

* NLP Pipeline *

* Break the problem down into sub. problems.
then try to develop a step-by-step procedure
to solve them.

① Data Acquisition →

- Available Data (CSV, Txt, PDF, XLS)
- Other Data (DB, Internal, API, scrapping)
- No data (create your own data).

Note → If you have less data then use
Data Augmentation.

→ Replace with synonyms.

Ex :- I am a Data scientist.
I am a AI Engineer.

→ Bigram Flip

Ex :- I am mamta.

Mamta is my name.

→ Back Translate

→ Add additional noise →

Ex :- I am a data scientist.

I love this job. (extra)

② Text Preparation →

① cleanup : HTML tags, emoji, spelling correction

② Basic preprocessing

③ Advance preprocessing.

① Basic Preprocessing →
Tokenization → Sentence
→ Word

Steps: →

① Stop words removal

② Stemming

③ Lemmatization

④ Removing punctuation (., !, ?, \$)

⑤ Lower case

⑥ Language Detection

② Advance Preprocessing :-

① Parts of speech tagging.

② Parsing

③ coreference resolution.

③ Feature Engineering →

① Text Vectorization

② TFIDF, Bag of words, OneHot,
Word2Vec, Encoding

④ Modeling →

(i) Heuristic

(ii) ML

(iii) DL

(iv) Cloud API (AWS, GCP, Azure)

⑤ Evaluation →

(i) Intrinsic (ML - classification test, accuracy, precision, confusion Matrix)

(ii) Extrinsic (How app to productivity uses, slow settings)

(iii) Deployment → (Monitoring, Networking)

* Co-occurrence Vectors →

The role of co-Occurrence Vectors →

① Statistical Insight →

Find insights about word frequency and distribution with co-occurrence vectors.

② Language Modeling → Improve your understanding of the language structure and underlying semantics.

③ Network Analysis → Uncover relationships between words and concepts with network analysis.

Co-Occurrence Vectors are an essential element of modern NLP and have numerous applications in improving text data analysis.

With the continued application of new techniques and methods, their use will only continue to grow.

* Constructing Co-Occurrence Matrices : →

① Definition → Understand the definition and mathematical representation of co-occurrence matrices.

② Window Size → choose an appropriate window size to capture the context of the words.

③ Contextual Weighting →

Explore different weighting schemes to emphasise the context that is most informative.

* Application of Co-Occurrence Vectors in NLP →

① Mobile Apps → Discover how co-occurrence vectors are used in NLP mobile apps like predictive text.

② Customer Support → Learn how customer support teams leverage co-occurrence vectors to identify themes and automate responses.

③ Book Recommendation Systems →

Find out how to build a book recommendation engine using co-occurrence vectors to capture user preferences.

* Doc2Vec →

- A novel technique for computing vector representations of entire documents.
- **How does it work?**
Doc2Vec learns vector representations of documents by predicting words within the document.
- **Why is it important?**
- Doc2Vec enables us to analyse documents as continuous vectors in multi-dimensional space, opening up new possibilities for next.

* Benefits of Doc2Vec in NLP ⇒

① Improved Text Classification →

Doc2Vec can be used to transform documents into continuous vectors that can be used as input to machine learning models for text classification.

② Efficient Document Retrieval →

Doc2Vec enables us to search for documents based on semantic similarity, rather than just keyword-based matching, leading to more accurate and efficient document retrieval.

③ Robust Information Extraction →

→ Doc2Vec can be used to extract meaningful information from unstructured text data, enabling more accurate and efficient information extraction for NLP tasks.

* Amplifying NLP With Doc2vec: Real-world Applications

→ ① Question Answering Systems →

Doc2vec can be used to create question answering systems that can automatically answer questions based on the semantic content of the document.

② Clinical NLP →

Doc2vec can be used to analyse and extract information from clinical reports, enabling more efficient and accurate analysis of patient medical data.

③ Sentiment Analysis →

Doc2vec can be used to analyse the sentiment of large volumes of text data, enabling more accurate and efficient sentiment analysis for market research and social media monitoring.

What are Contextualized Document Representations?

Ans → Contextualized document representations are document vectors that capture not only the meaning of the words in the document, but also their contextual relationships with other words in the document.

* ELMO : → (A Powerful Contextualized Encoder)

ELMO is a state-of-the-art contextual encoder that can generate contextualized embeddings for words and sentences, enabling more accurate and efficient NLP tasks.

BERT: (A Transformer-based Encoder)

BERT is a powerful transformer-based encoder that can generate contextualized embeddings for words and sentences, enabling highly accurate NLP tasks such as text classification and question answering.

* Unraveling the challenges of Implementing

① Data Preprocessing challenges → Doc2Vec

Preparing the large volume of text data for Doc2vec training can be challenging, as it requires extensive cleaning, normalization, and tokenization.

② Hyperparameter Tuning →

→ Finding the optimal combination of hyperparameters for the Doc2vec model can be challenging, as it requires extensive experimentation and Tuning.

③ Model Complexity →

The Doc2vec model can be computationally intensive to train, especially when dealing with very large volumes of text data.



* Text Blob → (Simplified Text Processing)

- Text Blob is a Python library for processing
- It provides a simple API for textual data.
common (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, classification, translation, and more.

* Text Blob is an exceptional library that makes NLP tasks a breeze. Its features allow for comprehensive text analysis with minimal code, making it a must-have for today's NLP workflow.

* Tokenization ⇒

Text Blob's tokenizer splits documents into sentences or words efficiently. Customize your tokenizer to suit the task at hand.

* Word Frequencies ⇒

Find out which words appear the most frequently. Create bar charts and word clouds to visualize your findings.

* N-gram Creation ⇒

Create word n-grams to uncover important patterns in your document. Text Blob allows for trigram and bigram generation.

Stop Words →

Remove noise from your documents. Text Blob makes it easy to remove stopwords and focus only on important words.

* Language Translation →

Translate your text into another language with ease. Carve out your niche in global markets easily and effectively.

* Language Detection →

Automatically detect which language your content is written in. Great for multi-language websites and document repositories.

* NLTK →

Natural Language Toolkit (NLTK) is a powerful open-source toolkit used to process human language data. It provides a range of functions that can help with text processing, manipulation and analysis.

NLTK's features →

- NLTK has a wide range of features that are useful for various language processing tasks.
- It also offers an extensive corpus of text data, ranging from Shakespeare's plays to the Brown Corpus of American English.

* Applications of NLTK →

- NLTK has been used in a vast array of applications, from chatbots and machine translation to sentiment analysis and spam detection.
- It is also widely used in academic research, as it provides a reliable and standardized toolkit for processing and analyzing human language data.

* Tokenization and Text Preprocessing with NLTK

* Tokenization → Tokenization refers to breaking up text into smaller units or tokens, such as words or phrases.

NLTK provides several tokenization tools, such as the `word_tokenize()` function and the `RegexpTokenizers` class.

* Text Preprocessing Practical

① How to perform text Preprocessing.

- HTML tags remove
- Punctuation Remove
- Emoji Handle
- Lemmatization
- Stemming
- Lower Case

② Text Representation & Word Embedding

- One Hot Encoding (\downarrow Text Vectorization)
- BOW (Bag of words)
- TF-IDF
- Word2Vec (convert text to vector (numbers))

③ Project \rightarrow Sentiment Analysis using ML

- ML Algorithm used. Approach.

Gensim →

Gensim is a Python library for topic modeling, document similarity analysis, and word embeddings. It provides a user-friendly interface for natural language processing tasks.

→ It is an open-source library for unsupervised semantic modeling of text.

* Purpose and Uses →

Gensim allows users to analyze and extract meaning from large collections of text documents.

(1) Word Embeddings →

Gensim provides tools for training and using word embeddings, such as Word2vec and FastText.

(2) Topic Modeling →

With Gensim, it's easy to discover latent topics in a collection of documents using techniques like Latent Dirichlet Allocation (LDA).

(3) Document Similarity →

Gensim allows you to measure the similarity between documents, which is useful for tasks like recommender systems.

Latent Dirichlet Allocation (LDA) →

Gensim supports LDA for discovering topics in a corpus and associating them with relevant documents.

Latent Semantic Analysis (LSA) →

LSA is another technique supported by Gensim for finding hidden semantic structures in a document collection.