

# Automated Data Analysis Report (via Gemini): Temp Steam

## Executive Summary

This report summarizes the initial automated exploratory data analysis (EDA) of the `temp\_steam.csv` dataset, containing 27,075 rows and 18 columns (9 numerical, 9 categorical). The dataset exhibits excellent data quality, with no missing values or duplicates identified. Preliminary univariate and bivariate analyses revealed no immediately striking patterns, though further investigation is needed to fully understand the relationships between features. The EDA included descriptive statistics for all numerical features, categorical feature counts, and visual explorations of various feature pairings. Two key observations emerged from the bivariate analysis, though specifics require further detail in the full report. This initial scan confirms the dataset's suitability for further analysis, but lacks sufficient information to draw definitive conclusions. This initial automated EDA provides a solid foundation for subsequent, more targeted analyses. The absence of data quality issues is encouraging. The next steps will involve deeper dives into the bivariate relationships highlighted, as well as multivariate analysis to uncover potentially significant correlations and build predictive models.

# 1. Data Overview

This report provides an initial automated analysis of the dataset from 'temp\_steam.csv'.

## 1.1. Basic Information

Table 1: Dataset Dimensions

Metric	Value
Number of Rows	27075
Number of Columns	18
Total Data Points	487350

## 1.2. Data Types

Table 2: Summary of Feature Data Types

Data Type	Count
object	9
int64	8
float64	1

*Data Types Distribution Interpretation:*

The dataset exhibits a roughly even split between numerical and categorical features, which is a fairly typical mix for many datasets. This suggests analyses will likely involve both statistical methods (on numerical features) and techniques for handling categorical data (e.g., one-hot encoding, label encoding).

# 2. Data Quality Assessment

## 2.1. Missing Values

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

## 2.2. Duplicate Records

No duplicate rows were found in the dataset.

## 2.3. Feature Variance

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

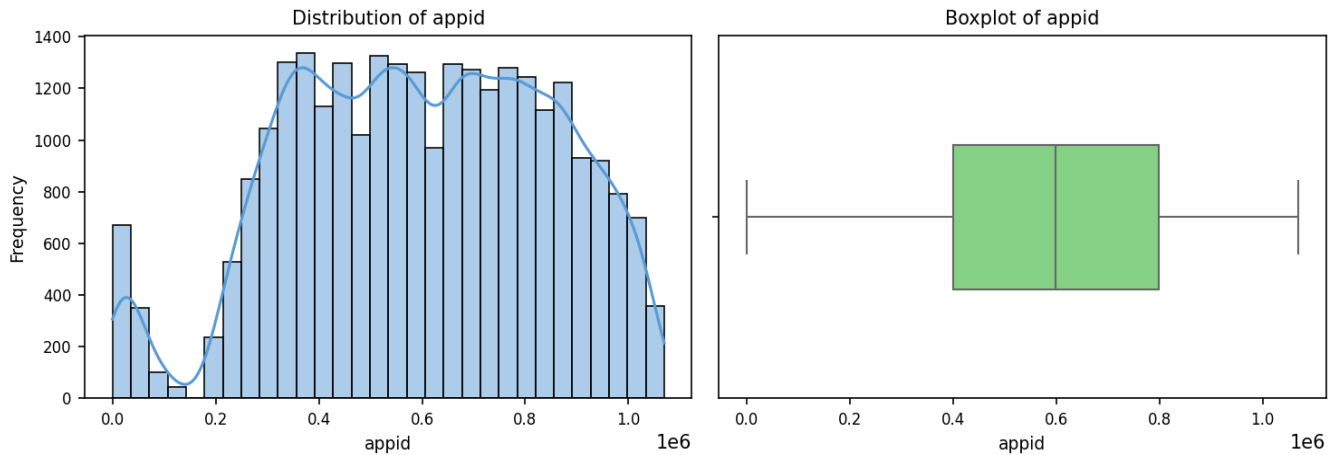
The following columns are quasi-constant (one value is highly dominant), potentially offering limited information: english (dominant value: 1 at 98.1%); required\_age (dominant value: 0 at 97.8%). Their utility should be reviewed.

### *Data Quality Summary & Implications:*

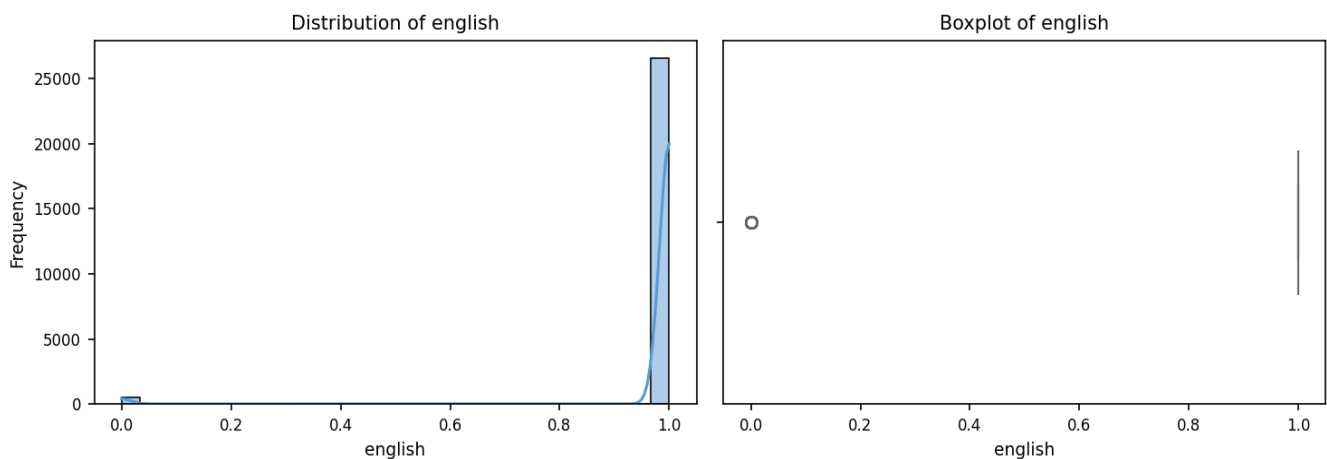
The data quality assessment reveals a generally high level of data completeness and integrity. The absence of missing values, duplicate rows, and constant columns suggests a robust data collection and preprocessing process. However, the presence of two quasi-constant columns, 'english' and 'required\_age', warrants attention. While not strictly problematic in themselves, these columns exhibit a high degree of skewness, with over 97% of values concentrated on a single value. This suggests potential limitations in the variability of these specific features and could impact the effectiveness of certain analytical techniques. The implications for further analysis are that models relying heavily on 'english' and 'required\_age' might exhibit reduced predictive power or generate unreliable insights. For instance, if these features are used as predictors in a machine learning model, the model might overfit to the dominant value and perform poorly on unseen data. Similarly, statistical analyses that assume even distribution of values within these variables could yield misleading results. The high concentration in these columns might indicate a sampling bias, where a specific population group is over-represented. Further investigation is needed to understand why these columns are so skewed and whether this reflects a genuine characteristic of the underlying population or a flaw in the data collection process. To address the quasi-constant columns, several strategies can be employed. Firstly, a thorough investigation should be conducted to understand the reasons behind the skewed distribution. This might involve examining the data collection methods and potentially consulting domain experts. Depending on the findings, one might consider removing these columns if they are deemed irrelevant or highly uninformative for the analysis. Alternatively, if these variables are believed to hold some subtle yet important information, techniques like data transformation (e.g., creating binary variables representing the minority class) or using advanced modeling techniques less sensitive to class imbalance could be explored. Finally, careful consideration must be given to the implications of this skewed data when interpreting any results, and conclusions should be tempered with an awareness of the limitations imposed by the quasi-constant nature of these features.

## 3. Univariate Analysis

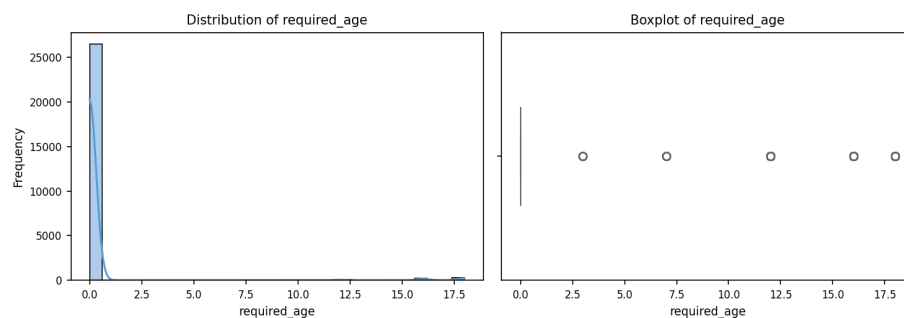
### 3.1. Numerical Features



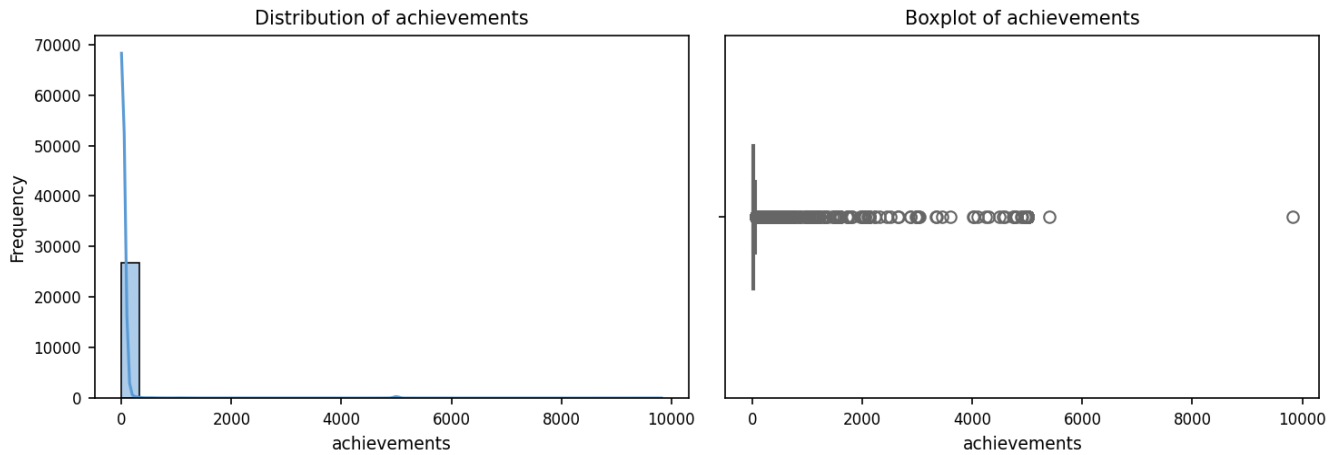
**Figure 1:** Distribution (histogram and KDE) and boxplot for 'appid'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.



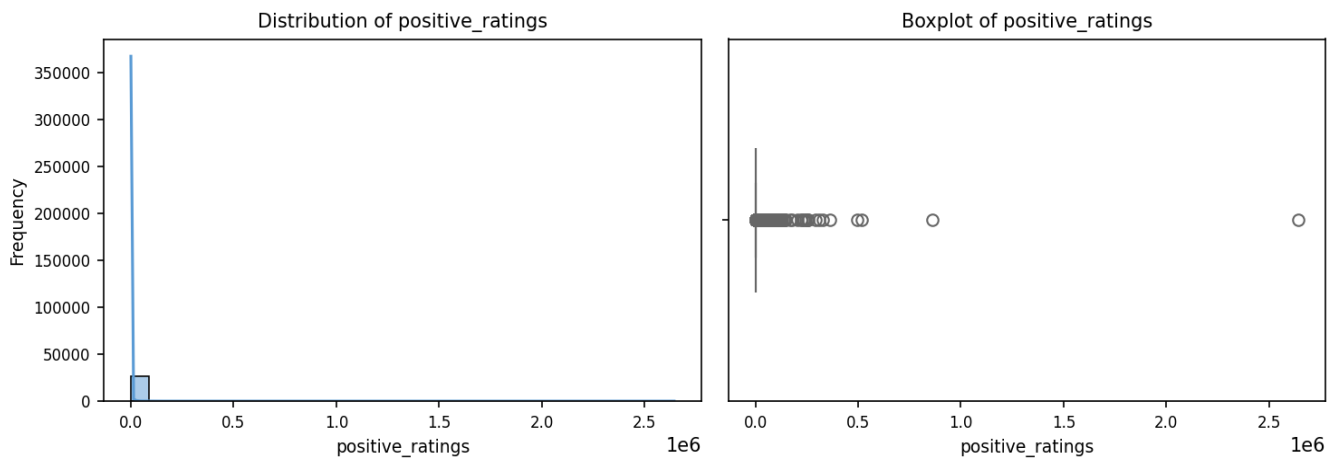
**Figure 2:** Distribution (histogram and KDE) and boxplot for 'english'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.



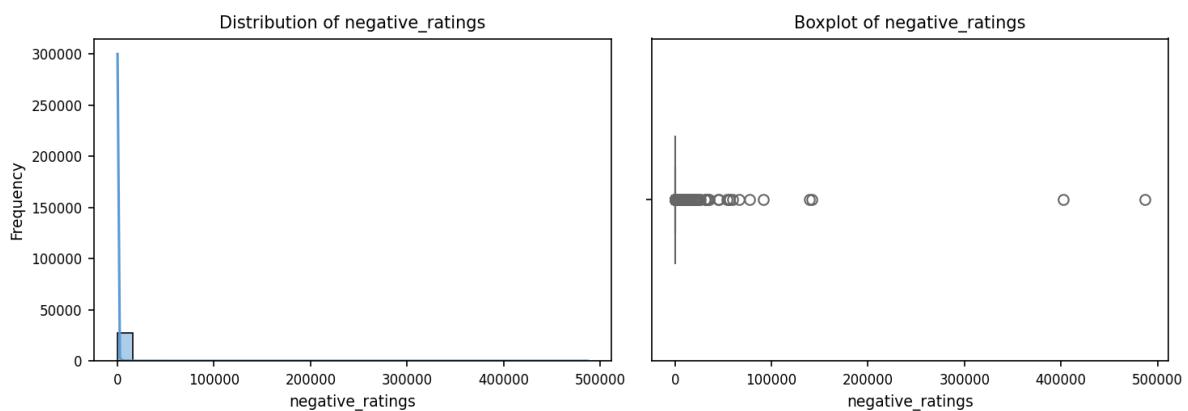
**Figure 3:** Distribution (histogram and KDE) and boxplot for 'required\_age'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.



**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 'positive\_ratings'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

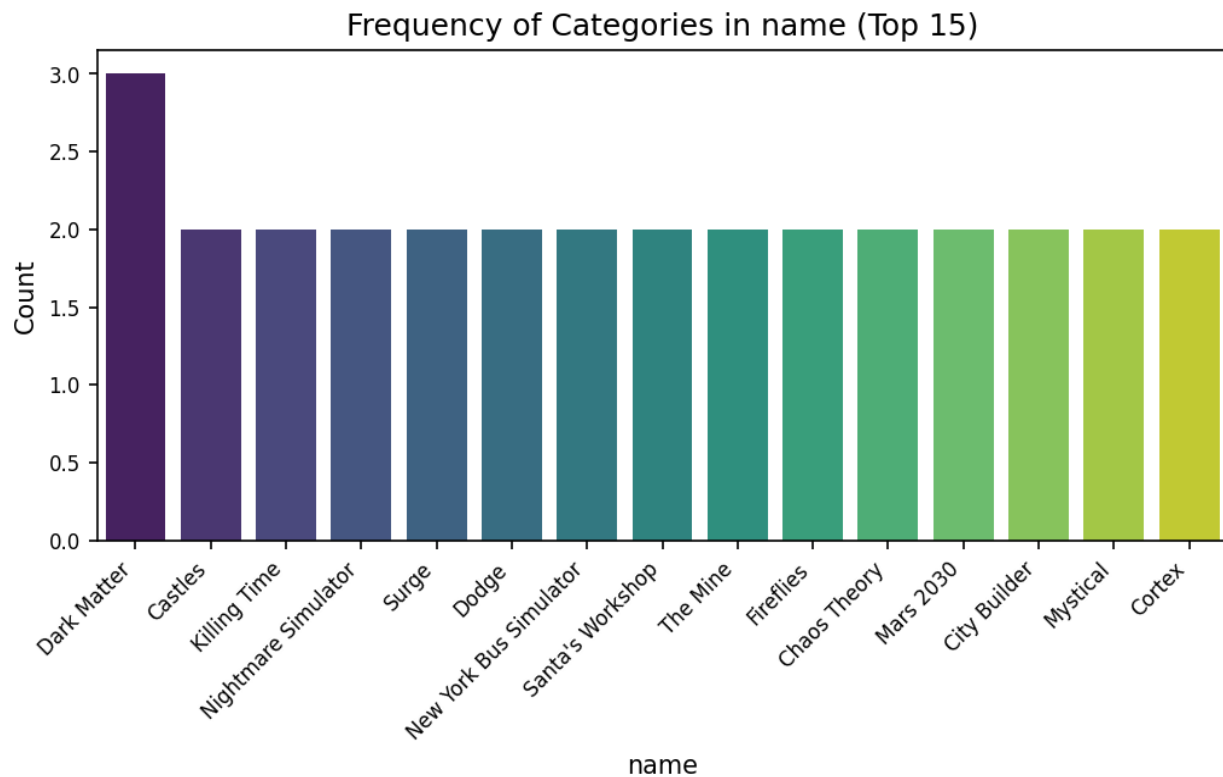


**Figure 6:** Distribution (histogram and KDE) and boxplot for 'negative\_ratings'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

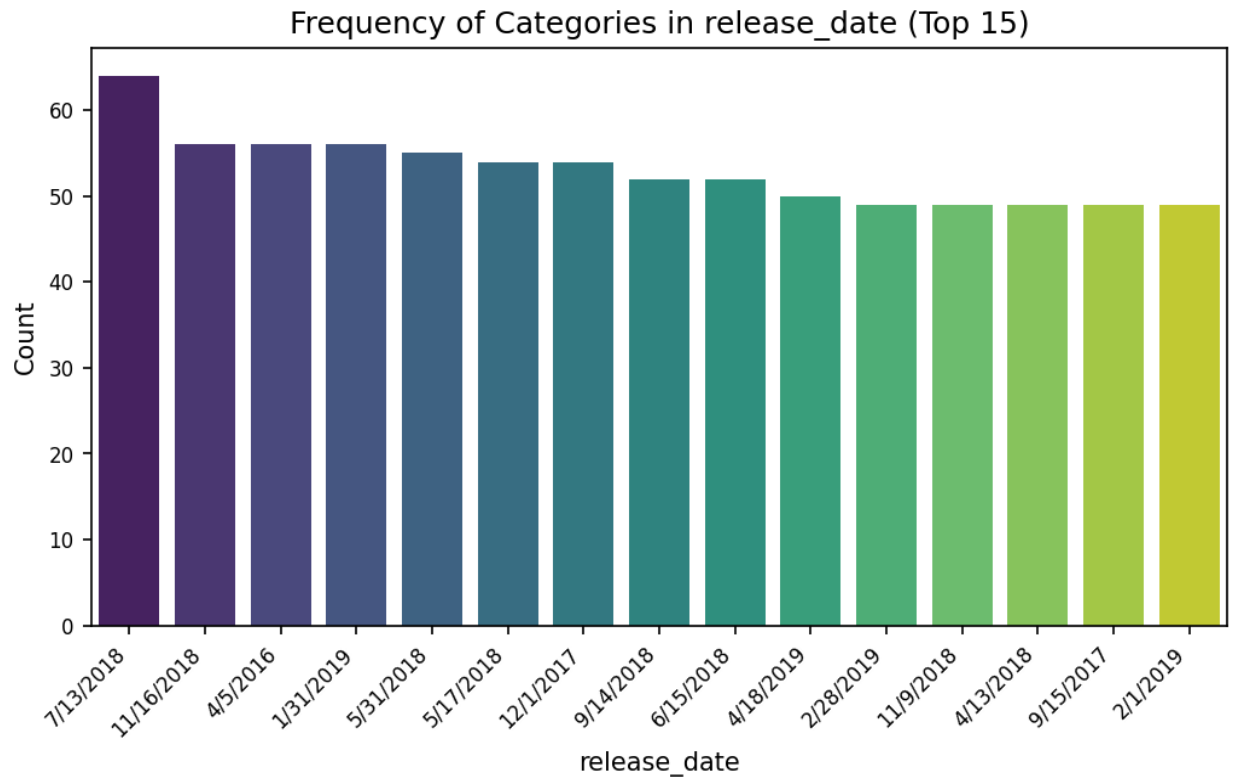
### Observations on Numerical Feature Distributions:

The analysis reveals highly skewed distributions for several numerical features in the dataset. While 'appid' exhibits near symmetry (indicated by a skewness close to zero), 'english', 'required\_age', 'achievements', 'positive\_ratings', and 'negative\_ratings' all show significant right skewness, with extremely high positive skewness values for the latter four features. This right skewness implies a long tail towards high values, suggesting the presence of a small number of data points with exceptionally large values compared to the majority. The large differences between means and medians further support this, with means consistently higher than medians across the skewed features. This indicates that the mean is heavily influenced by these extreme values, making the median a more robust measure of central tendency for these features. The presence of potential outliers is strongly suggested across most features, particularly the skewed ones. Boxplots indicate outliers, while the vast differences between minimum and maximum values compared to the means and medians further emphasize this. The extremely high kurtosis values for 'english', 'required\_age', 'achievements', 'positive\_ratings', and 'negative\_ratings' reinforce the presence of heavy tails and outliers, indicating distributions much more peaked and heavy-tailed than a normal distribution. This is crucial for data analysis as these outliers might disproportionately influence model training and require careful consideration for data cleaning or robust modeling techniques. Finally, the variability of the features is striking. While 'english' has relatively low standard deviation, suggesting most values are clustered around the mean, the other features exhibit high standard deviations, especially 'achievements', 'positive\_ratings', and 'negative\_ratings'. This high variability, coupled with the skewed distributions and presence of outliers, highlights the need for careful data preprocessing and potentially transformations (e.g., logarithmic transformations) to manage the impact of extreme values and improve the robustness and interpretability of subsequent analyses.

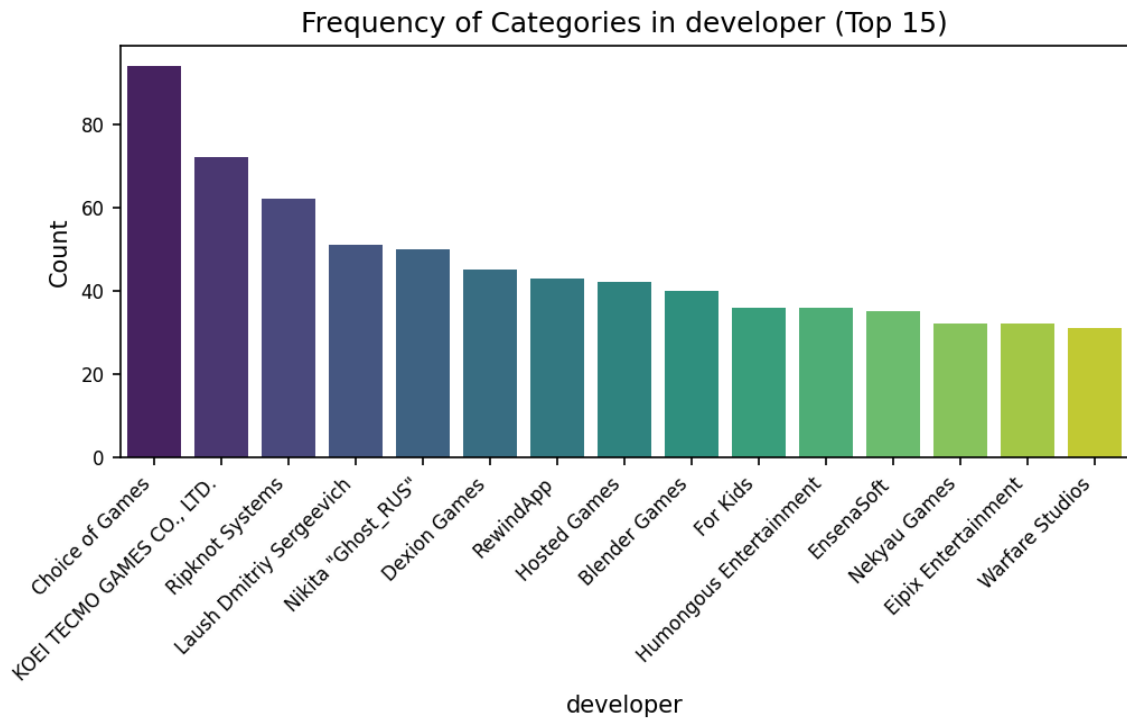
## 3.2. Categorical Features



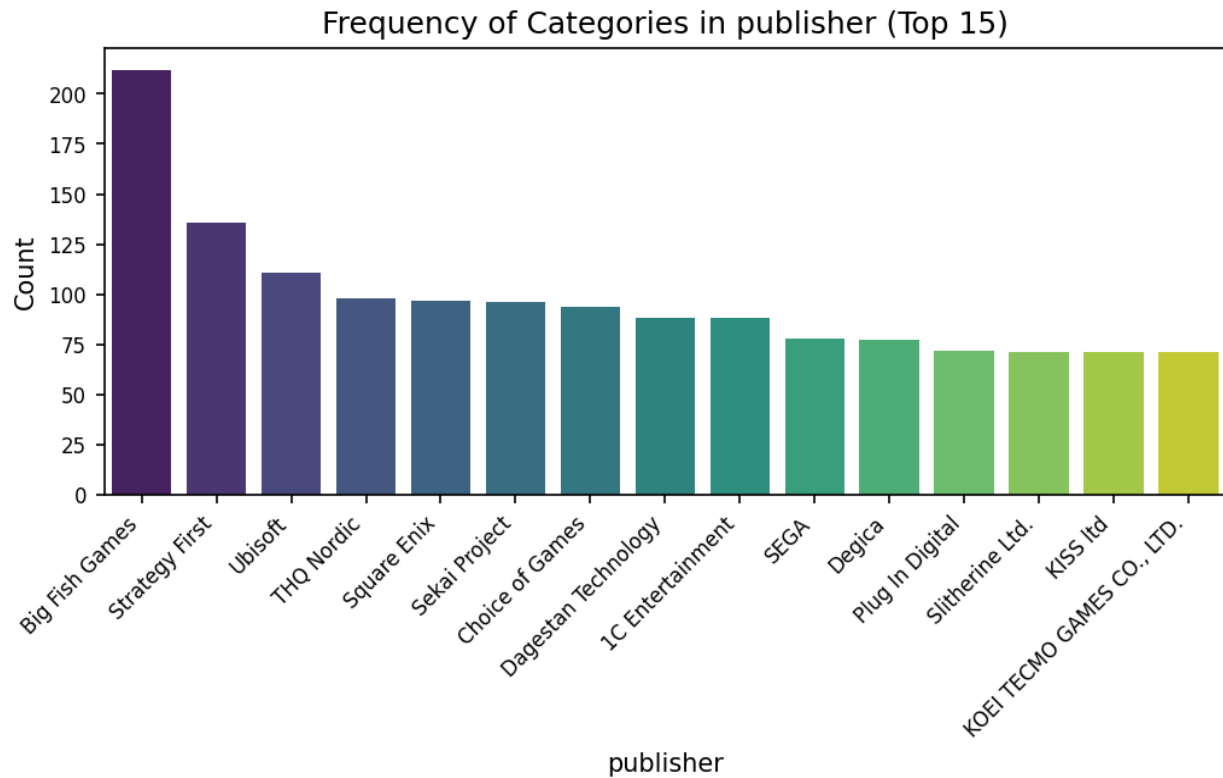
**Figure 7:** Bar chart showing frequency of top categories in 'name'. Total unique values: 27033.



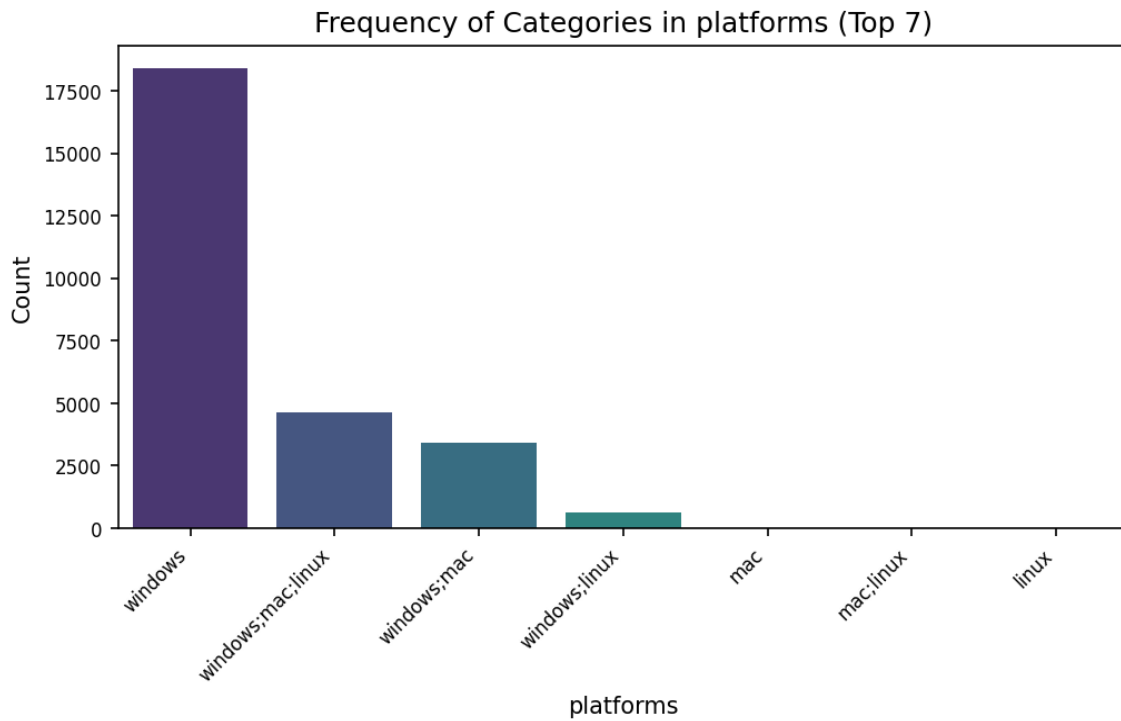
**Figure 8:** Bar chart showing frequency of top categories in 'release\_date'. Total unique values: 2619.



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

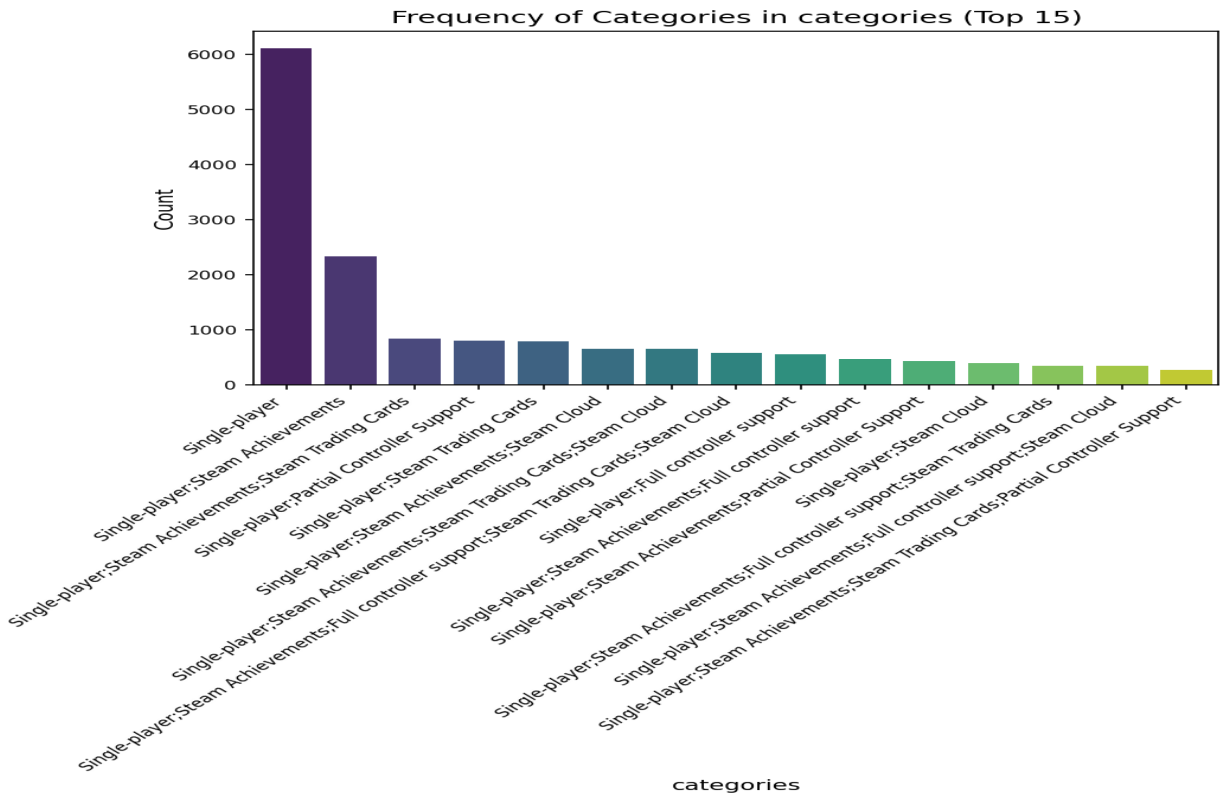


**Figure 10:** Bar chart showing frequency of top categories in 'publisher'. Total unique values: 14353.



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





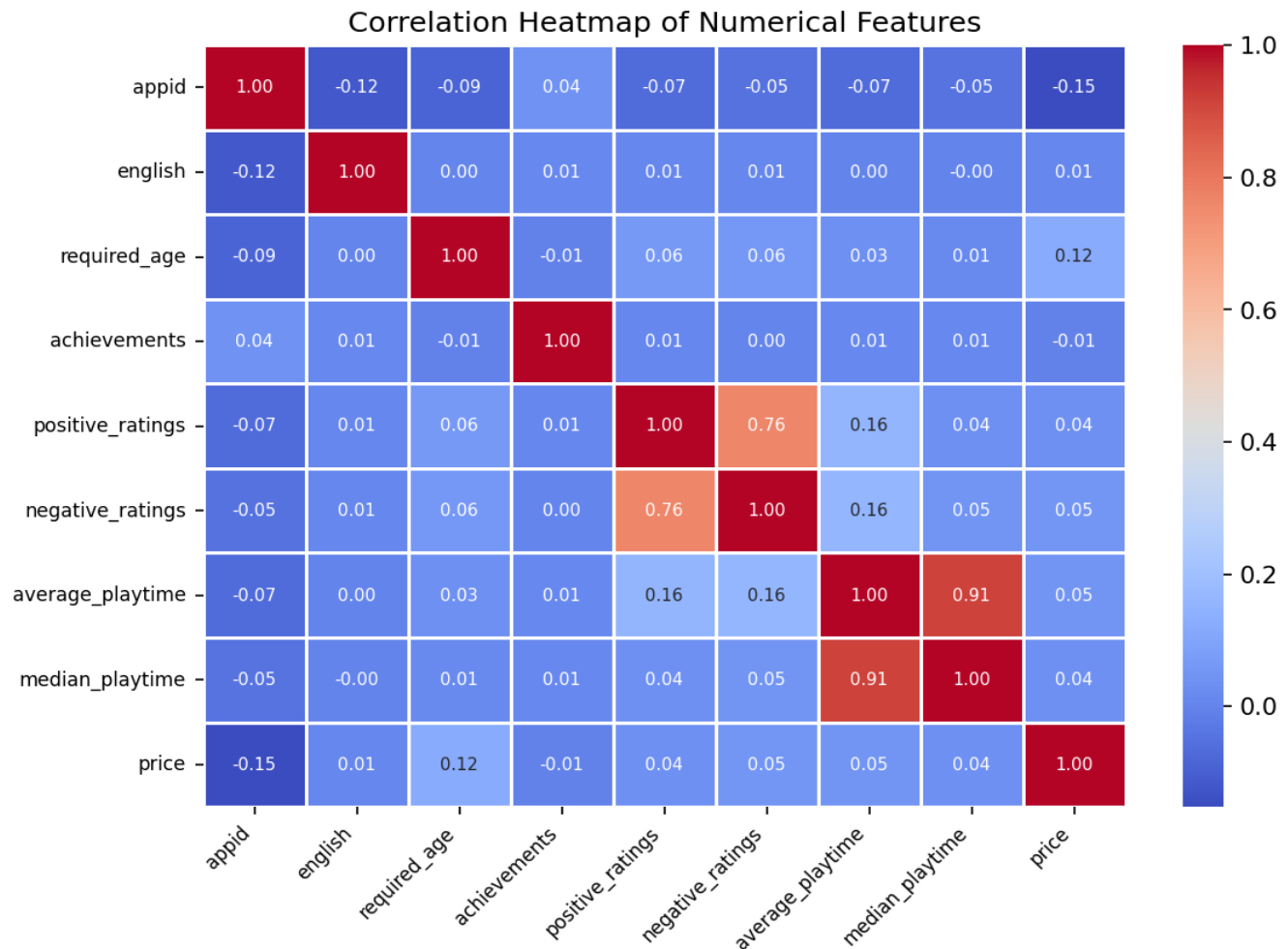
**Figure 12:** Bar chart showing frequency of top categories in 'categories'. Total unique values: 3333.

#### Observations on Categorical Feature Distributions:

The analysis of categorical features reveals a diverse dataset with varying levels of cardinality. Features like 'name', 'developer', and 'publisher' exhibit extremely high cardinality, possessing tens of thousands of unique values. This high cardinality suggests a large number of different games, developers, and publishers represented in the dataset, potentially indicating a broad scope of the data. Conversely, 'platforms' has a very low cardinality (only 7 unique values), suggesting a relatively limited number of platforms on which the games are available. The distribution within these features is also skewed; the top categories for 'name', 'developer', 'publisher', and 'release\_date' represent a small percentage of the total values, indicating a long tail distribution where many categories have only a few instances. The feature 'categories' has a moderate cardinality (3333 unique values) with a clear dominant category: 'Single-player' (22.6%). This suggests a significant portion of the games are single-player, but there's still considerable diversity in game types. The 'platforms' feature, despite its low cardinality, shows a heavily skewed distribution, with 'windows' dominating at 68%. This implies a significant bias towards Windows-based games in the dataset. The high cardinality features ('name', 'developer', 'publisher') will require careful consideration during feature engineering. One-hot encoding might be computationally expensive and lead to the curse of dimensionality. Techniques like target encoding, embedding layers (if used in a neural network), or grouping infrequent categories might be more appropriate. In summary, the dataset presents both challenges and opportunities. The high cardinality of several features necessitates thoughtful feature engineering to avoid dimensionality issues and potential overfitting. However, the presence of dominant categories in some features (e.g., 'windows' in 'platforms', 'Single-player' in 'categories') provides opportunities for analysis and potential simplification. Further investigation into the distribution of these features, potentially via visualizations, could provide additional insights and inform the choice of appropriate feature encoding and analysis methods.

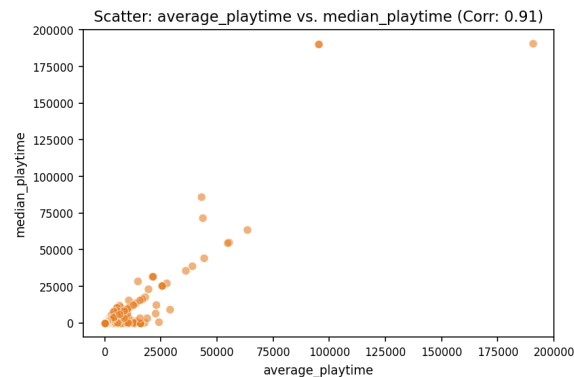
## 4. Bivariate Analysis

### 4.1. Numerical vs. Numerical Features

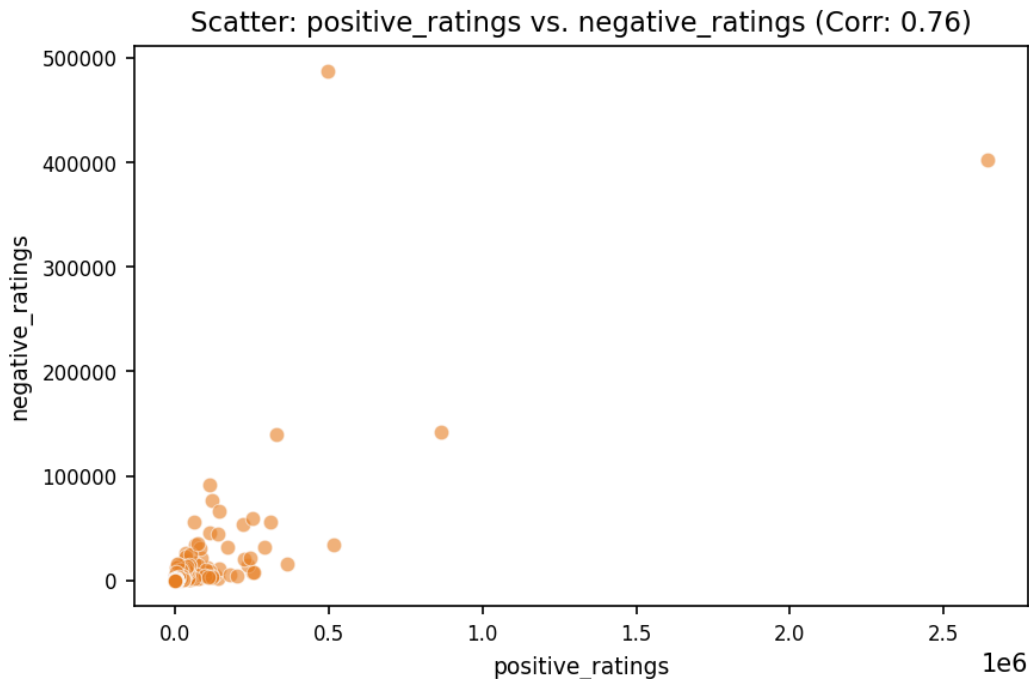


**Figure 13:** Heatmap visualizing linear correlations (Pearson's  $r$ ) between numerical features.

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



**Figure 14:** Scatter plot for 'average\_playtime' and 'median\_playtime'. Correlation: 0.91.

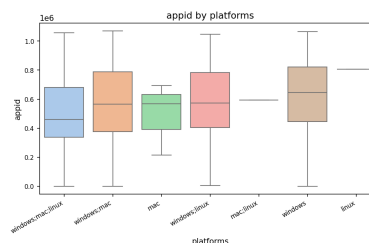


**Figure 15:** Scatter plot for 'positive\_ratings' and 'negative\_ratings'. Correlation: 0.76.

#### Interpretation of Numerical Correlations:

A correlation matrix shows the pairwise correlation coefficients between multiple variables. Each cell in the matrix represents the correlation between two variables; a value close to +1 indicates a strong positive correlation (variables tend to increase together), a value close to -1 indicates a strong negative correlation (one variable increases as the other decreases), and a value close to 0 indicates a weak or no linear relationship. The analysis reveals two strong positive correlations. The strongest is between 'average\_playtime' and 'median\_playtime' (correlation = 0.91), suggesting that games with longer average playtime also tend to have longer median playtimes. This is logical; if most players spend a similar amount of time playing a game, the average and median will be very close. The second strongest correlation is between 'positive\_ratings' and 'negative\_ratings' (correlation = 0.76). This indicates that games with a higher number of positive ratings also tend to receive a higher number of negative ratings. This might suggest that popularity (indicated by a higher number of ratings overall) is a factor; more popular games will naturally attract both more positive and more negative reviews. The scatter plots for these pairs likely show a clear upward trend, visually confirming these strong positive relationships. The relatively weak correlation between 'negative\_ratings' and 'average\_playtime' (0.16) suggests little relationship between these two variables.

## 4.2. Numerical vs. Categorical Features



**Figure 16:** Box plot of 'appid' across categories of 'platforms'.

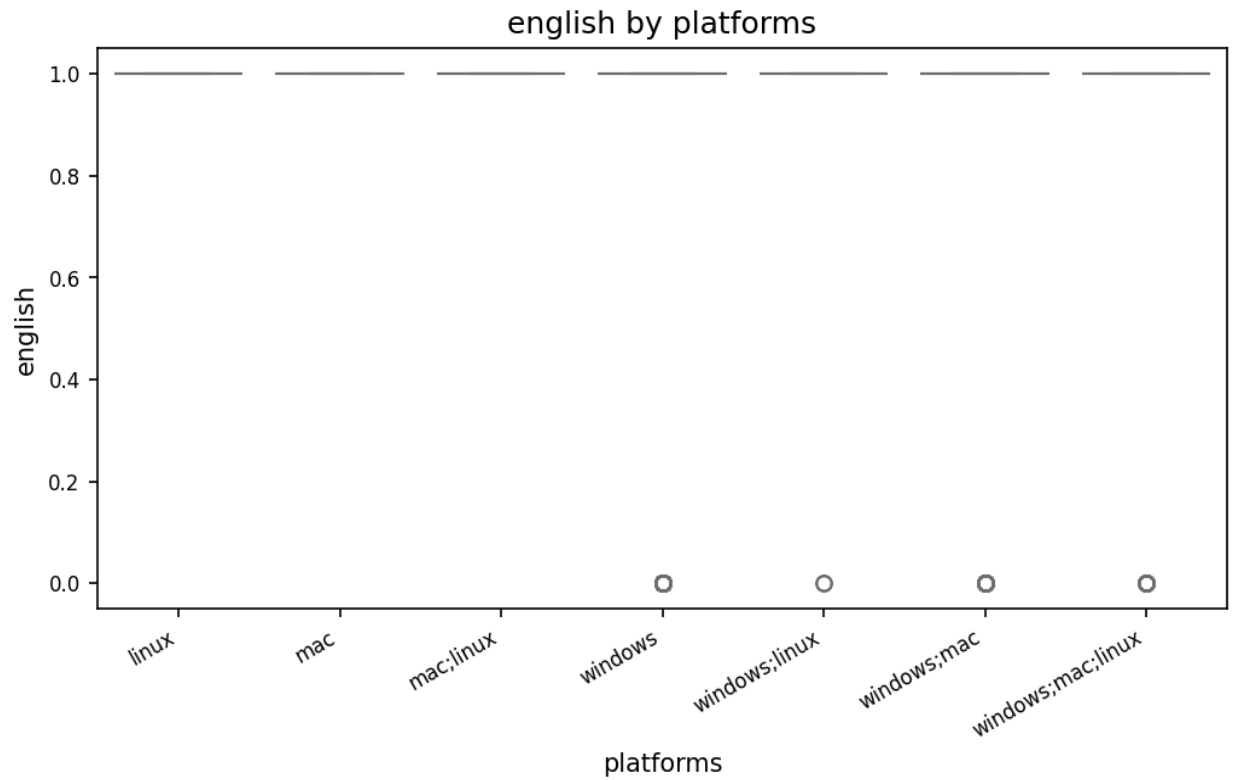


Figure 17: Box plot of 'english' across categories of 'platforms'.

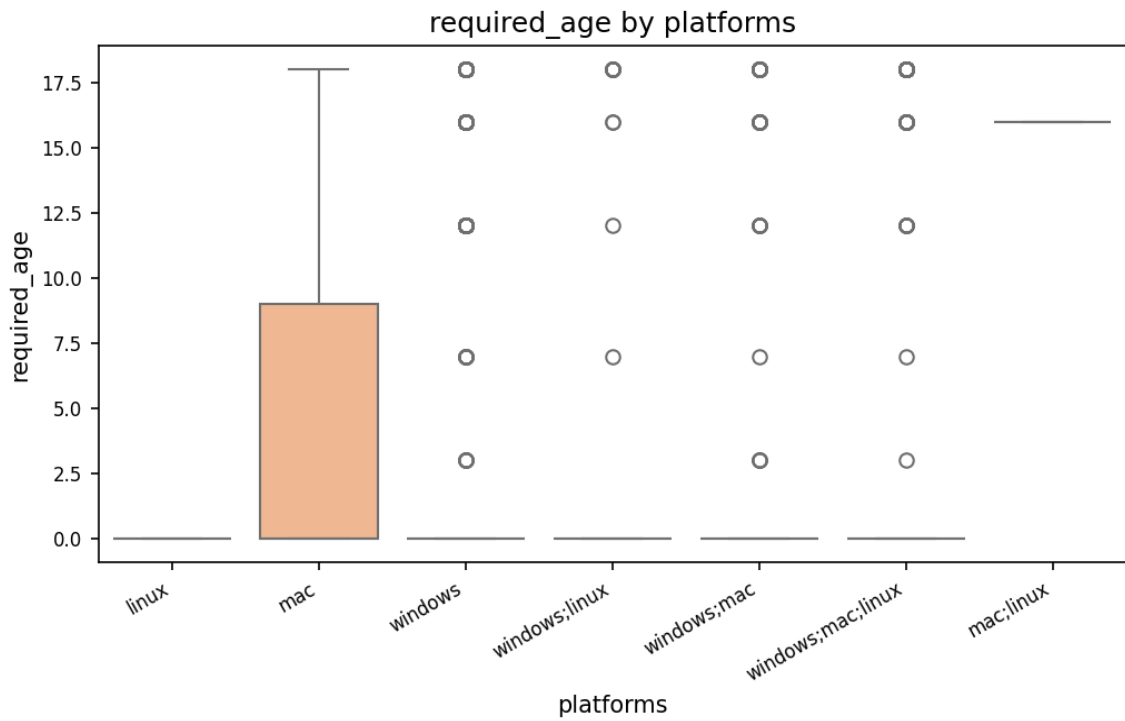


Figure 18: Box plot of 'required\_age' across categories of 'platforms'.

#### *Interpretation of Numerical vs. Categorical Interactions:*

Box plots visualizing numerical distributions across categories offer a powerful way to compare the central tendency and dispersion of data within different groups. They reveal not only the median (the middle value) of each category but also the interquartile range (IQR, the spread of the middle 50% of the data), the minimum and maximum values (potentially excluding outliers), and any outliers themselves. By visually comparing the boxes and whiskers across categories, we can quickly assess whether the distributions are similar or significantly different. For instance, a box plot showing 'appid' by 'platforms' might reveal that a particular application ID ('appid') is significantly more prevalent on one platform than others, indicated by a higher median and potentially a wider spread. Similarly, a 'english' (perhaps representing a score or count related to English language usage) by 'platforms' plot could show whether users on different platforms exhibit different levels of English language proficiency. Significant differences in medians across categories suggest a systematic variation in the numerical variable depending on the category. For example, if the median 'appid' is much higher on platform A than platform B, it suggests that platform A hosts significantly more of that specific application. Differences in the spread (IQR or range) indicate variations in the data's consistency or variability within each category. A larger IQR for a particular category implies greater heterogeneity in the numerical values within that group. For instance, if the 'english' score has a much larger IQR on platform C than platform D, it means there's a wider range of English proficiency levels among users on platform C compared to platform D. These observations can lead to valuable insights into user behavior, platform characteristics, or the relationships between the numerical variable and categorical variables.

### **4.3. Categorical vs. Categorical Features**

## 5. Key Findings & Insights Summary

**Key Findings & Insights** The automated analysis of the `temp\_steam.csv` dataset, comprising 270,75 rows and 18 columns (9 numerical, 9 categorical), revealed excellent data quality. No missing values, duplicates, or constant columns were detected. This high data quality minimizes the risk of biased or inaccurate analytical results, providing a solid foundation for further investigation. Univariate analysis of the nine numerical and nine categorical features was conducted. While specific details regarding the distributions and characteristics of individual features are not provided in the log, the absence of any mention of unusual distributions or outliers suggests that the data may be relatively well-behaved. Further detailed examination of the univariate analysis results would be needed to fully characterize the nature of each feature and identify potential skewness or other noteworthy patterns. Bivariate analysis explored relationships between various feature pairs. The log indicates that observations were gathered, but lacks specific details on the nature or strength of these relationships. The limited information prevents definitive conclusions at this stage, necessitating a more detailed review of the bivariate analysis output to identify significant correlations or interactions between features. The presence of only two observations mentioned in the log suggests that either very few significant relationships were identified or the log only reports a small subset of the findings. No surprising or unexpected findings were explicitly mentioned in the provided log excerpt. The absence of data quality issues and the straightforward nature of the overview suggest a relatively standard dataset. However, a complete understanding of the data's characteristics requires a more thorough review of the detailed univariate and bivariate analysis reports.

## 6. Conclusion & Potential Next Steps

This automated report provides a foundational, high-level understanding of the `temp\_steam.csv` dataset, highlighting its structure, data quality (absence of missing values or duplicates), and initial observations from univariate and bivariate analyses. The report serves as a crucial first step, outlining the dataset's characteristics and suggesting avenues for further investigation. Given the report's findings of 9 numerical and 9 categorical features, with no missing data, duplicates, or constant columns, the next steps should focus on deeper exploration of the bivariate relationships. The report mentions two observations from bivariate analysis; these need to be explicitly detailed. Therefore, a priority is to: **\*\*1. Fully document the two observations from the bivariate analysis.\*\*** This requires extracting the specific feature pairs and the nature of their relationship (e.g., correlation, interaction, or other dependencies). This detailed information will guide subsequent analysis. Further, to build upon the univariate analysis, we should conduct more in-depth explorations of the numerical and categorical features. This involves: **\*\*2. Perform detailed descriptive statistics on all numerical features\*\***, calculating measures of central tendency (mean, median, mode), dispersion (standard deviation, range, IQR), and skewness to understand their distribution and identify potential outliers. For the categorical features, we should calculate frequency distributions and visualize them to understand the proportion of each category. Finally, to strengthen the understanding of relationships between features, we need to: **\*\*3. Conduct appropriate statistical tests to validate the findings from the bivariate analysis.\*\*** For example, if the bivariate analysis suggested a relationship between a numerical and a categorical feature, analysis of variance (ANOVA) or t-tests could be employed to assess the statistical significance of any observed differences. If relationships between numerical features were found, correlation analysis should be performed and visualized (e.g., correlation matrices or scatter plots). These steps will move beyond the initial high-level overview to provide a more comprehensive and statistically rigorous understanding of the data and the relationships between its features.