

# Topic Modeling Report

**Introduction:** Topic modeling is a technique in natural language processing (NLP) that is used to discover abstract topics within a collection of documents. In this report, we applied topic modeling to an email dataset to identify key themes and topics present within the emails. The primary method used for topic modeling in this report is Latent Dirichlet Allocation (LDA).

## Approach

The overall approach consisted of the following steps:

### 1. Data Preprocessing:

**Loading Data:** The email data was loaded from an Excel file.

**Text Cleaning:** The email subjects were cleaned by removing stop words, punctuation, and non-alphabetic characters. Lemmatization was applied to convert words to their base forms using spaCy.

**Tokenization:** The cleaned subjects were split into lists of words (tokens).

### 2. Bag-of-Words Representation:

- A dictionary was created from the tokenized subjects.
- Each document was converted into a bag-of-words (BoW) representation.

### 3. Topic Modeling with LDA:

- An LDA model was trained on the corpus with a predefined number of topics (10 in this case).
- Parameters for the LDA model were fine-tuned to improve the coherence and interpretability of the topics.

### 4. Topic Interpretation and Labeling:

- The top words for each topic were extracted and manually reviewed.
- Based on the top words, each topic was assigned a descriptive label to reflect the underlying theme.

## Techniques Used

**Lemmatization:** Using spaCy, words were reduced to their base forms, which helps in standardizing the text and reducing dimensionality.

**Latent Dirichlet Allocation (LDA):** This probabilistic model was used to discover topics within the text. It assumes that documents are mixtures of topics and that topics are mixtures of words.

**Manual Topic Labeling:** After extracting the top words for each topic, a human-readable label was assigned to each topic based on the most frequent and relevant words.

## Why LDA?

### 1. Proven Effectiveness:

- LDA is a well-established method for topic modeling and has been widely used in both academia and industry. It is known for its ability to identify coherent and interpretable topics from large text corpora.

### 2. Probabilistic Approach:

- LDA uses a probabilistic approach to model the relationship between words, documents, and topics. This allows for a soft assignment of words to topics and documents to topics, providing a richer understanding of the data compared to hard clustering methods.

### 3. Flexibility:

- LDA is flexible and can handle different sizes of datasets effectively. It works well for a wide range of text sizes and can adapt to various types of textual data.

### 4. Interpretability:

- The topics generated by LDA are generally interpretable, as they are composed of sets of words that co-occur frequently. This makes it easier to label and understand the underlying themes in the data.

### 5. Implementation and Tuning:

- LDA is relatively straightforward to implement using libraries like Gensim, and it offers several parameters (e.g., number of topics, alpha, beta) that can be tuned to optimize performance for specific datasets.

**Conclusion:** The LDA model effectively identified distinct topics within the email dataset. Each topic was characterized by its top words, and based on these words, descriptive labels were assigned to provide a clear understanding of the themes. This topic modeling approach can be valuable for categorizing and summarizing large collections of text data, aiding in better information retrieval and data analysis. By implementing such techniques, organizations can gain insights into the key issues and themes present in their communication data, enabling them to address customer concerns more effectively and optimize their processes.