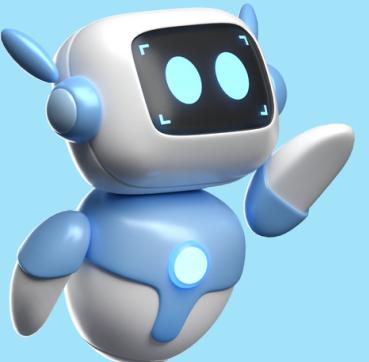


# FAKE NEWS DETECTOR

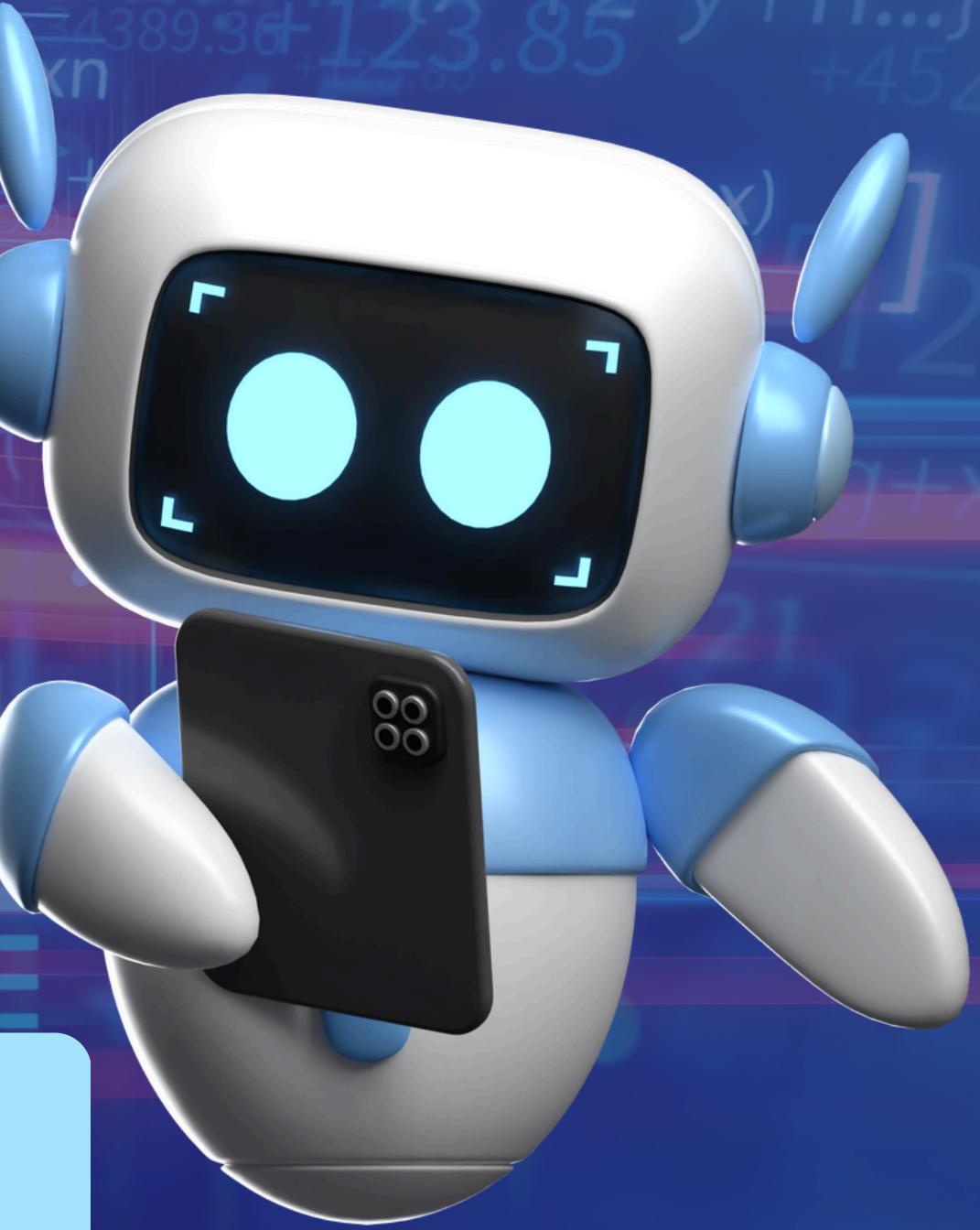
Our Group



**Aryan (U2220598L)**

**Maanya (U2323354L)**

**Akshay (U2323942B)**



# PROBLEM STATEMENT

**REAL LIFE  
PROBLEM**

**It is becoming increasingly difficult to discern fake news from real news in todays society, especially with the advent of Generative AI and LLMs such as ChatGPT and more recently, Google's Gemini**

**DATA SCIENCE  
PROBLEM**

**Classification:** We must classify a particular news as fake or real. Done by categorizing a particular piece of news' likelihood of being fake on a scale of 0 to 5

# SAMPLE COLLECTION

**USED 3 DATASETS, NAMELY LIAR, IFND AND FAKEEDIT**

Datasets from Kaggle, paperswithcode and huggingface.

3 datasets to get larger amount and variety of datapoints for more accurate modelling

# DATA CLEANING & PREPARATION

Unnamed: 0.2	Unnamed: 0	Unnamed: 0.1	Unnamed: 0.1.1	author	clean_title	created_utc	domain	hasImage	id	image_url	linked_submission_id	num_comments	score	subreddit	
0	0	0	NaN	NaN	buzzly6	virginia first lady criticized for handing	1.551316e+09	philly.com	False	avkxum	NaN	NaN	2.0	16	nottheonion han...

	Text	TruthRating	Country
0	Virginia first lady criticized for handing cot...	0	USA
3	Woman bites camel's testicles to save herself ...	0	USA

Deleted duplicates, null values

Also dropped columns to maintain conformity and uniformity amongst all the dataset columns

# Replacing MISLEADIND with MISLEADING

```
In [134]: df[df['Subject'] == 'MISLEADIND'].head(67)
```

```
Out[134]:
```

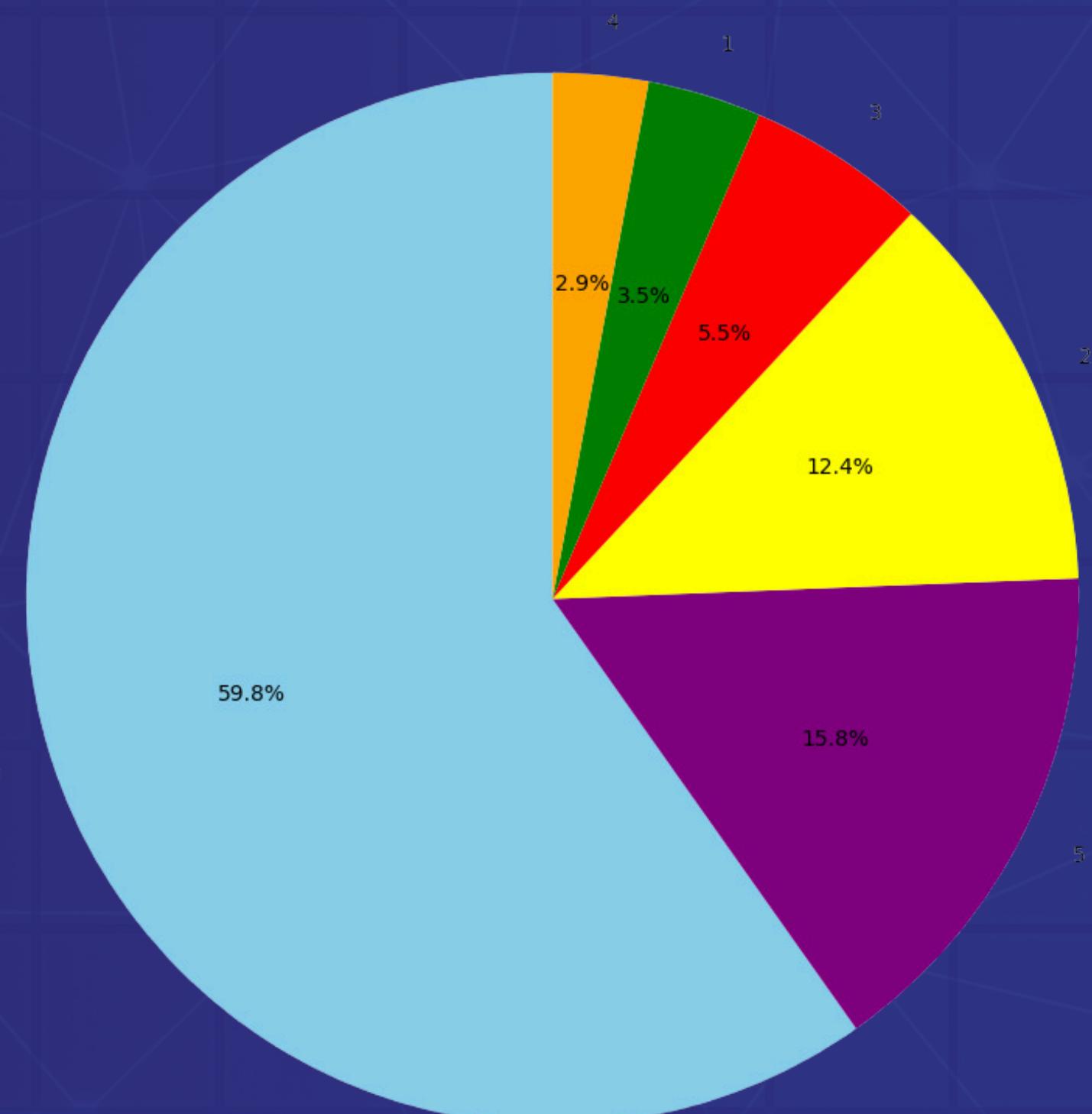
	Text	Subject	TruthRating	Country
54046	Fact Check: Old list of blood donors circulate...	MISLEADIND	Fake	India
54052	Fact Check: Picture of man injured in Mewat vi...	MISLEADIND	Fake	India
54053	Fact Check: This CISF officer was not injured ...	MISLEADIND	Fake	India
54054	Fact Check: Old picture from Prayagraj shared ...	MISLEADIND	Fake	India
54058	Fact Check: This picture of Kejriwal without m...	MISLEADIND	Fake	India
54063	Fact Check: Another lockdown in Delhi? No, thi...	MISLEADIND	Fake	India
54075	Fact Check: These pictures of farmers protest...	MISLEADIND	Fake	India
54077	Fact Check: This is not Amulya, the girl who r...	MISLEADIND	Fake	India
54083	Fact Check: Obama never warned Africans agains...	MISLEADIND	Fake	India
54089	Fact Check: Old images of pro-Khalistan demons...	MISLEADIND	Fake	India
54090	Fact Check: This is not a Wuhan lab that Obama...	MISLEADIND	Fake	India
54091	Fact Check: Old video of Khalistan supporters ...	MISLEADIND	Fake	India
54105	Fact Check: VFX video created by Russian artis...	MISLEADIND	Fake	India
54108	Fact Check: This is not an RSS man held for wa...	MISLEADIND	Fake	India
54116	Fact Check: Does this herbal medicine from Tan...	MISLEADIND	Fake	India
54118	Is this couple kissing during the ongoing prot...	MISLEADIND	Fake	India
54129	Fact Check: No, Canadian PM Trudeau is not pac...	MISLEADIND	Fake	India
54148	Fact Check: Pro-democracy clash in Hong Kong p...	MISLEADIND	Fake	India

Since there is no difference, we replace MISLEADIND to MISLEADING. Additionally, MISLEADING isn't an appropriate category, thus we change it to gossip/opinion (as it involved everything from factchecks, viral news, etc)

```
In [142]: df['Subject'] = df['Subject'].replace("MISLEADIND", "MISLEADING")
df['Subject'] = df['Subject'].replace("MISLEADING", "GOSSIP/OPINION")
df["TruthRating"] = df["TruthRating"].replace({'Fake' : 0})
df = df[df["TruthRating"]!="nan"]
df['TruthRating'] = pd.to_numeric(df['TruthRating'], errors='coerce', downcast='integer')
```

# **EXPLORATORY ANALYSIS & STATISTICAL DESCRIPTION**

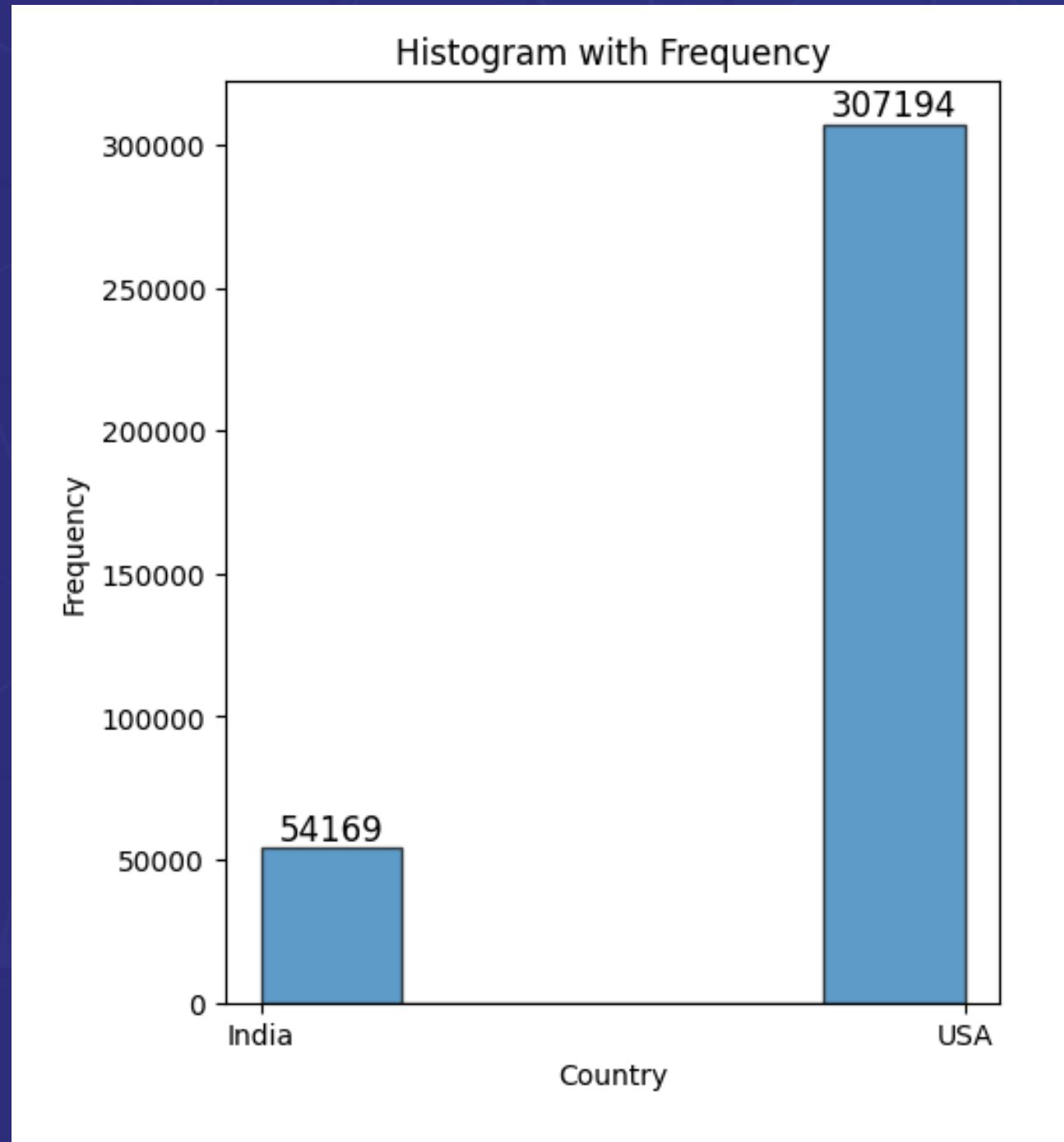
# Pie Chart of Percentage of Various Truth Ratings in Dataset:



## Our Insight

The training data majorly consists of fake news for better learning rate.

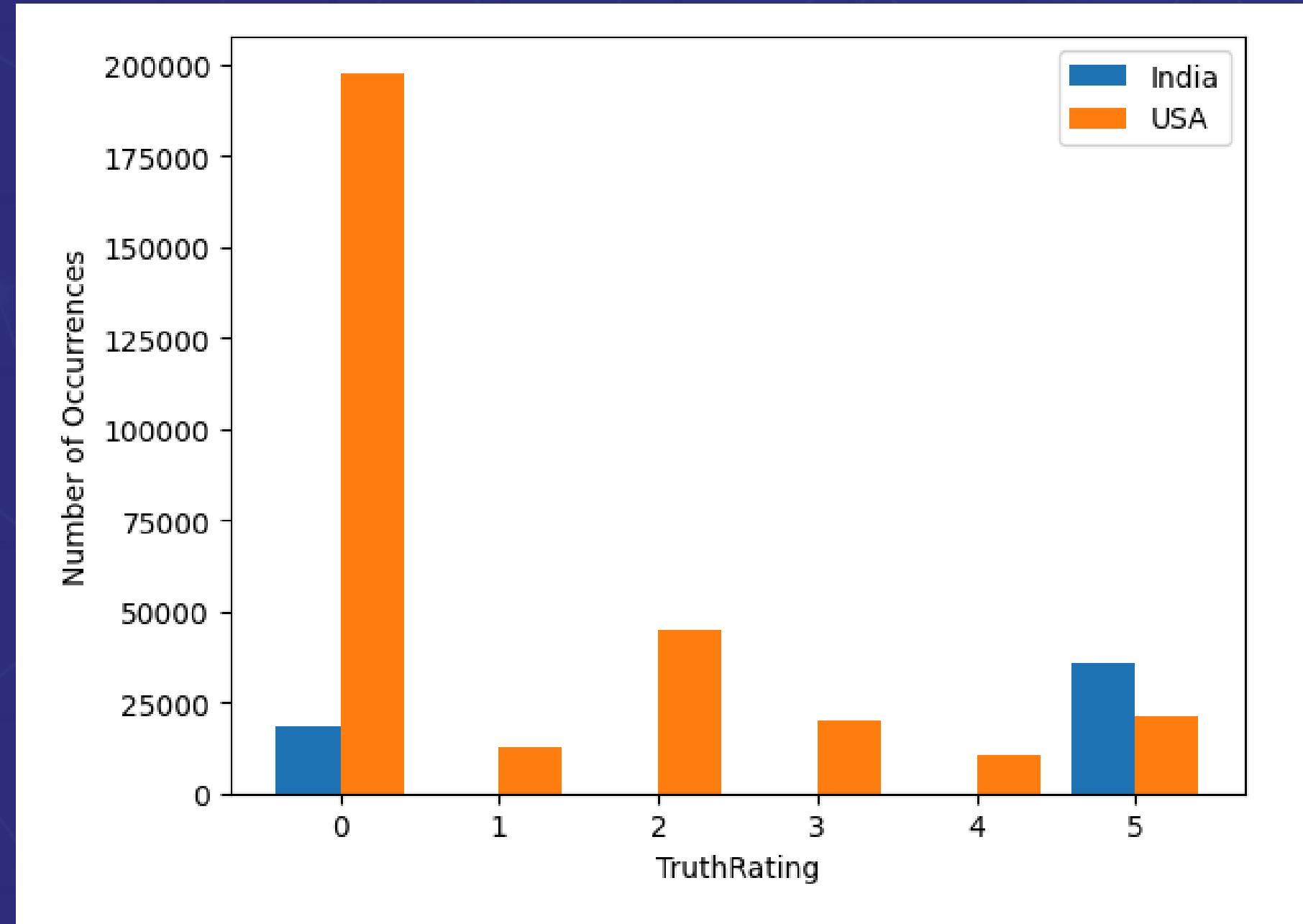
# Dual-variable Histogram of Number of Articles by TruthRating per Country



## Our Insight

American articles have a significantly higher proportion of articles

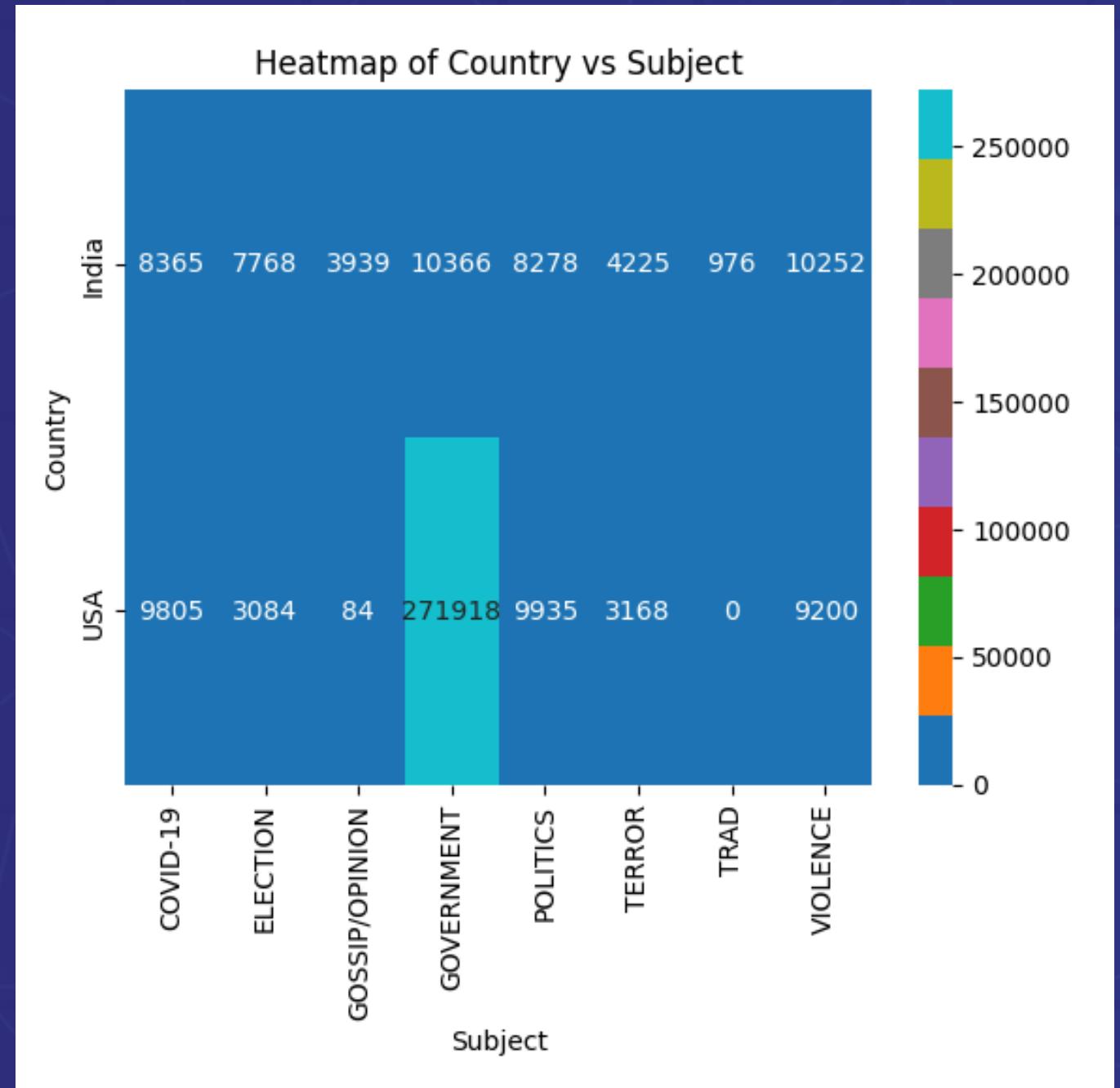
# Typical Histogram of Number of Articles per Country



## Our Insight

American articles have a significantly higher proportion of articles with TruthRating 0.

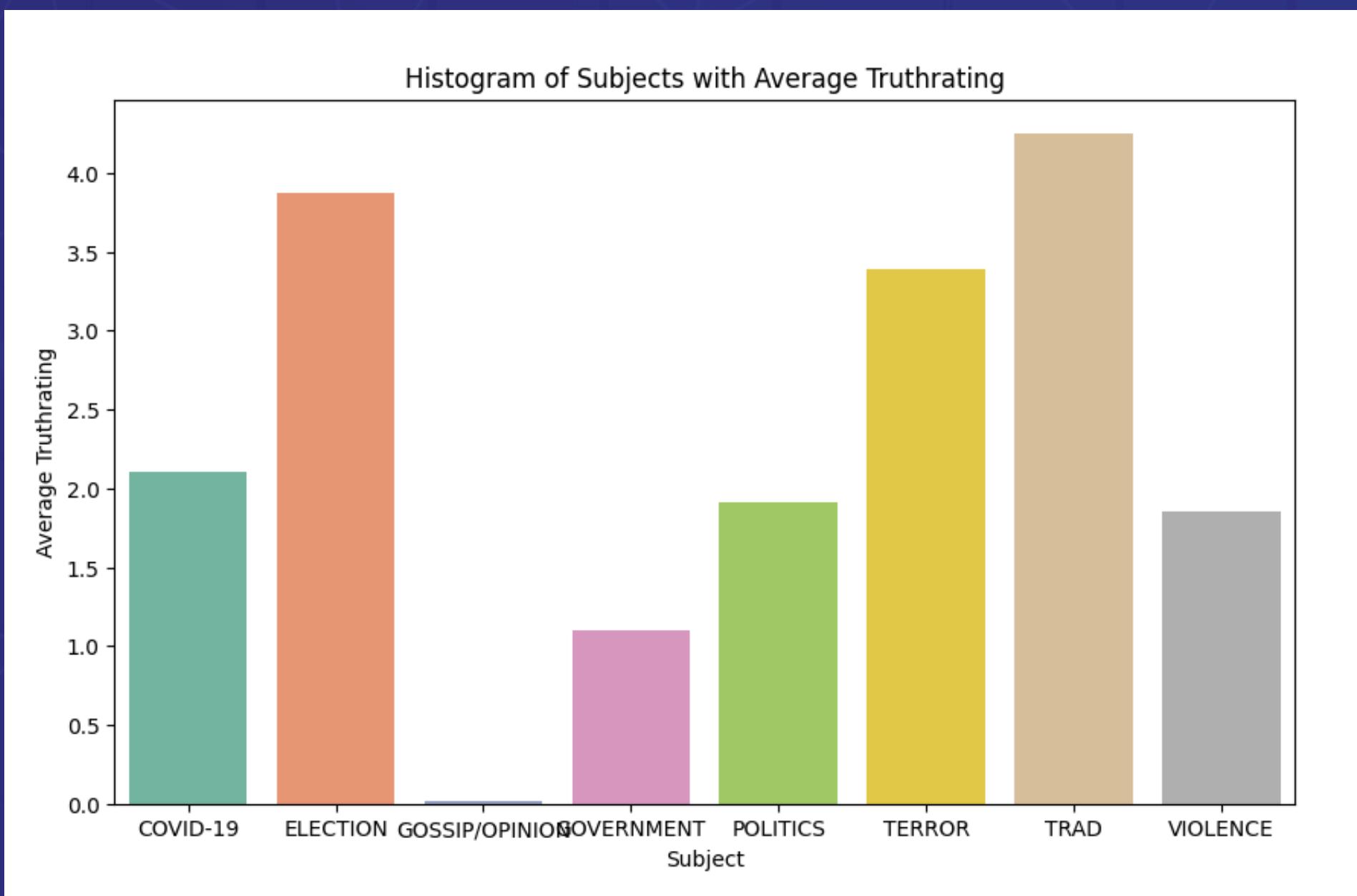
# Heatmap of Country vs Subject:



# Our Insight

US governmental articles far outnumber the rest of the country and subject combinations.

# Typical Histogram of Average TruthRating of Subject



Our Insight  
Average  
TruthRating of  
Government  
articles is about 1.0

# WORD CLOUD



# Our Insight

“Fact Check”, “new”, “say”, “look”, “people” and “found” are some of the most prominent words.

# APPLYING ML PRINCIPLES

- Cleaning Text - Removing Stop words like at, in, a, etc.
- Tokenization - Turn words into tokens of fixed length.
- Vectorization - Transform tokens into integer sequences(vectors) for ML fitting.
- Model Training - Training 2 models for comparison
- Metric Evaluation

# LINEAR REGRESSION

Linear regression is a foundational algorithm in the field of machine learning and statistics.

$$y = mx + b$$

## Libraries and Frameworks:

**scikit-learn** (for general-purpose ML)

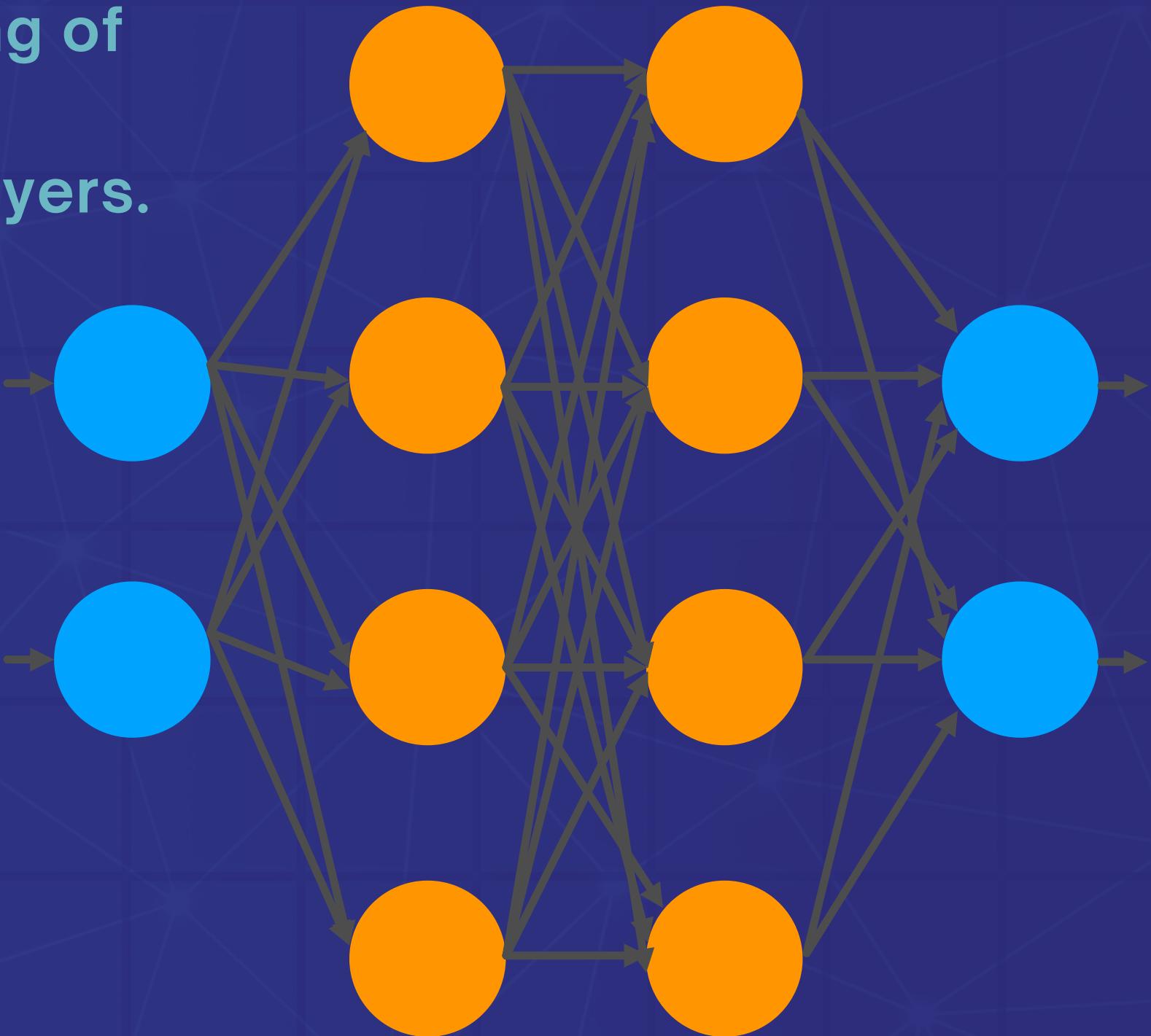
**statsmodels** (for more statistical-oriented tasks)

**TensorFlow and PyTorch** (for advanced ML applications, though often used for deep learning)

# SEQUENTIAL NEURAL NETWORKS

Neural networks are a type of machine learning algorithm inspired by the structure and functioning of the human brain. A neural network consists of interconnected nodes (neurons) organized into layers.

- Sequential
- Embedding layer
- LSTM (Long Short-Term Memory) layer with 128 units
- Dropout prevent overfitting.



# METRIC EVALUATION

Accuracy

F1 Score

Linear  
Regression



Neural  
Network



# DATA DRIVEN INSIGHTS

- Dataset skewed towards low Truth Ratings, suggesting prevalence of fake news.
- Distribution of articles biased towards US over Indian sources, possibly leading to US-centric bias in analyses.
- US governmental articles have lower TruthRatings, suggesting credibility issues.
- US articles show higher proportion of TruthRating 0, indicating higher fake news prevalence.



# RECOMENDATIONS

- Focus on monitoring and fact-checking American sources with high prevalence of low-truth articles.
- Balance dataset to prevent bias in machine learning models.
- Implement subject-specific fact-checking for government news.
- Use common words from word clouds to detect fake news.
- Include diverse sources for better generalization.
- Educate the public on spotting fake news and promote media literacy.



# THANK YOU!!

