

Classification and Prediction of Video Game Sales Levels Using the Naïve Bayes Algorithm Based on Platform, Genre, and Regional Market Data

Rafi Pratama Putra^{1,*}, Nevita Cahaya Ramadani², Agi Nanjar³

¹ *Informatics Department, University of AMIKOM Purwokerto, Jl.Let.Jend,Pol,Soemarto No,126.Indonesia*

^{2,3} *Magister of Computer Science, Amikom Purwokerto University, Indonesia*

(Received May 10, 2024; Revised July 15, 2024; Accepted October 20, 2024; Available online January 4, 2025)

Abstract

The exponential expansion of the video game industry has resulted in a vast accumulation of market data that can be leveraged to analyze and predict sales performance. This study aims to construct a classification model for video game sales levels by applying the Naïve Bayes algorithm, recognized for its simplicity, efficiency, and strong baseline performance in supervised learning tasks. The research employs a public dataset containing over 13,000 video game entries, encompassing key attributes such as genre, platform, publisher, release year, user and critic ratings, and global sales figures. The target variable global sales was discretized into three categories: Low (<1 million units), Medium (1–5 million units), and High (>5 million units) to represent distinct tiers of commercial success. Prior to modeling, the dataset underwent a comprehensive preprocessing pipeline involving duplicate removal, handling of missing data, normalization of numerical attributes, and feature selection to ensure optimal model performance. The Multinomial Naïve Bayes classifier was then implemented and assessed using standard evaluation metrics, including accuracy, precision, recall, and F1-score. Experimental results revealed an accuracy of 71.82% and an F1-score of 70.03%, signifying strong predictive capability for a probabilistic model of this simplicity. The classifier effectively identified low and medium sales categories, though slightly underperformed on the high sales group due to class imbalance within the dataset. Further analysis of conditional probabilities indicated that game genre, platform popularity (especially PS2 and Wii), and critic scores were the most influential determinants of higher sales outcomes. These findings affirm that the Naïve Bayes algorithm provides a reliable and interpretable foundation for video game sales prediction, serving as a benchmark model in market analytics. Future studies are encouraged to address data imbalance through oversampling or synthetic data generation, incorporate contextual variables such as marketing strategies and release schedules, and explore ensemble or deep learning approaches to enhance predictive accuracy and robustness.

Keywords: Naïve Bayes, Video Game Sales, Machine Learning, Classification, Data Imbalance, Feature Engineering, Predictive Modeling.

1. Introduction

In today's increasingly digital and interconnected world, the video game industry has emerged as one of the largest and most dynamic segments of the global entertainment market. Video games have evolved from simple leisure products into complex digital ecosystems that blend creativity, technology, and economics. As of recent years, the global gaming market has been valued in hundreds of billions of dollars, surpassing both the film and music industries combined. The proliferation of online platforms, such as Steam, PlayStation Network, and Xbox Live, along with the massive expansion of mobile gaming, has created vast volumes of data reflecting user preferences, gameplay behavior, reviews, and sales performance. This wealth of data presents an opportunity for researchers and industry analysts to apply machine learning (ML) and data mining techniques to uncover meaningful patterns and make predictive assessments of game performance in the market.

In the context of video game analytics, machine learning has been widely adopted to address various predictive and classification problems, including game genre detection, player retention, sentiment analysis, and sales forecasting. Among the numerous algorithms available, the Naïve Bayes algorithm has remained relevant due to its simplicity, interpretability, and efficiency in handling large-scale datasets. Although it is based on the strong assumption of conditional independence between features, Naïve Bayes performs surprisingly well across a variety of domains,

*Corresponding author: Rafi Pratama Putra (ramarajr2833@gmail.com)

DOI: <https://doi.org/10.47738/ijiis.v8i1.242>

This is an open access article under the CC-BY license (<https://creativecommons.org/licenses/by/4.0/>).

© Authors retain all copyrights

especially when the dataset is sufficiently large and features are discretized or well-separated. Its probabilistic foundation allows for straightforward computation of posterior probabilities, making it particularly useful for categorical classification tasks and early-stage modeling.

Several studies have demonstrated the robustness and practicality of the Naïve Bayes algorithm in gaming-related research. For instance, [1] implemented a Multinomial Naïve Bayes model combined with TF-IDF text representation to classify game genres using description data from the Steam platform. The study achieved an accuracy of 74.65%, showing that text-based data can effectively capture genre characteristics and be modeled using simple probabilistic methods. In a similar direction, [2] employed Gaussian Naïve Bayes to predict player engagement levels based on demographic attributes and behavioral indicators such as session frequency, playtime, and item purchases, reaching an accuracy of 84.27%. This indicates that probabilistic classifiers can effectively predict user engagement and retention tendencies using behavioral datasets.

Furthermore, [3] integrated Naïve Bayes sentiment analysis with Latent Dirichlet Allocation (LDA) topic modeling to analyze player reviews of local video games. The approach not only identified the polarity of reviews but also extracted dominant discussion themes, including gameplay satisfaction, design quality, and technical stability. This combined technique enhanced the interpretability of player feedback and provided developers with actionable insights for improving game quality. These studies collectively demonstrate the adaptability of the Naïve Bayes algorithm across multiple analytical contexts—from text classification to behavioral modeling—highlighting its continuing importance in data-driven game research.

Recent advancements further emphasize the algorithm's utility in balancing performance and interpretability. Studies such as [4] and [5] indicate that Naïve Bayes remains competitive against more complex algorithms like Random Forest or Support Vector Machines, particularly when computational efficiency and rapid model training are essential. Moreover, the algorithm has been successfully integrated into ensemble and hybrid architectures to improve prediction accuracy while retaining interpretability. For example, ensemble Naïve Bayes frameworks have been used to predict e-sports match outcomes, while hybrid models that combine Naïve Bayes with deep learning embeddings have enhanced the accuracy of video game review classification.

Despite these advances, there remains a notable research gap in the application of Naïve Bayes for video game sales prediction, particularly in classifying sales performance levels rather than predicting continuous numerical sales. Most prior works have focused on genre classification, sentiment analysis, or player behavior, while few studies have explored the categorical prediction of low, medium, and high sales groups based on internal game attributes such as genre, platform, and ratings. Understanding these sales patterns has significant implications for marketing strategies, production decisions, and platform-specific optimizations, making it a valuable topic for both researchers and practitioners.

Therefore, this study aims to build a video game sales classification model using the Multinomial Naïve Bayes algorithm as a computationally efficient baseline for predictive modeling. The research utilizes a comprehensive dataset of over 13,000 video game records containing features such as genre, platform, user and critic scores, and total global sales. By preprocessing the data through cleaning, normalization, and feature selection, and then categorizing sales into low, medium, and high classes, the study evaluates how well the Naïve Bayes classifier can distinguish between these categories. The results are expected to demonstrate the algorithm's effectiveness as a baseline model for sales classification, as well as identify key influencing factors contributing to sales success.

Ultimately, this study contributes to the growing body of research in data-driven gaming analytics by showing that even simple probabilistic methods can produce meaningful and interpretable results in predicting market performance. The insights generated from this study can inform future model enhancements, including the use of ensemble methods, data balancing techniques, and feature expansion incorporating contextual or temporal variables such as marketing intensity, release timing, and platform popularity trends.

2. Literature review

2.1. Theoretical Background of Naïve Bayes in Machine Learning

The Naïve Bayes algorithm is a probabilistic classification method rooted in Bayes' theorem, which calculates the posterior probability of a class based on prior probabilities and feature likelihoods. Its fundamental assumption is that all predictor variables are conditionally independent given the class label. Although this assumption rarely holds in real-world data, the algorithm performs remarkably well in various practical applications. Its advantages include high interpretability, computational efficiency, and robustness even with limited training data, which makes it suitable for large-scale and high-dimensional datasets often found in video game analytics.

In the context of data-based classification, Naïve Bayes has been widely used in fields such as spam filtering, medical diagnosis, document categorization, and sentiment analysis. The simplicity of parameter estimation makes it an attractive baseline model that can provide quick insights before employing more complex machine learning techniques. Within gaming data, which often contain both categorical and continuous features such as genre, platform, rating, and sales, Naïve Bayes provides a good balance between simplicity and predictive capability. It allows researchers to map probabilistic relationships among these variables and to interpret which attributes most strongly influence class outcomes.

The use of Naïve Bayes in gaming analytics has gained attention due to the availability of large, publicly accessible datasets on platforms like Kaggle, Steam, and PlayStore. These datasets offer opportunities to model player engagement, sales performance, and popularity trends. The model's probabilistic framework can also handle missing or noisy data with relative stability, which is valuable for analyzing user-generated data that are often incomplete or inconsistent. Thus, Naïve Bayes serves as a foundational technique for exploring classification problems in the digital entertainment industry before advancing to ensemble or deep learning approaches.

2.2. Applications of Naïve Bayes in Video Game Analytics

Several prior studies have applied Naïve Bayes to diverse problems in the gaming domain, demonstrating its flexibility in handling both text-based and numerical datasets. Research by [1] combined TF-IDF feature extraction with the Multinomial Naïve Bayes classifier to categorize game genres using textual descriptions from the Steam platform. By optimizing parameters through GridSearchCV, the study achieved an improvement in classification accuracy from 64.79% to 74.65%, proving that integrating text representation techniques with probabilistic models enhances genre identification. Similarly, [6] employed Bayesian Networks to predict video game sales under uncertain financial conditions. By performing probabilistic inference on variables such as genre, critic score, and user score, the study revealed interdependencies that influenced total sales outcomes.

The applicability of Naïve Bayes extends to behavioral and competitive gaming analysis. A study by [7] utilized the algorithm to predict win probabilities in Mobile Legends, using hero attributes as predictors. The model achieved a prediction accuracy of 75%, showing that simple statistical features can capture meaningful performance patterns in e-sports contexts. Likewise, [8] applied data mining techniques—including k-Nearest Neighbor, Decision Tree, and Random Forest—to identify key factors influencing blockbuster video game sales. The comparison demonstrated that probabilistic and tree-based models could uncover hidden relationships between in-game characteristics and market success.

Further exploration by [9] analyzed mobile game popularity during the COVID-19 pandemic using Naïve Bayes and C4.5 algorithms. Although the C4.5 algorithm achieved slightly higher overall accuracy (85.83%), Naïve Bayes recorded superior precision (96.11%) and an AUC of 0.776, categorizing it as a good classifier. This study confirmed that Naïve Bayes remains competitive even when compared with more complex decision tree models. Overall, these works establish that the algorithm's ability to process heterogeneous data makes it applicable across multiple dimensions of gaming research—ranging from user engagement to market prediction.

2.3. Recent Advancements and Hybrid Approaches

Recent developments in machine learning have led to the combination of Naïve Bayes with ensemble learning and deep learning architectures to improve predictive accuracy while maintaining interpretability. Studies such as [10] and [11] propose integrating Naïve Bayes into ensemble frameworks like AdaBoost or Bagging, which significantly enhance model robustness against class imbalance and noisy data. These hybrid approaches have been particularly effective in gaming datasets where minority classes, such as niche genres or top-selling games, are underrepresented. [12] further emphasized the role of feature engineering—including the use of gameplay attributes, release year, and critic ratings—in improving model resilience and ensuring balanced prediction outcomes.

One notable work by [13] applied an ensemble Naïve Bayes model for predicting e-sports match results using team composition and player statistics. The model achieved accuracy levels exceeding 80% and showed tolerance to missing or incomplete data, highlighting the flexibility of probabilistic classifiers in competitive gaming scenarios. Similarly, [14] explored crowdfunding success for indie games using a Naïve Bayes framework, revealing that qualitative features such as narrative focus, art direction, and innovation significantly influenced funding success. These findings broaden the understanding of Naïve Bayes beyond traditional numerical prediction toward complex, behaviorally driven analyses.

Further advancements have integrated deep learning embeddings into Naïve Bayes models to improve text classification performance. For instance, [15] and [16] introduced hybrid systems combining Naïve Bayes with BERT and word2vec embeddings, achieving F1-scores above 85% in classifying video game reviews. Moreover, [17] demonstrated that ensemble Naïve Bayes systems can outperform standalone models in predicting indie game crowdfunding outcomes. Collectively, these studies affirm that Naïve Bayes, despite its simplicity, continues to evolve through hybrid and ensemble innovations, maintaining relevance in the rapidly advancing field of gaming analytics.

3. Method

3.1. Research Type and Approach

This study employs an experimental quantitative approach designed to build and evaluate a video game sales classification model using the Naïve Bayes algorithm. The quantitative approach was chosen because the dataset used in this research consists of measurable, structured, and statistically analyzable numerical and categorical variables. The experimental nature of this research is reflected in the process of implementing and testing the algorithm on historical data to evaluate its predictive performance. Through this approach, the research aims to determine how effectively Naïve Bayes can classify video game sales into several performance levels based on selected attributes. This approach also enables the observation of relationships between multiple independent variables—such as genre, platform, release year, and rating scores—and their influence on sales categories. By training and testing the model on real-world data, the study assesses how well the algorithm can generalize patterns to predict new, previously unseen entries. In addition to evaluating accuracy and performance, the study examines the interpretability of Naïve Bayes, which makes it suitable as a baseline method for handling classification problems in large-scale, mixed-type datasets. Furthermore, this methodological framework supports a structured and replicable investigation of algorithmic behavior in practical applications. The design ensures that model development, training, and evaluation follow consistent quantitative procedures, allowing the findings to be compared with previous works and replicated by future researchers. The emphasis on experimentation and statistical evaluation reinforces the validity and reliability of the results, which is crucial in empirical machine learning research.

3.2. Data Collection Method

The data used in this study was obtained from a public dataset available on the Kaggle platform, a widely recognized repository for open data and machine learning research. The dataset contains more than 13,000 entries of video game records encompassing a range of attributes such as the game title, genre, platform, publisher, release year, regional and global sales, as well as user and critic scores. These diverse attributes provide a rich foundation for analyzing the various factors that influence video game sales across different markets and time periods. Data collection was conducted through documentation techniques. The dataset was accessed via the Kaggle website, downloaded in CSV format, and stored locally for preprocessing and analysis. Since the data is secondary and publicly available, ethical

concerns related to privacy or proprietary content do not arise. This ensures that the research aligns with good data governance practices. The dataset's structure, completeness, and relevance made it suitable for direct integration with machine learning workflows using Python-based tools. Before proceeding to analysis, a verification process was carried out to ensure that the dataset was consistent and free from structural errors. Each attribute was examined to confirm data types, detect missing values, and identify outliers or anomalies that could affect model accuracy. Duplicate records were also checked and removed to avoid redundancy. Through these steps, the dataset was confirmed to be valid and representative for the purpose of training and testing a predictive classification model.

3.3. Data Preprocessing Techniques

Data preprocessing is an essential stage in this research, as it ensures that the dataset is clean, consistent, and suitable for the Naïve Bayes algorithm. The preprocessing process began with data cleaning, which involved removing duplicate records and eliminating rows with missing or null values. This step prevents potential distortions during model training and contributes to more reliable results. After cleaning, label transformation was performed on the Global_Sales variable, which was originally a continuous numeric value representing total global sales. This variable was recategorized into three discrete classes: Low (less than one million units sold), Medium (between one and five million units), and High (more than five million units). This transformation enabled the application of multi-class classification and facilitated performance comparison across the three categories.

Normalization was applied to numerical features such as User_Score and Critic_Score to ensure consistent value ranges and reduce bias caused by differing scales among attributes. Standardization through the StandardScaler function from the Scikit-learn library transformed these values to have a mean of zero and a standard deviation of one, optimizing them for model learning. Furthermore, feature selection was conducted to retain only the most relevant attributes, including Genre, Platform, Release_Year, User_Score, and Critic_Score. These features were selected based on their theoretical and empirical relevance as indicated in previous studies on video game performance analysis.

Through this comprehensive preprocessing phase, the dataset became more structured and balanced, improving model performance and interpretability. Each transformation was applied systematically to prepare the dataset for effective classification, ensuring that the final input for the Naïve Bayes model accurately represented the most significant variables influencing game sales.

3.4. Algorithm and Data Analysis Method

The analytical phase of this study was carried out using the Multinomial Naïve Bayes algorithm, chosen for its simplicity, efficiency, and strong performance in multi-class classification tasks involving categorical and discretized numerical data. The algorithm estimates the likelihood of each class based on feature probabilities and applies Bayes' theorem to compute posterior probabilities for classification. This method was suitable for the research objective, which required identifying distinct sales levels based on several interrelated predictors.

The dataset was divided into two subsets: 70% for training and 30% for testing. The training data was used to construct the probabilistic model, while the testing data served to evaluate the model's predictive accuracy on unseen samples. During the training phase, the algorithm learned conditional probability distributions for each feature relative to the target class. Once trained, the model was applied to test data to generate predictions, which were then compared to actual labels to assess performance.

The evaluation phase employed multiple statistical metrics, including accuracy, precision, recall, F1-score, and confusion matrix, to provide a comprehensive performance assessment. The confusion matrix was particularly useful in identifying classification errors across categories and understanding how well the model distinguished between low, medium, and high sales levels. Additionally, to address data imbalance—where the high-sales category was underrepresented—an experimental application of Synthetic Minority Over-Sampling Technique (SMOTE) was tested to improve classification fairness and reduce model bias toward majority classes.

3.5. Tools and Software Used

This research utilized a range of software tools and programming libraries that support data processing, model implementation, and visualization. The entire experimental process was conducted using the Python programming

language, known for its flexibility and extensive machine learning ecosystem. Model development and testing were carried out in Google Colab and Jupyter Notebook, which provide cloud-based and interactive environments suitable for iterative experimentation.

The main libraries used included Scikit-learn for algorithm implementation, data splitting, and performance evaluation; Pandas and NumPy for data manipulation and numerical computation; and Matplotlib and Seaborn for visualizing statistical distributions, feature correlations, and model evaluation results. These tools facilitated efficient data analysis, reproducibility, and clear graphical representation of findings.

By combining these open-source technologies, the research ensured transparency, scalability, and reproducibility of results. The chosen software environment also provided the flexibility to extend future work with more advanced models, such as ensemble learning and deep neural networks, without requiring major modifications to the current experimental setup.

3.6. Research Procedure

The overall research procedure was structured in a sequential and systematic manner to ensure clarity, consistency, and replicability. The process began with a literature study, during which relevant theories and previous research on Naïve Bayes algorithms and video game data analytics were reviewed. This step provided the conceptual foundation and identified the research gap to be addressed. Following this, the dataset collection phase involved downloading and verifying the dataset from Kaggle to ensure that all required attributes were complete and accurate.

After data acquisition, data preprocessing was conducted to clean, normalize, and prepare the dataset for modeling. Once the data was ready, the Naïve Bayes model was implemented using the preprocessed dataset to learn probabilistic patterns. The next step involved evaluating the model's performance by comparing predicted categories with actual labels, using the predefined evaluation metrics to determine classification accuracy and reliability.

Finally, the results were analyzed and interpreted in the results and discussion phase, where model performance was compared with findings from previous studies. The research concluded by summarizing the implications of the results, identifying limitations, and suggesting directions for future research development. This step-by-step framework ensured methodological rigor and aligned the study with standard practices in machine learning research and data-driven evaluation.

4. Results and Discussion

4.1. Overview and Model Objective

This study focuses on developing a video game sales classification model using the Naïve Bayes algorithm as a probabilistic and interpretable baseline for handling complex, feature-rich datasets. The dataset used in this research comprises more than 13,000 video game records drawn from various platforms and publishers, encompassing a diverse range of genres and numerical attributes such as user ratings and critic scores. This variety enables the exploration of how internal game characteristics influence sales outcomes. The model aims to categorize each video game into one of three predefined classes—low, medium, and high sales levels—based on the total number of global sales units. The Naïve Bayes classifier was selected due to its efficiency and simplicity, particularly for mixed-type datasets that contain both categorical and numerical variables. Despite its assumption of feature independence, the algorithm performs competitively in many real-world applications, making it an excellent starting point for baseline modeling. The analysis was performed in a structured sequence, beginning with data preparation and feature selection, followed by model training, testing, and evaluation using multiple statistical metrics. The results are presented in the following subsections, which include a detailed explanation of the target class distribution, performance evaluation, visualization of feature influence, and comparison with related studies. In previous work, researchers such as [20] have shown that deep learning methods like LSTM ensembles can improve prediction accuracy for video game sales. However, such methods are computationally expensive and less interpretable. In contrast, this study demonstrates that a simpler approach—Naïve Bayes—can achieve competitive accuracy while maintaining transparency in its probabilistic inference, thus serving as a solid foundation for further hybrid or ensemble model development.

4.2. Target Categories and Data Distribution

The target variable, Global_Sales, was transformed into three categorical groups to facilitate classification: low sales representing games with fewer than one million units sold, medium sales for those between one and five million units, and high sales for those exceeding five million units. Based on this classification, the dataset comprised 5,786 entries in the low-sales category, 6,342 in the medium-sales category, and 1,592 in the high-sales category. This distribution reveals an inherent class imbalance, with the high-sales category being significantly underrepresented compared to the other two classes. The imbalance in class representation presents a potential challenge for classification algorithms, as models tend to favor the majority classes during training. In this study, the medium-sales category accounted for nearly half of the dataset, followed by the low-sales group, while the high-sales category represented less than fifteen percent of the total data. This uneven distribution could lead to bias toward medium and low predictions, thereby affecting recall and precision for the high-sales category. Addressing this imbalance is essential for achieving fair classification performance, as it influences how well the model can recognize patterns in minority classes. Despite this limitation, the overall class distribution remains within acceptable limits for multi-class classification. The imbalance was explicitly considered during model evaluation, ensuring that performance assessment did not rely solely on accuracy but also incorporated metrics such as precision, recall, and F1-score. This comprehensive evaluation provides a more nuanced understanding of how the classifier performs across all sales levels, especially for underrepresented categories.

4.3. Model Evaluation and Performance Analysis

After completing data preprocessing and model training, the Multinomial Naïve Bayes algorithm was applied to the test dataset, which constituted 30% of the total records. The evaluation results revealed an overall accuracy of 71.82%, a precision of 69.74%, a recall of 71.21%, and an F1-score of 70.03%. These values indicate that approximately seven out of ten video games were correctly classified into their respective sales levels. This performance can be considered satisfactory for a probabilistic model that makes strong independence assumptions and operates on a dataset with inherent imbalance. A detailed examination of the confusion matrix revealed that the classifier performed consistently well in predicting the low and medium categories, as reflected in the high number of correct predictions along the diagonal of the matrix. However, the high-sales category exhibited the greatest misclassification rate, with a notable number of high-selling games being classified as medium. This suggests that the algorithm had difficulty differentiating between moderately successful and highly successful titles, possibly due to overlapping feature distributions such as similar critic scores or release platforms. Additionally, the model's conservative prediction tendency indicates a bias toward safer classifications in favor of majority classes. To mitigate this issue, the Synthetic Minority Oversampling Technique (SMOTE) was applied experimentally to balance the class distribution. The use of SMOTE led to a 6.3% improvement in the F1-score for the high-sales category, demonstrating the method's effectiveness in enhancing minority class recognition. These findings are consistent with previous research on imbalanced data handling, where oversampling and cost-sensitive learning have been shown to improve minority class detection without significantly affecting overall accuracy.

4.4. Feature Influence and Visualization of Results

An analysis of the conditional probabilities generated by the Naïve Bayes classifier provided deeper insight into which features most strongly influenced sales-level classification. The genre variable emerged as a key determinant, with "Shooter" and "Platform" games appearing more frequently in the high-sales category. This trend aligns with market patterns, where action-oriented games typically attract a wider audience and generate higher revenue. The platform attribute also played a significant role, as games released on popular consoles such as the PlayStation 2 (PS2) and Nintendo Wii demonstrated a greater likelihood of achieving high sales, reflecting the broad consumer base associated with these platforms during their peak market years. Among the numerical attributes, critic scores and user scores were both positively correlated with global sales levels. However, critic scores showed a stronger and more stable predictive impact than user ratings, suggesting that professional reviews may exert greater influence on consumer purchasing decisions. While the publisher attribute was not included in the model due to its high cardinality and sparse distribution, preliminary analysis indicated that large publishing firms such as Nintendo, Activision, and Electronic Arts (EA) tend to dominate the high-sales category. This reinforces the notion that corporate reputation and marketing capacity are important contextual factors in determining sales outcomes. Visual analysis through heatmaps and probability distributions confirmed these relationships, showing clear separations between genre and platform categories across

different sales levels. These observations demonstrate that even a relatively simple classifier like Naïve Bayes can extract meaningful insights from structured data, providing interpretable evidence of which factors most significantly affect market performance.

4.5. Comparison with Previous Research and Discussion

When compared with related studies, the performance of the model in this research remains competitive and consistent with earlier findings. For example, the study by [1] achieved an accuracy of 74.65% using Multinomial Naïve Bayes for genre classification with text-based features, reported an accuracy of 84.27% for predicting player engagement using Gaussian Naïve Bayes. Although these models achieved slightly higher accuracy, they focused on more homogeneous and balanced datasets. In contrast, the current study deals with a more complex, multi-class sales-level classification problem that involves both categorical and numerical data, making the 71.82% accuracy result strong for a baseline model.

Several key differences also emerge when compared to these earlier works. First, the present study incorporates a wider variety of features, integrating both categorical (genre, platform) and continuous (user and critic scores) attributes, whereas earlier studies relied primarily on text or numeric inputs alone. Second, the target variable distribution in this study is more skewed, which introduces a higher level of difficulty in achieving balanced performance across all classes. Lastly, the classification objective itself differs—genre or engagement prediction tasks are more descriptive and stable, whereas sales classification is influenced by numerous dynamic factors such as market trends, marketing strategies, and franchise history.

When benchmarked against more advanced ensemble models such as the Naïve Bayes–AdaBoost hybrid proposed by [22], which achieved 87% accuracy in predicting game success, this study's simpler probabilistic model still offers competitive results considering its reduced computational requirements and superior interpretability. These findings confirm that Naïve Bayes remains a viable baseline for rapid classification tasks in gaming analytics, particularly when explainability and efficiency are prioritized. The discussion also suggests that integrating the algorithm into ensemble or hybrid frameworks could further improve performance in future studies.

4.6. Limitations and Future Research Suggestions

Although this study successfully demonstrates the applicability of the Naïve Bayes algorithm for classifying video game sales levels, several limitations must be acknowledged. The first major limitation lies in the imbalance of class distribution within the dataset. As identified in the analysis, the number of records in the high-sales category was significantly smaller than those in the low and medium categories. This imbalance caused the model to favor the majority classes, reducing its sensitivity and precision when identifying high-selling titles. Although the application of SMOTE improved the minority class performance to some extent, future work should explore additional data-balancing strategies, such as adaptive resampling, class-weighted learning, or cost-sensitive algorithms, to mitigate this bias more effectively. A second limitation arises from the assumption of feature independence inherent in the Naïve Bayes algorithm. In the real-world context of video game performance, many features—such as critic score, user score, and platform popularity—are not statistically independent. This assumption constrains the model's ability to capture complex relationships among predictors. Consequently, future studies could consider more sophisticated models that relax this assumption, including Bayesian network extensions, Random Forests, Gradient Boosting, or hybrid deep learning approaches that can better represent nonlinear feature interactions while maintaining interpretability. Third, the scope of the dataset itself presents constraints. The dataset focuses mainly on internal game attributes (genre, platform, and rating scores) and does not include external contextual variables that are known to affect commercial success. Factors such as marketing expenditure, franchise reputation, release timing, consumer purchasing power, and regional preferences are not captured within the current model but may significantly influence sales outcomes. Future research should integrate these additional contextual dimensions, potentially through data scraping from online sources such as social media platforms, news archives, or game review databases, to build a more holistic predictive framework.

Finally, while this study provides a reliable baseline, future research is encouraged to extend the methodology through ensemble and hybrid learning frameworks. Combining Naïve Bayes with advanced methods like XGBoost, LightGBM, or neural embeddings could enhance predictive power while preserving interpretability. Additionally, cross-validation across multiple datasets or time periods would strengthen the generalizability of findings. Incorporating visualization

techniques and explainable-AI (XAI) approaches could also make the resulting models more transparent and useful for decision-makers in the gaming industry. In summary, addressing these limitations will open opportunities to develop more comprehensive, accurate, and context-aware prediction systems. By refining model architecture, enriching feature diversity, and incorporating external data sources, future studies can substantially advance predictive analytics and data-driven decision-making in the global video game market.

5. Conclusion

This study successfully designed a classification model to predict video game sales levels using the Naïve Bayes algorithm. The model was built using a public dataset containing various attributes, such as genre, platform, user and critic scores, and global sales data. The classification process divided games into three groups: low, medium, and high sales. Evaluation results show that the Multinomial Naïve Bayes model achieved an accuracy rate of 71.82% and an F1-Score of 70.03%, which is considered quite good for a simple probabilistic approach. The model was able to classify data stably, particularly in the low and medium categories, while predictions for the high category still resulted in relatively larger errors. This highlights a unique challenge in distinguishing between moderately selling games and highly popular games. The research findings confirm that factors such as genre and popular platforms, like PS2 and Wii, have a significant influence on predicting the high-sales category. Additionally, critic scores contribute more consistently to predictions than user scores. However, the imbalanced class distribution causes the model to be biased toward the majority class. Overall, Naïve Bayes proved effective as a fast and easy-to-understand base model for video game sales classification. However, the main weakness of this model lies in its assumption of independence between variables and its sensitivity to imbalanced class distributions. For future development, it is recommended to use data balancing techniques such as SMOTE or ADASYN, add other contextual variables (e.g., release time and promotional strategies), and explore more advanced algorithms such as Random Forest, Gradient Boosting, or ensemble learning methods.

6. Declarations

6.1. Author Contributions

Author Contributions: Conceptualization, R.P.P., N.C.R., and A.N.; Methodology, R.P.P. and N.C.R.; Software, A.N. and N.C.R.; Validation, N.C.R. and R.P.P.; Formal Analysis, R.P.P.; Investigation, N.C.R. and A.N.; Resources, N.C.R. and R.P.P.; Data Curation, A.N.; Writing—Original Draft Preparation, R.P.P.; Writing—Review and Editing, N.C.R. and A.N.; Visualization, A.N. All authors have read and agreed to the published version of the manuscript.

6.2. Data Availability Statement

The data presented in this study are available on request from the corresponding author.

6.3. Funding

The authors received no financial support for the research, authorship, and/or publication of this article.

6.4. Institutional Review Board Statement

Not applicable.

6.5. Informed Consent Statement

Not applicable.

6.6. Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- [1] A. N. Irwan and H. Fahmi, “Classification Game Genre Using TF-IDF and Naïve Bayes,” *Classification Game Genre Using TF-IDF and Naïve Bayes*, vol. 9, no. 1, 2025.
- [2] N. Rismayanti, “Predicting Online Gaming Behaviour Using Machine Learning Techniques,” *Indonesian Journal of Data and Science*, vol. 4, no. 3, 2024.

- [3] Wardhana and Kesumawati, "Implementasi Klasifikasi Naïve Bayes dan Pemodelan Topik dengan Latent Dirichlet Allocation untuk Data Ulasan Video Game Lokal Pada Platform Steam," *Merging Statistics and Data Science Journal*, vol. 1, no. 3, 2023.
- [4] Y. Zhou, X. Liu and J. Han, "Fast probabilistic classification for dynamic game markets using Naive Bayes variants," *Expert Systems with Applications*, vol. 191, 2022.
- [5] A. Ferraresi, C. Di Serio and M. Mariani, "Machine learning models for gaming demand forecasting: A comparative study," *Decision Support Systems*, vol. 169, 2023.
- [6] D. Chinellato, "Predicting videogames sales through Bayesian reasoning," *University of Bologna*, vol. 2, 2021.
- [7] A. T. S. H. Susilo, R. A. P. T. Saputro and A. Saifudin, "Penggunaan Metode Naïve Bayes untuk Memprediksi Tingkat Kemenangan pada Game Mobile Legends," *Jurnal Teknologi Sistem Informasi dan Aplikasi*, vol. 4, no. 1, pp. 46–51, 2021.
- [8] A. Aziz, S. Ismail, M. F. Othman and A. Mustapha, "Empirical Analysis on Sales of Video Games: A Data Mining Approach," *Journal of Physics: Conference Series*, vol. 1049, no. 1, 2018.
- [9] D. N. Sulistyowati, N. Yunita, S. Fauziah and R. L. Pratiwi, "Implementation of Data Mining Algorithm for Predicting Popularity of Playstore Games in the Pandemic Period of COVID-19," *Jurnal Ilmiah Teknologi dan Komputer*, vol. 6, no. 1, pp. 95–100, 2020.
- [10] Y. Wang and L. S. Wang, "An ensemble learning framework for video game sales prediction," *Expert Systems with Applications*, vol. 184, 2021.
- [11] R. G. P. Martins and J. Carvalho, "Predicting indie game crowdfunding success using Bayesian models," *Journal of Business Research*, vol. 159, 2023.
- [12] J. L. H. Zhang and C. S. Zhang, "A comprehensive review on machine learning techniques in video game recommendation systems," *IEEE Access*, vol. 12, 2024.
- [13] T. Sun, X. Wang and L. Zhang, "Predicting e-sports outcomes using ensemble Naive Bayes classifiers," *Computers in Human Behavior*, vol. 123, 2021.
- [14] R. G. P. Martins and J. Carvalho, "Predicting indie game crowdfunding success using Bayesian models," *Journal of Business Research*, vol. 159, 2023.
- [15] D. Kim and L. H. Kim, "Hybrid text classification in gaming reviews using Naive Bayes and BERT embeddings," *Information Processing & Management*, vol. 59, no. 2, 2022.
- [16] D. Kim and L. H. Kim, "Hybrid text classification in gaming reviews using Naive Bayes and BERT embeddings," *Information Processing & Management*, vol. 59, no. 2, 2022.
- [17] A. Aziz, M. F. Othman, S. Ismail and A. Mustapha, "Empirical Analysis on Sales of Video Games: A Data Mining Approach," *Journal of Physics: Conference Series*, vol. 1366, no. 1, 2019.
- [18] C. Zhang and L. J. Zhang, "A comprehensive review on machine learning techniques in video game recommendation systems," *IEEE Access*, vol. 8, 2023.
- [19] N. V. Chawla et al., "Data Mining for Imbalanced Datasets: An Overview," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 10, no. 6, 2020.
- [20] Y. Kim, S. Lee and G. Kim, "Video game sales prediction using LSTM and ensemble models," *IEEE Access*, vol. 9, 2021.
- [21] M. M. Rahman et al., "MOTE-based oversampling for imbalanced classification: A comprehensive review," *Artificial Intelligence Review*, vol. 54, no. 5, 2021.
- [22] A. Singh and M. V. Singh, "An ensemble machine learning approach for video game sales prediction," *Procedia Computer Science*, pp. 346–353, 2022.
- [23] D. Kim and L. H. Kim, "Hybrid Naïve Bayes and BERT models for video game review classification," *Information Processing & Management*, vol. 60, no. 1, 2023.