

# English News Paper Analysis

This project focuses on news classification and analysis using Natural Language Processing (NLP) techniques. Leveraging a dataset from inshorts.com, the objective is to categorize news articles into predefined topics and extract actionable insights for efficient information retrieval and trend identification.

# Data Description

- Data Components:
  - Headline: Title of the news article.
  - Content: Textual content of the news article.
  - Category: Topic or subject matter of the news article.
  - Date: Date of scraping the news article.
- Source: Dataset sourced from inshorts.com, containing 199,706 English news articles.

# Background

- Dataset Origin: Scraped news articles from inshorts.com.
- Significance: Analyzing news articles aids in understanding public discourse and tracking emerging trends.
- Importance of NLP: Utilizing NLP techniques for efficient news classification and analysis.

# Objective

- Project Objective: Perform news classification and analysis using NLP techniques.
- Key Goal: Automate news categorization and extract insights for efficient information retrieval.
- Expected Impact: Facilitating trend identification and analysis in the rapidly evolving news landscape.

# Key Components

## **Data Collection and Preprocessing:**

- Preprocess textual content by removing noise and tokenizing for analysis.

## **Topic Classification:**

- Develop NLP models for categorizing news articles into predefined topics.

## **Model Training and Evaluation:**

- Train and evaluate classification models using appropriate metrics.

## **Topic Analysis and Insights:**

- Analyze classified news articles for insights into topic distribution and emerging trends.

## **Visualization and Reporting:**

- Create informative representations and generate comprehensive reports summarizing findings.

# NLP Experiments and Dataset Takeaways

- **Preprocessed Dataset:**
  - Contains 199,706 English news articles.
  - Preprocessing steps included lowercasing, removing noise, and tokenization.
- **Trained Classification Models:**
  - Multinomial Naive Bayes with Bag of Words achieved 0.89 accuracy.
  - Hyperparameter tuning improved accuracy to 0.893.
- **Evaluation Results:**
  - Best model correctly predicted 35,671 samples.

# Key Insights and Recommendations

- **Insights:**
  - Prevalent news categories: politics, education, sports, business, entertainment, and technology.
- **Recommendations:**
  - Explore other feature representations like TF-IDF or word embeddings.
  - Investigate alternative machine learning models for potentially better performance.
  - Analyze misclassified samples to identify patterns for improvement.
  - Incorporate additional features like publication source or author.
  - Develop a system to monitor changes in news categories over time.

# Conclusion

- Project Impact: Enhancing news classification and analysis through NLP techniques.
- Future Directions: Continued optimization and exploration of advanced methodologies for improved performance.
- Significance: Efficient dispute resolution and trend identification in the dynamic news landscape.