# FlipItNews: Natural Language Processing Business Case

❖ Topic: Natural Language Processing in News Article Categorization
❖ Duration: 1 week

---

## Why this case study?

### From the company's perspective:

- FlipItNews is transforming the Indian finance, business, and investment landscape using AI and ML.
- This case focuses on harnessing natural language processing to categorize news articles, a critical component for delivering relevant content to users.
- By analyzing its internal database of news articles, FlipItNews aims to intelligently categorize content into domains like politics, technology, sports, business, and entertainment.
- The insights from this analysis can significantly improve content discovery and user engagement, driving financial literacy and investment awareness among millennials and first-time investors.

### From the learner's perspective:

- This case offers a chance to delve into the burgeoning field of NLP, addressing a real-world problem in the media and content industry.
- NLP techniques, including text processing, stopwords removal, tokenization, lemmatization, and multi-class classification, are essential for understanding and categorizing textual data.
- Participants will gain hands-on experience in employing NLP methodologies and comparing different machine-learning models for text classification.
- The case promotes the development of practical, data-driven strategies, enhancing skills crucial for a career in data science and AI.

---

# Dataset Explanation: FlipItNews Data

Each row in this dataset corresponds to a unique news article, and each column represents features crucial for categorizing the content. The dataset 'FlipItNews Data' includes the following variables:

1. Article: The full text of the news article.
2. Category: The actual category of the news article (such as politics, technology, sports, business, and entertainment).

These features are self-explanatory yet pivotal for understanding and categorizing news content. The 'Article' feature requires extensive NLP processing to extract meaningful insights and patterns, while the 'Category' serves as the target variable for classification.

*Note: The primary objective of this study is to categorize each news article into its correct category based on the content using NLP techniques.*

_____


# What is Expected?

As a data scientist at FlipItNews, your task is to analyze the dataset of news articles to categorize them into their respective categories like politics, technology, sports, business, and entertainment. Your primary goal is to use natural language processing to create and compare at least three different models for this multi-class classification problem.

## Submission Process:

Once you've completed the case study...

- Compile your findings and the entire process in a Jupyter Notebook.

- In the notebook, ensure that you:

- ○ Display the Python code for all your analysis, model building, and evaluation.
- ○ Include visualizations such as frequency distributions, confusion matrices, and any other relevant plots that support your findings.
- ○ Provide valuable insights derived from the data analysis, and suggest actionable strategies for FlipItNews to enhance their content categorization process.
- Convert your Jupyter Notebook into a PDF (Save as PDF using the Chrome browser's Print command).

- Upload the PDF on the designated platform as per the submission guidelines.

*Note: Once submitted, you won't be able to edit your submission.*

## General Guidelines:

- Approach this as a real-world scenario, akin to the challenges faced by data scientists in the media and content industry.
- It's normal to encounter difficulties:
  - ○ Revisit the problem statement to ensure your analysis is aligned with the objectives.
  - ○ Break down complex tasks such as text preprocessing, model training, and evaluation into smaller, manageable segments.
  - ○ Utilize online resources, documentation, and forums for problem-solving – crucial skills in the field of data science and NLP.
  - ○ Collaborate with peers in discussion forums for different perspectives and solutions.
  - ○ Review course materials or external resources for a deeper understanding of NLP concepts and techniques.
  - ○ If you're stuck or need clarification on any aspect of the case study, reach out to your instructor or mentor for guidance.

_____

# What does 'good' look like?

## 1. Define the Problem Statement and perform Exploratory Data Analysis

|  | Hint | Approach |
|---|---|---|
| a. Definition of problem | Start with a clear objective. Why is categorizing news articles crucial for FlipItNews? | The goal is to accurately categorize news articles into various categories like politics, technology, sports, business, and entertainment to enhance user engagement and content relevancy. |
| b. Observations on Data | Gain a comprehensive understanding of the dataset's structure. | Examine the shape of the data, data types, presence of missing values, and statistical summary.<br><br>Use functions like data.info(), data.describe(), and data.shape in Python. Since the data primarily consists of text, focus on textual data characteristics. |
| c. Univariate Analysis | Analyze the 'Category' variable (distribution of news articles across categories). | Use bar plots or count plots to understand the distribution of articles across different categories. |
| d. Text Data Analysis | Examine a few news articles to understand the nature of the content. | Look for common themes, styles, or patterns in the text. Identify the average length of articles, common words, or unique terms. |
| e. Illustrate the insights based on EDA | Each analysis, whether it's the distribution of categories or text characteristics, should offer insights | Note patterns like predominant themes in certain categories, or the complexity of language used.<br><br>For example, "Bar plot shows a higher frequency of business-related articles, suggesting a focus on finance |

| | | | |
|---|---|---|---|
| | | and investment news." |
| f. | Comments on Text Data | Provide observations on text data such as common keywords in each category, diversity of topics within categories, or any unique textual patterns noticed. | Each observation should be accompanied by a brief comment, enhancing the understanding of the dataset's nature. |

## 2. Data Preprocessing

*Note: While there are several standard techniques, the approach should not be limited to these. Instead, consider these as starting points, and feel free to explore additional methods that might enhance your model*

| | | Hint | Approach |
|---|---|---|---|
| a. | Duplicate value check | In text data, exact duplicates are less common but still possible. Check for duplicates to ensure the integrity of the dataset. | Use Python's pandas library to identify and remove any duplicate news articles. |
| b. | Missing value treatment | Missing values in text data can significantly impact the analysis. | Identify if any articles or categories are missing and decide on a strategy like deletion or imputation. |
| c. | Text Data Preprocessing | Preprocessing text data is crucial for NLP tasks. | **Possible Steps:**<br><br>a. Remove non-letters and stopwords to reduce noise in the text.<br><br>b. Tokenize the text to convert sentences into words.<br><br>c. Perform lemmatization to bring words to their base or dictionary form. |
| d. | Feature engineering | Transforming text data into a format suitable for modeling is a key part of NLP. | **Possible Steps:**<br><br>a. Convert the textual data into numerical values using |

| | | techniques like Bag of Words or TF-IDF.

b. Consider creating features that capture the length of articles, frequency of certain words, or sentiment of the text. |
|---|---|---|
| e. Data preparation for modeling | The prepared text data needs to be split into training and testing sets. | Use train-test split ensuring a good representation of each category in both sets.

Apply suitable text vectorization techniques. Choose between Bag of Words and TF-IDF based on your model's requirement. |
| f. Identifying Patterns in Text | Look for patterns or distributions in the text data. | Analyze the frequency distribution of words or categories. Comment on any skewness or notable trends in the text data. |

## 3. Model building

| | Hint | Approach |
|---|---|---|
| a. Build the Neural Network Regression Model | Start by choosing the right model architecture for text classification. | Prepare your data for training and validation. Use train-test split to divide the dataset.

Fit models on the training data. Consider using models like Naive Bayes, Decision Trees, or Random Forests for a start. |
| b. Hyperparameter Tuning | Fine-tuning model parameters can significantly enhance performance in NLP tasks. | Experiment with different hyperparameters such as the number of trees in Random Forest or depth in Decision Trees.

Employ techniques like GridSearchCV to systematically |

| | | explore hyperparameter combinations. |
|---|---|---|
| c. Handling Text Data Specifics | Text data comes with unique challenges like feature representation and dimensionality. | Choose text representation methods (Bag of Words, TF-IDF) wisely based on the model and data characteristics.<br><br>Regularize your model to avoid overfitting, especially in high-dimensional text data. |
| d. Incorporating Advanced Models | After establishing a baseline with initial models, explore advanced NLP models. | **Possible Approaches:**<br><br>a. Experiment with custom transformer models or utilize pre-trained models like BERT (Bidirectional Encoder Representations from Transformers).<br><br>b. Explore how these advanced models can capture contextual nuances in the text better than traditional models.<br><br>**Scope for Experimentation:**<br><br>a. Tweak and fine-tune these advanced models to suit your specific dataset and classification goals.<br><br>b. Investigate the use of transfer learning techniques to adapt these models to your specific NLP task. |
| e. Model Evaluation Metrics | Select appropriate metrics for evaluating classification models. | Use metrics such as accuracy, precision, recall, and F1-score to evaluate the performance of your models. |
| f. Multiple Model Comparison | Compare different models to identify which performs best for the given text data. | Evaluate the performance of each model using the same metrics for a fair comparison.<br><br>Consider not just accuracy, but also how well the model handles different categories, |

| | | especially less represented ones. |
|---|---|---|

## 4. Results Interpretation & Stakeholder Presentation

| | **Hint** | **Approach** |
|---|---|---|
| a. Understand the Business Context | Deeply understand FlipItNews' goals and challenges in the context of content categorization and user engagement. | Frame your findings to align with key performance indicators such as user engagement, content relevance, and diversification of news categories. |
| b. Interpreting Model Results | Focus on how the models categorize different news articles and their accuracy in doing so. | Discuss how the model's categorization aligns with actual news content and the potential impact on user experience and content discovery. |
| c. Visual Representations | Use visualization tools to present your model's performance. | Create a confusion matrix, accuracy graphs, and category-wise performance charts to illustrate key findings and make the insights clear to stakeholders. |
| d. Trade-off Analysis | Discuss the balance between model accuracy and computational resources, or between generalization and overfitting. | Highlight any trade-offs made during feature engineering, model selection, or in choosing the text representation method. |
| e. Recommendations | Offer actionable strategies based on your findings, like enhancing the content categorization algorithm or diversifying news categories. | Use evidence from your analysis to support these recommendations, demonstrating how they can improve content relevance and user satisfaction. |
| f. Feedback Loop | Propose methods for continuous improvement of the categorization models. | Discuss the importance of regularly updating the model with new data and feedback to adapt to changing news |

| | | trends and user preferences.<br><br>Suggest a framework for incorporating user feedback and new content into the model for continuous enhancement. |
| --- | --- | --- |

_____

## Questionnaire (Answers should present in the text editor along with insights):

1. How many news articles are present in the dataset that we have?
2. Most of the news articles are from _____ category.
3. Only ____ no. of articles belong to the 'Technology' category.
4. What are Stop Words and why should they be removed from the text data?
5. Explain the difference between Stemming and lematization.
6. Which of the techniques Bag of Words or TF-IDF is considered to be more efficient than the other?
7. What's the shape of train & test data sets after performing a 75:25 split?
8. Which of the following is found to be the best-performing model..
   a. Random Forest b. Nearest Neighbors c. Naive Bayes

9. According to this particular use case, both precision and recall are equally important. (T/F)

_____