

Toxic Comment Classification



NaN-Prediction Pending

Aneri Dalwadi AU1940153

Mananshi Vyas AU1940289

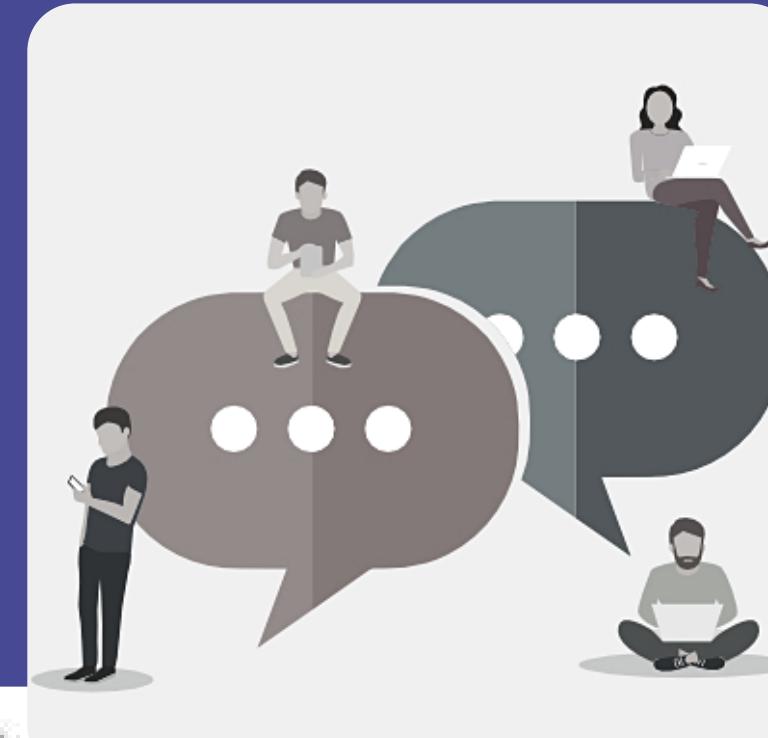
Nandini Bhatt AU1940283

Kairavi Shah AU1940177

...

Introduction

The data contains 1.5 million comments from Wikipedia's talk page and has to be classified into multiclass-classification - toxic, obscene, identity hate, severe, etc. And here we try to classify toxicity and calculate its severity using and comparing various Machine Learning algorithms.



...

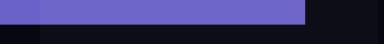
Problem Statement

Social media sites are one of the most popular websites on the internet today flooding with comments. It is vital to manage the user-generated offensive content on many of these sites that can make a user's online experience unpleasant.



...

GANTT CHART

Tasks	Week 1	Week 2	Week 3	Week 4	Week 5	Week 6	Week 7	Week 8	Week 9
Choosing a Problem Statement									
Visualizing and Pre-processing data									
Looking at various approaches									
Model Selection									
Model Training and Result Analysis									
Mid-Sem Presentation and Feedback									
Hyper-parameter fine tuning, new models train, test and result analysis									
Documentation and Presentation									

Existing Body Of Work



- WORK 1

Toxic Comment Classification on SocialMedia Using SVM and ChiSquare Feature Selection

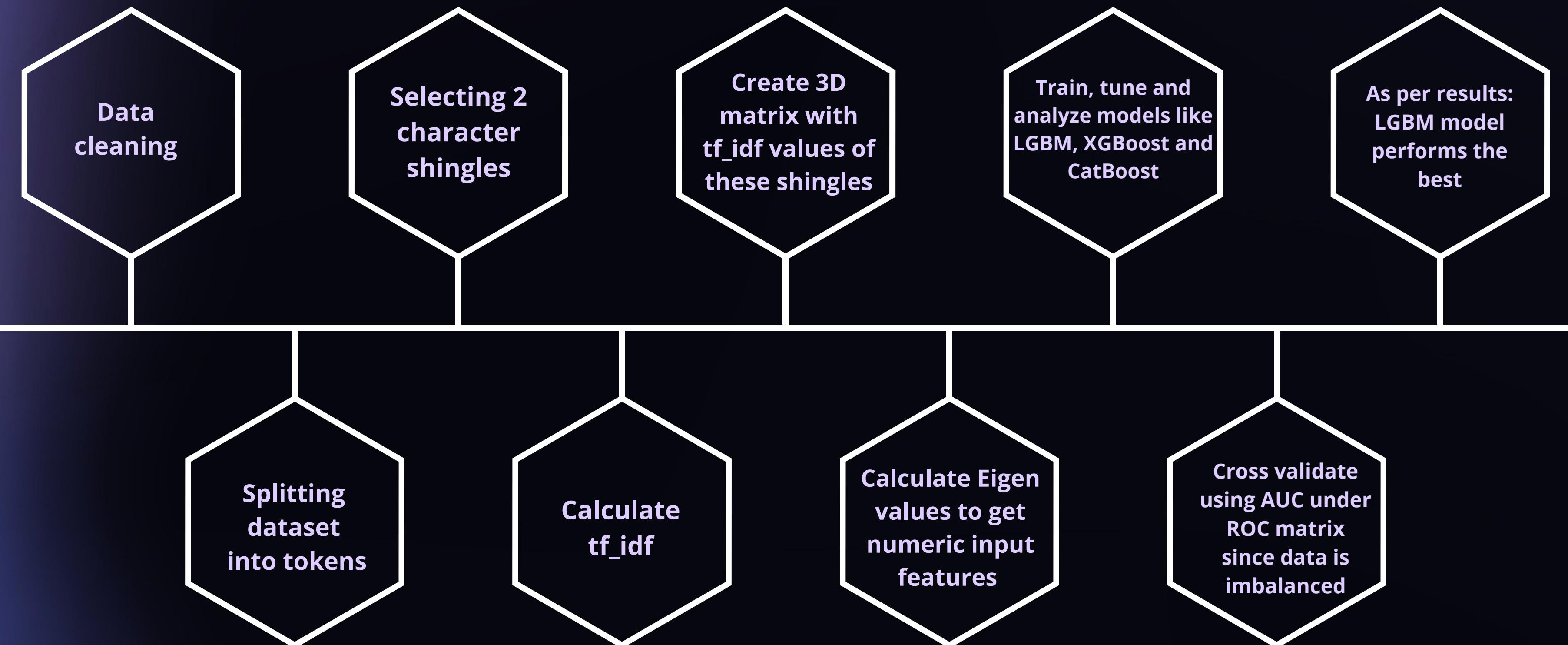
In this paper, authors have used SVM with TF-IDF as the feature extraction and Chi Square as the feature selection. The best performance obtained using the SVM model with a linear kernel, without implementing Chi Square, and using stemming and stopwords removal with the F 1 - Score equal to 76.57%.

- WORK 2

Detecting Offensive Tweets via Topical Feature Discovery over a Large Scale Twitter Corpus

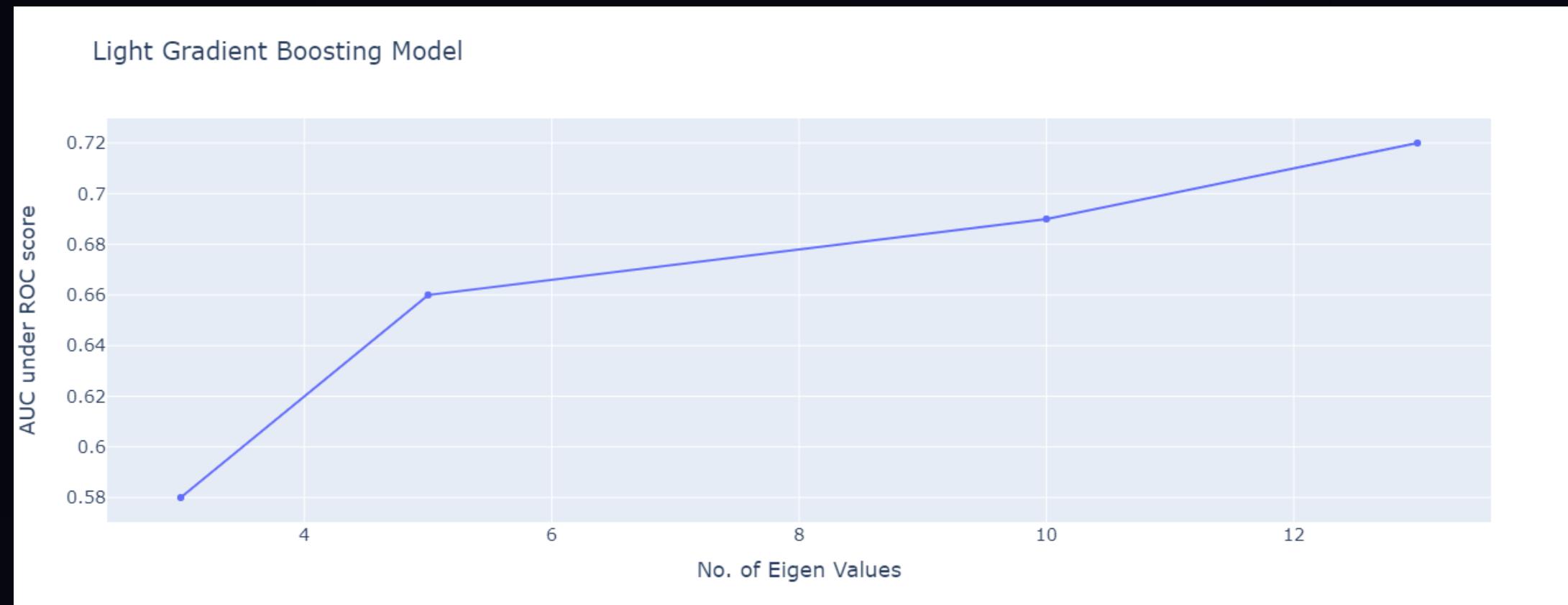
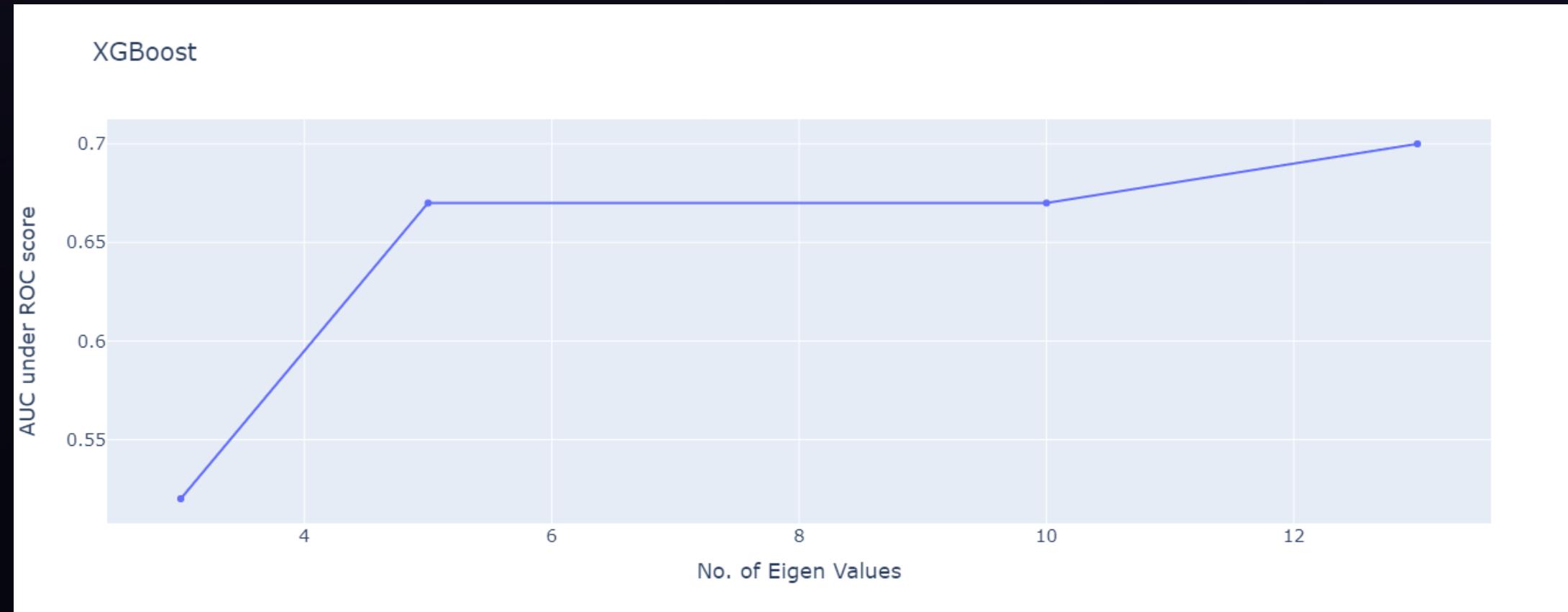
In this paper, authors have used a semi-supervised approach to detect profanity-related offensive content on Twitter. They achieved a 75.1% TP rate with Logistic Regression and a 69.7 % TP rate with popular keyword matching baseline. The false-positive rate was identical for both at about 3.77%.

Our Approach



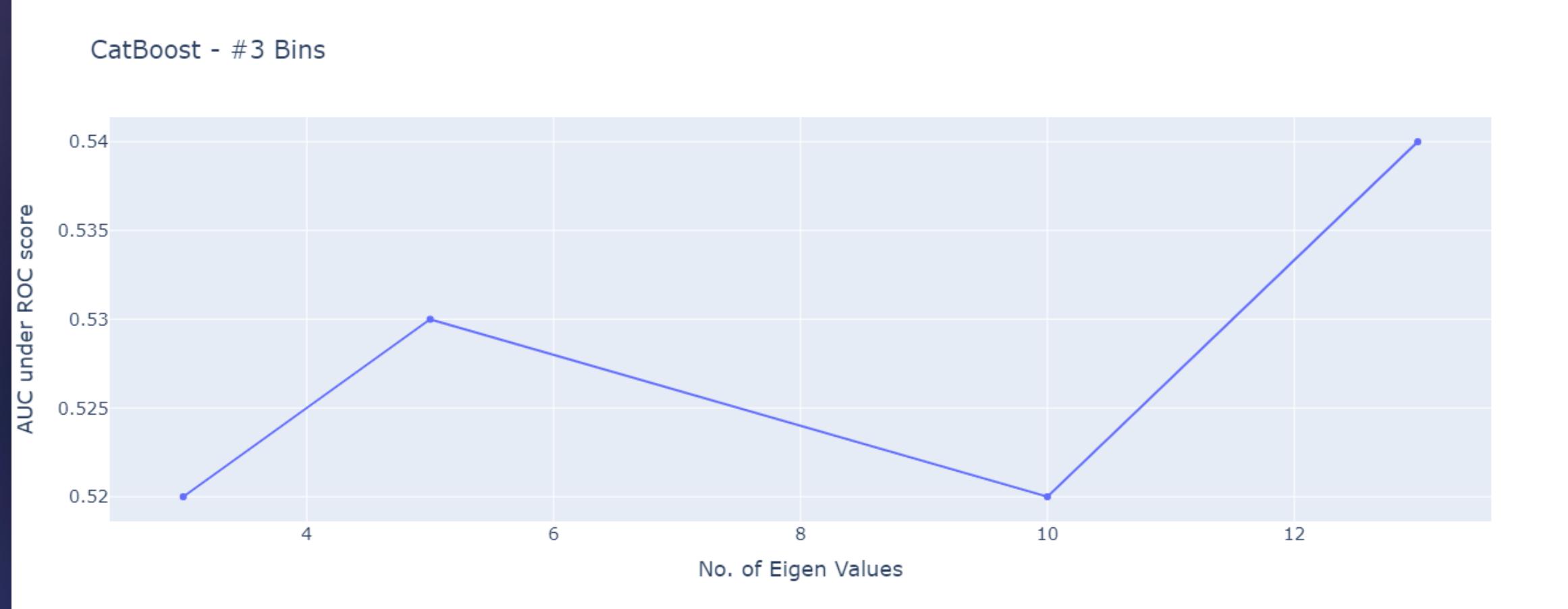
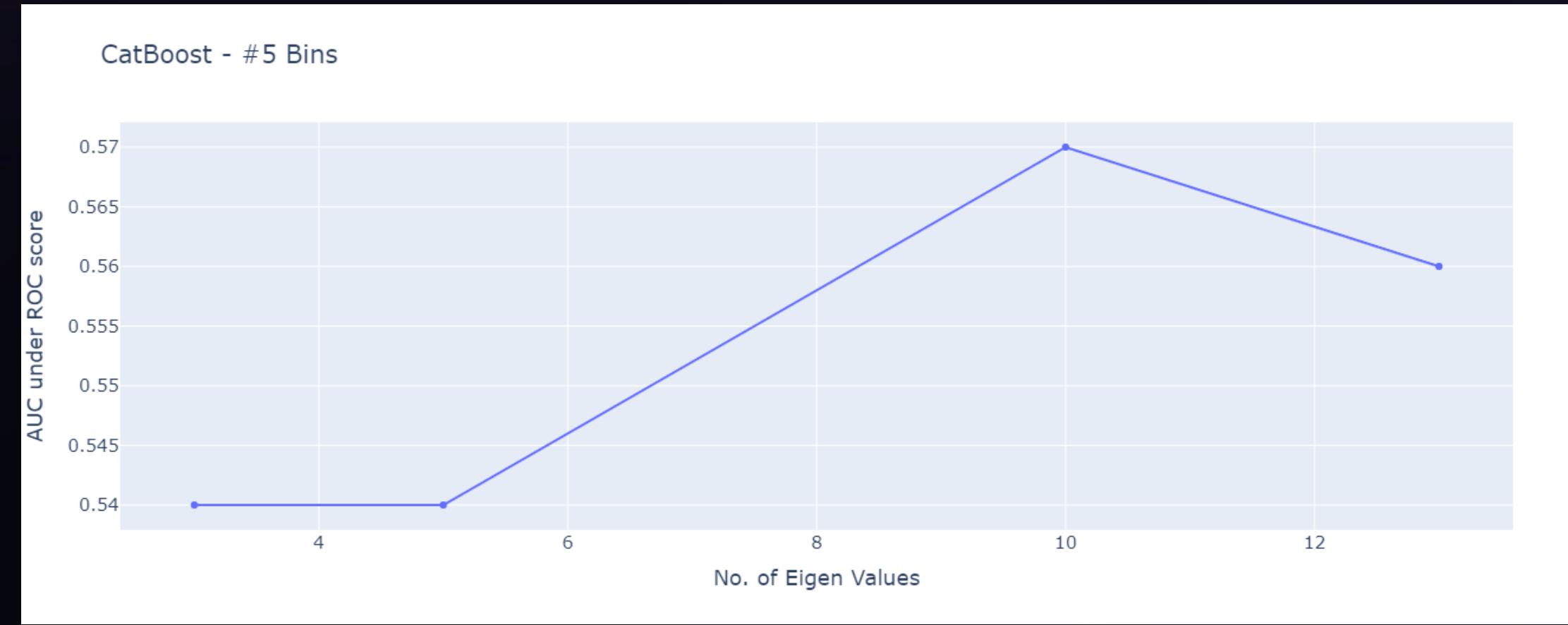
Final Results

- Applied 3 gradient boosting models amongst which LGBM outperforms with 13 eigen values with accuracy of 0.72 compared to XGBoost but underperforms in terms of training time .
- AUC Score drops on varying character shingles & n-gram for LGBM.
- Hyperparameter tuned are depth of tree, no. of leaves.



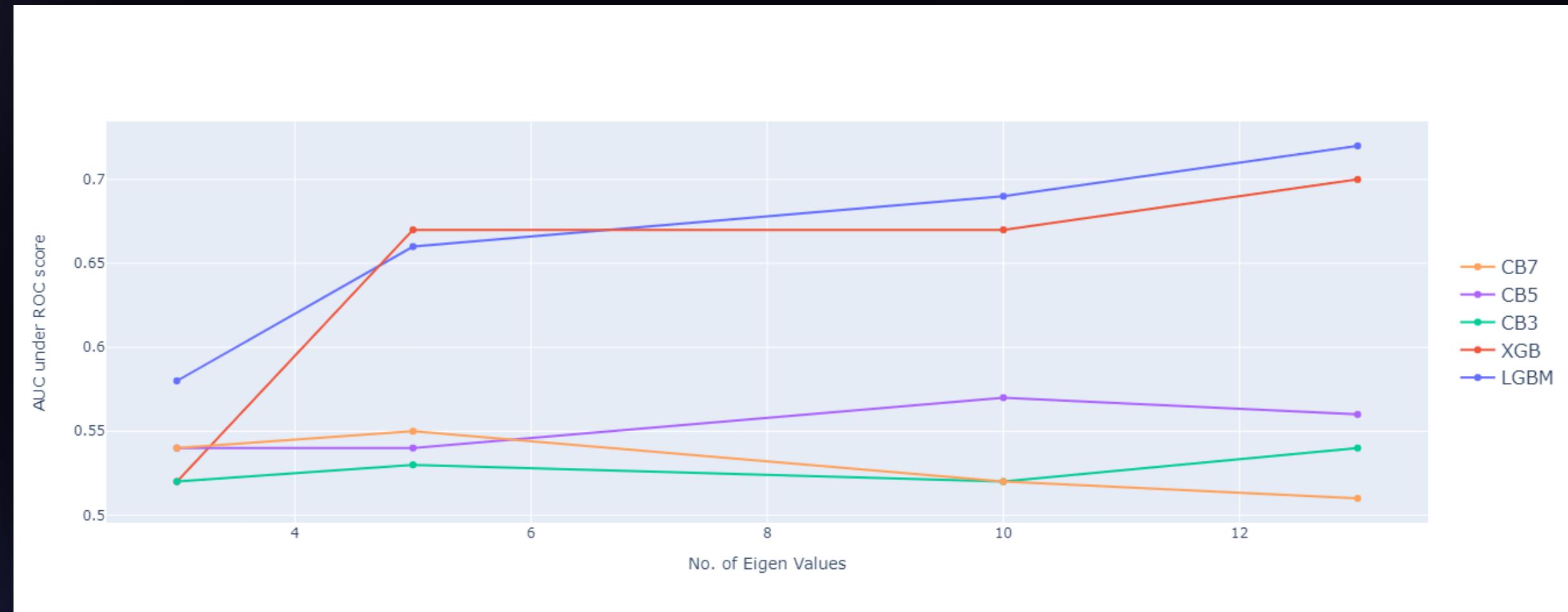
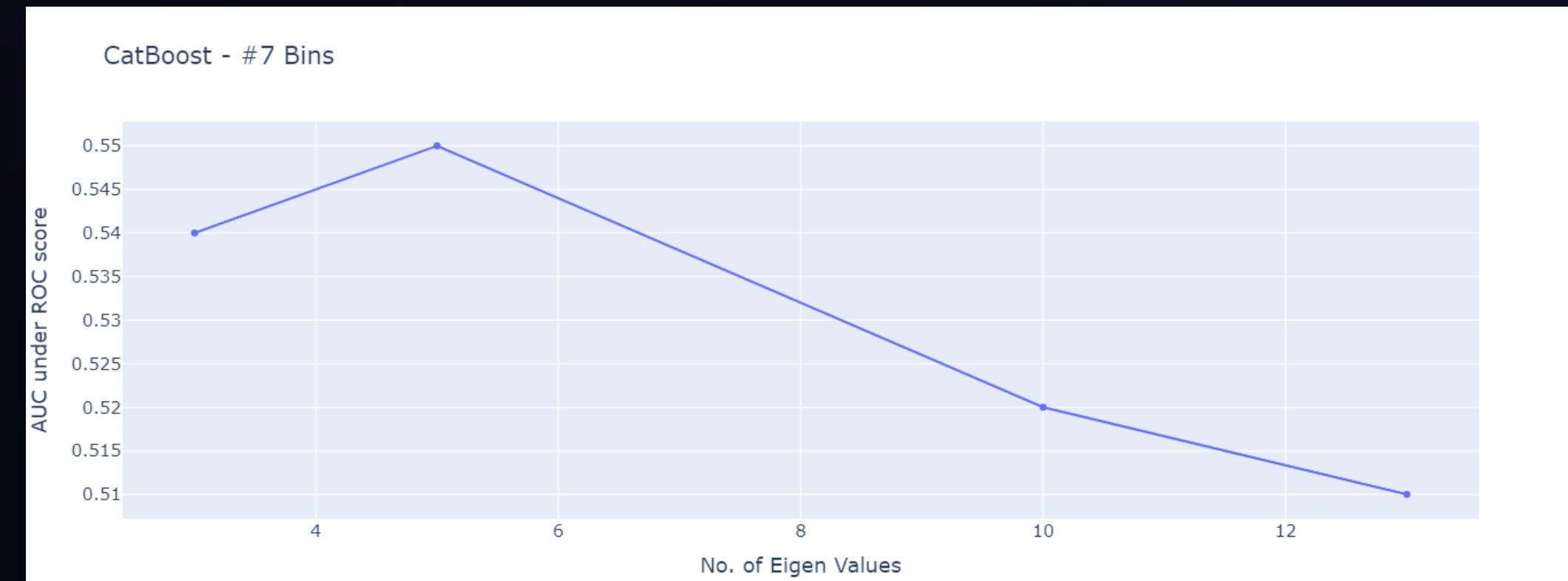
Final Results

- Increasing valuee of eignn values showed non-linear trend with AUC score around 0.5.
- Learning rate was fixed.



Final Results

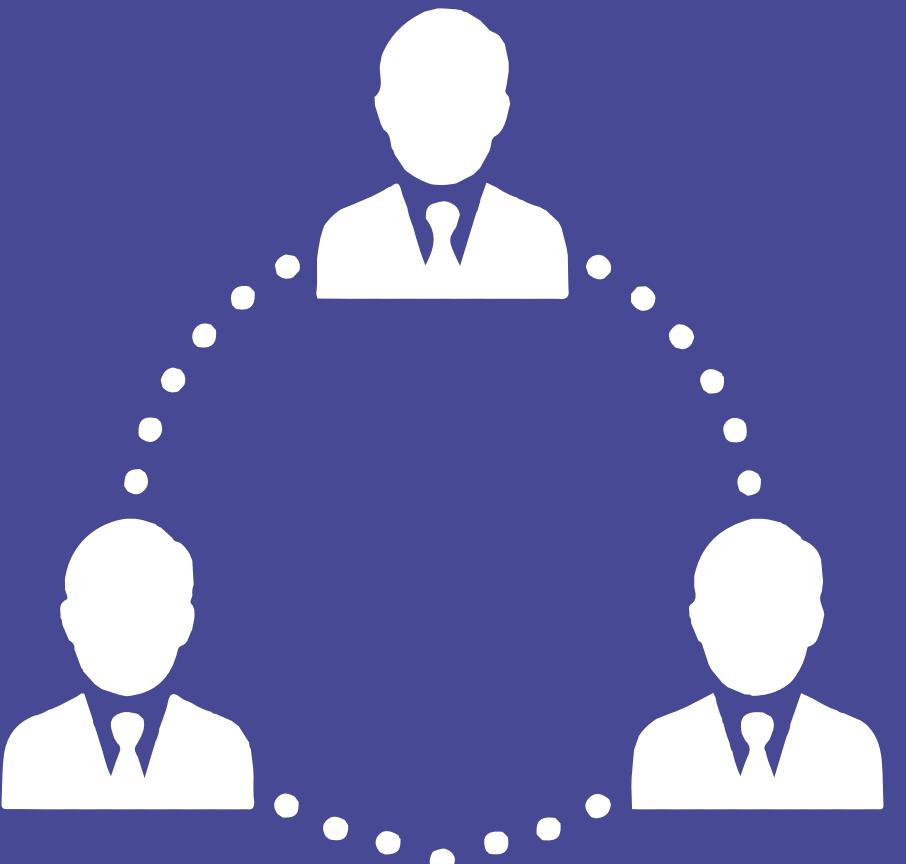
- LGBM perform best with 13 eigen vectors with less computation time of 6.5 hours & give AUC under ROC score of 0.72.
- Two character shingles works best for classification for our data.
- Threshold value for AUC curve is 0.7.



Conclusions

Model Used (Time Taken in sec & hrs)	Number of eigen values			
	3	5	10	13
LGBM (AUC)	0.58 (12600s [3.5 hrs])	0.66 (16200s [4.5 hrs])	0.69 (21600s [6 hrs])	0.72 (23400s [6.5 hrs])
XGBoost	0.52 (18000s [5 hrs])	0.67 (21600s [6 hrs])	0.67 (28800s [8 hrs])	0.70 (34200s [9.5 hrs])
CatBoost (No. of bins): 3	0.52 (10800s [3 hrs])	0.53 (16200s [4.5 hrs])	0.52 (25200s [7 hrs])	0.54 (28800s [8 hrs])
CatBoost (No. of bins): 5	0.54 (12600s [3.5 hrs])	0.54 (16200s [4.5 hrs])	0.57 (26280s [7.3 hrs])	0.56 (28080s [7.8 hrs])
CatBoost (No. of bins): 7	0.54 (10800s [3 hrs])	0.55 (14400s [4 hrs])	0.52 (25200s [7 hrs])	0.57 (27000s [7.5 hrs])

ROLE



ANERI

Implemented codes in R and ggplot for various inferences. Implementation and interpretation of base and gradient boosting models. Feature selection and hyper-parameter tuning.

MANANSHI

Implemented codes in Apache Spark for data cleaning and pre-processing. Understanding and interpreting models. Finding appropriate approach to solve the problem and formulating eigen value matrix.

NANDINI

Helped in implementing pre-processing with use of regex and apache spark. Interpretation of use of eigen values and fine tuning hyperparameters for model implementation.

KAIRAVI

Helped with implementation of models and conversion to csv files and reducing the data size. Interpretation of different eigen values and fine tuning hyperparameters for model implementation.

...

References

- "Toxic comment classification challenge," Kaggle. [Online]. Available: <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/data>. [Accessed: 20-Mar-2022].
- K. M. Lhaksmana, D. T. Murdiansyah, and N. S. Azzahra, "Toxic Comment Classification on SocialMedia Using Support Vector Machine and Chi Square Feature Selection," View of toxic comment classification on social media using support Vector Machine and Chi Square feature selection. [Online]. Available: <http://socj.telkomuniversity.ac.id/ojs/index.php/ijoict/article/view/552/316>. [Accessed: 20-Mar-2022].
- G. Xiang, B. Fan, L. Wang, J. I. Hong, and C. P. Rose, Detecting Offensive Tweets via Topical Feature Discovery over a Large Scale Twitter Corpus. [Online]. Available: <https://www.cs.cmu.edu/~lingwang/papers/sp250-xiang.pdf>. [Accessed: 20-Mar-2022].
- Van Rossum G, Drake Jr FL. Python tutorial. Centrum voor Wiskunde en Informatica Amsterdam, The Netherlands; 1995. Available: <https://www.python.org/>
- Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer- Verlag New York. ISBN 978-3-319-24277-4, 2016. Available: <https://ggplot2.tidyverse.org>
- R Core Team (2019). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. Available: <https://www.R-project.org/>.