# Predictive Modeling of App Ratings: A Review of Classical and Neural Network Methodologies.

Nivedita[1], Vivek Sharma[2]

[1]MTech Scholar, Department of Computer Science and Engineering, Technocrats Institute of Technology, Bhopal, India

[2]Professor, Department of Computer Science and Engineering, Technocrats Institute of Technology, Bhopal, India

**Abstract:** The Google Play Store, a dominant platform for Android applications, generates vast amounts of data including app descriptions, user reviews, and ratings. Analyzing this data is crucial for developers to understand user preferences, improve app quality, and predict future ratings. This review paper synthesizes recent research on Play Store app analysis and rating prediction, focusing on the application of classical machine learning (ML) models and artificial neural networks (ANNs). We discuss various feature engineering techniques, model selection strategies, and performance evaluation metrics, highlighting the strengths and limitations of different approaches. This paper provides a comprehensive overview of the current state of research and identifies potential avenues for future exploration.

## 1. Introduction

The digital landscape has been irrevocably transformed by the widespread adoption of smartphones, leading to an unprecedented surge in the availability of mobile applications. App stores, particularly the Google Play Store for Android devices, have become central hubs for accessing and distributing these applications. The sheer volume of apps available necessitates effective mechanisms for users to discern quality and relevance. Among these mechanisms, user ratings and reviews stand out as crucial indicators of app performance and user satisfaction [1]. These user-generated data points not only reflect the immediate experience of individuals but also collectively shape the perception and trajectory of an app's success. Consequently, the ability to accurately predict app ratings has emerged as a significant area of research, offering valuable insights for developers, marketers, and users alike.

The importance of accurate rating prediction stems from its potential to empower developers in numerous ways. By anticipating user feedback, developers can proactively identify areas for improvement, prioritize feature enhancements, and address potential pain points before they escalate into widespread dissatisfaction. This proactive approach can lead to higher user retention, increased app engagement, and ultimately, greater commercial success. Furthermore, accurate rating predictions can inform marketing strategies, allowing developers to target specific user segments with tailored campaigns and messaging.

In the realm of research, the vast and dynamic dataset provided by the Google Play Store has become a fertile ground for the application of machine learning (ML) and deep learning techniques [2]. These techniques enable the extraction of meaningful patterns and insights from the complex interplay of app descriptions, user reviews, installation data, and other relevant features. Classical ML models, such as linear regression, support vector regression (SVR), random forests (RF), and gradient boosting machines (GBMs), have demonstrated their efficacy in capturing linear and non-linear relationships within the data. These models offer advantages in terms of interpretability and computational efficiency, making them suitable for handling large datasets and providing clear insights into the factors influencing app ratings.

Alongside classical ML approaches, artificial neural networks (ANNs), particularly deep learning models, have gained prominence due to their ability to learn intricate, hierarchical representations from complex data. Multilayer perceptrons (MLPs), convolutional neural networks (CNNs), and recurrent neural networks (RNNs), including long short-term memory (LSTM) networks, have been employed to analyze textual data, extract sentiment, and model sequential dependencies within user reviews and app descriptions. The ability of deep learning models to capture nuanced semantic information and long-range dependencies has led to significant improvements in prediction accuracy.

This review paper aims to provide a comprehensive overview of recent research endeavors focused on Play Store app analysis and rating prediction. By synthesizing findings from various studies, this paper seeks to highlight the diverse approaches, methodologies, and techniques employed in this domain. A particular emphasis is placed on the application of classical ML models and ANNs, examining their respective strengths, limitations, and suitability for different

aspects of app rating prediction. Furthermore, this review will delve into the critical aspects of data preprocessing, feature engineering, and performance evaluation, providing a holistic understanding of the research landscape. By consolidating and analyzing the current state of research, this paper intends to serve as a valuable resource for researchers, developers, and practitioners seeking to leverage machine learning for app rating prediction and analysis. The paper will cover the preprocessing of text, and numerical data, the implementation of various classical and deep learning models, and the common metrics used to evaluate the performance of these models.

## 2. Data Preprocessing and Feature Engineering

The raw data extracted from the Google Play Store, while rich in information, is inherently unstructured and often noisy. To effectively train machine learning models for tasks like app rating prediction, user review analysis, or app category classification, a rigorous preprocessing and feature engineering phase is indispensable. This stage transforms the raw data into a clean, structured, and informative format, enabling the extraction of meaningful insights and the development of robust predictive models.

**Text Preprocessing:** A significant portion of Play Store data is textual, including app descriptions and user reviews. These texts carry valuable semantic information but require substantial cleaning and normalization. The initial step involves **tokenization**, which breaks down the text into individual words or sub-word units. Subsequently, **stop-word removal** eliminates common words like "the","a," and "is," which contribute little to the semantic meaning. Further normalization is achieved through **stemming** and **lemmatization**. Stemming reduces words to their root form by removing suffixes, while lemmatization considers the word's context and reduces it to its dictionary form. These techniques help standardize the vocabulary and reduce the dimensionality of the text data [3].

**Sentiment Analysis:** User reviews provide direct insights into user satisfaction and app performance. **Sentiment analysis** is crucial for extracting subjective opinions from these reviews. Techniques like **VADER (Valence Aware Dictionary and sEntiment Reasoner)** and **TextBlob** leverage pre-built lexicons to assign sentiment scores (positive, negative, or neutral) to

text segments. Machine learning-based sentiment classifiers, trained on labeled datasets, can also be employed for more nuanced sentiment detection. These classifiers can capture complex sentiment patterns and context-dependent expressions. The extracted sentiment scores become valuable numerical features, reflecting user perceptions and preferences [4].

**Categorical Feature Encoding:** The Play Store data contains various categorical features, such as app category (e.g., Games, Education, Productivity) and content rating (e.g., Everyone, Teen, Mature 17+). Machine learning models, in general, require numerical inputs. Therefore, these categorical features must be converted into numerical representations. **One-hot encoding** creates binary columns for each category, indicating the presence or absence of a particular category. While effective, it can lead to high dimensionality with numerous categories. **Label encoding** assigns a unique integer to each category. This approach is suitable for ordinal categorical features where there is an inherent order. However, it can introduce unintended ordinal relationships in nominal categorical features. The choice of encoding technique depends on the nature and characteristics of the categorical variables [5].

**Numerical Feature Scaling:** Numerical features, such as app size and the number of installs, often exhibit different scales and ranges. This can negatively impact the performance of machine learning algorithms, particularly those sensitive to feature scales, like gradient descent-based models. **Min-max scaling** transforms numerical features to a specific range, typically [0, 1]. **Standardization** centers the data around the mean and scales it to unit variance. Both techniques help normalize the numerical features, ensuring that no single feature dominates the learning process due to its magnitude. The selection of scaling technique depends on the data distribution and the requirements of the chosen machine learning algorithm [6].

**Feature Selection/Reduction:** High-dimensional datasets can lead to increased model complexity, overfitting, and computational inefficiency. **Feature selection** and **feature reduction** techniques address these challenges by identifying and retaining the most relevant features. **Principal Component Analysis (PCA)** is a dimensionality reduction technique that transforms the original features into a set of uncorrelated principal components, capturing the maximum variance in the data. **Feature importance** scores from tree-based models, such as Random Forests and Gradient Boosting Machines, can be used to identify the most influential

features. **Correlation analysis** helps identify highly correlated features, allowing for the removal of redundant information. These techniques reduce the dimensionality of the feature space, improve model generalization, and enhance computational efficiency [7].

By meticulously performing these preprocessing and feature engineering steps, we can transform the raw Play Store data into a structured and informative dataset, enabling the development of accurate and robust machine learning models. This rigorous preparation is crucial for extracting meaningful insights and building predictive systems that leverage the rich information contained within the Google Play Store.

**3. Classical Machine Learning Models in Play Store App Rating Prediction**

The prediction of Play Store app ratings has garnered significant attention from researchers and developers alike, aiming to provide users with reliable insights into app quality and popularity. Within this domain, classical machine learning (ML) models have played a pivotal role, offering a robust foundation for building predictive systems. These models, renowned for their interpretability and computational efficiency, are particularly advantageous when dealing with large datasets, a common characteristic of app store data.

**Linear Regression** [8], a foundational model in statistical learning, stands as a simple yet powerful tool for establishing a linear relationship between app features and user ratings. This model operates on the assumption that the rating can be expressed as a linear combination of input features, such as app size, number of downloads, or user reviews. While linear regression excels in its simplicity and interpretability, allowing for straightforward analysis of feature importance, it may struggle to capture complex, non-linear dependencies present in real-world app rating data. However, its efficiency and ease of implementation make it a valuable baseline model for comparison and preliminary analysis.

**Support Vector Regression (SVR)** [9] extends the principles of Support Vector Machines (SVMs) to regression tasks, enabling the modeling of non-linear relationships. By mapping the input data into a higher-dimensional space using kernel functions, SVR can effectively capture intricate patterns that linear models might miss. This capability is crucial for handling the inherent complexity of app rating prediction, where user preferences and app characteristics

often interact in non-linear ways. SVR's ability to handle high-dimensional data and its robustness against outliers make it a strong contender for achieving accurate predictions, though it may require careful parameter tuning to optimize performance.

**Random Forest (RF)** [10], an ensemble learning method, leverages the power of multiple decision trees to enhance prediction accuracy and mitigate overfitting. Each decision tree in the forest is trained on a random subset of the data and features, introducing diversity that reduces the risk of relying on specific data patterns. The final prediction is obtained by aggregating the predictions of individual trees, typically through averaging. RF's ability to handle both numerical and categorical features, its inherent feature importance estimation, and its resistance to overfitting contribute to its popularity in app rating prediction. Furthermore, the ensemble nature of RF allows it to capture complex interactions between features, leading to improved predictive performance compared to single decision tree models.

**Gradient Boosting Machines (GBMs)** [11] represent another powerful ensemble learning technique that builds decision trees sequentially. Unlike RF, where trees are built independently, GBMs construct trees in a stage-wise manner, with each subsequent tree focusing on correcting the errors made by the previous ones. This iterative approach allows GBMs to progressively refine the prediction model, leading to high accuracy. GBMs, such as XGBoost, LightGBM, and CatBoost, have demonstrated exceptional performance in various machine learning tasks, including regression problems like app rating prediction. Their ability to handle complex feature interactions, their robustness against outliers, and their efficient handling of large datasets make them highly effective for this application. However, GBMs often require careful hyperparameter tuning and can be prone to overfitting if not properly regularized.

These classical ML models collectively offer a diverse set of tools for tackling the challenge of Play Store app rating prediction. Their interpretability, computational efficiency, and proven performance on large datasets make them valuable assets in the development of robust and reliable predictive systems. While more advanced deep learning models have emerged, classical ML models continue to provide a strong foundation and often serve as essential benchmarks for evaluating the performance of more complex approaches. The selection of the appropriate model

depends on the specific characteristics of the dataset, the desired level of interpretability, and the computational resources available

The surge in mobile application usage has made the Google Play Store a crucial platform for developers and users alike. App ratings play a pivotal role in influencing user decisions, making accurate prediction of these ratings a valuable endeavor. In this context, Artificial Neural Networks (ANNs), particularly deep learning models, have emerged as powerful tools for capturing complex patterns within Play Store data.

## 4. Artificial Neural Networks (ANNs)

ANNs, inspired by the structure and function of the human brain, have demonstrated remarkable capabilities in various domains, including image recognition, natural language processing, and predictive modeling. Their ability to learn intricate non-linear relationships makes them well-suited for analyzing the complex and multifaceted data found in the Play Store.

**Multilayer Perceptrons (MLPs)**: These feedforward neural networks consist of multiple layers of interconnected nodes, or neurons. The input layer receives the raw data, while hidden layers perform non-linear transformations, and the output layer produces the predicted rating. The ability of MLPs to learn complex mappings between input features and target ratings has been extensively documented [12]. The flexibility offered by multiple hidden layers allows MLPs to approximate complex functions, making them suitable for capturing the nuanced relationships present in app rating data.

**Convolutional Neural Networks (CNNs)**: Originally designed for image processing, CNNs have proven effective in processing text data as well. By employing convolutional filters, CNNs can extract local features, such as word combinations and phrases, from app descriptions and user reviews [13]. These localized features can then be aggregated to form a comprehensive representation of the text, enabling the model to understand the sentiment and content of the textual data. This ability to capture contextual information from text is crucial for predicting app ratings, as user sentiment expressed in reviews often directly impacts the overall rating.

**Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks**: These architectures are specifically designed for sequential data, making them ideal for analyzing user reviews and app descriptions, which are essentially sequences of words [14]. RNNs maintain an internal memory that allows them to process information from previous time steps, enabling them to capture temporal dependencies within the data. However, standard RNNs suffer from the vanishing gradient problem, which limits their ability to learn long-range dependencies. LSTMs, a specialized type of RNN, address this issue by incorporating memory cells that can selectively store and retrieve information over extended periods. This makes LSTMs particularly effective for analyzing long reviews and capturing the evolution of user sentiment over time.

**Hybrid Models**: Combining the strengths of different neural network architectures can lead to improved performance. Hybrid models that integrate CNNs and RNNs, for instance, can leverage the local feature extraction capabilities of CNNs and the sequential processing abilities of RNNs [15]. This allows the model to capture both the local and global context of the data, resulting in more accurate and robust predictions. Such hybrid approaches are increasingly being explored to optimize the performance of app rating prediction models.

While ANNs can achieve higher prediction accuracy compared to classical machine learning models, they require more computational resources and are often less interpretable. The complexity of deep learning models can make it challenging to understand the factors that contribute to their predictions, which is a significant consideration in applications where interpretability is crucial.

## 5. Performance Evaluation Metrics

Evaluating the performance of rating prediction models is essential for ensuring their effectiveness. Several metrics are commonly used to assess the accuracy of these models.

**Mean Absolute Error (MAE)**: This metric calculates the average absolute difference between the predicted and actual ratings [16]. MAE provides a straightforward measure of the average magnitude of errors, making it easy to understand and interpret. It is particularly useful when the distribution of errors is not heavily skewed.

**Root Mean Squared Error (RMSE)**: RMSE calculates the square root of the average squared difference between predicted and actual ratings [17]. RMSE penalizes larger errors more heavily than MAE, making it more sensitive to outliers. It is commonly used when large errors are considered more significant than smaller errors.

**R-squared (R2) score**: This metric measures the proportion of variance in the dependent variable (actual ratings) that is predictable from the independent variables (features used by the model) [18]. R2 provides a measure of how well the model fits the data, with higher values indicating a better fit. However, R2 can be misleading in some cases, particularly when dealing with complex models or small datasets.

The choice of evaluation metric depends on the specific requirements of the application. For instance, if minimizing large errors is crucial, RMSE might be preferred. If a more robust measure of average error is desired, MAE might be more suitable. Understanding the strengths and weaknesses of each metric is essential for selecting the most appropriate one for a given task.

## 6. Recent Research Trends

Recent research in Play Store app rating prediction has focused on several key areas, aiming to improve the accuracy and robustness of prediction models.

**Incorporating user feedback dynamics and temporal aspects of review data**: User feedback and sentiment can evolve over time, making it crucial to consider the temporal dynamics of review data [19]. Researchers are exploring methods to incorporate time-dependent features, such as the recency and frequency of reviews, to capture the evolving nature of user sentiment.

**Utilizing advanced natural language processing (NLP) techniques**: Advanced NLP techniques, such as transformer models like BERT and RoBERTa, have shown significant improvements in sentiment analysis and feature extraction [20]. These models can capture contextual relationships between words and phrases, leading to more accurate representations of user reviews and app descriptions.

**Developing hybrid models**: Combining classical machine learning and deep learning techniques can leverage the strengths of both approaches [21]. Hybrid models can achieve optimal performance by integrating the interpretability of classical models with the accuracy of deep learning models.

**Analyzing the effect of different app categories**: The characteristics of app ratings can vary significantly across different app categories. Researchers are investigating the factors that influence rating prediction accuracy in different categories, aiming to develop more specialized models.

**Implementing explainable AI (XAI) techniques**: Improving the interpretability of deep learning models is crucial for gaining deeper insights into the factors that influence app ratings [22]. XAI techniques can help to understand the decision-making process of these models, making them more transparent and trustworthy.

## 7. Challenges and Future Directions

Despite significant progress, several challenges remain in the field of Play Store app rating prediction.

**Handling imbalanced datasets**: Rating distributions in the Play Store are often skewed, with a majority of apps receiving high ratings. This imbalance can lead to biased models that perform poorly on minority classes.

**Addressing the issue of spam and fake reviews**: The presence of spam and fake reviews can distort the accuracy of rating prediction models. Developing robust methods to detect and mitigate the impact of these reviews is crucial.

**Developing robust models that generalize well**: Models should be able to generalize well across different app categories and time periods. This requires developing models that are robust to variations in data distribution and user behavior.

**Improving model explainability**: Gaining deeper insights into the factors that influence app ratings requires improving the interpretability of prediction models.

**Exploring the use of graph neural networks (GNNs)**: GNNs can model the relationships between apps and users, providing a more comprehensive understanding of the factors that influence app ratings.

Future research should focus on addressing these challenges and exploring novel approaches to improve the accuracy and interpretability of Play Store app rating prediction models. By leveraging advancements in deep learning, NLP, and XAI, researchers can develop more effective and reliable models that benefit both developers and users.

**References**

[1]   A. Marcus, G. Couling, and S. Katevas, "App store reviews: Trends and challenges," in Proceedings of the 2013 International Conference on Mobile Software Engineering and Companion Workshop on Mobile Software Engineering, 2013, pp. 165–168.

[2]   S. Rana, S. Jain, and R. Kumar, "App rating prediction using machine learning techniques," in 2019 International Conference on Computing, Communication, and Intelligent Systems (ICCCIS), 2019, pp. 165–169.

[3]   D. Jurafsky and J. H. Martin, -Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition. Pearson Prentice Hall, 2009. 1

[4]   C. J. Hutto and E. Gilbert, "VADER: A parsimonious rule-based model for sentiment analysis of social media text," in Eighth international AAAI conference on weblogs and social media, 2014. 2

[5]   A. Géron, Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems. 3  O'Reilly Media, 2019.

[6]   I. Goodfellow, Y. Bengio, and A. Courville, Deep learning. MIT press, 2016.

[7] G. Chandrashekar and F. Sahin, "A survey on feature selection methods," Computers & Electrical Engineering, vol. 40, no. 1, pp. 16–35, 2014.

[8] N. R. Draper and H. Smith, Applied regression analysis. John Wiley & Sons, 2014.

[9] V. Vapnik, The nature of statistical learning theory. Springer science & business media, 1999.

[10] L. Breiman, "Random forests," Machine learning, vol. 45, no. 1, pp. 5–32, 2001.

[11] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," Annals of statistics, 4  vol. 29, no. 5, pp. 1189–1232, 2001. 5

[12] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," nature, vol. 323, no. 6088, pp. 533–536, 6  1986.

[13] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," Proceedings of the IEEE, vol. 86, no. 11, pp. 2278–2324, 1998. 7

[14] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural computation, vol. 9, no. 8, pp. 1735–1780, 1997. 8

[15] Z. Yang et al., "Hierarchical attention networks for document classification," in Proceedings of the 2016 conference of the North American chapter of the association for computational 9 linguistics: human language technologies, 10  2016, pp. 1480–1489.

[16] T. Chai and R. R. Draxler, "Root mean square error (RMSE) or mean absolute error (MAE)?—Arguments against avoiding RMSE in the presence of outliers," Geoscientific Model Development, vol. 7, no. 3, pp. 1247–1250, 2014.

[17] R. J. Hyndman and A. B. Koehler, "Another look at measures of forecast accuracy," International journal of forecasting, vol. 22, no. 4, pp. 679–688, 2006. 11

[18] C. D. Myers, "Classical and modern regression with applications," Boston, MA: PWS-Kent, 1990.

[19] J. Chen, Z. Yin, J. Tang, and S. Yang, "Modeling review dynamics for rating prediction," in Proceedings of the 24th ACM international on conference on information and knowledge management, 2015, pp. 1321–1330.

[20] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in Proceedings 12 of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), 2019, pp. 4171–4186. 13

[21] Y. Kim, "Convolutional neural networks for sentence classification," in Proceedings of the 2014 conference on empirical methods in natural language processing 14 (EMNLP), 2014, pp. 1746–1751.

[22] C. Molnar, Interpretable machine learning. A Guide for Making Black Box Models Explainable. Leanpub, 2020.