# Sentiment and Emotional Analysis of Indian Union Budget 2022-23 on Social Media

Aditya Dixit
aditya21004@iiitd.ac.in
IIIT Delhi
India

Mohnish Basarkar
mohnish21052@iiitd.ac.in
IIIT Delhi
India

Shivang Kanwar
shivang21084@iiitd.ac.in
IIIT Delhi
India

Shubham Aggarwal
shubham21090@iiitd.ac.in
IIIT Delhi
India

## ABSTRACT

This study aims at analyzing the sentiments and emotions of people toward the Indian Union Budget 2022-23 using the Tweets related to the Union Budget on Twitter and the comments on the live-streamed Budget Session videos on Youtube. Tweets and comments were taken in English and Hindi language. The analysis classified the tweets and comments into three Sentiment categories: Positive, Negative, and Neutral, and eight Emotion categories: Joy, Shame, Anger, Disgust, Neutral, Surprise, Sadness, and Fear. The study also compared the sentiments and emotions for the Union Budget 2022-23 with that of 2021-22.Random Forest classifier and TextBlob library of python used as baseline classifiers for emotions and sentiments respectively.Different machine learning classifiers were used like KNN, XGBoost, DecisionTree etc to find the accuracy of emotions depicted. 'XGBoost' came out to be the best classifier with an accuracy of 83% .After complete analysis, we found that the majority tweets/comments were classified into the emotion category 'Joy', indicating that people are happy with both budgets 2022-23 and 2021-22. However, the number of 'Joy' tweets/comments of budget 2022-23 was relatively less than that of the budget 2021-22, indicating that the people were more satisfied with the previous year budget as compared to the current budget.

## KEYWORDS

Union Budget,Sentiment Analysis,Twitter,YouTube,Emotions,Text Mining, TextBlob, Machine Learning, Classifiers

## 1 PROBLEM STATEMENT

Analyzing the emotions and sentiments of Indian citizens towards Union Budget for financial year 2022-23 using social media platforms like twitter and youtube.

## 2 INTRODUCTION

Social media not only facilitates conversations, creation, and sharing of information, it also offers a huge corpus for different types of analysis like Misinformation Analysis, Sentiment Analysis, and so on. Sentiment Analysis in social media can be useful to understand the opinion of people about a variety of topics, for example, Governments can analyze the public's opinion over Covid vaccines and can take necessary actions. Companies can study consumers' sentiment towards their product and develop business strategies, political parties can predict their winning chances in the upcoming elections by monitoring voters' opinions about them, and so on. Sentiment analysis on Twitter and Youtube can be used to monitor the public's opinion over the Indian Union Budget 2022-23. This will help Government get the feedback of the public on the Union Budget whether they liked the central theme of the budget or felt disappointed. It will also allow the government to make changes in the policies that will add to the growth of the country. This study aims at analyzing the sentiments and emotions of people toward the Indian Union Budget 2022-23 using the Tweets related to the Union Budget on Twitter and the comments on the live-streamed Budget Session videos on Youtube. Tweets and comments were taken in English and Hindi language. The analysis classified the tweets and comments into three Sentiment categories: Positive, Negative, and Neutral, and eight Emotion categories: Joy, Neutral, Anger, Disgust, Shame, Surprise, Sadness, and Fear. The study also compared the sentiments and emotions of the Union Budget 2022-23 with that of 2021-23. This paper presents a sentiment and emotion analysis of public opinions towards the Union Budget mined from Twitter and Youtube. We scraped the Tweets related to Union Budget from the microblogging site Twitter and the comments from the budget session live-streamed youtube videos. We scraped data in two phases: The pre-budget phase and Post-budget phase. Tweets and comments in the Hindi language were translated into English using the 'googletrans' python library. We performed the Sentiment analysis by classifying tweets into categories: Positive, Negative, and Neutral, followed by classifying them into eight emotion categories: Joy, Neutral, Anger, Disgust, Shame, Surprise, Sadness, and Fear. For classification,we used the following models: MNB,KNN,XGB,LR,SVC,LinearSVC,SGD,AdaBoost,DTC and GBC and later compared their accuracies. We compared the sentiments and emotions in the pre-budget phase with that of the post-budget phase. Finally, we also compared the sentiments and emotions towards Union Budget 2022-23 with that of the previous year i.e., Union Budget 2021-22. As many things in this world majorly depend upon perception as well as sentiments of people, hence sentiments of people towards any work done by the government can be a very strong driving force, thus this problem holds importance in designing policies according to the sentiment of the masses which would help in making efficient policies in future.

## 3 LITERATURE REVIEW

In 2015 [3], Satarupa Guha, Aditya Joshi, Vasudeva Varma in their paper Sentibase: Sentiment Analysis in Twitter on a budget, data

from Twitter is collected and operated upon. In preprocessing, authors have expanded all acronyms to their original form and change number to strings. The next step is vocabulary generation in which they have given unique IDs to all the words occurring in the data. All the hashtags are given single unique IDs as they are not important for sentiment analysis. Next, for the classification step, they have classified tweets into subjective(positive and negative) and objective(neutral). Various features of tweets were analyzed before classifying them. After this, a linear SVM classifier is trained using the python scikit-learn library for classification. The paper's accuracy in classifying tweets is checked by calculating the F1 score, precision, recall for the classifier.

In 2016 [6], Moonis Shakeel and Vikram Karwal have aimed to estimate the sentiment score of India's Union Budget document for 2016-17, whose central theme as anticipated, was rural India. In this, the Budget document was converted to a .txt file format for text data retrieval. R language was used for the analysis. The text file data was preprocessed and converted to a document term matrix displaying the frequency matrix of the text data and the words in it were sorted in decreasing order of their occurrence. Then the association between words that related to the words which were anticipated in the budget 2016-17 was carried out. The word cloud representing the most occurring words was created. The analysis pointed out that the words related to taxation had the highest frequency and deviated from the main focus of the Union Budget as anticipated. Inequality. The sentiment displayed by the budget document can be safely considered neutral. It mostly talks about taxation and less about other aspects of the economy. This paper suggested that the social and economic growth of India is regulated by regulating taxes rather than through other more relevant means.

In 2018 [7], the sentiment and emotional analysis were done by the students of Babasaheb Ambedkar University, Maharashtra. They aim to understand the relevance of public opinion towards the union budget of India 2018 through the large online data resource Twitter. They have used ML techniques, lexicon analysis such as Bing Liu and NRC, and statistical approaches for sentiment score calculation and detection of the sentiment using Twitter tweets. Studied strategies like supervised, unsupervised, lexicon, similarity metric, SVM, etc, which were used by other researchers. They have completed their study using R programming tools and packages. Around 10k tweets were collected using hashtags like 'Budget2018', 'UnionBudget2018', 'BudgetSession2018'. They have used 8 emotions to express their analysis i.e. joy, neutral, anger, disgust, shame, surprise, sadness, and fear. After performing sentence-level sentiment analysis they concluded that due to the positive polarity shown in text processing, the public inclination is towards government policy.

In 2018 [2], researchers accumulated 3 consecutive year tweets related to the Indian Budget for 2016, 2017 and 2018. Twitter is used to make a large corpus of dataset. This study aims to analyze the opinion of the Indian crowd and tries to infer some technical sense from that budget data. Out of document, sentence and feature level sentiment analysis, sentence level classifier is used to determine the polarity of sentiment i.e. positive, negative and neutral. To implement this sentiment analysis they have used R programming open source tools. Statistical, grammar and machine learning based methodologies used. They have used "Score = No. of positive words - No. of negative words". If score > 0 then it is positive sentiment , else if score < 0 then it is negative sentiment otherwise it is neutral sentiment. They have concluded that best polarity is observed in the tweets of budget 2018.

In 2020 [5] , Twitter sentiment analysis was done on Indian union budget 2020 by students of GNA University, Phagwara aims to study a public reaction on 6000 real-time tweets using the hashtag "Budget2020". For this, they have used tweepy which is a Twitter API to consolidate their study results. All the preprocessing stuff has been performed to attain higher accuracy. They finalize their sentimental score based on the concept of polarity and subjectivity with an overall positive score of +149.3387. To assign subjectivity and polarity to the most frequent words in the dataset they have used the Textblob library of python. Finally, they concluded that public reaction is in favor of government policy due to the high positive score achieved from the sentimental analysis.

In 2021, in the paper [4], Twitter was chosen to extract data from, all tweets were downloaded which were posted using the #unionbudget hashtag over a period of one month coinciding with the announcement of the budget. In this paper, data was captured using python and transferred to an excel file and then it was cleaned by removing punctuations, stop words, and neutral words. Words lesser than three letters, numbers, and special characters were removed. In this paper TermDocumentMatrix named function was used to make a document matrix table to depict the frequency of each word, the ten most frequent words are also mentioned in the results. Here sentiment analysis was done using EmoLex ( list of words and their association with 8 basic emotions ). And two sentiments ie negative and positive, here a dataframe was created where each row of dataframe had a sentence whose sentiment analysis needs to be done and there are ten columns that represent eight emotions and two sentiments. In this, a word cloud was generated using the above methodologies which depicted different words with different sizes depict their relevancy, it was observed that the word "economy" had the highest frequency with maximum use in 427 tweets.

In 2022 in the paper [1], the authors tried to corroborate the results of sentiment analysis with event analysis for the Indian Budget 2022 on cryptocurrencies. They tried to collect the tweets based on keywords like 'Indian Budget 2022' and 'Bitcoin'. They have used supervised Machine Learning algorithms. Tweets were collected 120 days prior to the release of the budget and 10 days after its announcement as the estimation window. They have used "Abnormal normal returns were measured for the constant mean model". For feature extraction, they have used the TF-IDF strategy and for finding the polarity of tweets they have used the TextBlob library of python. Models like Support Vector Machine, Bernoulli NaÃŕve Bayes, Logistic Regression for comparing accuracies. Logistic regression recorded with the best figures around 81.4 percent. They find that "negative sentiments which were expressed in tweets on BTC-INR were slightly higher than the positive tweets and statistics revealed that the positive sentiments were negated by the negative sentiments."

| Year | Before Budget | After Budget |
|------|---------------|--------------|
| 2021 | 2025 Tweets | 5312 Tweets, 2655 Youtube comments |
| 2022 | 1631 Tweets | 4194 Tweets, 1109 Youtube comments |

**Table 1: Number of Tweets Scraped for Each Year**

## 4 METHODOLOGY

### 4.1 Design and Dataset

As explained in the fig., 1. Dataset for each year 2021 and 2022 was created from Twitter and Youtube video's comment and replies. Tweets related to budget 2022-23 and budget 2021-22 for the months of January and February were scrapped from Twitter using Snscrape library and comments from Youtube videos that live streamed Budget session were scraped. Pre-phase data refers to tweets that were tweeted before the announcement of budget for the month of January and Post-phase data consisted of tweets and youtube comment generated after the budget announcement for month of February. Tweets were scraped in two languages, Hindi and English from twitter.

Query used for scraping tweets : (#unionbudget2022-23 OR #budget2022 OR #unionbudget OR #budget OR #unionbudget2022 OR #unionbudget22-23 OR "corporate tax" OR "income tax").

Pre-processing of tweets was carried by converting into lowercase text,removing whitespaces,punctuations,removing all URLs and username,removing # from words,emoticons.Later preprocessed tweets were followed by stemming.

### 4.2 Classifying pre-labeled into Emotions

We used a pre-labeled dataset to classify the tweets into one of the 8 emotion categories: 'neutral', 'joy', 'sadness', 'fear', 'surprise', 'anger', 'shame' and 'disgust'.

### 4.3 Pre-budget phase analysis

(1) **Sentiment Analysis:** We found the sentiments of each tweet as 'Positive' or 'Negative' or 'Neutral' using the subjectivity and polarity score of each tweet calculated using the Textblob library.

(2) **Emotional Analysis:** We used the pre-labeled dataset to train our Random-Forest classification model. Later, we used this model to classify the tweets and comments scraped from twitter and youtube into one of the 8 emotional categories.

(3) **Visualization using Word cloud and Frequency distribution:** We created the word cloud to visualize the most frequently occurring words in our tweets dataset. We also created a frequency distribution chart representing the top 30 most frequent words on the x-axis and their frequency on the y-axis. Refer to figure 2.

(4) **Classification through various models and finding their accuracies:** Considering the classification output of the Random-Forest model as the baseline result, we splitted the output into 20-80% and used it to train and test on two other classification models: Naive Bayes and XGB classifier.
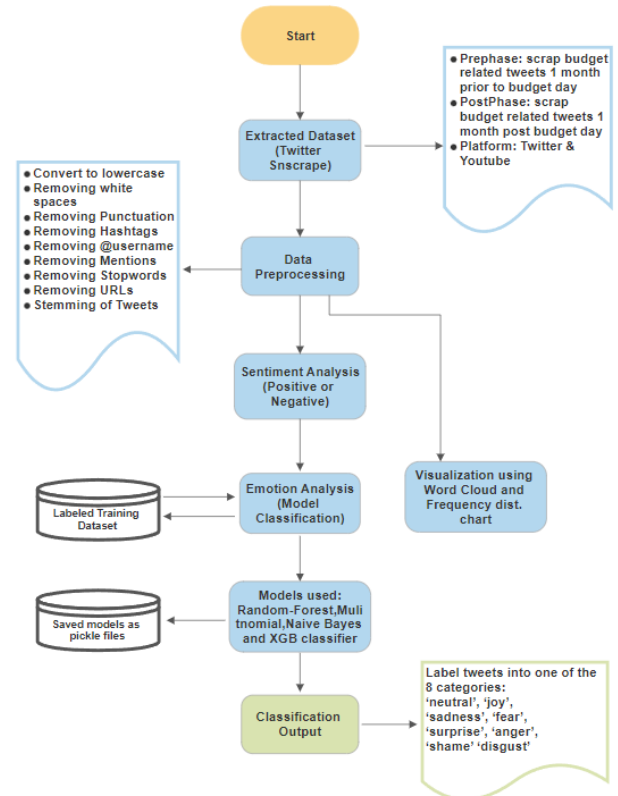
**Figure 1: Workflow**

### 4.4 Post-budget phase analysis

Similar to Pre-budget phase, sentiment analysis, emotional analysis, visualization using word cloud and frequency distribution and classification through other models was done and calculate corresponding accuracies. Refer to figure 3.

## 5 EVALUATION

- **Pre-budget Phase:** On performing Sentiment analysis on the tweets prior to the Budget, there were approximately 671 'Positive' tweets, 340 'Negative' tweets and 620 'Neutral' tweets. The comparison between the different sentiments for

| Sentiments | PreBudget 2022-23 | PostBudget 2022-23 | PreBudget 2021-22 | PostBudget 2021-22 |
|------------|---------|---------|---------|---------|
| Positive | 671 | 1669 | 808 | 3571 |
| Negative | 340 | 857 | 412 | 1051 |
| Neutral | 620 | 2066 | 805 | 4335 |

**Table 2: Count of Tweets/comments Scraped for each phase and each year**

both phases for the financial year 2021 and 2022 are shown

**Figure 2: Word cloud for pre budget phase**



**Figure 3: Word cloud for post budget phase**

in the figure. 8

Emotional classification using the Random-Forest model resulted in approximately 1200 tweets that were labeled as 'joy',190 tweets as 'neutral', 150 as 'sadness' followed by 'fear', 'anger', and 'surprise' that were less than 100.However, tweets with labeling as 'shame' and 'disgust' were negligible. Refer to figure 4

The Word Cloud and the Frequency distribution chart for the top 30 most occurring words showed that the words like 'budget', 'tax', 'inflation', 'income' 'defense' were most mentioned in the tweets.

On considering the above classification result as a baseline, splitting and using it for training and testing a Multinomial Naive Bayes classifier model, the results gave an accuracy of 78%. However, using the XGB classifier resulted in an accuracy of 83%. Refer to figure 5.

- **Post-budget Phase:** On performing sentiment analysis on tweets and comments replies from the youtube videos post the Budget, there were approximately 1669 'Positive' tweets, 857 'Negative' tweets and 2066 'Neutral' tweets.

  Doing emotion classification using the Random forest model resulted in approximately 500 as 'neutral', 3000 as 'joy', 490 as 'sadness', 50 as 'fear', 100 as 'surprise',120 as 'anger', and around 10 as 'disgust' emotions. Refer to figure 6.

  On considering the above classification result as a baseline, splitting and using it for training and testing a Multinomial Naive Bayes classifier model, the results gave an accuracy of 74%. However, using the XGB classifier resulted in an accuracy of 87%. Refer to figure 7.

For the Budget session 2022-23, percentage of negative tweets remained the same. However, the percentage of positive tweets dropped from 41.1% during pre-budget phase to 35.4% after the declaration of the budget.

On performing the emotion analysis, maximum number of tweets were classified as 'Joy' with 79% and 69.6% of the total tweets and comments for the financial year 2022-23 and 2021-23 respectively. However, the number of tweets and comments from youtube dropped for the year 2022 as compared to that of year 2021.

We started with Random Forest model as the baseline, and then we applied other models including MultinomialNB, KNeighborsClassifier, XGBClassifier, LogisticRegression,SVC, LinearSVC,SGDClassifier,AdaBoostClassifier, DecisionTreeClassifier,GradientBoostingClassifier.The accuracy achieved by these classifiers for both the year can be seen in the table below. Out of all these, XGBoost gave the

| Classifiers | Accuracy(year 2022) | Accuracy(year 2023) |
|---|---|---|
| MNB | 72.730656 | 71.946288 |
| KNN | 31.962293 | 39.770844 |
| XGBoost | 82.318900 | 84.643285 |
| LR | 76.547908 | 80.948627 |
| SVC | 76.873487 | 81.148660 |
| LSVC | 79.743024 | 82.457978 |
| SGD | 79.273583 | 82.373468 |
| ADABoost | 67.190344 | 68.339026 |
| DT | 79.071567 | 82.994568 |
| GRADBoost | 80.827219 | 83.137014 |

**Table 3: Accuracy achieved by different classifiers for year 2022 2021**

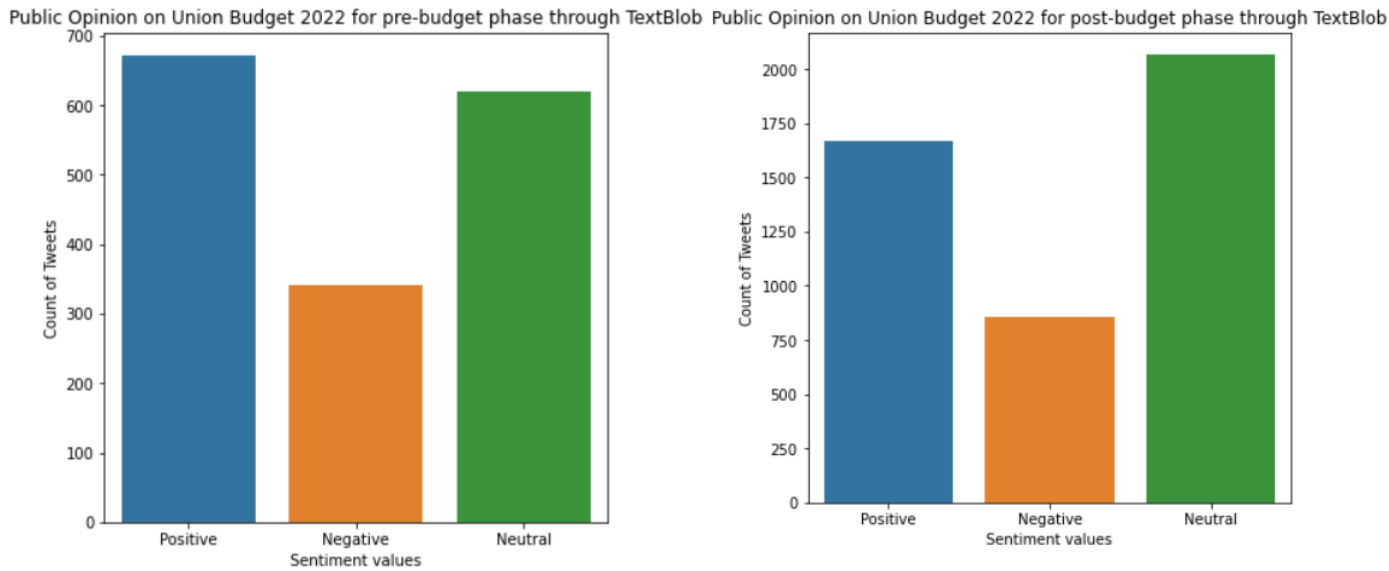highest accuracy of 84.5% and 83.10% for year 2021 and 2022

Figure 4: Public Opinion on Union Budget 2022 for pre-budget phase



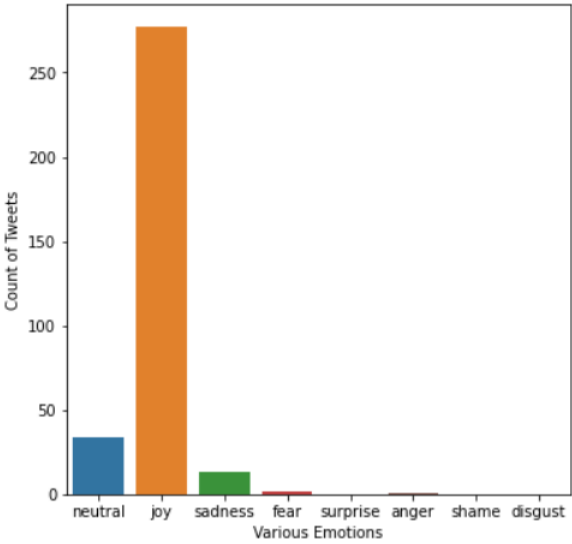Figure 6: Public Opinion on Union Budget 2022 for post-budget phase



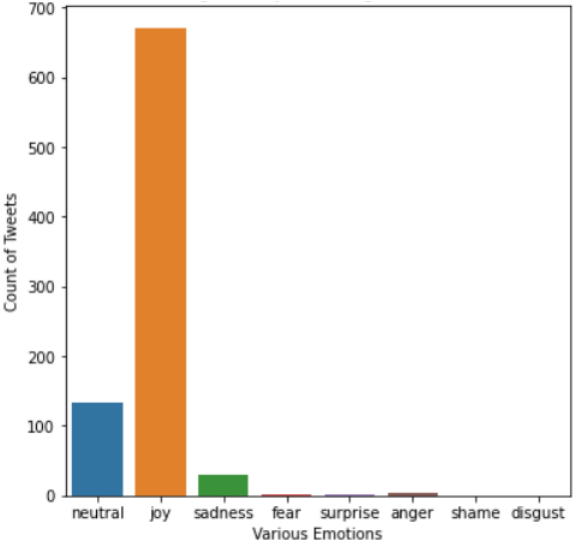Figure 5: Pre-Budget analysis through XGB classifier



Figure 7: Post-Budget analysis through XGB classifier

respectively.Comparison between the accuracies by different classfiers for both the finacial years i.e. 2021 and 2022 is shown in the figure.9 Also, the comparison between different emotions through XGBoost classifier for both financial years is shown in figure.10

## 6 NOVELTY

Our work is different from others in the sense that:

- Dataset corpus is made using social media platforms like twitter and youtube.

- Analyze the tweets tweeted in multiple language across the country by converting them to english with the help of 'googletrans' library.
- Classifying tweets in different sets of emotions for example joy, anger, sadness, trust, etc. Along with that sentiments were grouped into 3 categories i.e. positive, neutral, and negative.
- Try to compare the tweets of budget 2021-22 with 2022-23 and try to analyze the difference that has come in Indian society in a gap of 1 year.
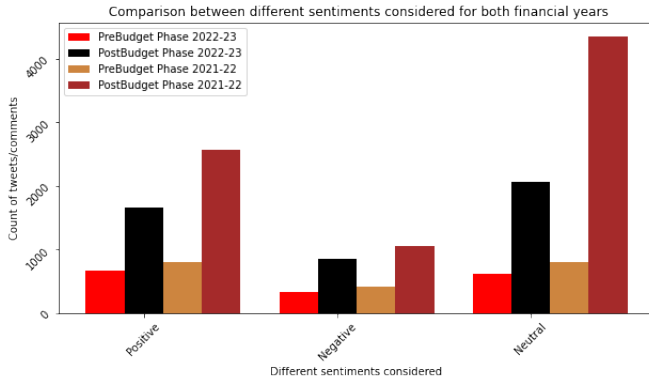
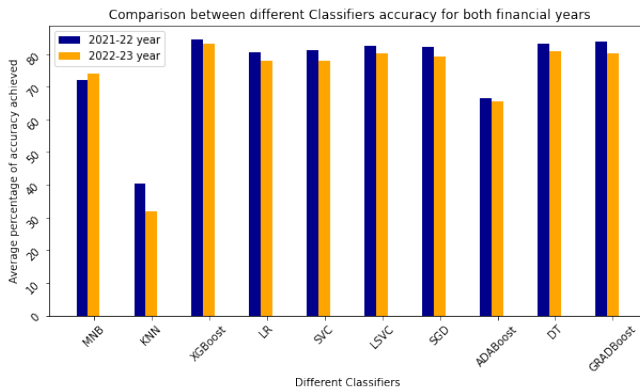**Figure 8: Sentiment comparison for both phases for the years 2021 and 2022**



**Figure 9: Accuracy comparison of different classifiers for both financial years**
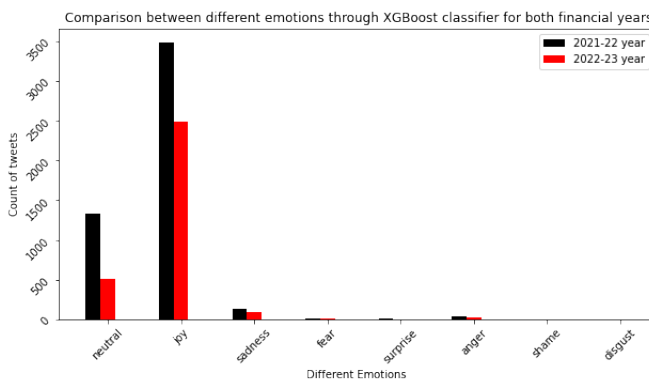


**Figure 10: Comparison between different emotions through XGBoost classifier for both financial years**

- Analysis performed in 2 phases i.e. Pre-phase(before release of budget) and post-phase(after release of budget).

# 7 FUTURE SCOPE

- Compare tweets of more than 1 previous year based on data availability.
- Analyze the responses in the format of GIFS, videos etc.
- Sentimental Analysis is done after classifying the tweets into major sectors of budget i.e. finance, education, agriculture, service, defense, etc, and correspondingly analyzing the budget impact on the above-mentioned sectors.

# 8 CONCLUSION

Our study provides research and practical implications in public policy making. Our findings offer insights into how people's emotions changes on social media before and after the Union Budget is declared. This can help the government and decision-makers navigate better policies in the future that benefit the country's citizens. The future work can be attributed in the direction of sector-wise analysis of the Union Budget to understand the people's emotions within different sectors of the economy and the contributing factors behind it. Also, other social media platforms, along with non-text data in multiple languages can be explored.

# REFERENCES

[1] Dr V Uma Maheswari Abhinand G. 2022. Corroboration of Twitter Sentiment Analysis and Event Analysis of Indian Budget 2022 on Bitcoin Market. (April 2022), 12. https://doi.org/10.21203/rs.3.rs-1515523/v1
[2] Seema S. Kawathekar Bharat Naiknaware. 2018. Peoples Opinion on Indian Budget Using Sentiment Analysis Techniques. (2018), 5.
[3] Satarupa Guha, Aditya Joshi, and Vasudeva Varma. 2015. Sentibase: Sentiment Analysis in Twitter on a Budget. (07 2015), 590–594. https://doi.org/10.18653/v1/S15-2098
[4] Sheetal Mahendher, Toshith Sastry, Yashus Gopal, and Rohith M S. 2021. Sentimental Analysis on the Union Budget, India-2020. 10 (02 2021), 14–21.
[5] Manpreet Singh3 Rupinder Kaur1, Rajvir Kaur2 and Dr. Sandeep Ranjan4. 2020. Twitter Sentiment Analysis of the Indian Union Budget 2020. 29 (2020), 8. http://www.ijsrcsams.com/images/stories/Past_Issue_Docs/ijsrcsamsv7i4p149.pdf
[6] Moonis Shakeel and Vikram Karwal. 2016. Lexicon-based sentiment analysis of Indian Union Budget 2016−17. (july 2016), 299–302. https://doi.org/10.1109/ICSPCom.2016.7980595
[7] Monali Waghmare and Sachin Deshmukh. 2018. Sentiment and Emotion analysis on Indian Budget. 7 (july 2018), 4.