

# **Artificial Intelligence and Machine Learning**

## **Project Report**

**Semester-IV (Batch-2022)**

## **Text Summariser**



**Supervised By:**

Dr. Jatin Arora

**Submitted By:**

Divya Bhoria 2210990293

Geetansh Sood 2210990323

Ojas Gupta 2210990627

**Department of Computer Science and Engineering  
Chitkara University Institute of Engineering & Technology,  
Chitkara University, Punjab**

## Abstract

The exponential growth of digital news content has led to significant challenges in information management, creating a pressing need for efficient summarization tools. This project focuses on developing a machine learning-based text summarizer to effectively condense news articles from major news websites. The primary goal is to address the issue of information overload, which affects both individuals and organizations, by providing a solution that enables quick and easy access to essential information.

Our summarizer utilizes advanced natural language processing (NLP) techniques to automatically generate concise and coherent summaries of lengthy news articles. The project involved collecting and processing a large dataset of archived news articles, which served as the basis for training our machine learning model. The model was designed to identify key points and main ideas within the articles, ensuring that the summaries retain the most important information while discarding extraneous details.

Through rigorous testing and evaluation, the summarizer demonstrated a significant reduction in the time required to comprehend news content, thereby enhancing productivity and decision-making capabilities for users. The results indicate that our approach not only improves information retrieval but also facilitates better retention and understanding of the summarized content.

The implications of this project are far-reaching, with potential applications in various fields such as journalism, academia, finance, and policy-making, where timely and accurate information is crucial. By leveraging the power of machine learning and NLP, this project contributes to the broader field of AI-driven information management solutions, paving the way for future advancements in automated content summarization.

In conclusion, this research underscores the potential of machine learning technologies in transforming how we manage and consume vast amounts of information. The developed text summarizer offers a practical and efficient tool for navigating the digital information landscape, ultimately helping users to stay informed and make better-informed decisions in an increasingly data-saturated world.

# CONTENTS

S.NO.	TABLE CONTENTS
1.	INTRODUCTION
2.	PROBLEM DEFINITION AND REQUIREMENTS
3.	PROPOSED DIAGRAM AND DEFINITION
4.	RESULTS

## Introduction

In today's digital age, the relentless flow of information from various news sources has become both a boon and a burden. While access to a vast array of news articles enables individuals and organizations to stay informed about current events and trends, the sheer volume of content can be overwhelming. The challenge lies in efficiently sifting through this information to extract relevant and important insights without investing excessive time and effort. This issue of information overload is particularly acute for professionals, researchers, and decision-makers who rely on timely and accurate data to guide their actions. To address this pressing need, our project aims to develop a machine learning-based text summarizer capable of automatically condensing lengthy news articles into concise and coherent summaries. By leveraging advanced natural language processing (NLP) techniques, this tool seeks to enhance productivity, improve information management, and support better decision-making in a world inundated with data.

## Background

The digital revolution has dramatically transformed the way news is produced, distributed, and consumed. With the advent of the internet and mobile technologies, news organizations publish a continuous stream of articles, covering an extensive range of topics. Major news websites update their content multiple times a day, contributing to an overwhelming abundance of information. This phenomenon, known as information overload, presents a significant challenge to individuals and organizations alike. For professionals in fields such as finance, marketing, policy-making, and academia, staying informed about the latest developments is crucial. However, the traditional method of manually reading through numerous articles is not only time-consuming but also inefficient.

Moreover, organizations that archive news articles for reference, analysis, or research purposes face difficulties in managing and utilizing this vast repository of data. The valuable insights embedded within these articles often remain untapped due to the impracticality of manually sifting through such large volumes of text. This situation highlights the urgent need for automated solutions that can streamline the process of information extraction and utilization.

Advances in machine learning and natural language processing (NLP) have opened new possibilities for tackling these challenges. Machine learning algorithms can be trained to understand and process human language, enabling the development of tools that can automatically summarize large bodies of text. NLP techniques allow for the extraction of key points and essential information from lengthy articles, creating concise summaries that retain the core message and important details.

This project leverages these technological advancements to create a machine learning-based text summarizer tailored to the needs of users overwhelmed by the influx of news articles. By transforming how news content is consumed, our summarizer aims to make it easier for

individuals and organizations to stay informed, make timely decisions, and efficiently manage their information resources. The development of this tool is not only a response to the current challenges posed by information overload but also a step towards more intelligent and automated information management solutions for the future.

## Objectives

**Develop an Automated Summarization Tool:** Create a machine learning-based text summarizer capable of automatically generating concise summaries of lengthy news articles from major news websites.

**Improve Efficiency in Information Consumption:** Reduce the time and effort required for individuals and professionals to stay informed by providing quick and accurate summaries of news articles.

**Enhance Information Management:** Facilitate better organization and utilization of archived news articles by transforming large volumes of text into manageable summaries.

**Support Decision-Making:** Provide timely and relevant summaries that aid decision-makers in fields such as finance, marketing, policy-making, and academia, ensuring they have access to critical information without delay.

**Leverage Advanced NLP Techniques:** Utilize state-of-the-art natural language processing methods to accurately identify and extract key points and main ideas from news articles.

**Ensure Coherent and Accurate Summaries:** Develop algorithms that generate summaries which are not only concise but also maintain the coherence and accuracy of the original content.

**Test and Validate the Summarizer:** Conduct rigorous testing and evaluation to ensure the effectiveness and reliability of the summarizer in various real-world scenarios and with different types of news content.

**Contribute to Technological Advancements:** Advance the field of machine learning and NLP by applying these technologies to the practical problem of news summarization, demonstrating their potential and effectiveness.

**Provide a User-Friendly Solution:** Design the summarizer to be easily accessible and user-friendly, catering to the needs of both individual users and organizations.

**Promote Future Research and Development:** Lay the groundwork for future innovations in automated content summarization and intelligent information management solutions.

## Significance

The development of a machine learning-based text summarizer addresses the pressing issue of information overload, providing a crucial tool for navigating the vast amounts of news content produced daily. By generating concise summaries, the summarizer significantly enhances productivity for professionals, researchers, and decision-makers, allowing them to access critical information quickly and efficiently. This tool supports better decision-making by ensuring that users have timely and accurate information at their fingertips, which is particularly valuable in fields such as finance, marketing, policy-making, and academia. Moreover, the summarizer improves information management by transforming large repositories of archived news articles into manageable and useful summaries, facilitating better organization, retrieval, and utilization of valuable data.

The project also advances the technological frontiers of machine learning and natural language processing, demonstrating their practical applications and encouraging further research and development. By reducing the cognitive load on users and helping them absorb and retain important information more effectively, the summarizer enhances overall comprehension and reduces cognitive fatigue. Additionally, the tool's ability to deliver personalized content aligned with individual preferences improves user engagement and satisfaction.

This project exemplifies the practical applications of artificial intelligence in solving everyday challenges, promoting wider adoption of AI-driven solutions across various sectors. Furthermore, the successful development of this summarizer can inspire new approaches to information management and content consumption, fostering innovation in the field. Ultimately, by addressing current challenges and setting a precedent for effective summarization tools, this project lays a solid foundation for future advancements in automated content summarization and intelligent information systems.

## Problem Statement

In the digital age, the constant stream of news articles overwhelms individuals and organizations, leading to information overload. Manually summarizing these articles is time-consuming and inefficient, hindering productivity and decision-making. Despite extensive news archives, the lack of an automated summarization tool results in underutilized data and cognitive fatigue. Thus, there's a critical need for a machine learning-based text summarizer to swiftly condense news articles, addressing information overload and enhancing information management and usability.

## Software Requirements

1. **Python (Version 3.6 or higher):** Python is required as the primary programming language for developing the text summarization tool.
2. **NLTK (Natural Language Toolkit):** NLTK will be used for natural language processing tasks such as tokenization, stemming, and lemmatization.
3. **TextBlob:** TextBlob, a Python library built on top of NLTK and Pattern libraries, will be utilized for sentiment analysis and other text processing functionalities.
4. **Newspaper3k:** The Newspaper3k library will be used for web scraping news articles from major news websites for training and testing the summarization model.
5. **Pandas:** Pandas, a powerful data manipulation library, will be used for data preprocessing, analysis, and manipulation tasks.
6. **NumPy:** NumPy, a fundamental package for scientific computing with Python, will be used for numerical computations and data manipulation tasks.
7. **Matplotlib:** Matplotlib, a comprehensive library for creating static, interactive, and animated visualizations in Python, will be used for data visualization tasks.
8. **Seaborn:** Seaborn, a Python visualization library based on Matplotlib, will be used for creating visually appealing and informative statistical graphics.
9. **Plotly Express:** Plotly Express, a high-level interface for creating expressive and interactive visualizations, will be used for advanced data visualization tasks, such as interactive plots and dashboards.

These software requirements will provide the necessary tools and libraries for developing a robust and effective text summarization application.

## Hardware Requirements

The hardware requirements for running the GitHub profile viewer are relatively modest, as the primary processing will be handled on the client-side and the backend server:

1. Computer or Server:
2. Internet Connection

## Proposed Design / Methodology

### 1. Data Collection:

- Scrape news articles using the Newspaper3k library.
- Clean the data by removing HTML tags and non-textual content.

### 2. Text Preprocessing:

- Tokenize text into words or sentences with NLTK or spaCy.
- Remove stop words and punctuation, and apply stemming or lemmatization.

### 3. Feature Extraction:

- Extract features such as word frequency, TF-IDF, or word embeddings.

### 4. Model Development:

- Implement and train extractive (e.g., TextRank) or abstractive (e.g., seq2seq with attention) summarization models.

### 5. Evaluation:

- Use ROUGE or BLEU scores to evaluate model performance.
- Refine models based on evaluation results and user feedback.

### 6. Deployment:

- Package the tool into an executable or web service for deployment.

### 7. Testing and Validation:



- Perform thorough testing to ensure functionality and reliability.
- Validate with real users to gather feedback for improvements.

#### 8. Documentation and Maintenance:

- Document the tool's design, implementation, and usage.
- Provide updates and maintenance for bug fixes and enhancements.

#### 9. Scalability and Extensibility:

- Ensure the tool is scalable and extensible for future enhancements and larger datasets.

## **Schematic structure:**

### **1. Data Loading and Preprocessing:**

- The process begins by loading the data from a CSV file containing URLs of news articles using the ``Load_data`` class.
- Unnecessary columns such as "Id", "Author", and "Date" are dropped from the DataFrame to clean the data and prepare it for analysis.

### **2. Sentiment Analysis:**

- The ``sentiment_EDA()`` method of the ``Load_data`` class performs sentiment analysis on the news articles.
- For each article, the sentiment polarity and subjectivity scores are calculated using the TextBlob library.
- The results are stored in a DataFrame, which is then used for exploratory data analysis (EDA).

### **3. Exploratory Data Analysis (EDA):**

- EDA is performed to gain insights into the sentiment distribution of the news articles.
- Two types of plots are generated:
  - Scatter plot: Polarity vs Subjectivity to visualize the relationship between sentiment polarity and subjectivity.
  - Histograms: Distribution of polarity and subjectivity scores to understand their frequency distribution.

#### **4. Visualization:**

- Matplotlib is used to create visualizations for EDA, including scatter plots and histograms.
- These visualizations provide a graphical representation of the sentiment analysis results, enabling easy interpretation and analysis.

#### **5. Individual Article Analysis:**

- The `analyze_data()` method of the `Load_data` class analyzes individual news articles.
- For each article, the title, summary, and sentiment analysis results (polarity and subjectivity) are printed.
- This analysis provides insights into the sentiment of each article and helps understand the overall sentiment distribution across the dataset.

#### **6. Integration and Execution:**

- The code snippets for data loading, preprocessing, sentiment analysis, EDA, and individual article analysis are integrated into a single script.
- The script is executed to perform the entire analysis pipeline, from loading the data to visualizing the sentiment distribution and analyzing individual articles.

Overall, the schematic structure outlines the sequential flow of operations, starting from data loading and preprocessing to sentiment analysis, EDA, and individual article analysis. It demonstrates a systematic approach to analyzing sentiment in news articles, facilitating data-driven insights and decision-making.

## Algorithms Used:

Algorithm Type: Lexicon-based sentiment analysis algorithm.

Library Used: TextBlob library in Python.

Sentiment Analysis Process:

1. Text Preprocessing: The algorithm preprocesses the text data, including tokenization and removal of stopwords and punctuation.
2. Lexicon-based Scoring: Each word in the text is assigned a pre-defined sentiment score based on a lexicon or dictionary of sentiment words.
3. Calculation of Sentiment Polarity: The sentiment polarity score of the text is computed by aggregating the sentiment scores of individual words and normalizing the result. The polarity score ranges from -1 (negative sentiment) to 1 (positive sentiment), with 0 representing neutral sentiment.
4. Calculation of Subjectivity: The subjectivity score of the text is calculated by averaging the subjectivity scores of individual words. The subjectivity score ranges from 0 (objective/factual) to 1 (subjective/opinionated).

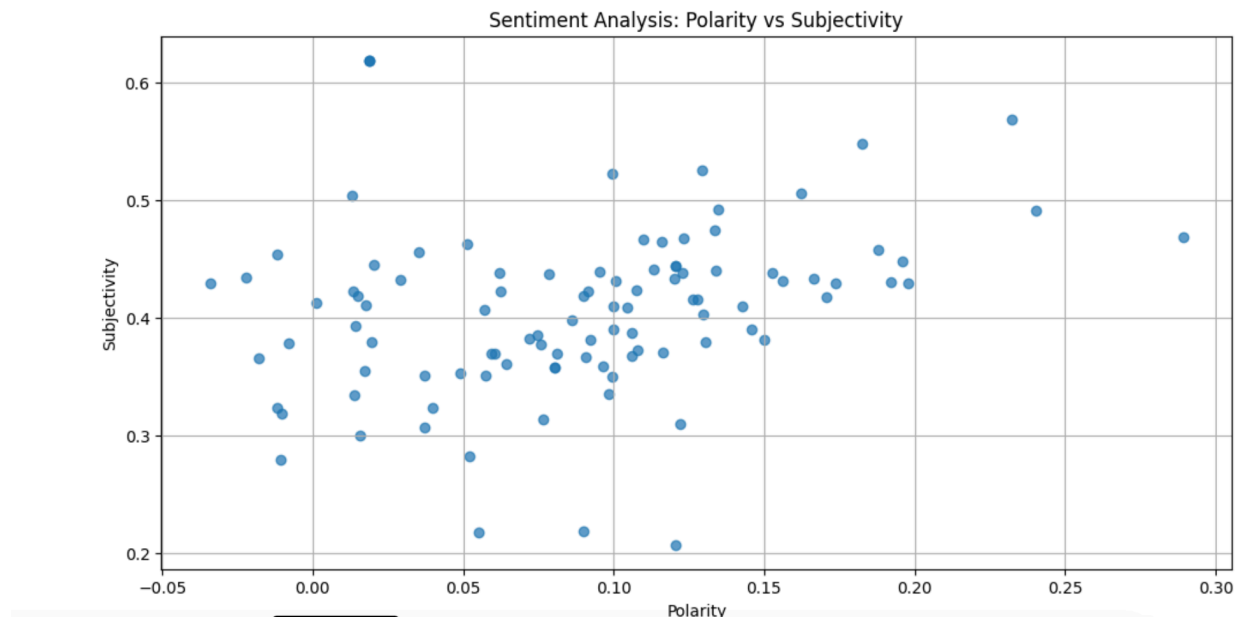
Output: The algorithm generates sentiment analysis results, including the sentiment polarity and subjectivity scores for each news article.

Purpose: The algorithm is utilized for exploratory data analysis (EDA) to analyze the sentiment distribution across the dataset, visualize sentiment trends, and gain insights into the overall sentiment of the news articles.

## Summary:

The sentiment analysis algorithm utilized for the analysis of news articles employs a lexicon-based approach facilitated by the TextBlob library in Python. Through this method, the algorithm systematically preprocesses the text data, assigning sentiment scores to individual words based on pre-defined lexicons or dictionaries. By aggregating these scores, the algorithm computes sentiment polarity and subjectivity scores for each article. This process enables the algorithm to generate comprehensive sentiment analysis results, offering insights into the emotional tone and subjective nature of the news content. Despite its simplicity and efficiency, the algorithm's reliance on pre-defined lexicons may limit its ability to capture nuanced sentiments accurately. However, it remains a valuable tool for conducting exploratory data analysis, visualizing sentiment trends, and gaining a deeper understanding of the overall sentiment landscape within the dataset.

## Result

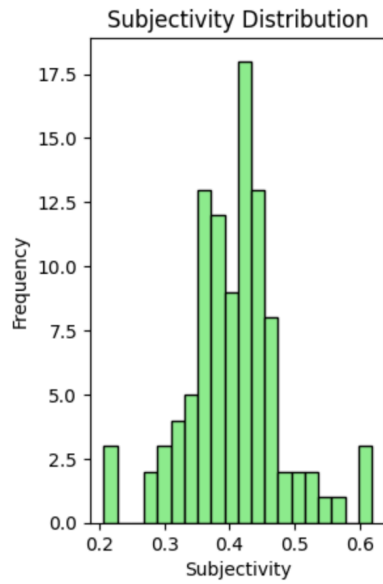


The sentiment analysis graph shows the comparison between polarity and subjectivity , which tells us how biased or unbiased , positive or negative the news article is . News articles have been denoted as blue markers with polarity being on the x-axis and subjectivity being on the y-axis.

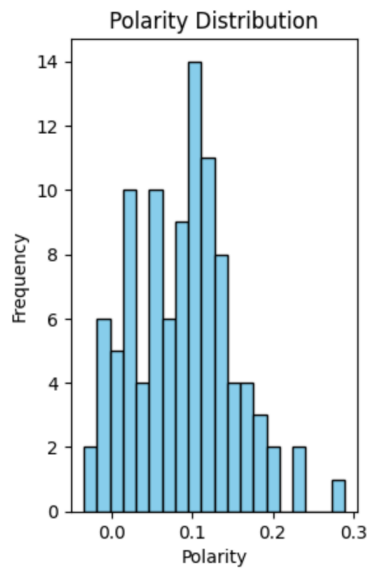
```
In [50]: obj.analyze_data()

Title: Biden is up against nostalgia for Trump's first term
Summary: More than half, 55%, of all Americans say they see Trump's presidency as a success, while 44% see it as a failure.
Four years ago, he was able to assail Trump's time in office from his position as a challenger.
Biden won all three in 2020 after Trump won them in 2016 in his victory over Democratic nominee Hillary Clinton.
In the poll, 92% of Republicans view Trump's time in office a success, while just 73% of Democrats say Biden's has been a success.
And while 85% of Democrats polled say they back Biden, 91% of Republicans say they support Trump.
-----
```

Here the title and summarised news articles have been shown , it is summarised in nearly 60-70 words each . All 101 news articles will be represented as shown above.



This is a subjectivity distribution graph which compares the biases of an article with its frequency.



This is a polarity distribution graph which compares the positivity/negativity of an article with its frequency.

```

-----
Sentiment(polarity=0.01494917826024697, subjectivity=0.4184030325060858)
0.01494917826024697
Polarity: positive
Subjectivity: Biased And Less Factual
-----

```

Here we can see the polarity and subjectivity numeric values and character values.

Out[6]:

	Id	Author	Date	URL
0	1	Stephen Collinson	April 29, 2024	<a href="https://edition.cnn.com/2024/04/28/politics/bi...">https://edition.cnn.com/2024/04/28/politics/bi...</a>
1	2	John Towfighi	April 28, 2024	<a href="https://edition.cnn.com/2023/10/11/business/ha...">https://edition.cnn.com/2023/10/11/business/ha...</a>
2	3	Andrew Carey and Olga Voitovych	April 28, 2024	<a href="https://edition.cnn.com/2024/04/28/europe/russ...">https://edition.cnn.com/2024/04/28/europe/russ...</a>
3	4	Andrew Carey and Olga Voitovych	April 28, 2024	<a href="https://edition.cnn.com/2024/03/13/travel/aust...">https://edition.cnn.com/2024/03/13/travel/aust...</a>
4	5	Silvia Marchetti	April 28, 2024	<a href="https://edition.cnn.com/travel/italy-house-bou...">https://edition.cnn.com/travel/italy-house-bou...</a>


ID , author , date and url have been listed after retrieving the following from the CNN dataset articles .

Out[46]:

	URL
0	<a href="https://edition.cnn.com/2024/04/28/politics/bi...">https://edition.cnn.com/2024/04/28/politics/bi...</a>
1	<a href="https://edition.cnn.com/2023/10/11/business/ha...">https://edition.cnn.com/2023/10/11/business/ha...</a>
2	<a href="https://edition.cnn.com/2024/04/28/europe/russ...">https://edition.cnn.com/2024/04/28/europe/russ...</a>
3	<a href="https://edition.cnn.com/2024/03/13/travel/aust...">https://edition.cnn.com/2024/03/13/travel/aust...</a>
4	<a href="https://edition.cnn.com/travel/italy-house-bou...">https://edition.cnn.com/travel/italy-house-bou...</a>

URL has been displayed after dropping the id , author and date .

## Reference

 [Summarize News Articles with Machine Learning in Python](#)