



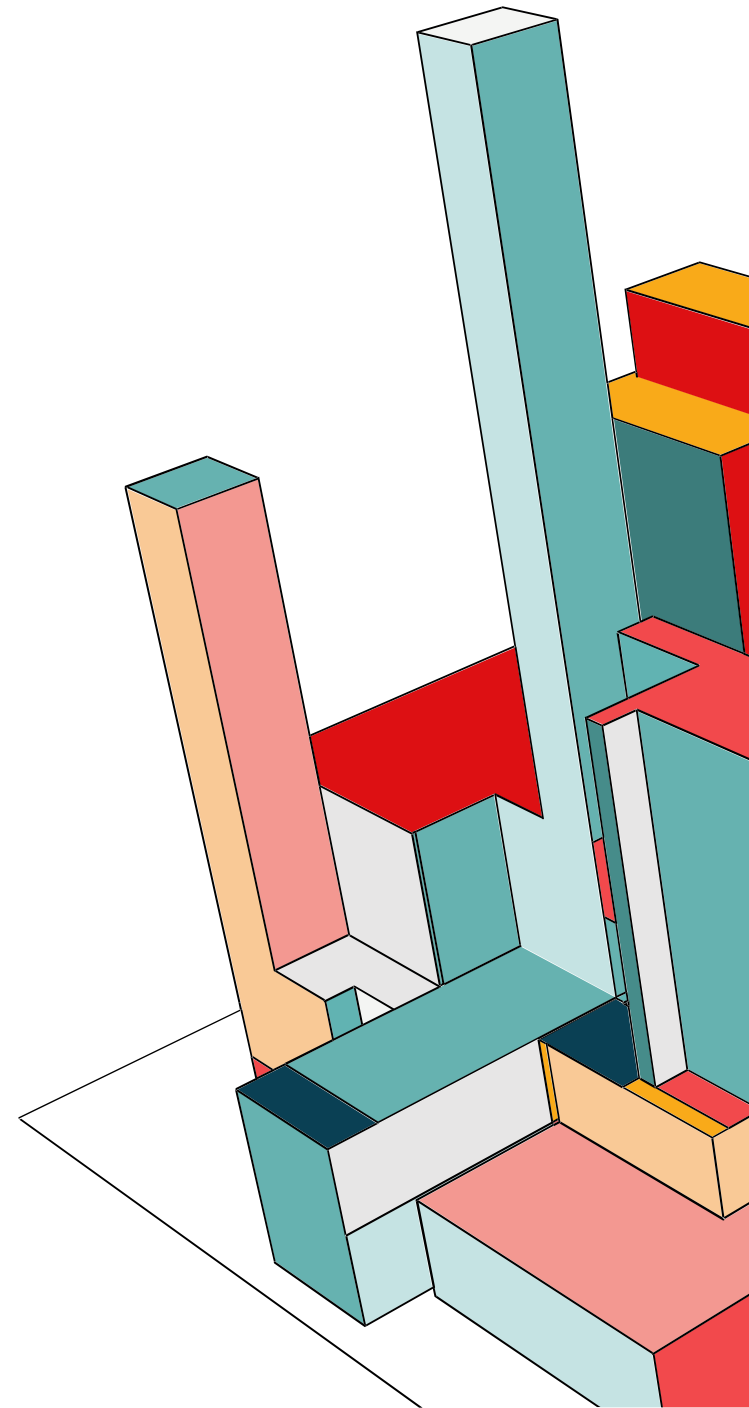
# **FAKE OR FACT?**

Zeenat Zahoor

21BDA61

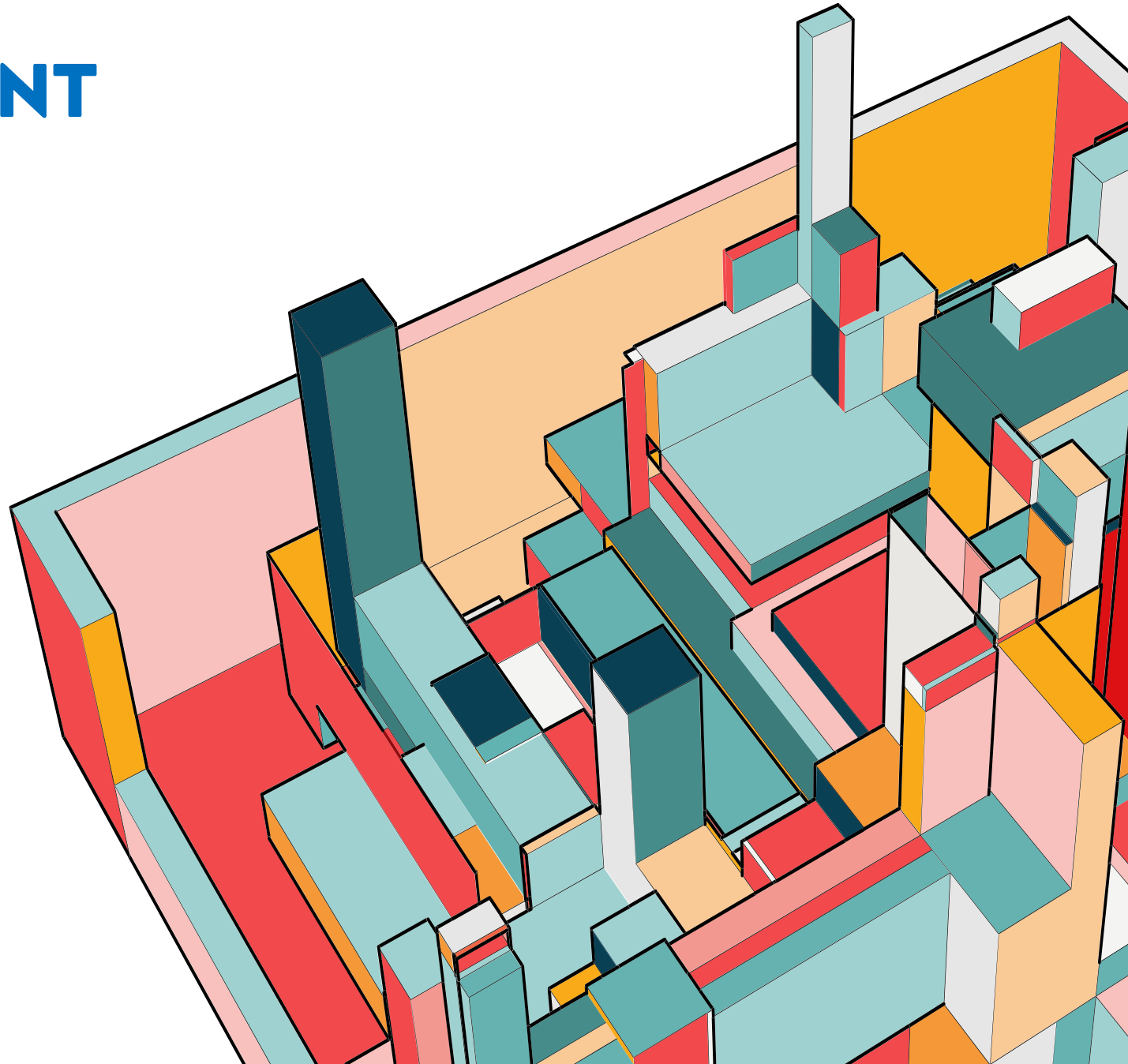
# WHY THIS PROJECT?

**We have all seen fake news forwards on our WhatsApp messages. Generally, these articles are generated by bots and internet trolls and are used with an intent to intrigue the audience and mislead them. Fake news can be very dangerous as it can spread misinformation and inflict rage in public. It is now becoming a serious problem in India due to more and more people using social media and lower levels of digital awareness.**



# PROBLEM STATEMENT

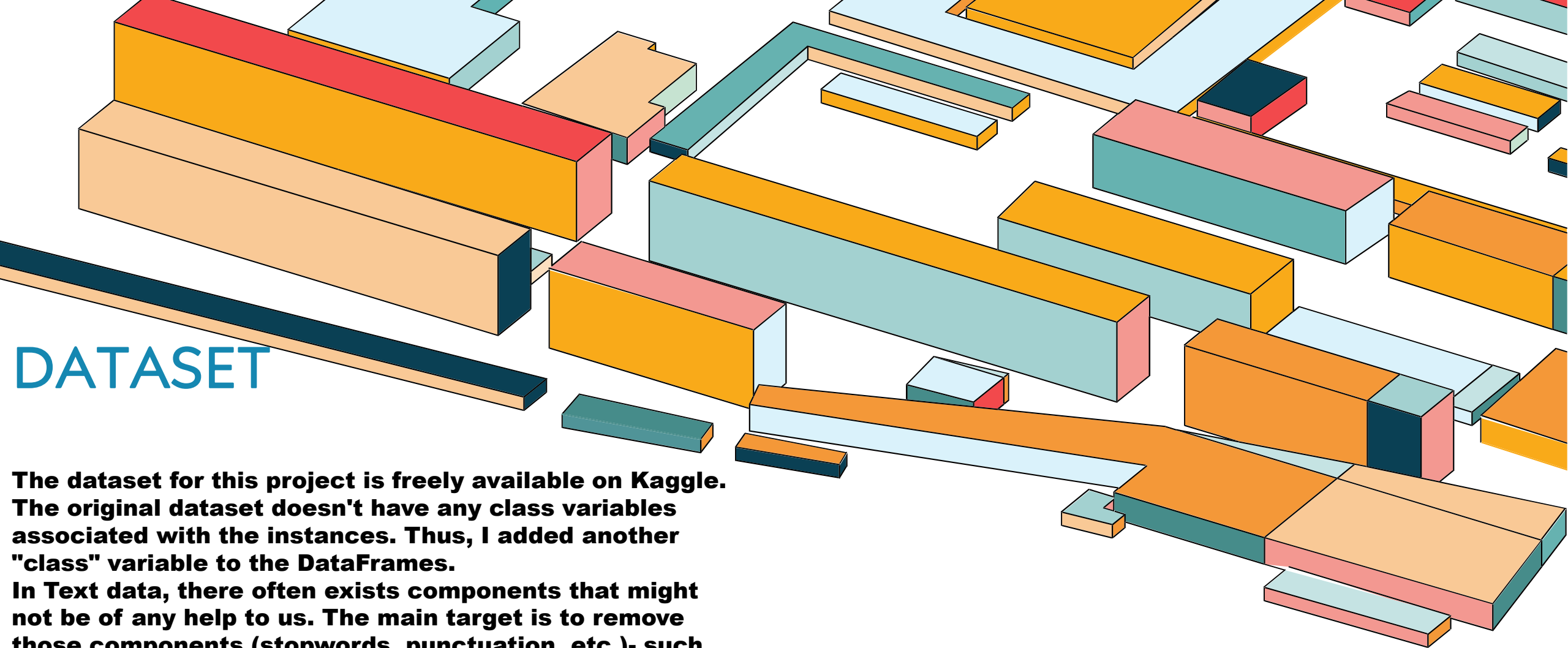
**The aim of this project is to make a Fake News Classification Web Application that can identify fake or true(fact) news.**





# SOLUTION

**To build a streamlit web app that can differentiate between fake and real news when an article is given to it.**



# DATASET

**The dataset for this project is freely available on Kaggle. The original dataset doesn't have any class variables associated with the instances. Thus, I added another "class" variable to the DataFrames.**

**In Text data, there often exists components that might not be of any help to us. The main target is to remove those components (stopwords, punctuation, etc.)- such that we can have a smoother analysis. Also, to get a reliable and authentic score for classification I concatenate the "text" and "title" columns. Then drop the redundant columns from both the DataFrames. Then, I just make a single DataFrame out of both the DataFrames.**

# METHODOLOGY

- **The first step was data pre-processing and cleaning.**
- **Form pre-trained word embeddings for text datasets using a pre-trained language model with these embeddings. From all of the visualizations, each technique produced satisfactory results in terms of separation between real and fake articles.**
- **Following these methods, some basic classical machine learning models are trained on the dataset and compared in terms of accuracy. For the best results, an LSTM/Advanced RNN model was used. The methods in this project were just some quick exploration and analysis methods of the given dataset.**
- **After attempting to use statistical models such as Decision Trees, Random Forest Regression, Support Vector Machines, etc., and deep learning models such as Recurrent Neural Networks, Long Short Term Memory Units, and Transformers, Bidirectional LSTMs proved to perform the best in terms of Accuracy, Evaluation and Generalisation.**

# CONCLUSION

**The real and fake datasets have unique 'subject' categories between them. Therefore, it is extremely important that we do not use the 'subject' as a feature for making predictions; this will introduce data leakage that will make it extremely easy for our models to achieve 100% accuracy. With this data leakage in our models, we will have extremely high performance on our training and validation sets, but very poor generalisation performance on unseen real world data. This problem is common occurrence throughout data science, and something that we need to carefully look out for throughout projects.**

**It seems in general, fake news articles tend to be longer than true articles. This is especially true for the article titles, and in fact the stark difference noted above suggests that we could classify articles to a reasonable accuracy using only title length.**

**The user is supposed to paste an article in the input box and then click on the Submit Button to run inference. After some time, the model returns a prediction whether the article is supposed to be Fake or Fact.**

**THANK YOU**

