# National University
## of computer and emerging sciences

# Data Visualization and Analysis

## Instructor: Ma'am Eisha tur Razia

## Project Report

## Section: BDS 5B

| Roll No | Name |
|---------|------|
| 22L-7478 | Abdullah Maqsood |
| 22L-7465 | Omar Bandial |
| 22L-7541 | Eman Atiq |
| 22L-7555 | Ameer Tufail |

# Comparative Linguistic Analysis on Narratives of Mainstream International Media Outlets

## Introduction

News articles play a vital role in shaping public opinions and understanding of global events. They are complex and often carry subtle biases, either openly or in hidden ways. These biases and the tone of reporting can influence how events are perceived by readers. Articles use nuanced language, specific word choices, and framing techniques that may shape narratives in ways that are not immediately obvious. Understanding and analyzing these elements require a careful and systematic approach.

This project aimed to analyze the sentiments expressed in news articles from five major global media outlets: BBC, CNN, TRT, Al Jazeera, and Fox News. These sources were chosen to ensure a wide range of perspectives, representing different regions and viewpoints. The project focused on five significant topics: the Israel War, Ukraine War, Islamophobia, US Presidential Elections, and China News. These topics were selected due to their global relevance and the likelihood of generating diverse media coverage.

A dataset of over 5,000 articles was collected by scraping these news sources. The dataset included features such as the article source, headline, description, timestamp, author, region, and full article content. These features provided rich context for each article, allowing for a deeper exploration of how sentiment is conveyed across topics and sources.

The project used advanced data processing and machine learning techniques to classify the sentiment of articles into Positive, Negative, or Neutral categories. Beyond simply classifying sentiment, the project sought to identify key patterns in the language used, uncover keywords influencing sentiment, and explore how sentiments vary across sources and topics. Tools like SHAP were used to make the analysis more transparent and explainable, helping to show why certain articles were classified with specific sentiments.

The significance of this work lies in its ability to uncover the hidden layers of media narratives. Different news outlets may report the same event in ways that reflect their own cultural, political, or ideological perspectives. By combining data analysis with visualizations, this project provides a clear and structured way to understand these dynamics, empowering readers to approach media with a more critical and informed perspective.

## Data Collection through Web Scraping

The project began with collecting news articles from five major online media sources: BBC, CNN, TRT, Al Jazeera, and Fox News. The primary objective was to gather a diverse dataset that represented a variety of perspectives and regional viewpoints. Web scraping was conducted using Python and libraries such as **BeautifulSoup** and **Selenium**. The scraping process targeted multiple attributes for each article, including:

- **Source**: The media outlet that published the article.
- **Headline**: The title of the article, often reflecting the main focus or angle of reporting.
- **Description**: A brief summary or introductory text about the article content.
- **Timestamp**: The publication date and time, useful for tracking trends over time.
- **Author**: Where available, the journalist or author of the article.
- **Region**: Geographical information linked to the event or article focus.
- **Full Content**: The complete text of the article, which served as the primary input for sentiment analysis.

Over 5,000 articles were successfully gathered. The diversity of sources ensured a dataset that captured a wide range of reporting styles, tones, and biases.

## Data Preparation and Transformation

Raw data scraped from the web often contains inconsistencies and noise. Preparing the data for analysis was a crucial step to ensure accuracy and reliability. This phase involved several cleaning and transformation tasks, including:

- **Removing duplicates**: Articles with identical content or metadata were eliminated to avoid redundancy.
- **Standardizing formats**: Text fields such as dates and headlines were standardized for uniformity.
- **Handling missing values**: Missing data in fields like author or region was either imputed or marked as unavailable.
- **Feature creation**: New columns, such as word counts or sentiment placeholders, were added to support further analysis.

This structured dataset ensured consistency and quality, forming a strong foundation for the subsequent steps.

## Sentiment Classification Model

To determine the sentiment of each article, a **transformers-based sentiment analysis pipeline** was employed. The model was pre-trained on a vast corpus of textual data and fine-tuned using Twitter sentiment analysis parameters to better understand shorter and contextually rich text. The sentiment model categorized each article into one of three classes:

- **Positive**: Articles conveying optimism, praise, or favorable tones.
- **Negative**: Articles reflecting criticism, pessimism, or adverse tones.
- **Neutral**: Articles that maintained a balanced or impartial tone.

The sentiment classification added a new column to the dataset labeled "Sentiment." This enabled a detailed exploration of how sentiment was distributed across topics, sources, and time.

## Keyword Extraction

Understanding why an article was classified with a specific sentiment required transparency in the model's decisions. For this, **SHAP (SHapley Additive exPlanations)** was utilized. SHAP identified key words or phrases within each article that significantly influenced its sentiment classification. For example:

- In a **Positive** article, words like "success," "peace," or "cooperation" might play a dominant role.
- In a **Negative** article, terms such as "violence," "failure," or "controversy" might be influential.
- **Neutral** articles often contained descriptive or factual terms that lacked emotional weight.

This step provided deeper insights into the patterns of language usage that shaped sentiment classification and helped identify potential biases in how topics were discussed.

## Machine Learning Model Testing

To validate the sentiment classification process, traditional machine learning models were tested on the same dataset. These included:

- **Naive Bayes**: A probabilistic model well-suited for text classification tasks.
- **Logistic Regression**: A simple yet effective algorithm for binary and multiclass classification problems.
- **Decision Tree**: A model that splits data based on feature importance to make predictions.

Each model was trained using the sentiment labels generated by the transformers pipeline, and their performance was evaluated using metrics such as accuracy, precision, and recall. Comparing these models ensured that the chosen method was both reliable and robust.


## Visualizations

To make the analysis understandable and engaging, data visualization was a critical step. Various tools, including **Seaborn**, **Plotly**, **Streamlit**, and **Power BI**, were used to create informative and interactive visualizations:

- **Bar Charts**: Displayed sentiment distribution across topics and media outlets.
- **Pie Charts**: Illustrated the proportions of Positive, Negative, and Neutral sentiments within the dataset.
- **Word Clouds**: Highlighted frequently occurring keywords in articles associated with each sentiment.
- **Line Charts**: Showed trends in sentiment over time, revealing how media tone changed during significant events.
- **Bivariate Dot Charts**: Explored relationships between variables, such as region and sentiment.
- **Map Charts**: Visualized geographic patterns in the dataset, showing where specific topics were more prominently covered.

These visualizations were not only valuable for analysis but also played a key role in presenting findings to stakeholders.

## Results and Outcomes

The analysis produced several important findings:

- Sentiments varied significantly by topic and source. For instance, articles on the **US Presidential Elections** had a higher proportion of Negative sentiments, while **China News** showed a more balanced distribution.
- Keyword analysis revealed recurring terms that shaped sentiments, providing transparency into the classification process.
- Media outlets differed in their tone and focus when covering the same topics, reflecting potential biases.
- Sentiment trends over time showed shifts in tone, especially during major developments in global events.

These insights were deployed through an interactive **Streamlit application** and a polished **Power BI dashboard**, making them accessible to a broader audience.

## Future Improvements

The project can be further enhanced by expanding the dataset to include additional sources and languages.
Incorporating advanced natural language processing techniques and real-time analysis could also provide
richer insights. Additionally, exploring context analysis and framing in articles could offer a deeper
understanding of media narratives.

## Conclusion

This project demonstrated a comprehensive approach to analyzing sentiments in news articles, offering valuable insights into how media outlets frame and communicate major global topics. By collecting over 5,000 articles from diverse sources, the analysis captured variations in reporting styles, tones, and potential biases across outlets like BBC, CNN, TRT, Al Jazeera, and Fox News. The thorough data cleaning and structuring ensured the dataset's quality, forming a solid foundation for analysis.
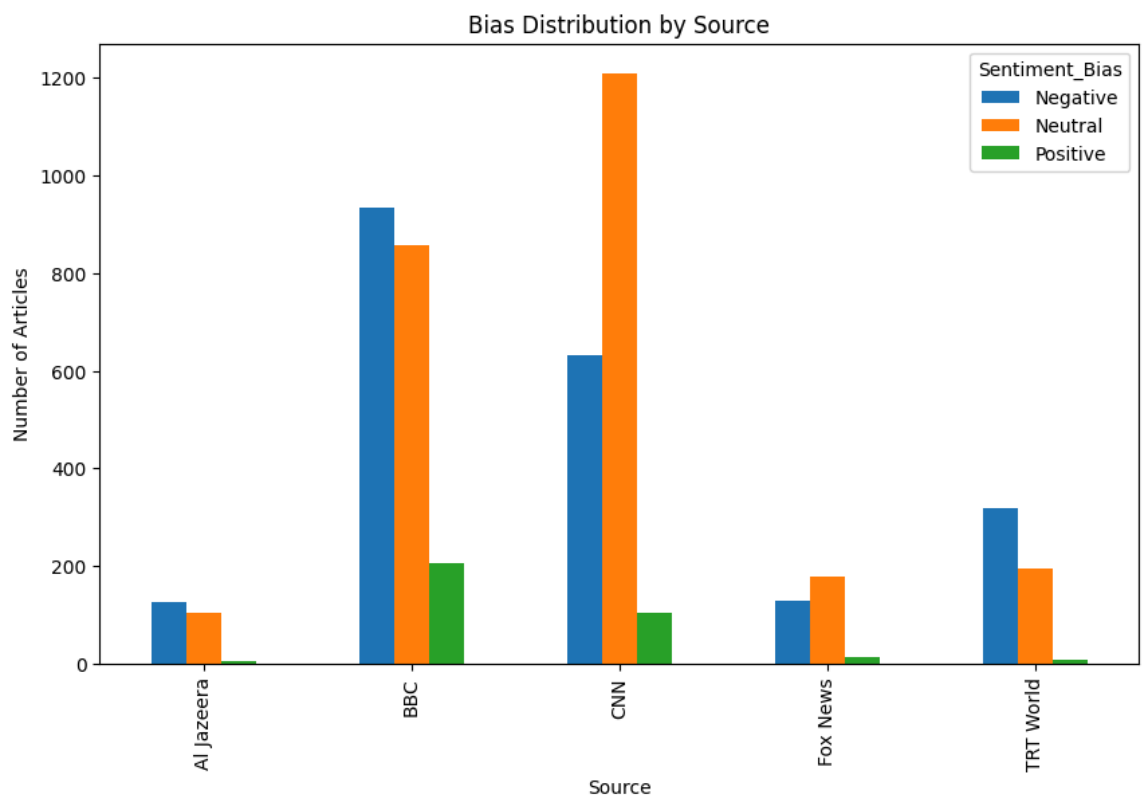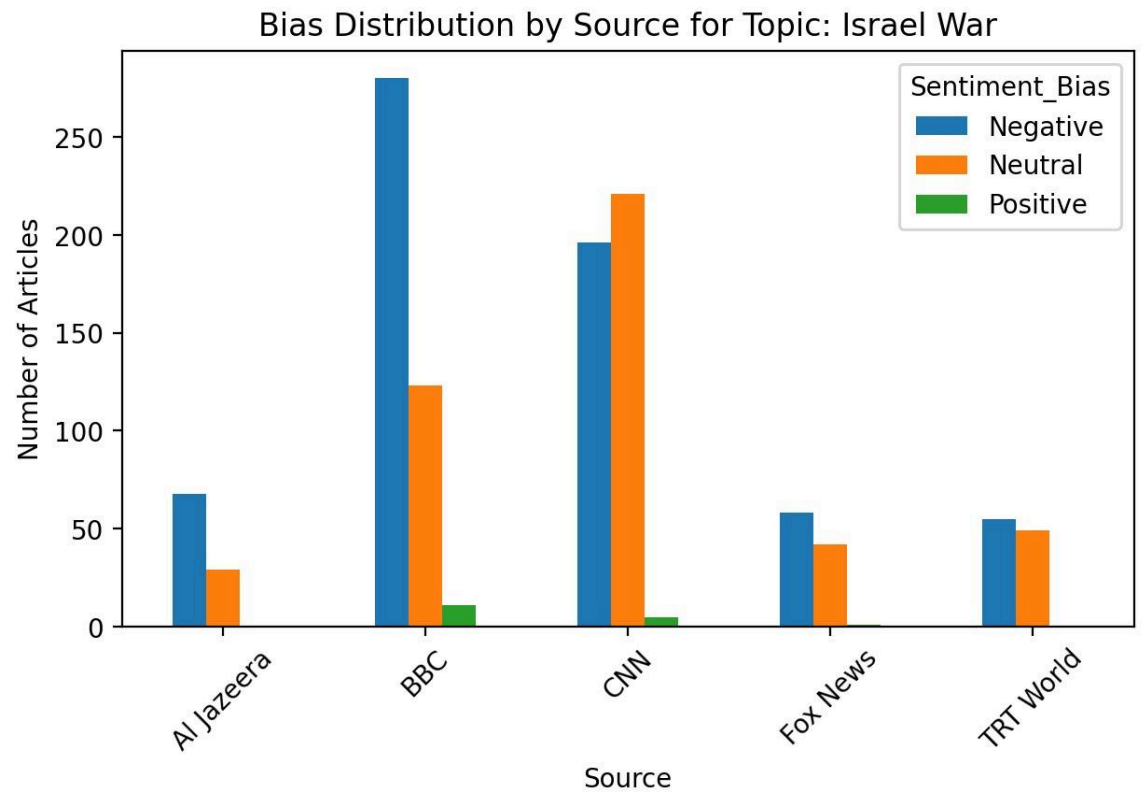
The sentiment classification, powered by a transformers-based model, accurately labeled articles as Positive, Negative, or Neutral. SHAP further enhanced transparency by identifying keywords influencing sentiment, providing a deeper understanding of language patterns. Comparing traditional machine learning models like Naive Bayes and Logistic Regression validated the robustness of the classification approach.
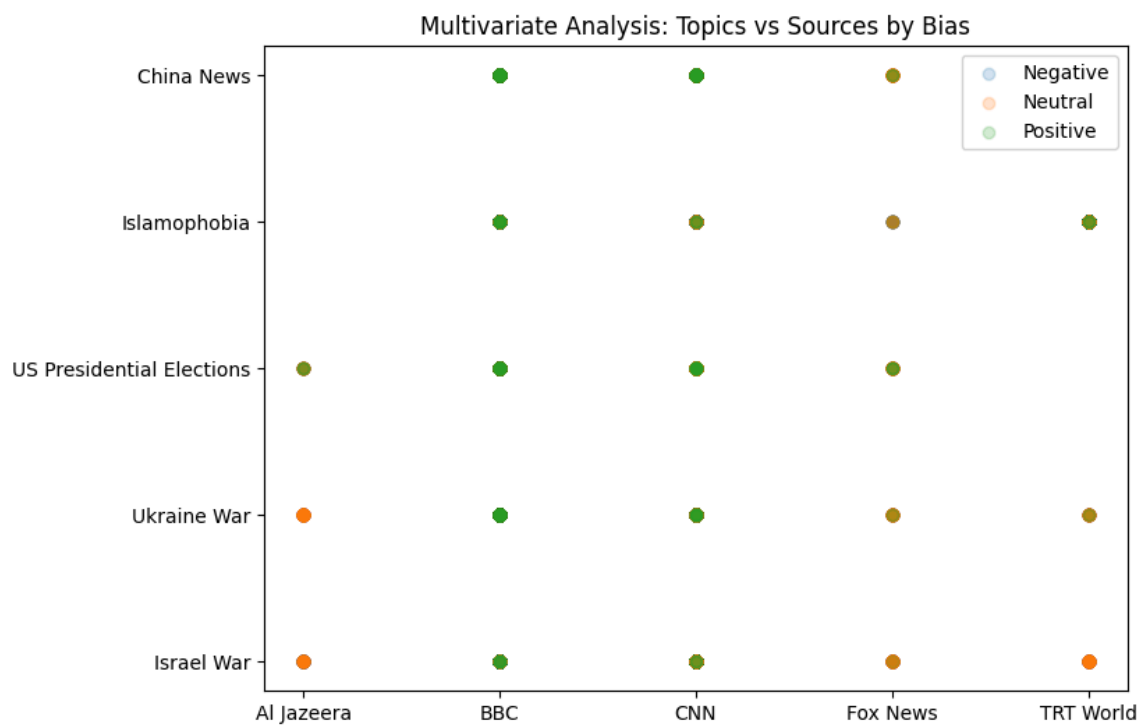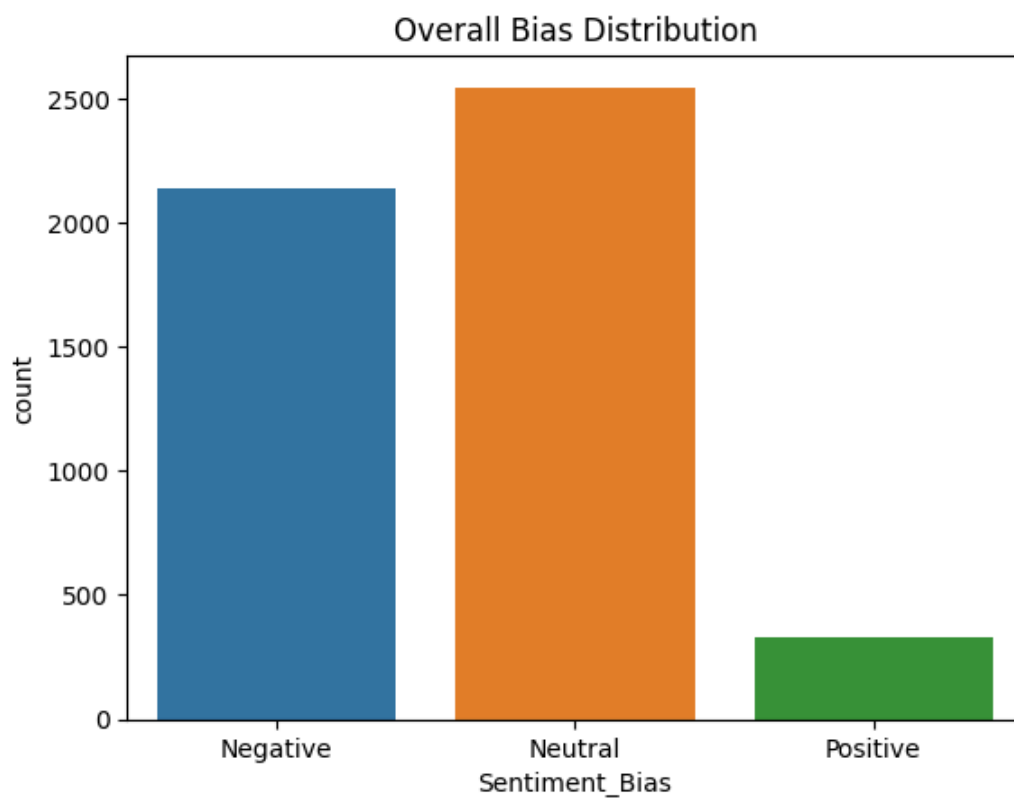
Visualizations created using Seaborn, Plotly, Streamlit, and Power BI brought the findings to life, showcasing sentiment distributions, keyword trends, and regional patterns. Interactive tools made the insights accessible and actionable for deeper exploration. Key findings highlighted significant differences in how media outlets report on topics such as the Israel War and US Presidential Elections, revealing shifts in tone and framing across sources.

This project underscores the importance of critically engaging with media content to understand biases and nuanced language. Future enhancements could include expanding data coverage, incorporating real-time analysis, and exploring advanced techniques to uncover hidden biases.

In conclusion, this project successfully provided a framework for sentiment analysis and media narrative exploration, highlighting the role of media in shaping perceptions. With further development, it has the potential to become a powerful tool for monitoring and analyzing media reporting on a global scale.

# Visualizations

## Bias Distribution by Source for Topic: Israel War



## Bias Distribution by Source

Overall Bias Distribution

Multivariate Analysis: Topics vs Sources by Bias

## Heatmap: Topics vs Sources

| Topic | Al Jazeera | BBC | CNN | Fox News | TRT World |
|---|---|---|---|---|---|
| China News | 0 | 4.2e+02 | 4e+02 | 67 | 0 |
| Islamophobia | 0 | 3.6e+02 | 2.6e+02 | 23 | 3.7e+02 |
| Israel War | 97 | 4.1e+02 | 4.2e+02 | 1e+02 | 1e+02 |
| US Presidential Elections | 38 | 3.7e+02 | 4.1e+02 | 68 | 0 |
| Ukraine War | 99 | 4.3e+02 | 4.5e+02 | 60 | 48 |

## Word Cloud of Articles Keywords

# Streamlit App Dashboard