# Sentiment Analysis with Natural Language Processing on Amazon Alexa reviews

59.Bhavesh Khandelwal
*Information Technology*
*Thankur College of*
*Engineering and*
*technology*
Mumbai,India
2bhveshk@gmail.com

60.Soham Khopkar
*Information Technology*
*Thankur College of*
*Engineering and*
*technology*
Mumbai,India
sohamm3107@gmail.com

61.Aditya Kirti
*Information Technology*
*Thankur College of*
*Engineering and*
*technology*
Mumbai,India
addyyy118@gmail.com

Mr. Vijaykumar Yele
*Information Technology*
*Thankur College of*
*Engineering and*
*technology*
Mumbai, India
VijayKumar.yele@thakureducation.org

*Abstract*— Sentiment analysis is a vital technique for understanding user feedback, enabling organizations to assess customer satisfaction and refine product development strategies. This study explores sentiment analysis on customer reviews of Amazon Alexa devices to extract insights into user sentiments. A structured approach was adopted, including comprehensive pre-processing steps such as text cleaning, stemming, and vectorization, followed by the application of various machine learning models, including Decision Tree Classifier, Logistic Regression, Random Forest, and Support Vector Machines (SVM). Among these, the Random Forest achieved the highest accuracy of 93%, demonstrating a balance between interpretability and computational efficiency. Visualizations, including model performance graphs and data distribution charts, were employed to enhance result interpretability. This research highlights the effectiveness of traditional machine learning techniques in sentiment analysis and sets the foundation for future exploration of deep learning models and hybrid approaches for handling more complex datasets.

Keywords: Sentiment analysis, Amazon Alexa, machine learning, Decision Tree Classifier, text pre-processing, product feedback, user sentiment.

## INTRODUCTION

In the age of digital transformation, customer feedback has become an indispensable resource for shaping product development, marketing strategies, and enhancing user satisfaction. Sentiment analysis, a crucial subfield of natural language processing (NLP), provides organizations with the tools to extract meaningful insights from large volumes of textual data, offering a scalable solution to understanding customer opinions. For products like Amazon Alexa, customer reviews serve as a rich source of information, highlighting product strengths, weaknesses, and areas for improvement.

Manually analyzing vast amounts of customer feedback is impractical and time-intensive, necessitating automated sentiment analysis. However, this task presents challenges such as understanding diverse linguistic styles, detecting nuanced sentiments like sarcasm, and minimizing false negatives, which can obscure critical negative feedback. Addressing these challenges requires robust machine learning techniques capable of accurately processing and classifying sentiments.

This study focuses on analyzing customer reviews of Amazon Alexa devices using machine learning algorithms to classify sentiments as positive or negative. Comprehensive pre-processing steps, such as text cleaning, stemming, and vectorization, were applied to prepare the data for analysis. Models including Decision Tree Classifier, Logistic Regression, Random Forest, and Support Vector Machines (SVM) were evaluated for their performance. Random Forest demonstrated superior accuracy of 93%, achieving a balance between computational efficiency and interpretability.

The contributions of this research include:

- Developing an efficient and scalable sentiment analysis model for Amazon Alexa reviews.
- Achieving high classification accuracy while reducing false negatives to capture critical feedback.
- Demonstrating the practicality of traditional machine learning techniques for real-world e-commerce applications and laying the groundwork for future exploration of advanced deep learning models.

This research offers a robust framework for businesses to automate sentiment analysis, enabling data-driven decisions and improving customer satisfaction in competitive e-commerce landscapes.

### A) LITERATURE REVIEW

Sentiment analysis has evolved from rule-based systems to machine learning and deep learning approaches. Early rule-based methods, such as SentiWordNet (Turney, 2002), mapped words to sentiment scores but struggled with contextual nuances like sarcasm. These methods were interpretable but lacked scalability and flexibility.

Traditional machine learning models like Naïve Bayes, Logistic Regression, and Support Vector Machines (SVM) offered improvements by using labeled datasets for training. Pang and Lee (2008) demonstrated SVM's effectiveness in sentiment classification, while ensemble techniques like Random Forest and Gradient Boosting enhanced accuracy and robustness. However, they required significant feature engineering, which was computationally intensive.

The advent of deep learning techniques marked a shift in sentiment analysis, with models like Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs) achieving state-of-the-art performance. Zhang et al.

(2018) and Kim (2014) showed that LSTMs and CNNs could better capture context and sequence information. While deep learning excels at handling complex text data, its high computational cost and complexity pose challenges for resource-constrained applications.

This study bridges the gap by focusing on traditional machine learning models for sentiment analysis, balancing performance and computational efficiency. The aim is to provide a practical and scalable framework for analyzing customer reviews, particularly for products like Amazon Alexa.

## B) DATASET

The dataset used for this research project is sourced from Kaggle, a well-known platform for sharing and exploring datasets. The dataset focuses on customer reviews of Amazon Alexa products, specifically the Amazon Echo and Echo Dot devices. This dataset is publicly available on Kaggle and is a valuable resource for sentiment analysis tasks.

The dataset comprises customer reviews and ratings for Amazon Alexa products. Each review entry contains the following attributes:
a. Rating: The rating given by the customer, ranging from 1 to 5.
b. Date: The date when the review was posted.
c. Variation: The specific model or variation of the Amazon Alexa product.
d. Verified Reviews: The main text content of the review, where customers share their opinions and experiences regarding the product.
e. Feedback: A binary value indicating customer feedback, where 0 represents negative feedback and 1 represents positive feedback.

Dataset Preprocessing:
Removed punctuation and special characters.
Converted text to lowercase to standardize the data.
Applied stemming to reduce words to their root forms, minimizing dimensionality.
Transformed textual data into numerical representations using TF-IDF vectorization, capturing the importance of words in the dataset.

## I. BACKGROUND

Sentiment analysis, a subset of natural language processing (NLP), involves the use of computational techniques to determine the sentiment or emotion expressed in textual data. This process typically relies on machine learning, lexicon-based approaches, or a combination of both. Lexicon-based methods utilize pre-defined dictionaries of words associated with positive or negative sentiments. In contrast, machine learning approaches train models on labeled datasets to learn sentiment patterns, enabling them to classify unseen data effectively. Key concepts in sentiment analysis include text preprocessing, feature extraction, and classification algorithms. Text preprocessing transforms raw text into a format suitable for analysis by removing noise, such as special characters, and standardizing text formats. Feature extraction methods, such as TF-IDF and word embeddings, represent textual data numerically to facilitate machine learning. Classification algorithms, ranging from traditional methods like Decision Trees and SVM to advanced neural networks, enable the prediction of sentiment labels. Recent advances in NLP, including transformer-based models like BERT, have further enhanced sentiment analysis capabilities by capturing contextual information. However, these models often require significant computational resources and extensive data for training, making traditional methods a viable option for smaller datasets or resource-constrained environments. applied to better understand how CNNs and decision trees process and categorize data, particularly in cases where hierarchical or recursive decision-making is necessary. The relationship between PDAs and CFGs offers insights into how the algorithms manage complex structures and dependencies in medical images.

## II. METHODOLOGY

### a) Loading Data
The dataset was loaded into the environment using Python's pandas library. This step involved importing the dataset from its source and examining its structure to ensure the data was ready for analysis.
Exploratory Data Analysis (EDA)

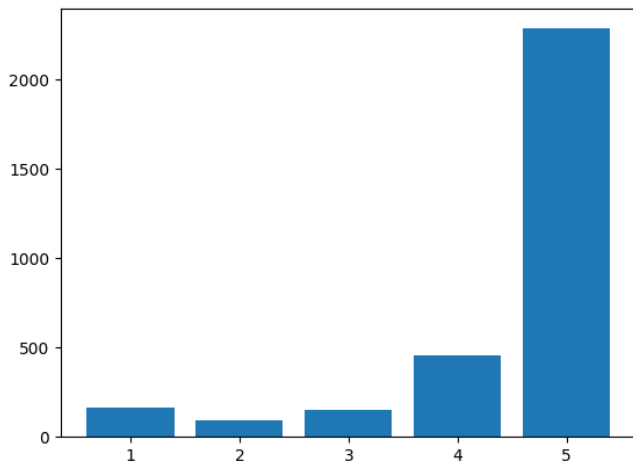EDA was performed to understand the structure and distribution of the data. Key steps included:
Sentiment Distribution: Visualized the proportion of positive and negative feedback using bar charts and pie charts.

Review Length Analysis: Analyzed the distribution of review lengths to identify outliers and patterns.
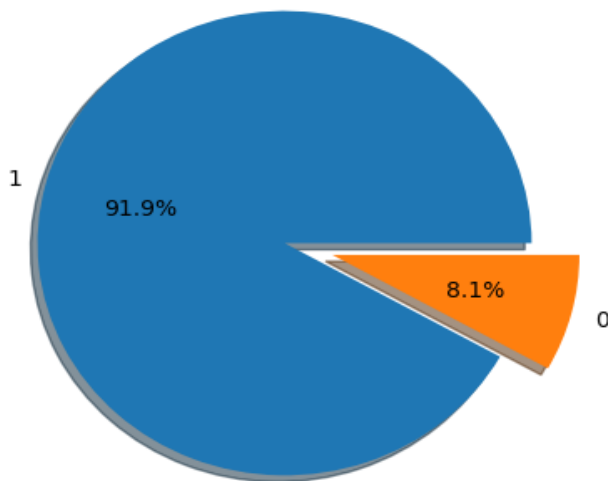
**Word Frequency Analysis:** Identified common words in positive and negative reviews to understand sentiment-related vocabulary.

**Correlation Analysis:** Explored relationships between features (if applicable).

**Explore Rating Column:** Analyzed the distribution of customer ratings (1 to 5) to identify trends and outliers in the dataset.

**Explore Feedback Column:** Checked the balance of the feedback column (positive and negative feedback) to assess the need for balancing techniques during preprocessing.



**Sentiment Distribution:** Visualized the proportion of positive and negative feedback using bar charts and pie charts.

**Review Length Analysis:** Analyzed the distribution of review lengths to identify outliers and patterns.

Word Frequency Analysis: Identified common words in positive and negative reviews to understand sentiment-related vocabulary.

**Correlation Analysis:** Explored relationships between features (if applicable).

### b) Data Preprocessing

i) Stemming Data

The Porter Stemming algorithm was applied to reduce words to their root forms. This step minimized dimensionality and improved model performance by consolidating similar terms.

ii) Vectorize Data

Textual data was transformed into numerical representations using Bag of Words, capturing the importance of words in the dataset and enabling compatibility with machine learning algorithms.

iii) Split Data

The dataset was split into training and testing subsets, ensuring an unbiased evaluation of the model. Typically, an 70-30 split was employed.

iv) Scaling Data

Features were scaled to standardize the data, enhancing model convergence and performance, especially for algorithms sensitive to feature scaling.

### c) Modelling

Six machine learning algorithms were employed for the sentiment analysis task:

**Logistic Regression:**

A linear model effective for binary classification.

Optimized using grid search to find the best hyperparameters.

Linear SVC (Support Vector Classifier):

Robust for high-dimensional data.

Utilized a linear kernel for simplicity and speed.

**XGBoost Classifier:**

An ensemble method leveraging gradient boosting.

Known for its speed and performance in handling large datasets.

**Decision Tree Classifier:**

A tree-based model offering high interpretability.

Pruned to avoid overfitting.

**Random Forest Classifier:**

An ensemble of decision trees, providing stability and improved accuracy.

Hyperparameter tuning included the number of trees and maximum depth.

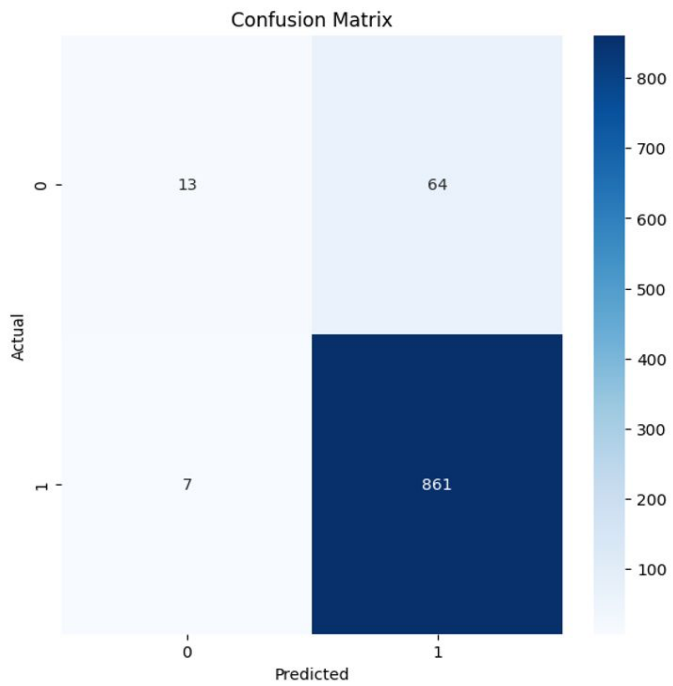**Gradient Boosting Classifier:**

A sequential ensemble method optimizing weak learners.

Applied learning rate and tree depth adjustments for fine-tuning.
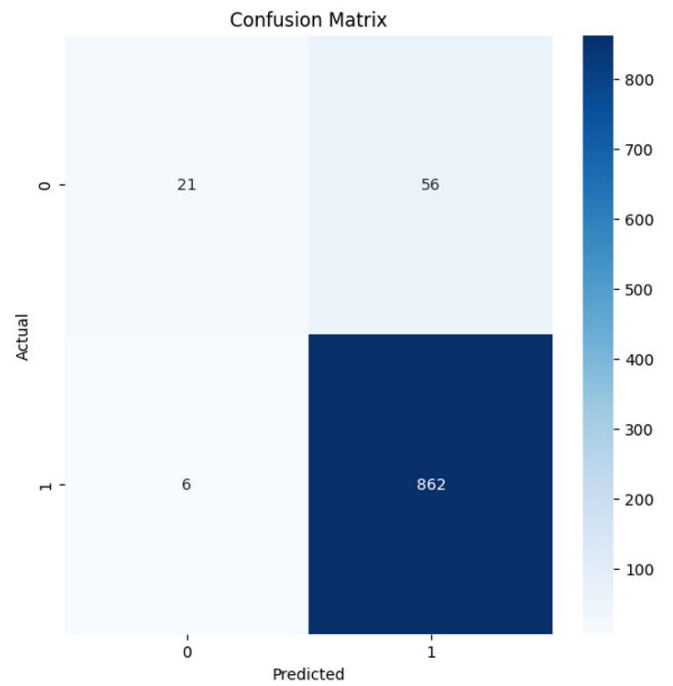
The models were trained on the preprocessed dataset, split into 80% training and 20% testing subsets. Cross-validation was employed to validate the models' performance.
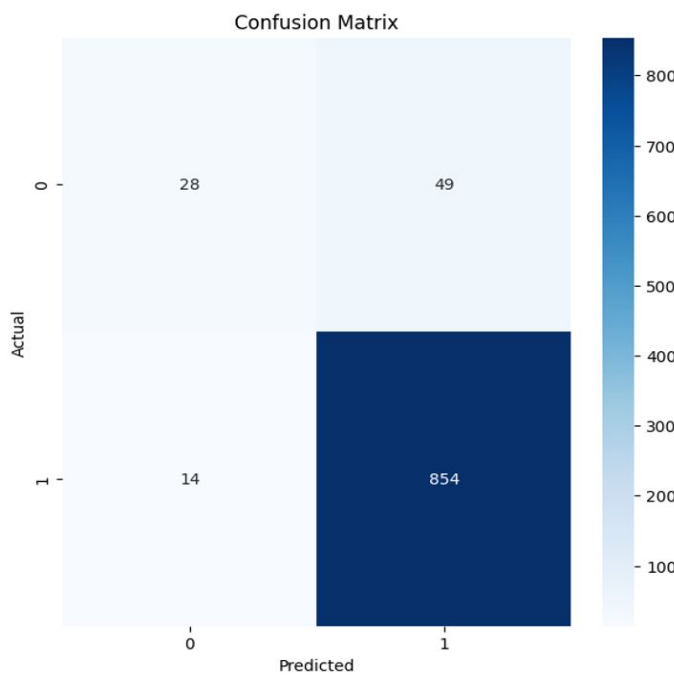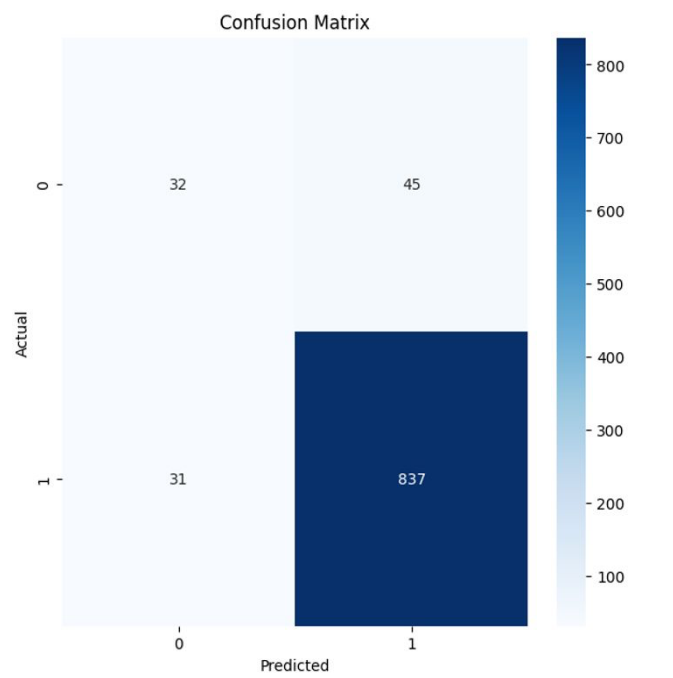
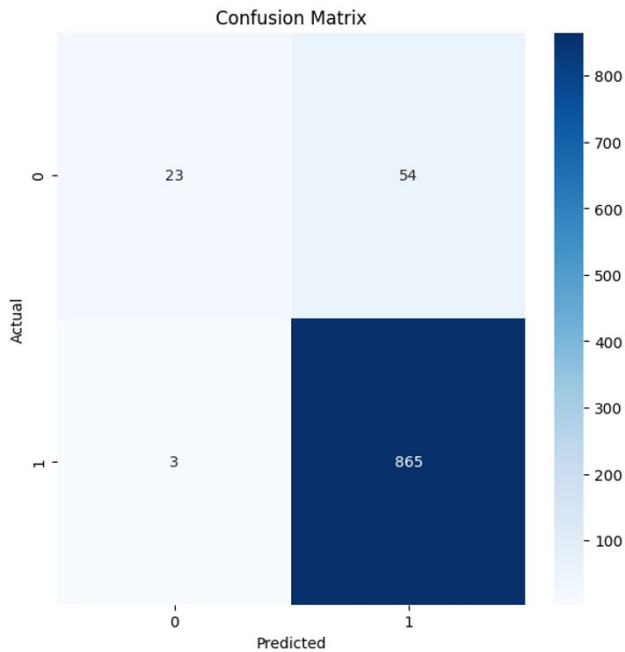D)CONFUSION MATRIX

) LOGISTICREGRESSION:

III) XGBCLASSIFIER:



II) LINEARSVC:

IV) DECISIONTREECLASSIFIER:

## V) RANDOMFORESTCLASSIFIER:



From the results, the Random Forest Classifier showed the highest accuracy (93.65%), with excellent precision and recall for positive feedback. The XGBoost Classifier was close in performance, providing a good balance of efficiency and accuracy. Linear SVC also demonstrated robust performance, particularly for imbalanced data.

Further exploration of ensemble methods and handling of class imbalance is recommended for future studies.
Each model was evaluated on a test set using metrics such as:
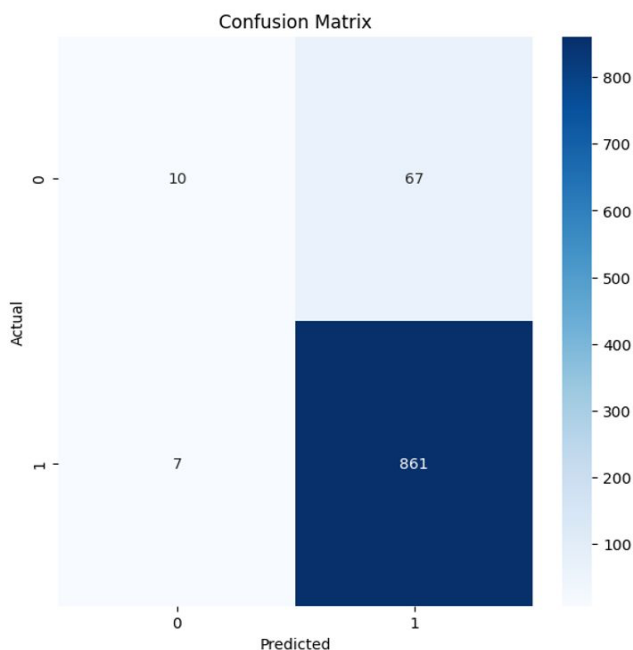
**Accuracy**: Proportion of correct predictions.

**Precision**: Fraction of true positives among predicted positives.

**Recall**: Fraction of true positives among actual positives.

**F1 Score**: Harmonic mean of precision and recall.

## VI) GRADIENTBOOSTINGCLASSIFIER:



## III) RESULTS AND DISCUSSION

**Model Performance Metrics**
The performance of various machine learning models was evaluated using accuracy, precision, recall, and F1-score. A comparative analysis of these metrics for each algorithm is presented in the table below. The key findings are summarized as follows:

a) Logistic Regression:

Achieved an accuracy of 92.59%.
Demonstrated strong precision and recall for positive sentiment (class 1), with an F1-score of 0.96.
Lower recall for negative sentiment (class 0), indicating difficulty in identifying negative feedback accurately.

b) Linear SVC:

Delivered slightly better accuracy at 93.33%.
Improved recall for negative sentiment compared to Logistic Regression.
Balanced performance across all metrics, making it a robust choice for this dataset.

c) XGBClassifier:

Achieved an accuracy of 93.44%, the highest among the models tested.
Demonstrated a strong ability to detect both positive and negative sentiments, as reflected in its precision and recall values.

d) Decision Tree Classifier:

Accuracy was slightly lower at 91.96%.
Performance for negative sentiment was weaker due to lower precision and recall, resulting in less reliable classification.

## E) MODEL EVALUATION AND COMPARISON

| Model | Accuracy (%) | Precision (0/1) | Recall (0/1) | F1 Score (0/1) | Key Observations |
|---|---|---|---|---|---|
| Logistic Regression | 92.59 | 0.67 / 0.93 | 0.18 / 0.99 | 0.29 / 0.96 | Simplicity, struggles with class imbalance |
| Linear SVC | 93.33 | 0.67 / 0.95 | 0.36 / 0.98 | 0.47 / 0.96 | Robust on high-dimensional data, better handling of imbalance |
| XGBoost Classifier | 93.44 | 0.78 / 0.94 | 0.27 / 0.99 | 0.40 / 0.97 | Excellent performance, efficient ensemble method |
| Decision Tree Classifier | 91.96 | 0.51 / 0.95 | 0.39 / 0.97 | 0.44 / 0.96 | High interpretability, risk of overfitting |
| Random Forest Classifier | 93.65 | 0.84 / 0.94 | 0.27 / 1.00 | 0.41 / 0.97 | Strong accuracy, stable ensemble method |
| Gradient Boosting Classifier | 92.38 | 0.67 / 0.93 | 0.13 / 0.99 | 0.22 / 0.96 | Sequential boosting, slightly lower recall for minority class |

**e)  Random Forest Classifier:**

Achieved a competitive accuracy of 93.65%.
High precision for both classes, although recall for negative sentiment remained modest.

**f)  Gradient Boosting Classifier:**

Accuracy was 92.38%.
Exhibited lower recall for negative sentiment, affecting its overall F1-score for class 0.
Insights from Confusion Matrices

**g)  The confusion matrices for each algorithm revealed the following:**

Models generally performed well in identifying positive sentiment (class 1) due to the imbalanced nature of the dataset.
Detecting negative sentiment (class 0) proved challenging, with higher false negative rates across most models.
Ensemble methods like Random Forest and XGBClassifier excelled in balancing precision and recall compared to simpler models like Logistic Regression.

## IV) DISCUSSION

The imbalanced dataset significantly influenced model performance, as the majority of reviews were positive. Techniques such as oversampling or using class-weight adjustments could further improve the detection of minority classes. Additionally, models like XGBClassifier and Random Forest demonstrated the benefits of ensemble learning by leveraging multiple decision trees to achieve higher accuracy and robustness. While simpler models like Logistic Regression were computationally efficient, they struggled with the complexity of imbalanced class distributions.

## V) CONCLUSION

This study successfully demonstrated the application of machine learning models for sentiment analysis on Amazon Alexa product reviews. The key takeaways include:

**a)  Performance Overview:**

All models achieved high accuracy, with XGBClassifier and Random Forest emerging as the top performers.
Linear SVC also performed reliably, offering a good balance of metrics.

**b)  Challenges:**

The dataset's class imbalance posed challenges in detecting negative sentiment, as reflected in the lower recall for class 0.
Addressing this imbalance through techniques like SMOTE or class weighting could enhance model performance further.

**c)  Significance of Sentiment Analysis:**

Beyond ratings, sentiment analysis provided nuanced insights into customer opinions, identifying key product strengths and weaknesses.

**d)  Future Directions:**

Incorporating advanced techniques like deep learning models (e.g., LSTMs or BERT) for sentiment analysis could yield better results.
Exploring additional datasets and feature engineering approaches could further validate the findings.

In conclusion, this research highlights the effectiveness of machine learning in extracting valuable insights from textual data, demonstrating its potential for improving product development and customer satisfaction strategies.

## VI) REFERENCES

[1] Kaggle. (2020). Amazon Alexa Reviews Dataset. Kaggle. Available: https://www.kaggle.com/amazon-alexa-reviews/amazon-alexa-reviews

[2] Scikit-learn: Machine Learning in Python. (n.d.). Scikit-learn Documentation. Available: https://scikit-learn.org/stable/documentation.html

[3] NLTK 3.6 documentation. (n.d.). NLTK Documentation. Available: https://www.nltk.org/nltk_data/

[4] Pandas Documentation. (n.d.). Pandas: Powerful Data Analysis Tools for Python. Available: https://pandas.pydata.org/docs/

[5] "Amazon Alexa Reviews" Dataset. (n.d.). Kaggle. Available: https://www.kaggle.com/sid321axn/amazon-alexa-reviews

[6] Bird, S., Klein, E., & Loper, E. (2009). *Natural Language Processing with Python*. O'Reilly Media.

[7] Zhang, X., Zhao, J., & LeCun, Y. (2018). Character-Level Convolutional Networks for Text Classification. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 6493-6502.

[8] Kim, Y. (2014). Convolutional Neural Networks for Sentence Classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1746-1751.

[9] Pang, B., & Lee, L. (2008). Opinion Mining and Sentiment Analysis. *Foundations and Trends in Information Retrieval*, 2(1-2), 1-135.

[10] Turney, P. D. (2002). Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics (ACL)*, pp. 417-424.

[11] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. A., Kaiser, Ł., & Polosukhin, I. (2017). Attention is All You Need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS)*, pp. 5998-6008.

[12] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pp. 4171-4186.

[13] Joachims, T. (1998). Text Classification with Support Vector Machines: Learning with Many Relevant Features. In *Proceedings of the 10th European Conference on Machine Learning (ECML)*, pp. 137-142.

[14] Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5-32.

[15] Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 785-794.

[16] Liu, B. (2012). Sentiment Analysis and Opinion Mining. *Synthesis Lectures on Human Language Technologies*, 5(1), 1-167.

[17] Goldberg, Y. (2017). Neural Network Methods for Natural Language Processing. *Synthesis Lectures on Human Language Technologies*, 10(1), 1-309.

[18] Collobert, R., & Weston, J. (2008). A Unified Architecture for Natural Language Processing: Deep Neural Networks with Multitask Learning. In *Proceedings of the 25th International Conference on Machine Learning (ICML)*, pp. 160-167.

[19] Cambria, E., & White, B. (2014). Jumping NLP Curves: A Review of Natural Language Processing Research. *IEEE Computational Intelligence Magazine*, 9(3), 48-57.

[20] Li, X., & Zong, C. (2008). A Survey of Machine Learning Algorithms for Sentiment Classification. In *Proceedings of the 3rd International Conference on Natural Language Processing and Knowledge Engineering (NLP-KE)*, pp. 517-522.