

# Predicting Emotional Well-Being from Social Media Usage: A Machine Learning Approach

Muntasir Maruf  
Faculty of Science and Technology  
American International University-Bangladesh  
Dhaka, Bangladesh  
[muntasir.maruf26@gmail.com](mailto:muntasir.maruf26@gmail.com)

Nahid Hasan Nobil  
Faculty of Science and Technology  
American International University-Bangladesh  
Dhaka, Bangladesh  
[nobil921@gmail.com](mailto:nobil921@gmail.com)

Shusmita Anjum Aziz  
Faculty of Science and Technology  
American International University-Bangladesh  
Dhaka, Bangladesh  
[azizshusmita@gmail.com](mailto:azizshusmita@gmail.com)

Md. Nishaduzzaman Khandoker  
Faculty of Science and Technology  
American International University-Bangladesh  
Dhaka, Bangladesh  
[nishaduzzaman11031@gmail.com](mailto:nishaduzzaman11031@gmail.com)

**Abstract**— The rapid expansion of social media has significantly influenced emotional well-being, impacting users both positively and negatively. This paper aims to investigate the relationship between social media usage patterns (Posts, Likes, Comments, etc.) and emotional well-being using various machine learning models. The study leverages the "Social Media Usage and Emotional Well-Being" dataset and applies multiple Machine Learning algorithms, including Decision Tree, Random Forest, LightGBM, CatBoost, K-Nearest Neighbors, to predict emotional states based on user activity. The results demonstrate that machine learning models like Random Forest, Extra Trees and LightGBM achieved F1-Scores of 97.98%, 96.99% and 96.97%, respectively. The findings highlight how advanced machine learning techniques are crucial to determine complex relationship between social media behaviors and its impact on emotional health. This research contributes valuable insights into the data-driven approaches to enhance emotional well-being through targeted interventions and improved user experiences.

**Keywords**— *Social Media Usage, Emotional Well-Being, Machine Learning, Predictive Modeling, Mental Health.*

## INTRODUCTION

The rapid expansion of social media use has reshaped the way of communication, sharing information, and engagement with the world [1]. Platforms such as Facebook, Twitter, Instagram, and TikTok have become integral parts of everyday life, affecting not just how people relate to others but also how they perceive and express emotions [2]. Social media usage has both positive and negative impacts on psychological well-being, with excessive use linked to anxiety, loneliness, and addiction, while moderate use can foster social connectedness and support [1]. As the digital landscape continues to evolve, understanding the complicated relationship between the usage of social media and its impact on emotional well-being is becoming increasingly important [3]. This paper investigates this relationship using several machine learning approaches to predict emotional well-being from social media usage patterns, offering new insights into how online behaviors can influence mental health.

Consequently, concerns have grown over the potential negative effects of social media addiction on psychological well-being, such as anxiety, loneliness, and depression. Research has shown that excessive use, often driven by the

fear of missing out, can lead to smartphone addiction and behaviors like "phubbing," where users are distracted by their phones during in-person interactions. These behaviors can harm social connections and increase feelings of social isolation. However, some studies highlight that social media can also reduce isolation by fostering a sense of connectedness and facilitating interactions with both close and distant ties, emphasizing that the impact on well-being depends on how these platforms are used [1]. The varying results highlights that the relationship between social media usage and emotional well-being is not straightforward and is likely influenced by several features such as usage patterns, the nature of the content, and individual differences among users [4].

The background research problem addressed in this paper is predicting and understanding how different social media usage patterns affect emotional well-being. This study addresses employing machine learning models to analyze large datasets and identify the nature of them by finding patterns that contribute to users' emotional well-being. Unlike traditional methods that may struggle to interpret complex datasets, machine learning techniques can efficiently analyze large volumes of data and uncover hidden patterns and relationships, making them increasingly valuable for extracting meaningful insights from the abundance of available data [5].

This research contributes to society by providing machine learning-based models to predict emotional well-being from social media usage patterns. This paper provides a deeper understanding of how social media usage impacts mental health [6]. It explores the prediction of machine learning models, the results and discussion of findings, and concludes with suggestions, limitations, and future research directions.

## LITERATURE REVIEW

Social media has become an inseparable part of the life of this generation. Intensive use of social media in many cases is impacting the very way of living. The posts and blogs of Facebook, Instagram, Twitter, Reddit etc. are influencing the users.

A systematic review of studies which applies machine learning approaches to data in text form from social media for detecting depressive symptoms has been done summarizing

findings from previous research. Which predicts emotions using various machine learning models. Bayesian Classifier (Mean absolute error = 0.186), Random Forest (Post classification accuracy = 0.898), SVM (Post classification accuracy = 0.8) models showed promising results [7].

An overview of various sentiment analysis methods and machine learning for emotion detection from social media data aims to compare techniques and discuss limitations. Datasets were used from ISEAR, SemEval and AWS. Where most of the models are for classification problems. CNN, BiLSTM are deep learning models and performed better in large and complex datasets. Traditional Machine Learning Models such as Naïve Bayes and KNN showed decent result. CNN achieved 92.09% accuracy outperforming other ML and Deep learning models [8].

Machine Learning models were used to recognize speech in emotion using Linguistic Data Consortium (LDC) Emotional Prosody Speech Corpus. Which contains approximately 2300 utterance from 7 actor (3 males and 4 females) expressing 14 emotional states and a neutral state. Support Vector Machine, KNN and hybrid models were used to where K Nearest Neighbor performed with the best accuracy (79.5%) [9].

In a study to develop a Chatbot that can recognize and analyze human emotions through textual conversations. Using mainly conversations data from various messaging platforms. The goal is to classify emotions using machine learning model Naïve Bayes Classifier, and tools such as NLTK Vader, TextBlob, Flair and DeepMoj. The Naïve Bayes Classifier achieved emotion detection accuracy of 76% [10].

Research was conducted to create a rule-based algorithm to collect and annotate data from Twitter posts which focuses on classification of specific emotions using Plutchik's eight core emotions. Data sets used from twitter posts. Decision Tree, Random Forest and Neural Networks were used, where Neural Networks Performed best with 82% to 90% accuracy [11].

A review was performed on text-based data analyzing sentiment and detecting sentiment which used collected data from various datasets crawled from social media platforms, blogs, and e-commerce websites. Stanford Sentiment Treebank (SST), emEval Task, and other data sets were used for training and testing multiple machines and deep learning models such as Support Vector Machines (SVM), Random Forest (RF), Convolution Neural Networks (CNN), Long Short-Term Memory (LSTM), K-Nearest Neighbors, Logistic Regression, Naive Bayes, Hybrid models combining SVM and RF etc. Comparing their performances to detect and classify emotions such as joy, sadness, fear, and others effectively [12].

Ref.	Findings		
	Objectives	Models	Results
[7]	Using text from social media to detect symptoms of depression. It summarizes findings from previous research and suggests directions for future work, focusing on the challenges of sampling, prediction optimization, generalizability, privacy, and ethical issues	Bayesian Classifier, Lasso, Random Forest, SVM, CNN, LSTM	Bayesian Classifier (Mean absolute error = 0.186), Random Forest (Post classification accuracy = 0.898), SVM (Post classification accuracy = 0.8)

Ref.	Findings		
	Objectives	Models	Results
			models showed promising results.
[8]	The paper provides an overview of various sentiment analysis methodologies and machine learning approaches for emotion detection from social media data. The aim is to compare techniques and discuss limitations and future research directions.	CNN, Bi-LSTM, Naive Bayes, K-Nearest Neighbors, SVM.	, CNN outperformed others with an accuracy of 92.09% in the movie domain and 91.19% in the agriculture domain.
[9]	ML techniques for recognizing emotions from human speech. The focus is on extracting acoustic features, reducing dimensionality, and using machine learning classifiers to identify emotions conveyed in speech.	Hybrid classification methods, K-NN, Neural network, SVM	K-Nearest Neighbors (K-NN) achieved the highest accuracy of 79.5% and Neural Networks (NN) performed the worst with an accuracy of 50%
[10]	Develop a Chatbot that can recognize and analyze human emotions through textual conversations. The system is aimed at enhancing human-machine interactions by detecting emotions like joy, sorrow, irritation, and anger, and responding accordingly. This application has potential uses in social networking, business applications, and healthcare.	NLTK, TextBlob, Flair, Naive Bayes, Deepmoji	Naive Bayes classifier achieved an emotion detection accuracy of 76%, making it the best-performing model
[11]	Create a rule-based algorithm to collect and annotate data from Twitter posts, focusing on identifying specific emotions using Plutchik's eight core emotions.	Decision Tree, Random Forest, Neural Networks	Best performing models are Random Forest (accuracy: 75%-85%) and Neural Networks (accuracy: 82%-90%)
[12]	Review of machine learning algorithms analyzing sentiment from text.	Several Traditional and deep learning models and hybrid models	Naïve Bayes, decision tree, and SVM Accuracy = 81.16%, Hybrid, Accuracy = 96.75%

## METHODOLOGY

From the systematic literature review it can be noticed that most of the research conducted is speech or text based. Not much research has been done using social media behavior. Such as, daily usage time, daily post count, like count, comments count etc.

To implement models to predict emotional well-being using features were extracted from the "Social Media Usage and Emotional Well-Being" dataset available on Kaggle. Which includes these types of features. The primary goal was to evaluation and comparison of the performance of multiple machine learning models to identify which one could provide the highest accuracy for emotion detection. Several essential Python libraries were employed throughout the analysis, including Pandas and NumPy for efficient data manipulation, and Matplotlib and Seaborn for data visualization. Scikit-learn library and its preprocessing tools were used to prepare the data before training machine learning models.

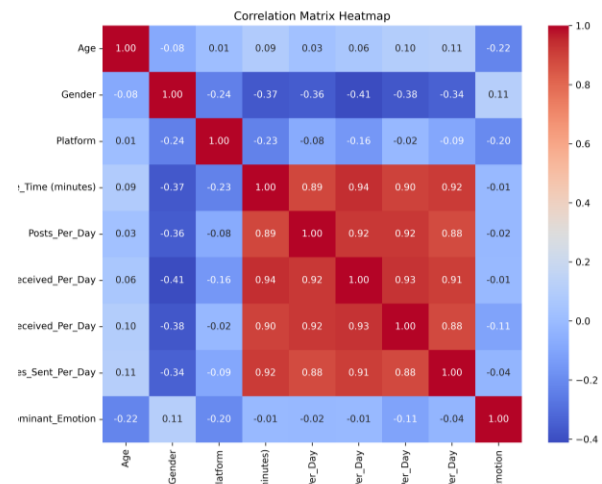
The initial step in the data analysis process involved comprehensive data preprocessing. This stage included identifying and treating missing values, as data quality issues might affect the accuracy and overall performance of the models. Mislabelled entries, particularly in critical columns such as "Age" and "Gender," were corrected to ensure consistency and accuracy. The dataset contained categorical (String) variables. To use that data to train/test machine learning algorithms, it needed to be encoded into numerical formats. Label encoding was applied for the conversion making the data compatible with model training.

The dominant emotion was encoded to, 0 – Anger, 1 – Anxiety, 2 – Boredom, 3 – Happiness, 4 – Neutral, 5 – Sadness. For gender, 0 – Female, 1 – Male, 2 – Non-Binary. And for the social media, 0 – Facebook, 1 – Instagram 2 – LinkedIn, 3 – Snapchat, 4 – Telegram, 5 – WhatsApp, 6 – Twitter.

After preprocessing, Exploratory Data Analysis (EDA) was conducted to investigate the data to find relationships between the features and the targeted classes, "Dominant Emotion." Several visualization techniques were used, including pairplots, histograms, and boxplots. Pairplots helped visualize the relationships between features, providing insight into potential patterns and correlations between variables and emotion labels. Histograms and boxplots allowed for the exploration of the distribution and spread of numerical features, while a correlation matrix was used to detect feature dependencies and collinearity. This step was essential to identify redundant or highly correlated features, which could negatively impact model performance by introducing noise or unnecessary complexity.

For the model selection phase, fifteen machine learning models were chosen to evaluate: Logistic Regression (LR), K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), Gradient Boosting (GB), AdaBoost (AB), Extra Trees (ET), Naive Bayes (NB), Quadratic Discriminant Analysis (QDA), a neural network-based Multi-Layer Perceptron (MLP), LightGBM, XG Boost (XGB), Linear Discriminant Analysis (LDA) and Cat Boost (CB). Each of these models represents different algorithmic approaches to classification, ranging from simple linear models to more complex ensemble methods. The selected dataset was divided into three files separately for train, test and validation. Due to some issues with validation dataset the training dataset provided for training was divided with 70% of the data allocated for training and 30% for validation. Testing was done by the test data. The evaluation metric used was accuracy, precision, recall, f1-score, AUC value and ROC curve, assessing how

well each model predicted the dominant emotion in the test and validation set.



## ANALYSIS

In this Research, we evaluated performance of various machine learning models, highlighting how their ability to handle non-linearity influenced their results. It is essential to first understand the nature of the dataset used for predicting emotional well-being based on social media usage.

Figure 1: Heatmap to show correlation of variables

From the heatmap, we can observe strong positive correlations between certain features, such as the daily usages (minutes) and posts per day, messages sent per day, and messages received per day. These variables exhibit highly correlated behavior. On the other hand, we see weaker correlations between features like dominant emotion and most other variables. This suggests that emotional well-being, as measured here, is more complex and less dependent on simple social media usage metrics, thus making it challenging for linear models to predict.

Given these non-linear relationships, models that excel in handling non-linear data performed better, while linear models struggled. Linear models like Logistic Regression works best for binary classification and probabilistic models like Naive Bayes assume relationships between features and the target are linearly separable and there is a conditional independence between every pair of features. Since the data exhibits clear non-linearity and overlapping class relationships, and the classification is not binary, these models may perform poorly. These models are not suitable to capture the complex relationships between features, leading to suboptimal performance [13], [14].

Decision Trees, Random Forest, and Extra Trees are tree-based models, and might perform well due to their design to manage non-linear patterns and complex interactions among features. These models utilize recursive data splitting based on feature values, enabling them to adeptly navigate the non-linear relationships between social media usage metrics and emotional well-being [15], [16].

Boosting models such as LightGBM, CatBoost, XGBoost, Gradient Boosting further enhanced performance due to their ability to combine multiple weak learners into a strong learner. Also, these models can handle categorical features and manage overfitting, also yielding high performance. These boosting models effectively handle non-linearities and complex feature interactions, making them suitable for this classification task [17], [18], [19], [20].

## RESULT & DISCUSSION

The study revealed significant differences in the performance of the models. The Random Forest model emerged as the top performer, achieving an accuracy of 97.98%, precision of 98.11%, 97.98% for recall and f1-score for test dataset. For validation it also performed best achieving an accuracy of 98.67%, precision of 98.7%, 98.67% for recall and 98.66% for f1-score.

In contrast, simpler models like Logistic Regression and the MLP classifier performed poorly, with both achieving lower accuracies. These models were less capable of capturing the intricate patterns and complex feature interactions necessary for accurate emotion detection. The relatively low performance of these models underscores their limitations in addressing non-linear relationships and feature dependencies in the dataset.

Table 1: Test Results

Models	Accuracy	Precision	Recall	F1	AUC
LR	53.54%	50.48%	53.54%	50.37%	0.84
KNN	92.93%	93.18%	92.93%	92.95%	0.97
SVM	47.47%	42.84%	47.47%	41.63%	N/A
DT	93.94%	94.22%	93.94%	93.93%	0.96
RF	97.98%	98.11%	97.98%	97.98%	1.00
GB	95.96%	96.13%	95.96%	95.96%	0.99
AB	38.38%	16.37%	38.38%	22.76%	0.75
ET	96.97%	97.14%	96.97%	96.99%	1.00
NB	39.39%	31.53%	39.39%	32.57%	0.80
LDA	56.57%	57.42%	56.57%	55.28%	0.86
QDA	63.64%	66.55%	63.64%	63.19%	0.90
XGB	94.95%	95.25%	94.95%	94.99%	1.00
LGBM	96.97%	97.04%	96.97%	96.97%	1.00
CB	95.96%	96.06%	95.96%	95.97%	1.00
MLP	55.56%	60.97%	55.56%	54.05%	0.90

Table 2: Validation Results

Models	Accuracy	Precision	Recall	F1	AUC
LR	55.81%	53.42%	55.81%	53.25%	0.85
KNN	98.01%	98.03%	98.01%	98.0%	0.99
SVM	53.16%	46.31%	53.16%	45.87%	N/A
DT	97.01%	97.07%	97.01%	97.01%	0.98
RF	98.67%	98.7%	98.67%	98.66%	1.00
GB	97.67%	97.75%	97.67%	97.660/0	1.00
AB	39.87%	30.65%	39.87%	26.19%	0.78
ET	99.0%	99.02%	99.0%	98.99%	1.00
NB	47.84%	48.96%	47.84%	41.88%	0.83
LDA	59.47%	59.02%	59.47%	58.34%	0.87
QDA	82.06%	81.98%	82.06%	81.76%	0.97
XGB	97.67%	97.75%	97.67%	97.68%	1.00
LGBM	98.01%	98.14%	98.01%	98.02%	1.00
CB	98.34%	98.38%	98.34%	98.32%	1.00
MLP	55.81%	61.91%	61.46%	60.24%	0.91

The findings suggest that gradient boosting algorithms, particularly LightGBM and CatBoost, are highly effective for emotion detection tasks involving structured data, significantly outperforming traditional classification methods. Their superior performance can be attributed to their ability to handle complex feature interactions, non-linear patterns, and structured data more efficiently than simpler models. These results highlight the importance of selecting advanced algorithms for emotion prediction, where the nature of the data involves intricate relationships between features and the target classes. Random Forest, Extra Trees and Decision Tree, are also ahead in terms of every margin. This could be because they are more suitable for our selected data set which is non liner with mixed feature type. They can also be better at handling outliers, high dimensional and complex decision boundaries further improving the result.

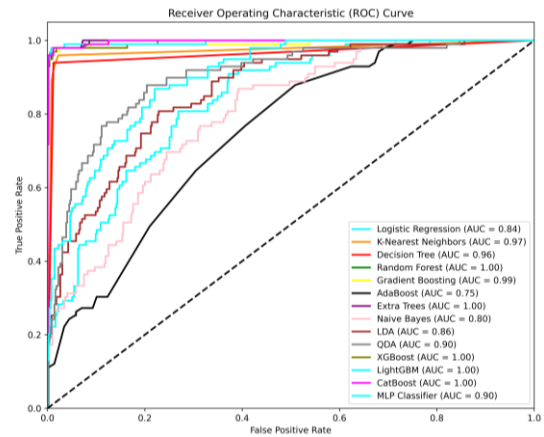


Figure 2: ROC Curve for Test Results

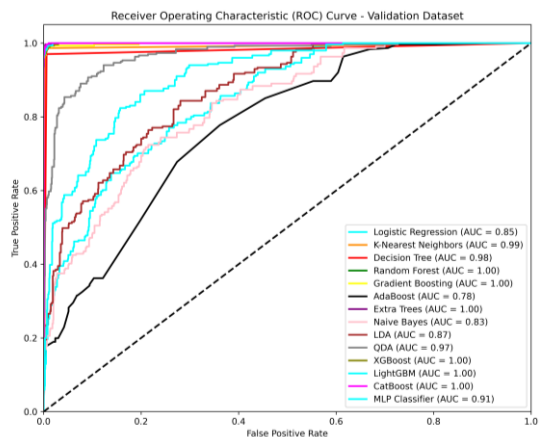


Figure 3: ROC Curve for Validation Results

From analyzing the ROC curve for both test and validation results it is noticeable that Random Forest, Extra Trees, XG Boost, Light GBM, and Cat Boost all have very high AUC scores, and their ROC curves indicates exceptional performance or possible overfitting. Which requires further research to analyze. The size of the data set could be a limiting factor. Using a larger data set might also help with the results to be more refined.

Future research could explore more advanced deep learning techniques, such as Recurrent Neural Networks (RNNs) or Transformers, which are particularly adept at capturing

temporal dependencies. Such models may prove beneficial when dealing with social media usage data, as user activity often follows temporal patterns. Additionally, incorporating new features, such as sentiment analysis or textual features extracted from social media posts, could further enhance model performance and accuracy in predicting emotional well-being. By integrating such approaches, future work could be continued improving overall performance measures and robustness of the implemented models.

## CONCLUSION

This study highlights the efficacy of machine learning models in predicting emotional well-being from social media usage patterns. The research demonstrates ML models, particularly Random Forest and other tree-based models are highly effective in handling the correlated and non-linear relationships present in social media data. These models significantly outperformed other models, providing robust predictions that can help in identifying potential emotional distress in users. By understanding the intricate patterns of social media interaction, this research opens pathways for proactive mental health interventions and improved social media designs that promote emotional well-being. However, the study is limited by the absence of temporal data and deeper sentiment analysis from user-generated content. Future work could be the exploration of several deep learning models and incorporate sentiment analysis of textual data to enhance the predictive power further. Additionally, more diverse datasets and real-time analysis could deepen our understanding of how specific social media behaviors influence psychological outcomes.

## REFERENCES

- [1] D. Ostic *et al.*, "Effects of Social Media Use on Psychological Well-Being: A Mediated Model," *Front Psychol*, vol. 12, Jun. 2021, doi: 10.3389/fpsyg.2021.678766.
- [2] D. G. Gracyal and D. Viswam, "Social Media and Emotional Well-being: Pursuit of Happiness or Pleasure," *Asia Pacific Media Educator*, vol. 31, no. 1, pp. 99–115, Jun. 2021, doi: 10.1177/1326365X211003737.
- [3] D. Smith, T. Leonis, and S. Anandavalli, "Belonging and loneliness in cyberspace: impacts of social media on adolescents' well-being," *Aust J Psychol*, vol. 73, no. 1, pp. 12–23, Jan. 2021, doi: 10.1080/00049530.2021.1898914.
- [4] E. Bulut, "Social Media Usage and Emotional Well-Being." Accessed: Sep. 10, 2024. [Online]. Available: <https://www.kaggle.com/datasets/emirhanai/social-media-usage-and-emotional-well-being>
- [5] B. Mahesh, "Machine Learning Algorithms - A Review," *International Journal of Science and Research (IJSR)*, vol. 9, no. 1, pp. 381–386, Jan. 2020, doi: 10.21275/ART20203995.
- [6] L. Braghieri, R. Levy, and A. Makarin, "Social Media and Mental Health," *American Economic Review*, vol. 112, no. 11, pp. 3660–3693, Nov. 2022, doi: 10.1257/aer.20211218.
- [7] D. Liu, X. L. Feng, F. Ahmed, M. Shahid, and J. Guo, "Detecting and Measuring Depression on Social Media Using a Machine Learning Approach: Systematic Review," *JMIR Ment Health*, vol. 9, no. 3, p. e27244, Mar. 2022, doi: 10.2196/27244.
- [8] V. Ahire and S. Borse, "Emotion Detection from Social Media Using Machine Learning Techniques: A Survey," 2022, pp. 83–92. doi: 10.1007/978-981-16-2008-9\_8.
- [9] L. Cen, M. Dong, H. L. Z. Liang Yu, and P. Ch, "Machine Learning Methods in the Application of Speech Emotion Recognition," in *Application of Machine Learning*, InTech, 2010. doi: 10.5772/8613.
- [10] Ch. Sekhar, M. S. Rao, A. S. K. Nayani, and D. Bhattacharyya, "Emotion Recognition Through Human Conversation Using Machine Learning Techniques," 2021, pp. 113–122. doi: 10.1007/978-981-15-9516-5\_10.
- [11] M. Krommyda, A. Rigos, K. Bouklas, and A. Amditis, "Emotion detection in Twitter posts: a rule-based algorithm for annotated data acquisition," in *2020 International Conference on Computational Science and Computational Intelligence (CSCI)*, IEEE, Dec. 2020, pp. 257–262. doi: 10.1109/CSCI51800.2020.00050.
- [12] P. Nandwani and R. Verma, "A review on sentiment analysis and emotion detection from text," *Soc Netw Anal Min*, vol. 11, no. 1, p. 81, Dec. 2021, doi: 10.1007/s13278-021-00776-6.
- [13] A. DeMaris and S. H. Selman, "Logistic Regression," in *Converting Data into Evidence*, New York, NY: Springer New York, 2013, pp. 115–136. doi: 10.1007/978-1-4614-7792-1\_7.
- [14] "Naive Bayes," Scikit learn. Accessed: Sep. 29, 2024. [Online]. Available: [https://scikit-learn.org/stable/modules/naive\\_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html)
- [15] "What is random forest?," IBM. Accessed: Sep. 29, 2024. [Online]. Available: <https://www.ibm.com/topics/random-forest#:~:text=Random%20forest%20is%20a%20commonly,both%20classification%20and%20regression%20problems.>
- [16] "Decision Trees," Scikit learn. Accessed: Sep. 29, 2024. [Online]. Available: <https://scikit-learn.org/stable/modules/tree.html>
- [17] "LightGBM (Light Gradient Boosting Machine)," Geeksforgeeks. Accessed: Sep. 29, 2024. [Online]. Available: <https://www.geeksforgeeks.org/lightgbm-light-gradient-boosting-machine/>
- [18] "CatBoost in Machine Learning," Javapoint. Accessed: Sep. 29, 2024. [Online]. Available: <https://www.javatpoint.com/catboost-in-machine-learning>
- [19] "XGBoost," Geeksforgeeks. Accessed: Sep. 29, 2024. [Online]. Available: <https://www.geeksforgeeks.org/xgboost/>
- [20] "Gradient Boosting in ML," Geeksforgeeks. Accessed: Sep. 29, 2024. [Online]. Available: <https://www.geeksforgeeks.org/ml-gradient-boosting/>