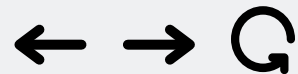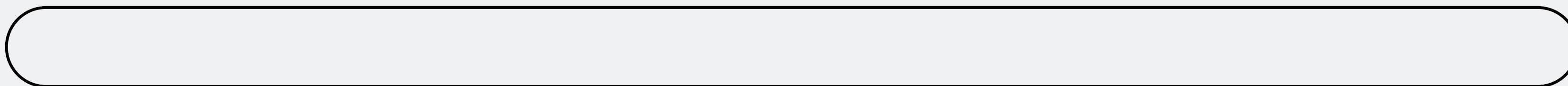Fusemachines

# Youtube Sentimental Analysis
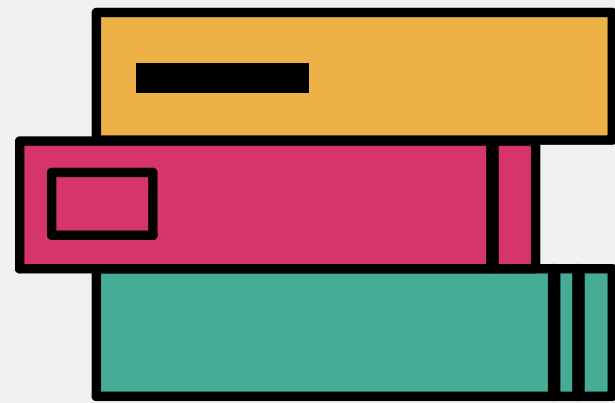
# Introduction

Sentimental Analysis involves the classification of human language sentences into predefined categoriees.

Our youtube comment sentimental analysis involves three categories namely positive , neutral and negative.

Steps

# Introduction

## Steps done

1. Scraped comment data via youtube api
2. Preprocessed the scraped data
3. Labelled the dataset with textBlob package

4. Naiive Bayes Model was trained
5. Model performance measured
6. Retraining and reevaluating model with some changes

Scrape

# Scraping

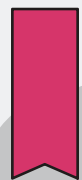## The First Attempt

Initial attempt was done with selenium which was very slow and unfeasible

## The Second Attempt

Involved using the youtube api to get json of requested number of comments (16333)

Preprocess

# Preprocessing



**Html tags and links removed**

**Emoji replaced with words**

**Removal of non english alphabets**

# Labelling TextBlob

## Positive Comments

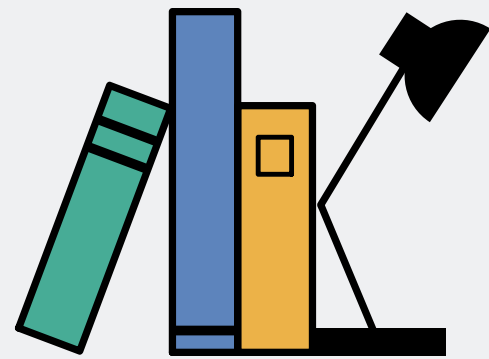4609 comments were labelled as positive

## Neutral Comments

10865 comments were labelled as neutral

## Negative Comments

859 comments were labelled as negative

# The dataset is imbalanced

# Model training

Remove special characters

Remove stop words

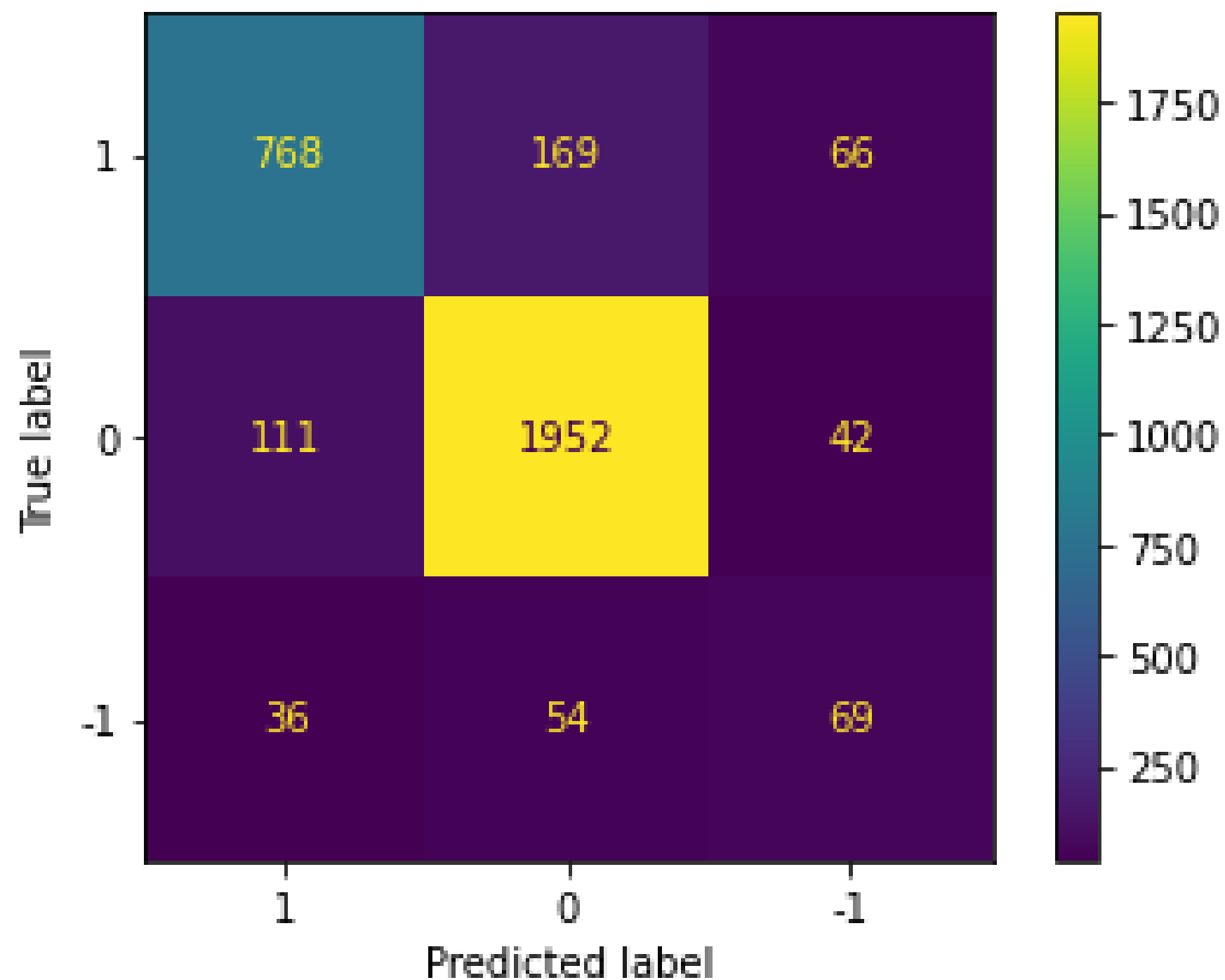Lemmatize instead of stemming

Vectorize by bag of word model

**Multinomial Naiive bayes classifier model used**

Insert your topic here

# Results 1

## No consideration for the skewed dataset taken

# Results 2

## Number of features limited to only 25k

Insert your topic here

# Results 3

## Sampalled similar number of each category

Insert your topic here

# Conclusion



The model with equally sampled training data performed the best for the classificatiion task. This emphasises the analysis of dataset skewness before the model training steps