjacent.
- phi-coefficient:- This is a measure of association for binary data.
 - It quantifies how much more likely two words are to co-occur, offering insights into significant word pairings within a corpus.
- phi coefficient is equivalent pearson correlation.

2. tf-idf for bigrams

- tf-idf can be applied to n-grams to identify contextually significant word pairs or phrases in document, similar to how it highlights important single words. This helps identify word pairs unique to specific documents or sections within a text collection.

3. Negation handling in sentiment analysis:-

- It is an approach to ~~handle~~ adjust sentiment analysis by identifying words that change the sentiment of following terms. (like "not", "never")
- By using bigrams we can capture phrases such as "not happy" to reverse or reverse or adjust the sentiment score, providing a more accurate sentiment analysis.

4. Markov chains:-

- A markov chain is a model where the occurrence of each word depends only on the preceding word , making it relevant ~~for~~ for bigram analysis.
- Network graphs built from bigrams can be seen as visualizations of a Markov chain , providing insights into sequences and transitions b/w words,

5. Conceptual Co-Occurrence Analysis:-

- Conceptual co-occurrence analysis uses pair of words that appear within the same section, chapter or document.
- This can reveal thematic patterns or topics within a text by identifying words that frequently appear together.

Chapter 5. Converting to and from non-tidy formats:-

1. Document-Term Matrix & Sparse Matrix:-

- DTM's are usually sparse matrices (most values are zeroes) because not all terms appear in all documents.
- Sparse matrix is a matrix with many zero values, used to represent DTM's efficiently. Rather than storing all values, sparse matrices save space by storing only non-zero values.

2. tidy()

- tidy() function from broom package converts non-tidy objects like DTM into a tidy dataformat.
- This allows users to use the tidy data

structure to conduct analysis, such as sentiment analysis or tf-IDF on datasets initially stored as DTM's.

3. last_* Functions:

- last_dtm() :- converts tidy data into a DTM object from tm package
- last_dfm() :- converts tidy data into a dfm (document feature matrix) from the quanteda package.
- last_sparse() :- converts tidy data into a sparse matrix, which is useful for ML libraries that require sparse matrices as input.
- last() turns a tidy one-term-per-row data frame into a matrix.

H. Quanteda package:-

- quanteda is a text mining package in R that focuses on quantitative analysis of textual data.
- It provides an alternative format for DTM known as dtm's.

5. tm package:-

- tm (Text mining) is another package in R for NLP that works with DTM and corpus objects.
- It's widely used for text processing, including tokenization and creating DTM's.

E. stopom package:-

6. Loughran and McDonald dictionary:-

- This dictionary is a sentiment lexicon developed for analyzing financial text.
- It categorizes words into sentiments like positive, negative, litigious, uncertain, containing and superfluous.
- It is designed to prevent false positive or negative readings in finance contexts where words like "growth" or "stares" might have neutral meanings.

Chapter 6: Topic Modelling

1. Topic Modelling:-

- This is an unsupervised learning technique used to discover the underlying themes or "topics" within a collection of documents.
- It automatically groups related words and documents into topics based on patterns in text data.
- It's useful for organizing large sets of documents, like news article, books or tweets, where each topic may contain multiple overlapping topics.

2. Latent Dirichlet Allocation (LDA) :-

LDA assumes:-

- Each document is a mixture of several topic's and.
- Each topic is a mixture of specific words.

Understanding topics and words:-

- In LDA, topics are generally probability distribution over words.
- This means each topic is represented by a set of words, where each word has the probability of belonging to that topic.
- Similarly, documents are probability distributions over topics.
- Each topic is thus made up of certain proportion of words from each topic.

Generative process:-

- LDA is a generative model, which means it assumes that the documents were generated based on probabilistic process. A LDA essentially tries to reverse this process to discover the topics:-

1. Choose the distribution of topics for each document.
2. Choose a distribution of words for each topic.
3. Generate words for each document.

parameters in LDA :-

LDA relies on 3 parameters :-

1. No. of topics (k) :- This is usually often chosen in advance based on the no. of distinct topics you expect to find.
 - For instance, if analyzing new articles, you might choose $k=10$ if you think there are 10 distinct subjects like politics, business and sports, etc.

2. Alpha (α) :- Controls the per document topic distribution.

- A higher alpha means documents will be

composed of more topics in similar proportions, while lower alpha means documents ~~are~~ will lean more heavily on few specific topics.

Beta (β):- Controls the per-topic word distribution

- A higher beta means topics will include a broader range of words, while lower beta means topics will be more concentrated with a smaller number of frequent words.

augment():

- The augment() function assigns model results to each observation in the original data.
- In topic modelling, augment() assigns each word in each document to a topic creating a dataset that shows which topics are most associated with specific words within document.

Words within

Twitter data project - Covid-19 tweets:-

1. Data prepossessing

* Text cleaning:-

- Remove URL's
- Remove User mention (@username)
- Remove Hashtags (#).
- Remove punctuation, special characters, and Numbers.
- Convert to lowercase.
- Remove emoji's.

* Tokenization:-

- Breakdown each tweet into individual words (tokens) for easier processing

* Stopword Removal

* Stemming and Lemmatization

* Handling Negation.

2. Feature Engineering:-

- * Bag-of-words (BoW)

- * TF-IDF

- * Word embeddings.

- use pretrained embedding's like Word2Vec,

- Glove or FastText

- BERT embeddings.

- * Sentiment Specific Features:-

- Emoji Counts.

- Exclamation marks Count.

- Sentiment lexicon scores.

3. EDA:-

- * Sentiment distribution.

- * Word cloud for the and are sweets.

- * Top hashtags.

- * Top n-grams.