# *Words unleashed: python-powered text analysis adventure*

**[1]G. Sai Chand, [2]K. Karthik Reddy**
[1,2]Student, Jayaprakash Narayan College of Engineering

**[3]M. Bharathi, [4]T. Aditya Sai Srinivas**
[3,4]Assistant Professor
Jayaprakash Narayan College of Engineering

*Corresponding Author: T. Aditya Sai Srinivas*
*Email: taditya1033@gmail.com*

## Abstract:

Text Analysis, encompassing techniques such as Text Mining and Natural Language Processing (NLP), is a transformative method for deriving valuable insights from unstructured text, spanning documents, emails, social media, and customer reviews. This article serves as a gateway to Text Analysis proficiency, focusing on Python as the primary tool. Offering a comprehensive guide, it navigates readers through the intricacies of Text Analysis, equipping them with the skills to extract meaningful information from diverse textual sources. Ideal for those eager to delve into the realm of linguistic data, this article promises a hands-on exploration of Text Analysis, empowering readers to unravel the potential of Python in this domain.

## Keywords:

Text Analysis, Text Mining, Natural Language Processing (NLP), unstructured text data, insights.

# 1. Introduction:

In an era defined by an unprecedented abundance of textual information, the ability to derive meaningful insights from unstructured data has become a paramount skill. Text Analysis, a dynamic field encompassing methodologies such as Text Mining and Natural Language Processing (NLP), stands as the key to unraveling the latent narratives embedded within vast corpora of text. This article embarks on a journey into the realm of Text Analysis, delving into the techniques that empower us to distill valuable information from diverse textual sources. As we navigate this landscape, Python emerges as a powerful and versatile ally, enabling a hands-on exploration of the intricacies involved in extracting knowledge from the written word. Join us on this enlightening exploration where language meets analytics, and the untapped potential of unstructured text data is revealed.

# 2. Dataset overview: text analysis on articles and titles:

The dataset under consideration is tailored for comprehensive text analysis, a field encompassing text mining, natural language processing (NLP), and information extraction from unstructured textual data. This dataset includes articles accompanied by corresponding titles, offering a rich source of information on diverse topics within the realms of data analysis, machine learning, and related fields.

## 2.1. Key features:

1) Text Content: Each article contains extensive text, delving into various aspects of data analysis, machine learning, and related domains.

2) Title Summaries: Titles serve as succinct summaries, encapsulating the main subject or theme of the respective articles.

Overall Goal: The primary goal is to conduct a thorough and insightful text analysis on this dataset, unraveling patterns, sentiments, and latent topics. The outcomes aim to enhance understanding and contribute valuable insights to the domains covered in the articles.

# 3. Text analysis:

Initiating Text Analysis with Python: Commence the task by importing essential Python libraries and loading the dataset. Install essential Python libraries for data analysis and text processing using the command: `pip install pandas nltk wordcloud matplotlib scikit-learn`.
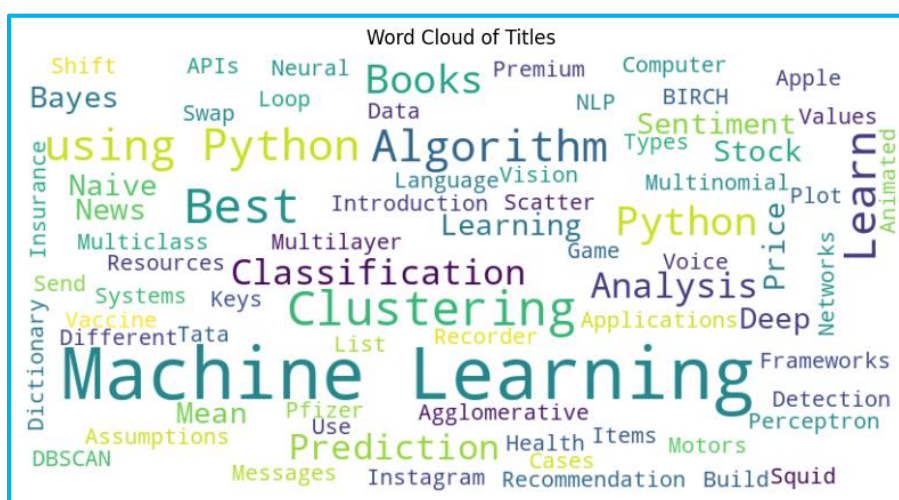
```python
# Import necessary libraries
import pandas as pd
import plotly.express as px
from wordcloud import WordCloud
import matplotlib.pyplot as plt
from textblob import TextBlob
import spacy
from collections import defaultdict
from collections import Counter
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.decomposition import LatentDirichletAllocation
# Load English language model for spaCy
nlp = spacy.load('en_core_web_sm')
# Read the CSV file into a pandas DataFrame
data = pd.read_csv("articles.csv", encoding='latin-1')
# Display the first few rows of the dataset
print(data.head())
```

```
                                             Article  \
0  Data analysis is the process of inspecting and...
1  The performance of a machine learning algorith...
2  You must have seen the news divided into categ...
3  When there are only two classes in a classific...
4  The Multinomial Naive Bayes is one of the vari...

                                               Title
0              Best Books to Learn Data Analysis
1         Assumptions of Machine Learning Algorithms
2            News Classification with Machine Learning
3  Multiclass Classification Algorithms in Machin...
4        Multinomial Naive Bayes in Machine Learning
```
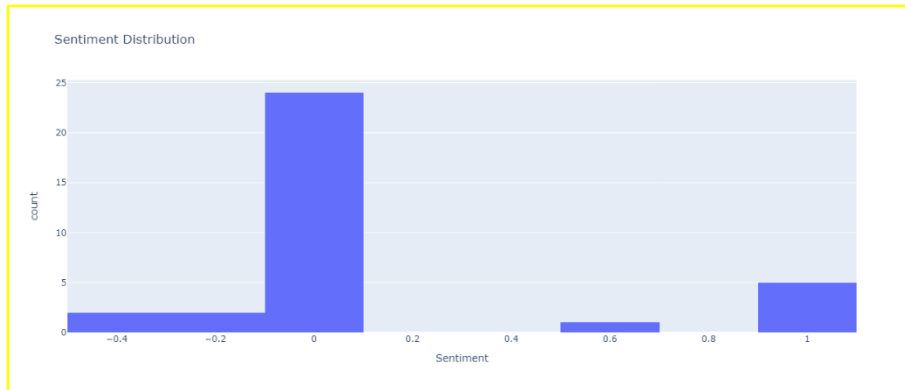
Now, proceed to visualize title word clouds.



Word Cloud of Titles

Here, a word cloud visually encapsulates article titles. The process starts by merging individual titles into a cohesive string, `titles_text`, using the `join` method. Next, a `WordCloud` object is created, incorporating parameters like width, height, and background color to influence
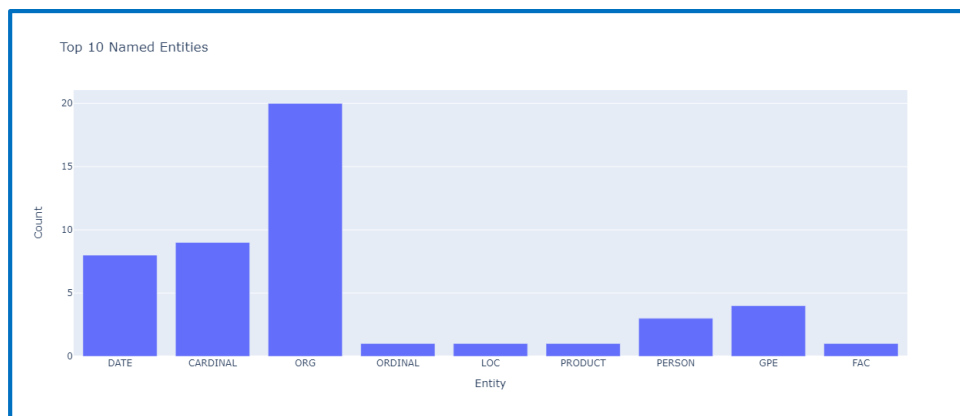
visual aesthetics. This object is then utilized to generate the word cloud, where word size correlates with frequency. This visualization highlights the most prevalent terms within the titles, offering a concise yet impactful representation of the overarching themes present in the articles.

## 3.1. Next, analyze the sentiment distribution in the data:



In sentiment analysis, we utilize the TextBlob library to evaluate emotional tones in the dataset's articles. Calculating sentiment polarity using numerical scores, positive values convey positivity, negatives indicate negativity, and scores near zero suggest neutrality. Subsequently, a histogram visually illustrates the distribution of sentiment scores, providing insights into the overall sentiment patterns within the dataset.
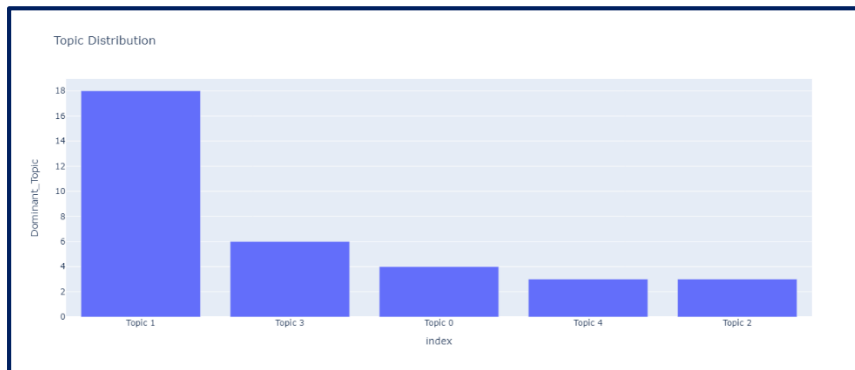
## 3.2. Next, initiate named entity recognition (NER):



This phase focuses on executing Named Entity Recognition (NER), a crucial natural language processing technique for identifying entities like organizations, locations, names, and dates from text. The `extract_named_entities` function, powered by spaCy, systematically analyzes each article, categorizing entities based on labels such as "ORG" for organizations and "LOC" for locations. The NER outcomes are stored in a new 'Named_Entities' dataset column. A concise visualization showcases the top 10 entities and their counts, providing a swift overview

of the predominant entities within the textual data. This graphical representation enhances comprehension of key entities extracted through NER.

### 3.3. Next, initiate topic modeling:



In Topic Modeling, we employ Latent Dirichlet Allocation (LDA) to unveil latent topics in text documents. We begin by vectorizing text using CountVectorizer, converting it into a numerical format. Parameters, including document frequencies, feature limits, and stopwords, are specified. LDA is then applied with five topics. The resulting matrix shows article distribution across topics. Dominant topics are assigned based on the highest probability, visualized in a bar chart. This method offers a comprehensive overview of prevalent themes in articles, showcasing proficient Text Analysis with Python.

## 4. Conclusion:

Text Analysis journey using Python unraveled profound insights from a dataset of articles and titles. We traversed the realms of sentiment analysis, named entity recognition, and topic modeling. Through word clouds, sentiment histograms, and visualizing prominent named entities, we gained a nuanced understanding of the dataset's linguistic landscape. The culmination of Latent Dirichlet Allocation offered a glimpse into latent topics, revealing prevalent themes. This comprehensive exploration showcases the power of Python in deciphering and extracting valuable information from textual data, underscoring its significance in the broader domain of data science and analysis.

## 5. References:

(1)     https://www.datacamp.com/tutorial/text-analytics-beginners-nltk

(2)     https://thecleverprogrammer.com/2023/10/02/text-analysis-using-python/

(3)    https://www.geeksforgeeks.org/text-analysis-in-python-3/

(4)    https://towardsdatascience.com/getting-started-with-text-analysis-in-python-ca13590eb4f7

(5)    https://towardsdatascience.com/getting-started-with-text-analysis-in-python-ca13590eb4f7