

Automated Data Analysis Report: Steam

Executive Summary

****Executive Summary: Automated Data Analysis Report**** This report presents the preliminary findings of an initial automated exploration of the 'steam.csv' dataset, comprising 27,085 rows and 18 columns. The primary objective of this exercise is to provide a high-level overview of the dataset's structure, quality, and characteristics. Notably, the dataset exhibits a robust size, with no missing values, and only 10 duplicates, indicating a relatively clean and comprehensive data collection. The analysis reveals a balanced mix of numerical and categorical features, with 9 columns in each category. Univariate and bivariate analyses have been performed, including histogram and box plot visualizations for numerical features, bar charts for categorical features, and correlation assessments between numerical columns. Three numerical-categorical plot analyses have also been conducted. While no datetime columns are present, the dataset's overall quality and diversity offer a promising foundation for further investigation. This initial scan has provided valuable insights into the dataset's composition and potential patterns. The descriptive statistics, quality checks, and visualizations performed have laid the groundwork for more in-depth analysis and exploration. As we move forward, this preliminary assessment will serve as a useful starting point for identifying key trends, correlations, and areas of interest. The next steps will involve drilling deeper into the data to uncover more nuanced insights, ultimately informing strategic decision-making and driving business value.

1. Data Overview

This report provides an initial automated analysis of the dataset from 'steam.csv'.

1.1. Basic Information

Table 1: Dataset Dimensions

Metric	Value
Number of Rows	27085
Number of Columns	18
Total Data Points	487530

1.2. Data Types

Table 2: Summary of Feature Data Types

Data Type	Count
object	9
int64	8
float64	1

Data Types Distribution:

The distribution of data types in this dataset suggests that it is predominantly composed of numerical and categorical features, with a notable absence of datetime features. The presence of 9 numerical features indicates that the dataset contains a significant amount of quantitative information, which can be easily analyzed and modeled using statistical and machine learning techniques. The 9 categorical features, on the other hand, provide context and descriptive information about the data, which can be used to identify patterns, relationships, and trends. The lack of datetime features implies that the data may not have a strong temporal component, and analysis may focus more on cross-sectional relationships rather than time-series patterns. Overall, this distribution of data types suggests that the analysis will likely involve a mix of statistical modeling, feature engineering, and data transformation to effectively harness the numerical and categorical information, with a focus on understanding the relationships and patterns within the data.

2. Data Quality Assessment

This section evaluates common data quality aspects such as missing values, duplicates, and feature variance.

2.1. Missing Values

No missing values were found in the dataset. This is excellent for data completeness.

2.2. Duplicate Records

The dataset contains 10 duplicate rows (representing 0.04% of the data). These may need to be investigated or removed depending on the analysis context, as they can skew results.

2.3. Feature Variance

No constant columns (columns with only one unique value) were identified.

The following columns are quasi-constant (one value is highly dominant), potentially offering limited information: english (dominant value: 1 at 98.1%); requiredage (dominant value: 0 at 97.8%). Their utility should be reviewed.

Data Quality Summary & Implications:

The data quality assessment reveals a generally clean dataset with some minor issues. The absence of missing values is a positive finding, indicating that the data is complete and ready for analysis. The small percentage of duplicate rows (0.04%) is unlikely to have a significant impact on the analysis, but it's still worth considering removing them to ensure the accuracy of the results. The lack of constant columns suggests that the data is not plagued by redundant or unnecessary features. However, the presence of quasi-constant columns, such as "english" and "requiredage", may pose some challenges. The dominant values in these columns (98.1% and 97.8%, respectively) indicate that they may not be highly informative or useful for modeling purposes. This could lead to issues with model performance, as the algorithms may struggle to learn from features that are largely constant. Furthermore, the quasi-constant columns may also lead to biased or ungeneralizable results, as the models may overfit to the dominant values. For instance, if a model is trained on a dataset where 98.1% of the values are "1" for the "english" column, it may not perform well on new, unseen data where this proportion is different. The potential implications of these findings for further analysis are that the quasi-constant columns may impact the reliability of insights and model performance. To address these issues, general strategies can be employed, such as removing or transforming the quasi-constant columns to create more informative features. For example, the "english" column could be transformed into a more nuanced feature, such as a language proficiency score, to capture more variation. Additionally, techniques like feature engineering or dimensionality reduction (e.g., PCA) can be used to create new features that are more informative and less prone to dominance by a single value. By addressing these issues, the data can be better prepared for analysis, and the resulting insights are likely to be more reliable and generalizable.

3. Univariate Analysis

This section examines individual features to understand their distributions, central tendencies, spread, and potential outliers.

3.1. Numerical Features

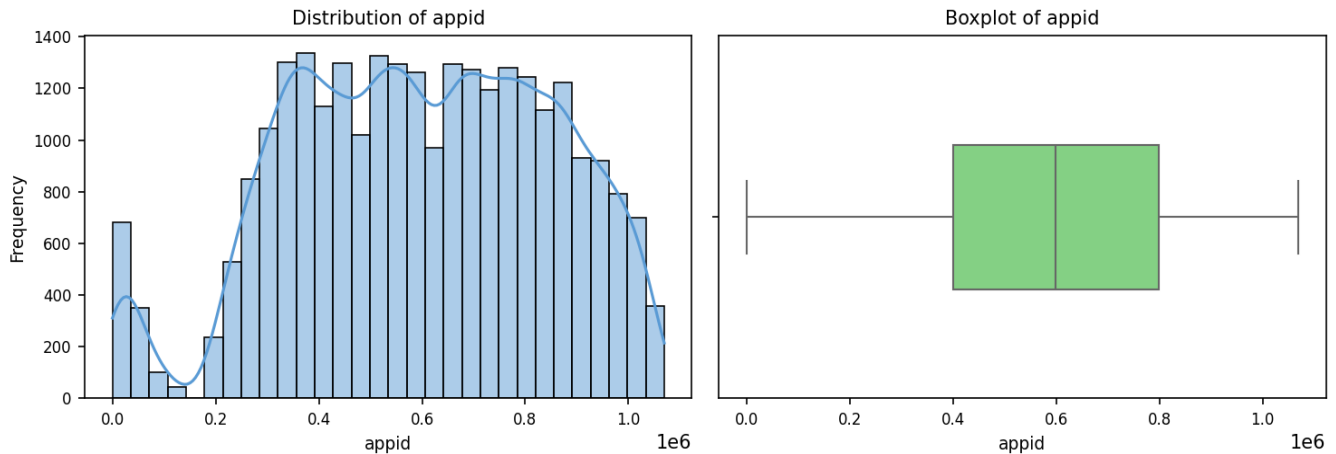


Figure 1: Distribution (histogram and KDE) and boxplot for 'appid'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

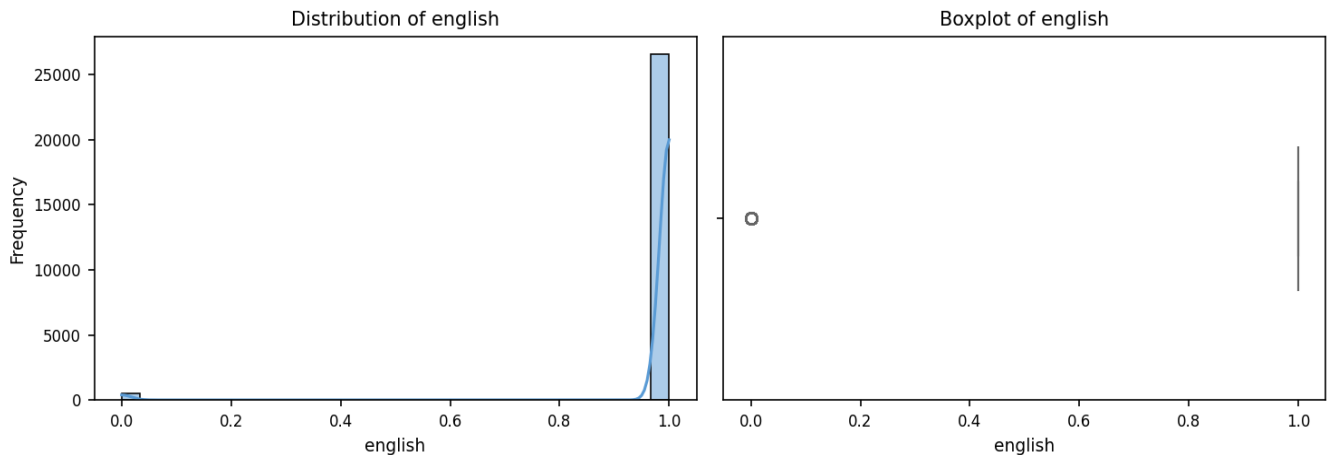


Figure 2: Distribution (histogram and KDE) and boxplot for 'english'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

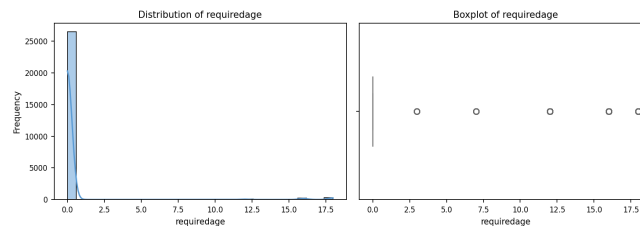


Figure 3: Distribution (histogram and KDE) and boxplot for 'requiredage'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

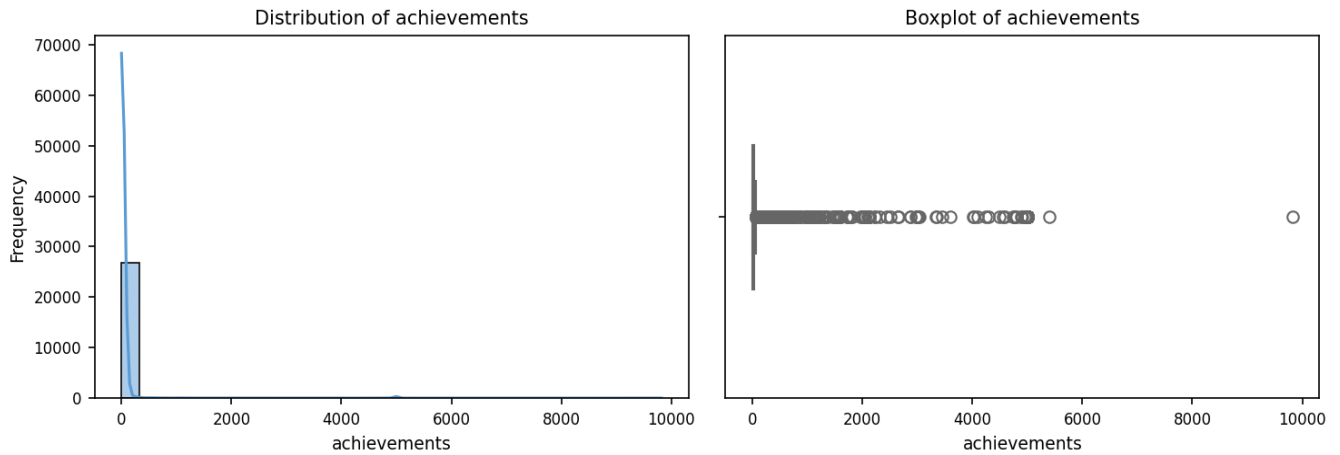


Figure 4: Distribution (histogram and KDE) and boxplot for 'achievements'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

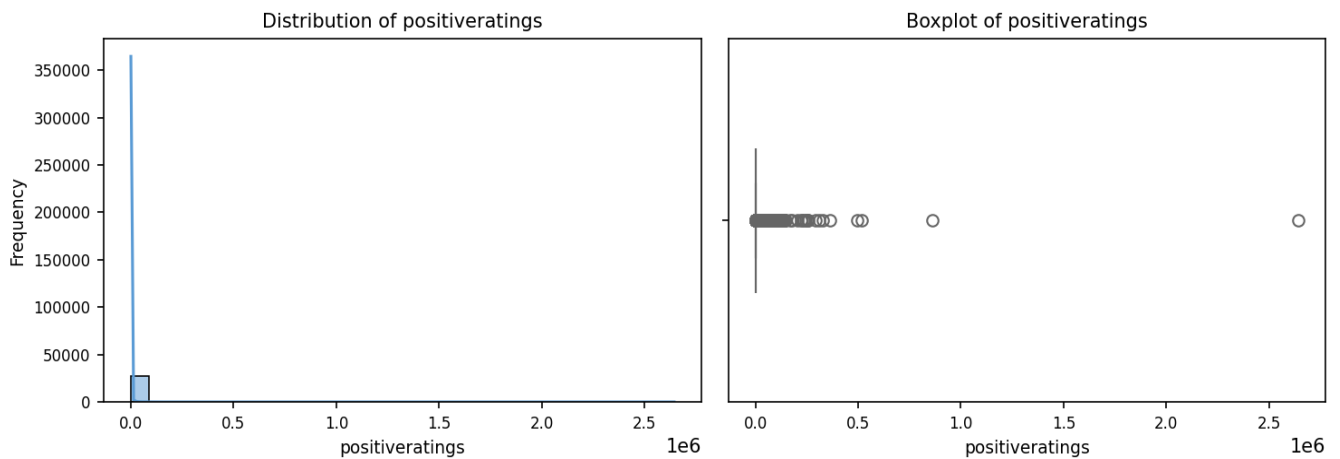


Figure 5: Distribution (histogram and KDE) and boxplot for 'positiveratings'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

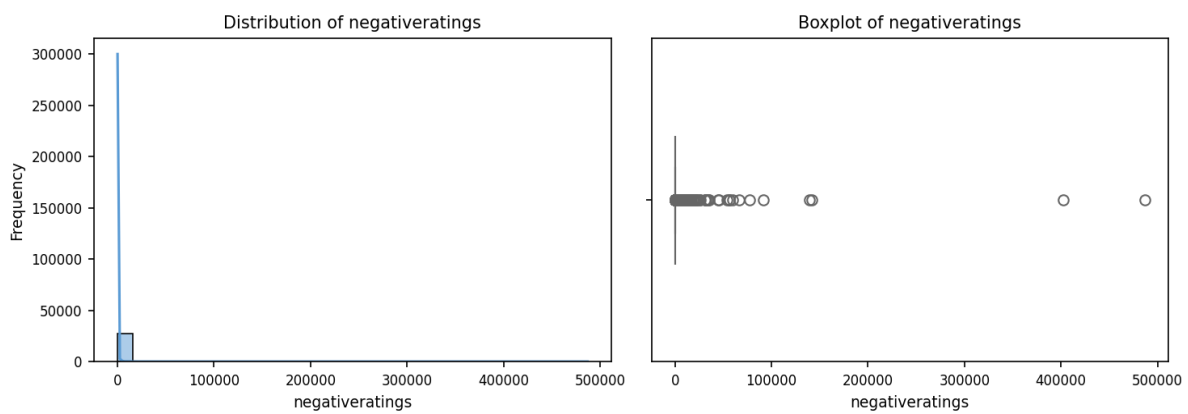


Figure 6: Distribution (histogram and KDE) and boxplot for 'negativeratings'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

Observations on Numerical Feature Distributions:

The numerical features analyzed exhibit a range of distribution shapes and characteristics. Notably, most features are highly skewed, with 'achievements', 'positiveratings', and 'negativeratings' displaying extreme positive skewness, indicating that the majority of values are concentrated at the lower end of the scale, with a long tail of extreme values at the higher end. In contrast, 'appid' is mildly skewed to the left, while 'english' is heavily skewed to the left due to its binary nature. The 'requiredage' feature is also positively skewed, but to a lesser extent. These skewed distributions suggest that the data may not be normally distributed, and alternative statistical methods or transformations may be necessary for analysis. The presence of potential outliers is a striking characteristic of these distributions. Boxplots for several features, including 'appid', 'achievements', 'positiveratings', and 'negativeratings', may indicate outliers, which is supported by the extreme min/max values relative to the mean/median. For example, the 'positiveratings' feature has a maximum value of 2644404, which is vastly larger than its median of 24. These extreme values may be influential in modeling and analysis, and may require special consideration or handling to prevent biased results. The 'english' feature, being a binary variable, does not exhibit outliers in the classical sense, but its skewness is notable. The spread or variability of the features is also noteworthy. 'appid', 'requiredage', and 'english' have relatively low standard deviations, indicating that these features are relatively consistent or stable. In contrast, 'achievements', 'positiveratings', and 'negativeratings' have extremely high standard deviations, reflecting their highly skewed and variable distributions. The 'achievements' feature, for instance, has a standard deviation of 352.61, which is more than 50 times its median value. This high variability suggests that these features may be more challenging to model or analyze, and may require specialized techniques or transformations to capture their underlying patterns. Overall, the distributions of these numerical features exhibit a range of patterns and characteristics, highlighting the importance of careful data exploration and preprocessing in analysis.

3.2. Categorical Features

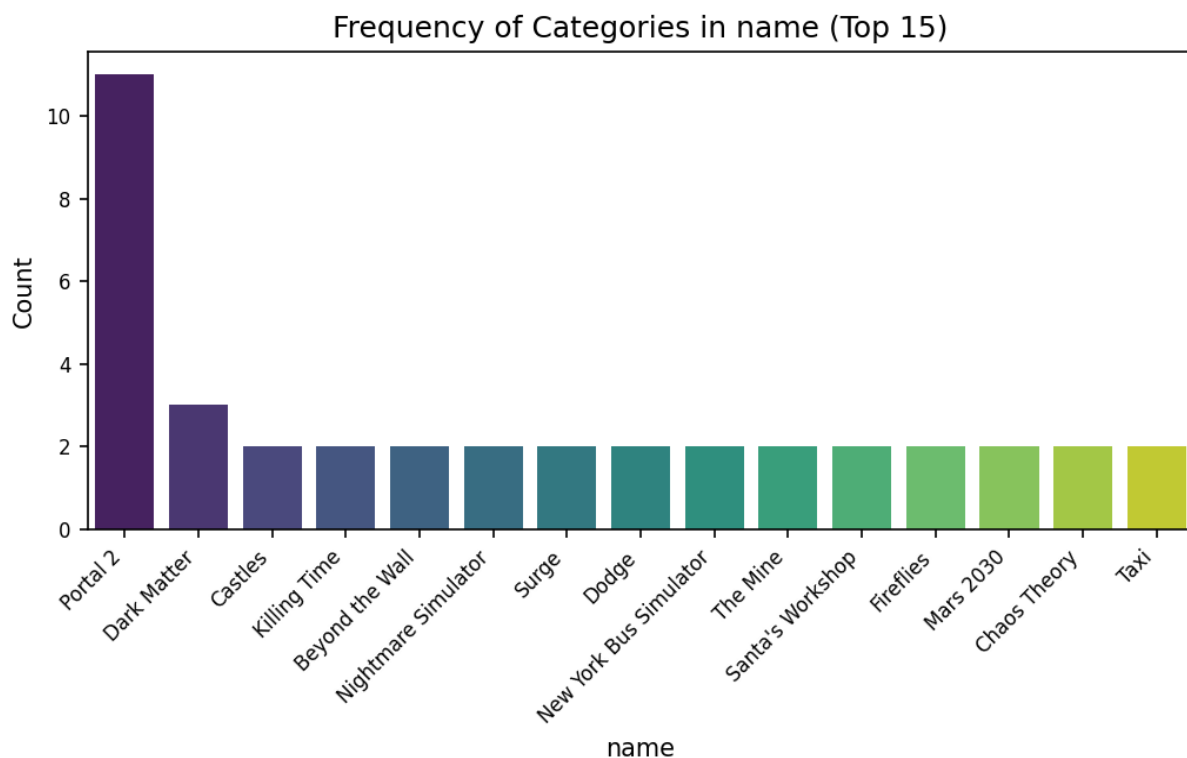


Figure 7: Bar chart showing frequency of top categories in 'name'. Total unique values: 27033.

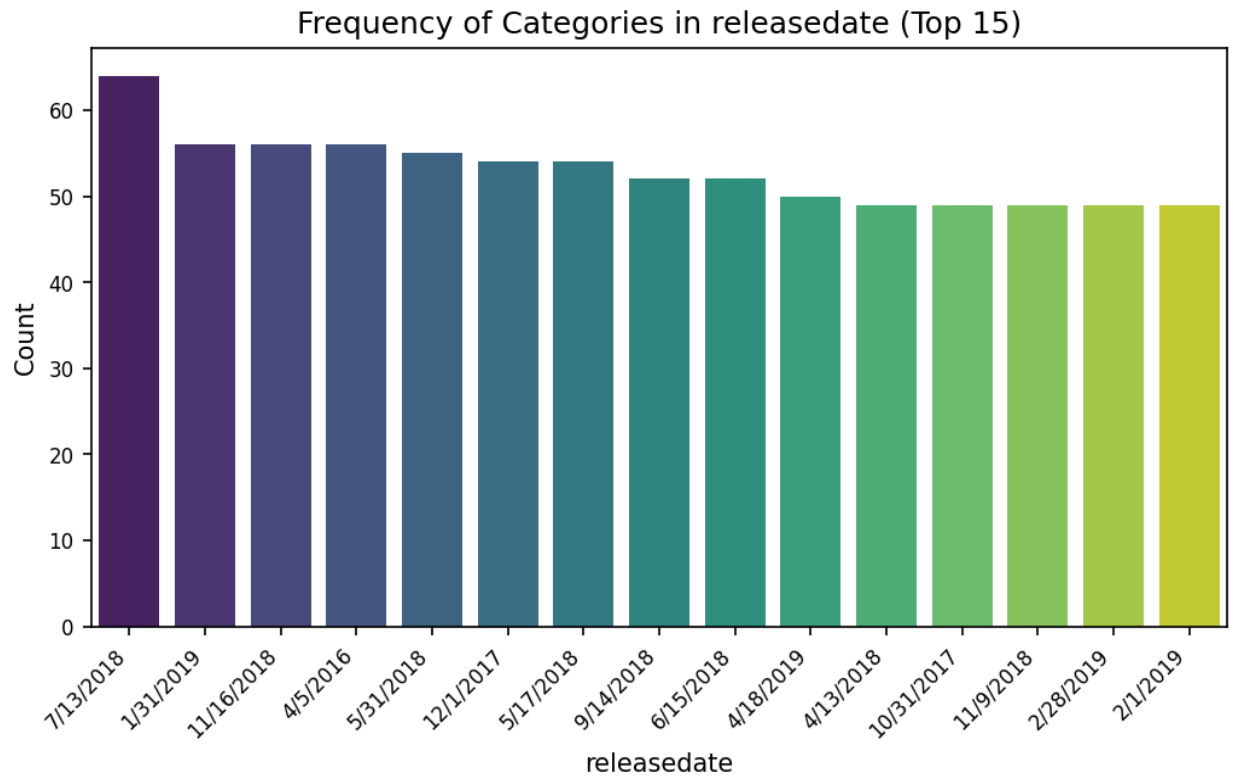


Figure 8: Bar chart showing frequency of top categories in 'releasedate'. Total unique values: 2619.

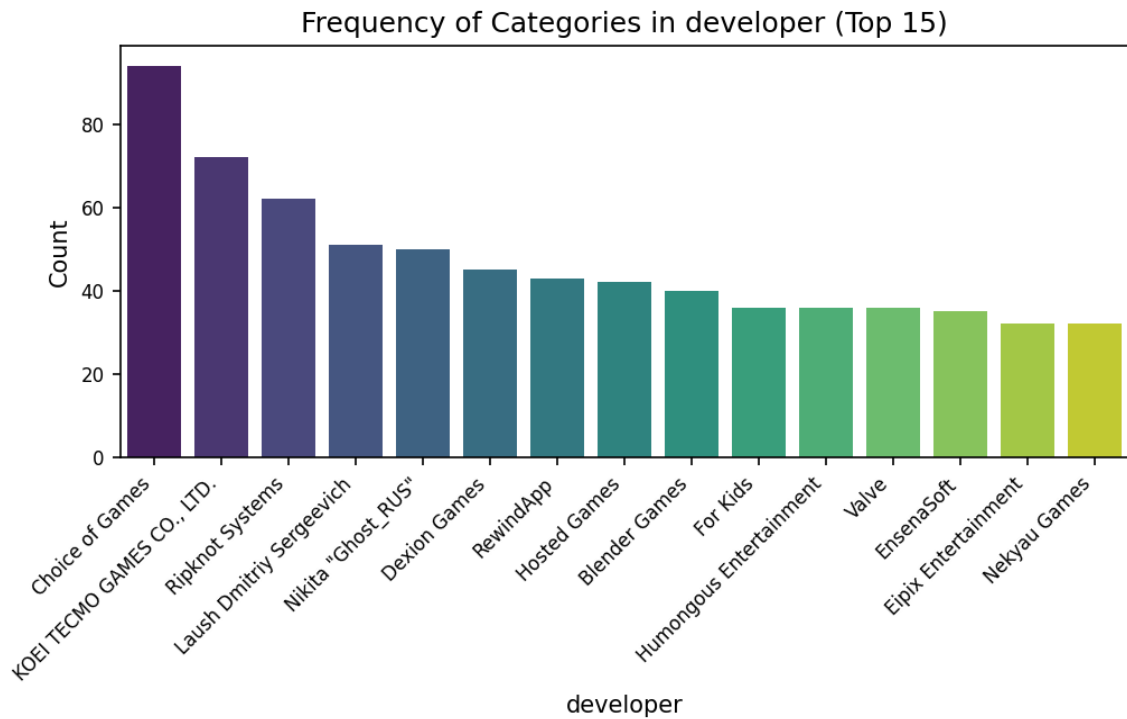


Figure 9: Bar chart showing frequency of top categories in 'developer'. Total unique values: 17112.

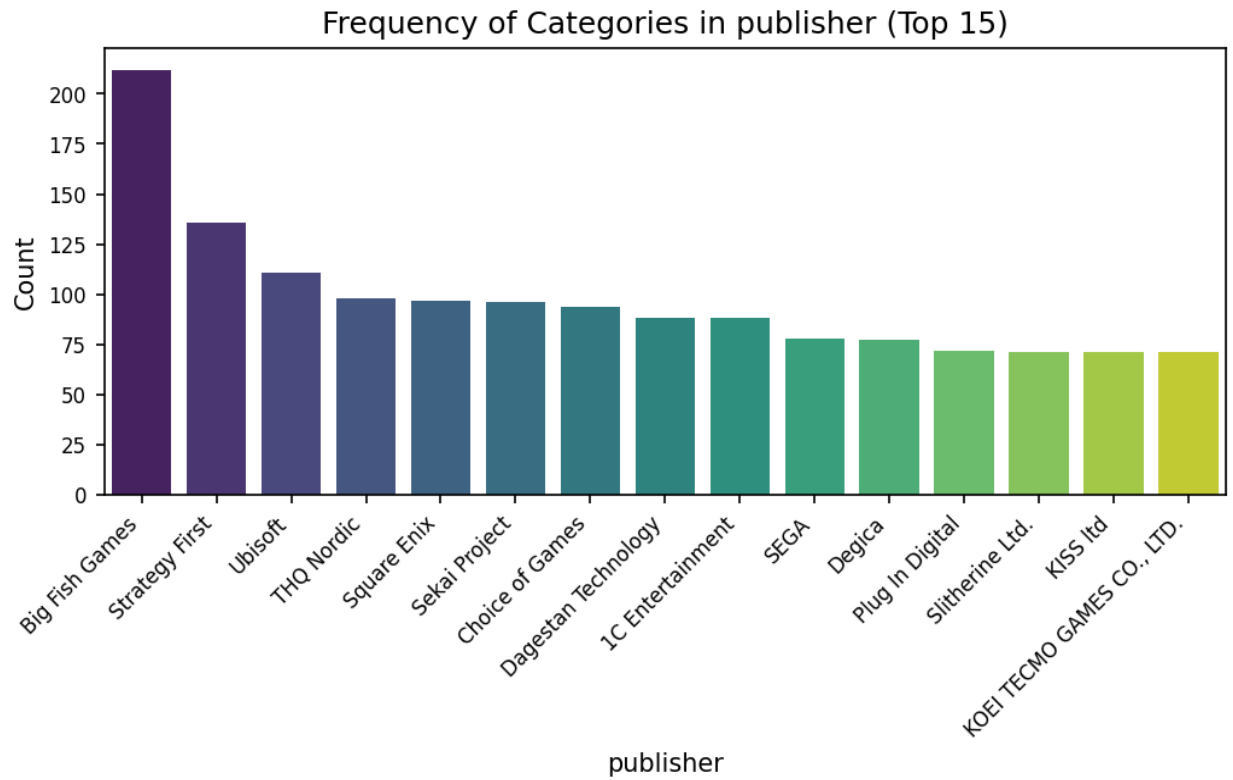


Figure 10: Bar chart showing frequency of top categories in 'publisher'. Total unique values: 14353.

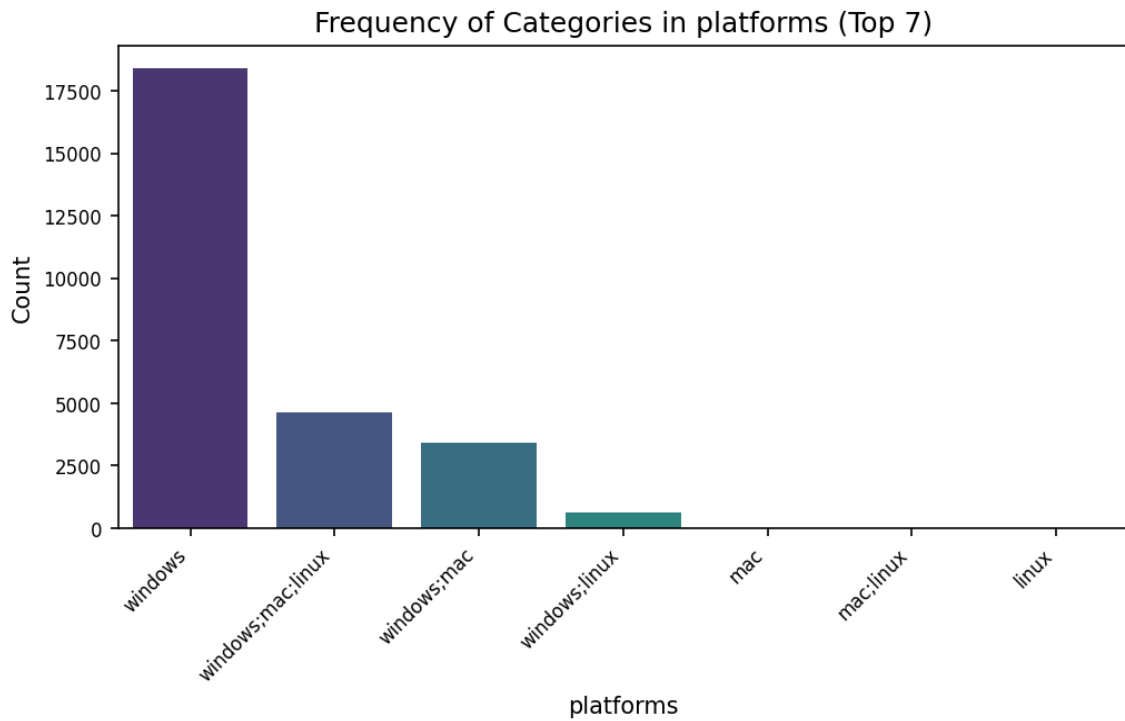


Figure 11: Bar chart showing frequency of top categories in 'platforms'. Total unique values: 7.

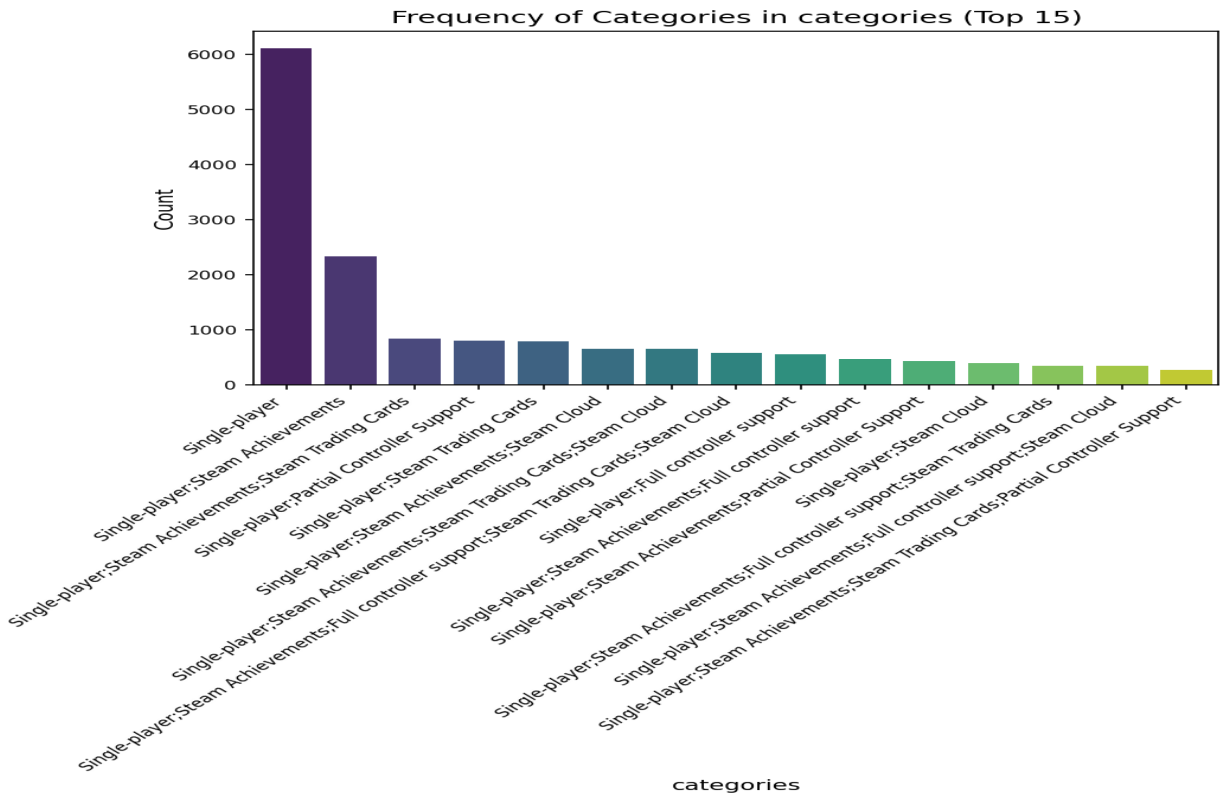


Figure 12: Bar chart showing frequency of top categories in 'categories'. Total unique values: 3333.

Observations on Categorical Feature Distributions:

The analysis of the categorical features reveals a range of observations. First, we notice that the features 'name', 'developer', and 'publisher' exhibit high cardinality, with 27033, 17112, and 14353 unique values, respectively. These features have a large number of distinct categories, which may require special handling during feature encoding or analysis. On the other hand, features like 'platforms' and 'categories' have relatively low cardinality, with 7 and 3333 unique values, respectively. The feature 'releasedate' has a moderate number of unique values, with 2619 distinct categories. The distribution of categories within each feature also provides valuable insights. The feature 'platforms' has a dominant category, 'windows', which accounts for 67.9% of the observations. Similarly, 'categories' has a top category, 'Single-player', which makes up 22.6% of the data. In contrast, features like 'name', 'developer', and 'publisher' do not have dominant categories, as indicated by the low percentages of their top categories (0.0%, 0.3%, and 0.8%, respectively). The feature 'releasedate' also has a relatively low percentage for its top category, with 0.2%. These observations suggest that features with dominant categories may be more amenable to analysis, while those with more even distributions may require more careful consideration. The implications of these observations are significant for feature encoding and analysis. Features with high cardinality, such as 'name', 'developer', and 'publisher', may require techniques like one-hot encoding, label encoding, or hashing to reduce their dimensionality. Alternatively, these features could be treated as text data and analyzed using natural language processing techniques. Features with dominant categories, like 'platforms', may be more suitable for analysis using traditional categorical encoding methods. The relatively even distributions of some features may also suggest the use of clustering or dimensionality reduction techniques to identify patterns or relationships in the data. Overall, the analysis of these categorical features highlights the importance of careful consideration and handling of high-cardinality features and uneven distributions to ensure effective and accurate analysis.

4. Bivariate Analysis

This section explores relationships between pairs of features, which can reveal correlations, dependencies, and interactions.

4.1. Numerical vs. Numerical Features

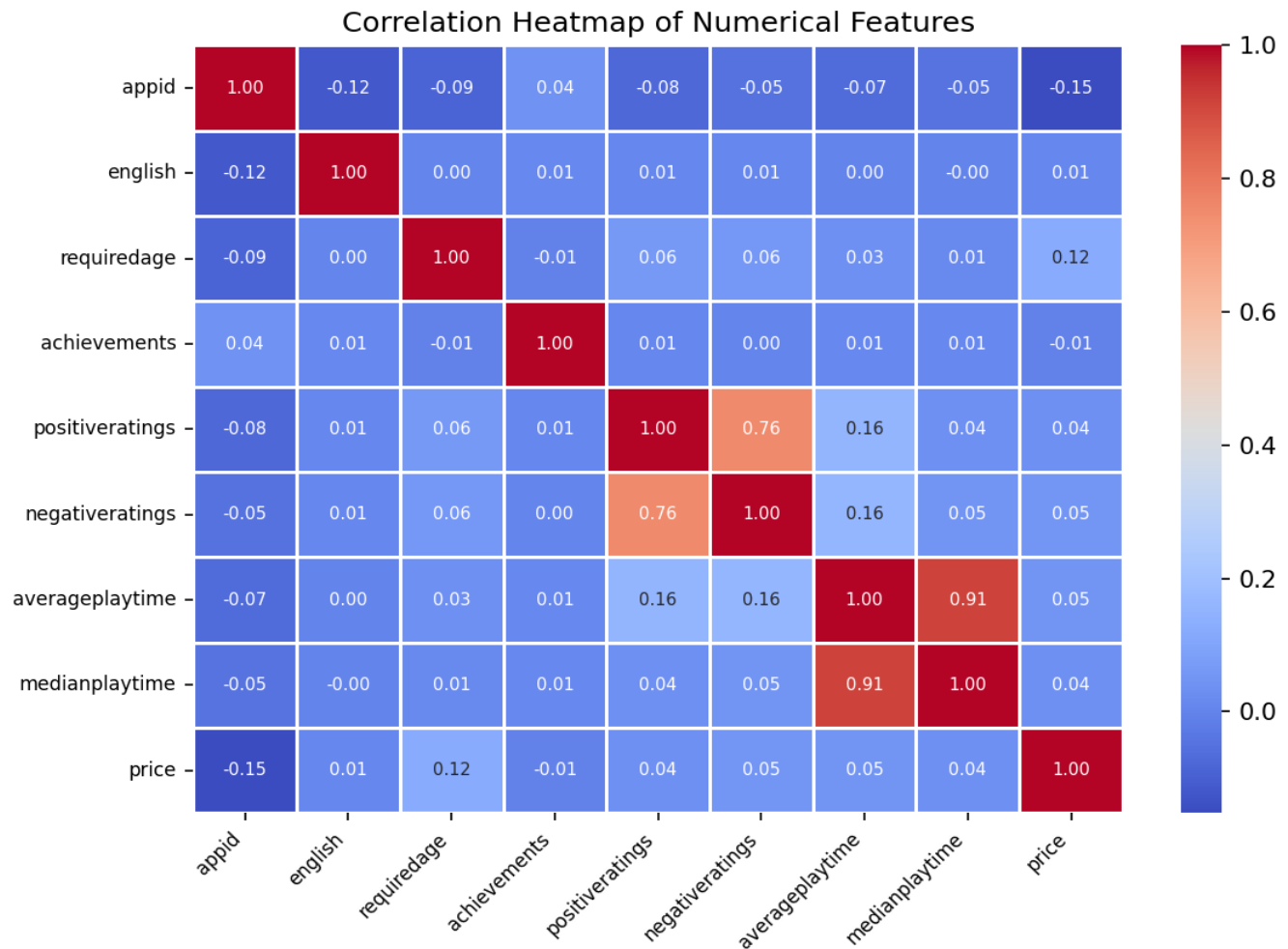


Figure 13: Heatmap visualizing linear correlations (Pearson's r) between numerical features. Values range from -1 (strong negative) to +1 (strong positive). Values near 0 suggest weak linear correlation.

Scatter plots for up to 3 most correlated pairs (absolute value > 0.3):

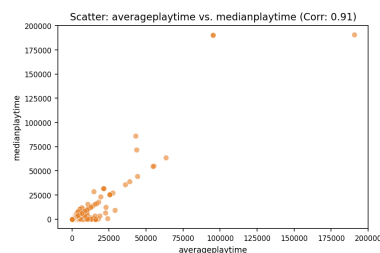


Figure 14: Scatter plot for 'averageplaytime' and 'medianplaytime'. Correlation: 0.91.

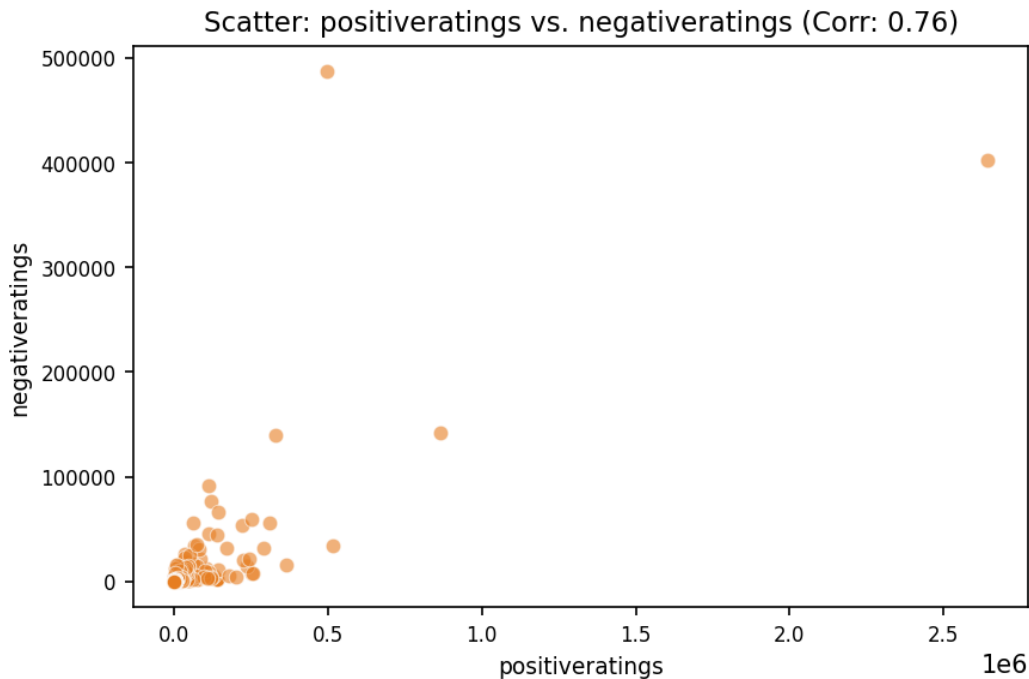


Figure 15: Scatter plot for 'positiveratings' and 'negativeratings'. Correlation: 0.76.

Interpretation of Numerical Correlations:

A correlation matrix is a statistical tool that displays the correlation coefficients between different variables in a dataset. It provides a snapshot of the strength and direction of the relationships between each pair of variables, with values ranging from -1 (perfect negative correlation) to 1 (perfect positive correlation). The correlation matrix can help identify patterns, relationships, and potential interactions between variables, which can be useful in understanding the underlying structure of the data. From the provided list, the strongest correlations observed are between "averageplaytime" and "medianplaytime" (0.914881) and between "positiveratings" and "negativeratings" (0.756570). These strong positive correlations imply that as the average playtime increases, the median playtime also tends to increase, suggesting a consistent relationship between these two variables. Similarly, the correlation between positive and negative ratings suggests that games with more positive ratings tend to have more negative ratings as well, possibly indicating that games with strong opinions (either positive or negative) tend to elicit more ratings overall. If scatter plots were available, it would be interesting to visualize these relationships to see if they appear linear or if there are any outliers or non-linear patterns. For example, a scatter plot of "averageplaytime" vs "medianplaytime" might show a clear linear relationship, while a plot of "positiveratings" vs "negativeratings" might reveal a more complex pattern, such as a curve or a cluster of points. Unfortunately, without more information, we can only speculate about the patterns that might be observed in the scatter plots.

4.2. Numerical vs. Categorical Features

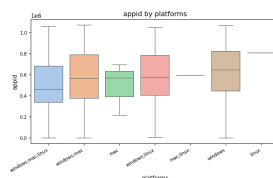


Figure 16: Box plot showing distribution of 'appet' across categories of 'platforms'. This helps visualize differences in central tendency and spread of 'appet' for each category.

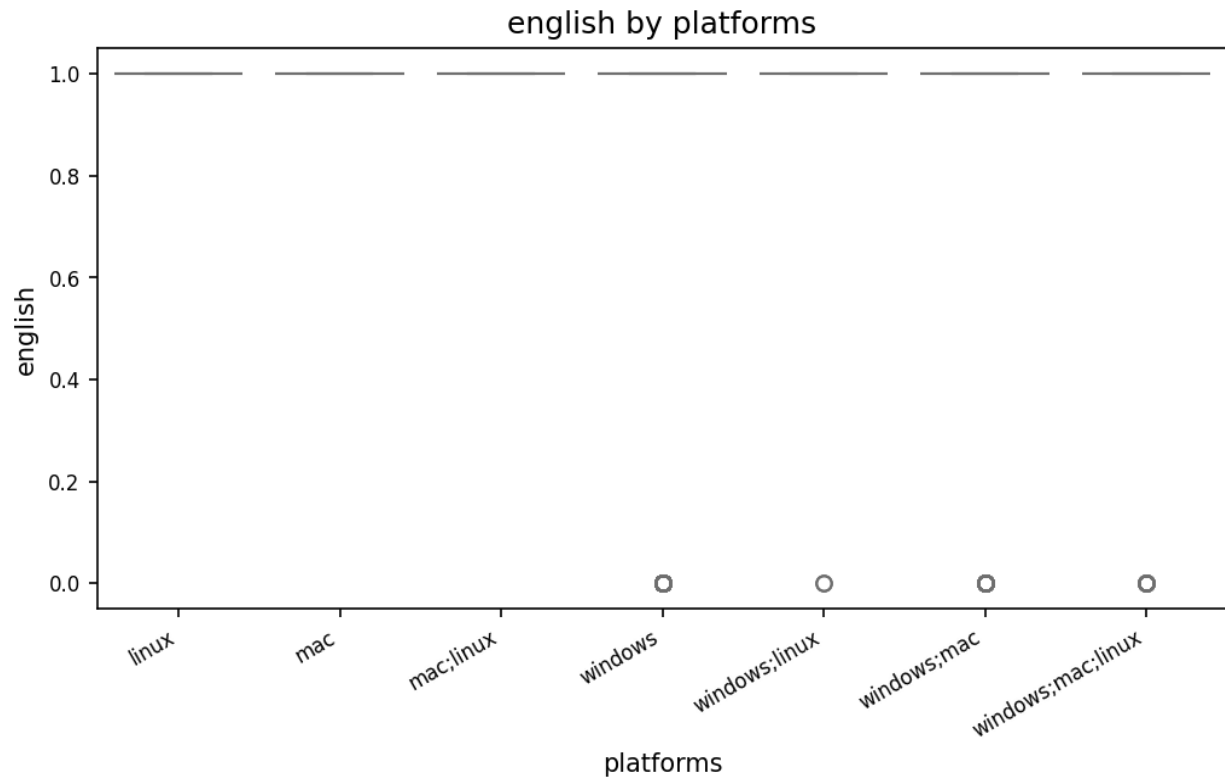


Figure 17: Box plot showing distribution of 'english' across categories of 'platforms'. This helps visualize differences in central tendency and spread of 'english' for each category.

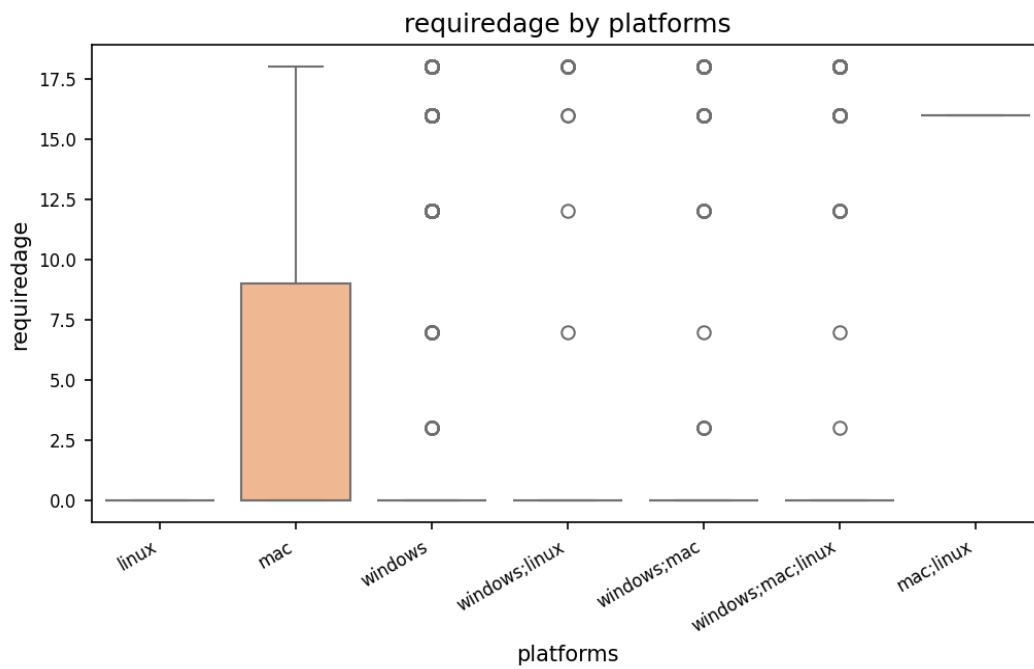


Figure 18: Box plot showing distribution of 'requiredage' across categories of 'platforms'. This helps visualize differences in central tendency and spread of 'requiredage' for each category.

Interpretation of Numerical vs. Categorical Interactions:

Box plots showing numerical distributions per category provide valuable insights into the relationships between categorical features and numerical features in a dataset. By plotting numerical features against categorical features, these box plots reveal the distribution of numerical values within each category. The plots display the median, quartiles, and outliers of the numerical feature for each category, allowing for visual comparisons across categories. For instance, in the case of 'appid' vs 'platforms', the plot might show that the distribution of 'appid' values varies significantly across different 'platforms', with some platforms having higher or lower median 'appid' values than others. If significant differences are observed in medians or spreads across categories, it can indicate meaningful relationships between the categorical feature and the numerical feature. For example, if the median 'requiredage' is higher for games on one platform compared to another, it may suggest that the games on the first platform tend to be more mature or complex. Similarly, if the spread of 'appid' values is wider for one platform than another, it could indicate that the first platform has a more diverse range of applications. These insights can be useful in various applications, such as identifying trends in game development, informing marketing strategies, or optimizing platform-specific features. By examining these plots, one can draw conclusions like "Category A tends to have higher values of Numeric X than Category B", which can inform decisions or guide further analysis to explore the underlying causes of these differences.

4.3. Categorical vs. Categorical Features

5. Key Findings & Insights Summary

Key Findings & Insights The initial analysis of the dataset has revealed several critical data quality issues that could significantly impact the reliability and accuracy of subsequent analyses. Notably, the presence of missing values is reported as "None", indicating a clean dataset in this regard. However, the duplicate count is 10, suggesting that there might be redundancy in the data that could affect analysis outcomes if not addressed. Furthermore, there are no constant columns, which is a positive indicator of data variability. Despite these mixed signals, it is essential to monitor data quality closely to ensure the integrity of the analysis. Univariate analysis has provided insights into the distribution of features within the dataset. The analysis focused on 9 numerical (hist/box) and 9 categorical (bar) features. This suggests a balanced dataset with a variety of feature types, which can be useful for exploring different aspects of the data. The diversity in feature types indicates that the dataset might be suitable for a range of analytical approaches, from identifying patterns and trends in numerical data to understanding category distributions and potential correlations within categorical data. Bivariate analysis aimed to uncover relationships between features. Although the specific details of these relationships are not extensively reported, the completion of bivariate analysis, including num-num correlation and cat-cat analysis, implies that there are findings worth exploring further. The fact that 3 num-cat plots were generated suggests that there are interesting intersections between numerical and categorical features that could reveal valuable insights into how different variables interact. However, without more specific details on the correlations found, the exact nature and strength of these relationships remain to be further elucidated. Considering the summaries provided, there are no explicitly stated surprising or unexpected findings. However, the mere presence of certain types of analyses (like bivariate analysis) implies that the data might hold some nuances or patterns worth deeper investigation. The absence of constant columns and the presence of duplicates might also hint at underlying complexities within the dataset that require careful consideration. Overall, these initial findings offer a foundation for further exploration and a more detailed understanding of the dataset's characteristics and potential insights it may yield.

6. Conclusion & Potential Next Steps

The automated analysis of the 'steam.csv' dataset provides a foundational understanding of the data's characteristics, quality, and potential relationships. This high-level overview reveals key insights into the dataset's structure, including the presence of numerical and categorical features, the absence of missing values, and the presence of a small number of duplicates. These findings synthesize the quality, univariate, and bivariate analyses, offering a solid base for further exploration. Based on the provided summaries, the lack of missing values and constant columns suggests that the data is relatively clean and ready for more in-depth analysis. However, the presence of 10 duplicates may warrant further investigation to determine if these duplicates are due to errors in data collection or if they represent genuine repeats. Additionally, the completion of numerical-numerical correlation analysis and numerical-categorical plots provides a starting point for exploring relationships between variables. Potential next steps for a more in-depth analysis include developing a strategy for handling the 10 duplicate rows to ensure that they do not skew the results of subsequent analyses. Further investigation into the relationships between numerical features, as hinted at by the completion of numerical-numerical correlation analysis, could provide valuable insights into the underlying structure of the data. Examining the results of the numerical-categorical plots to identify if certain categories are associated with distinct numerical distributions could also be a fruitful area of exploration. Lastly, exploring the distribution of values in the numerical features using the histograms and box plots generated during the univariate analysis could help to identify potential outliers or skewness that may impact the results of later analyses. These next steps are designed to build upon the foundational understanding provided by the automated analysis, allowing for a more nuanced and detailed exploration of the 'steam.csv' dataset. By addressing the duplicates, exploring relationships between variables, and examining the distribution of numerical features, a more comprehensive understanding of the data can be developed, ultimately informing future analyses and modeling efforts.