# 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 27088 rows and 18 columns (9 numerical, 9 categorical). Preliminary analysis revealed a relatively clean dataset with minimal issues: only 13 duplicate entries were identified, and no missing values or constant columns were detected. Initial univariate and bivariate analyses, including descriptive statistics and visualizations, have been conducted. Two key observations emerged from the bivariate analysis, though further investigation is warranted. The dataset's size and apparent cleanliness suggest ample potential for meaningful insights. The absence of significant data quality problems facilitates efficient progression to more advanced analytical techniques. This initial EDA provides a strong foundation for subsequent, more in-depth analyses. The findings suggest the data is suitable for further modeling and predictive analysis, and the next steps will focus on exploring the identified bivariate relationships and developing more specific hypotheses.

# 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 | 27088 |
| Number of Columns | 18 |
| Total Data Points | 487584 |

## 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 shows a roughly even split between numerical and categorical features, which is a typical mix for many datasets. The absence of datetime features might limit the ability to perform time-series analysis or incorporate temporal trends into the analysis.

# 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

The dataset contains 13 duplicate rows (representing 0.05% 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%); required_age (dominant value: 0 at 97.8%). Their utility should be reviewed.

*Data Quality Summary & Implications:*

The data quality assessment reveals a dataset with generally high quality. The absence of missing values is a significant positive, indicating a complete dataset ready for analysis. The extremely low rate of duplicate rows (0.05%) is also negligible and can likely be addressed with simple deduplication. The presence of quasi-constant columns, 'english' and 'required_age', however, warrants further investigation. While not strictly problematic, their high dominance in a single value (98.1% and 97.8% respectively) suggests these variables may offer limited predictive power in many modeling scenarios and could potentially be redundant. Their inclusion might unnecessarily inflate the dimensionality of the data, potentially leading to less efficient or less interpretable models. The implications for further analysis are relatively minor given the overall high quality. The negligible number of duplicate rows can be easily cleaned. The quasi-constant columns, however, may impact model performance depending on the analytical goals. If these variables are not crucial for the research question, they could be removed to simplify the dataset and potentially improve model efficiency. If they are deemed important, their limited variability should be considered when interpreting results and building models. Over-reliance on insights derived from these variables should be avoided. To address the identified issues, a simple deduplication process should be implemented to remove the 13 duplicate rows. For the quasi-constant columns, a thorough examination of their relevance to the research question is necessary. If they prove to be unhelpful or redundant, they can be removed from the dataset. Alternatively, if maintaining these columns is deemed essential, their low variance should be explicitly acknowledged and accounted for during analysis and model interpretation. Further investigation into the few cases where 'english' is 0 or 'required_age' is 1 could provide valuable context and potentially reveal underlying patterns.

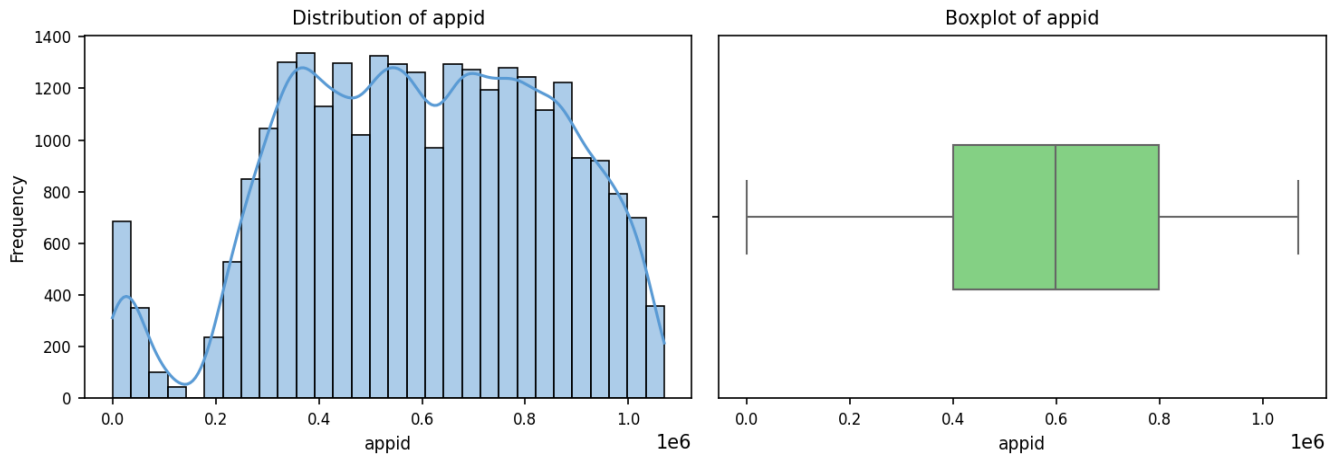# 3. Univariate Analysis

## 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.*
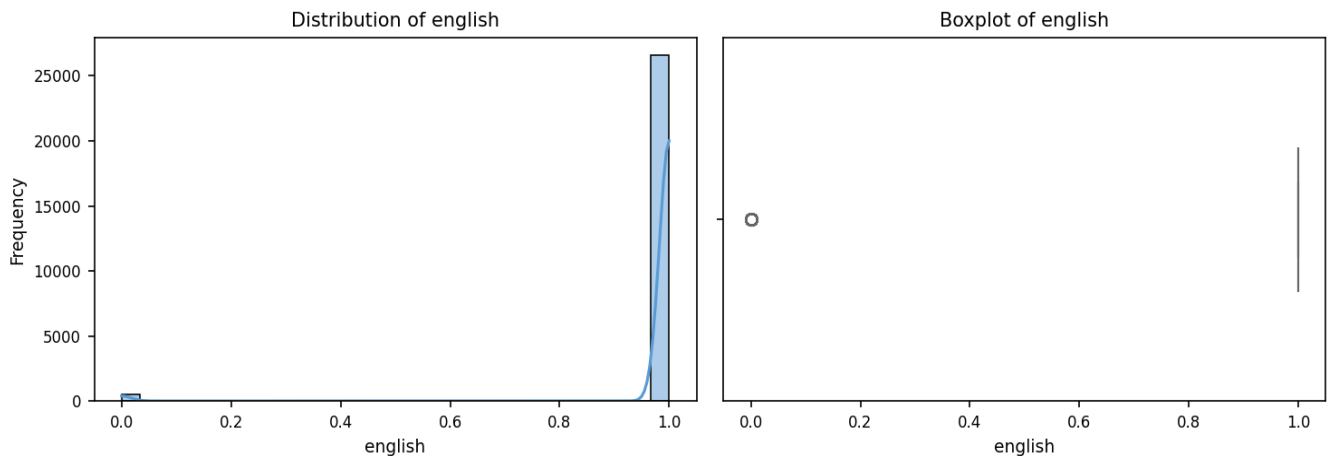


*Figure 2:* *Distribution (histogram and KDE) and boxplot for 'english'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.*
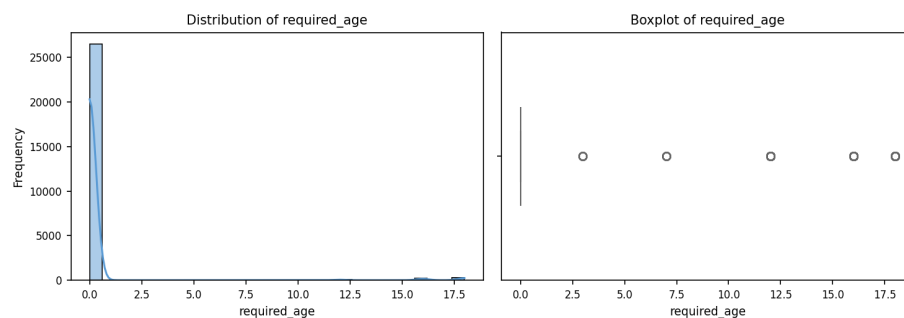


*Figure 3:* *Distribution (histogram and KDE) and boxplot for 'required_age'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.*
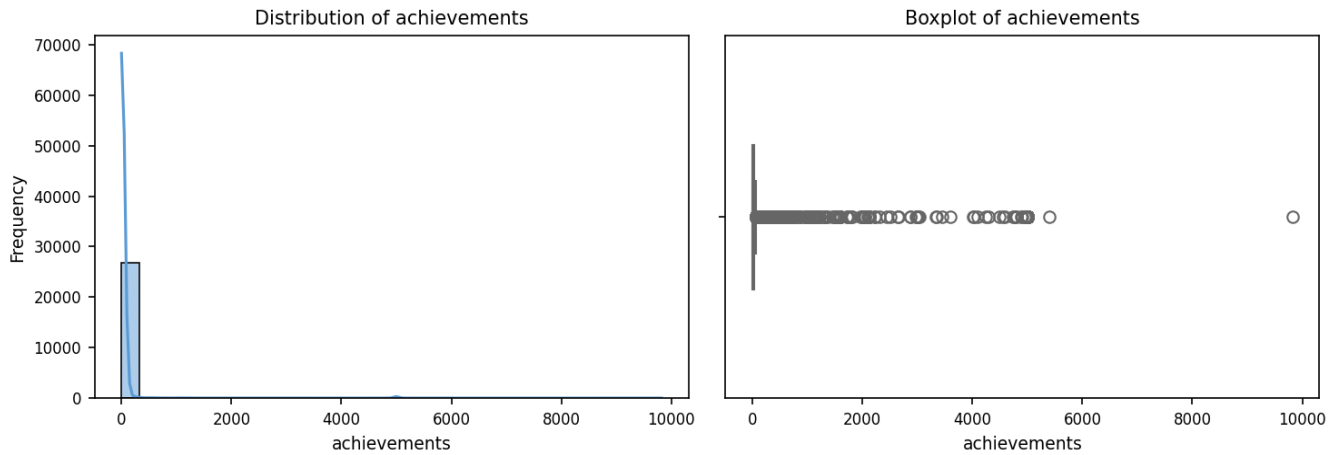
***Figure 4:*** *Distribution (histogram and KDE) and boxplot for 'achievements'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.*
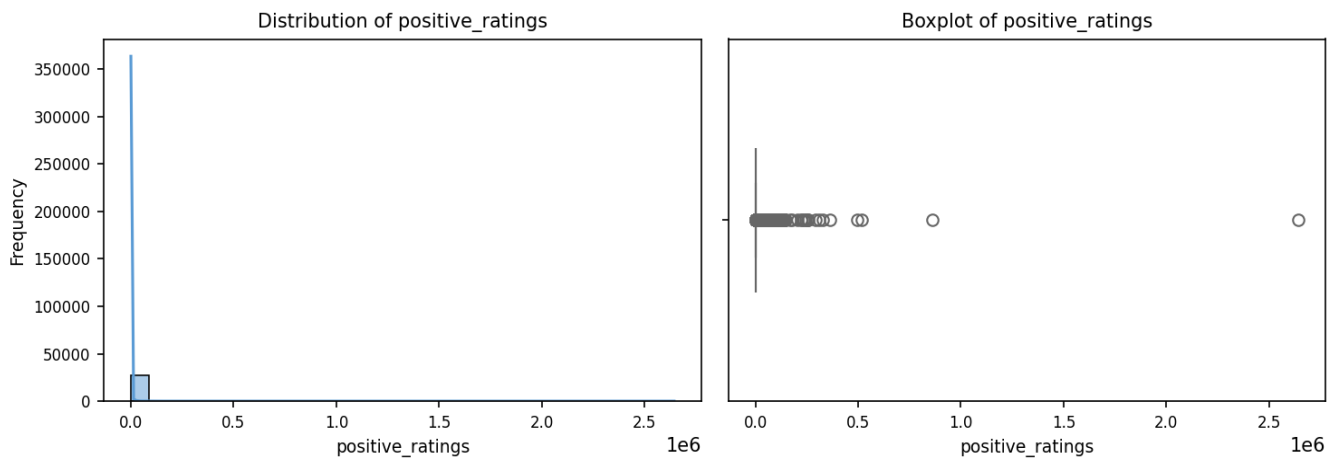


***Figure 5:*** *Distribution (histogram and KDE) and boxplot for 'positive_ratings'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.*
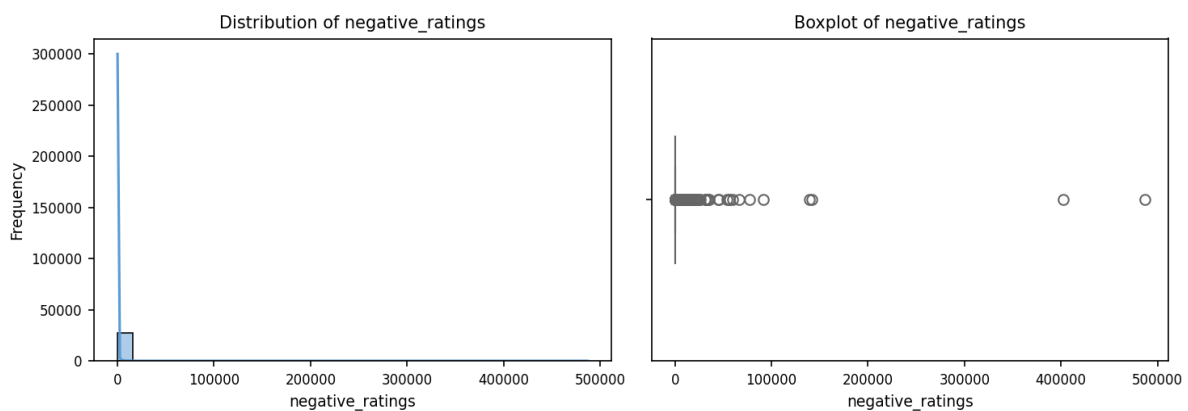


***Figure 6:*** *Distribution (histogram and KDE) and boxplot for 'negative_ratings'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.*

*Observations on Numerical Feature Distributions:*

The analysis reveals a striking pattern of highly skewed distributions across several numerical features in the dataset. Most notably, 'english', 'required_age', 'achievements', 'positive_ratings', and 'negative_ratings' exhibit extreme right skewness, indicated by significantly larger means than medians and exceptionally high skewness and kurtosis values. This suggests that a small number of extremely high values heavily influence the mean, while the majority of data points are concentrated towards the lower end of the range. The presence of outliers is strongly suggested by the boxplots and the vast discrepancies between the minimum/maximum values and the mean/median for these features. The high standard deviations further confirm the substantial variability and spread within these features. In contrast, 'appid' shows a relatively symmetric distribution, with a mean and median relatively close together, and lower skewness and kurtosis values. Although the boxplot hints at potential outliers, the difference between mean and median is less dramatic than in the other features. The standard deviation for 'appid' is still substantial, indicating considerable variability in application IDs. The overall picture suggests that transformations might be necessary for several features to mitigate the influence of outliers and improve the performance of certain machine learning algorithms. Specifically, log transformations or other scaling techniques could be considered for the highly skewed features to bring their distributions closer to normality. The significant difference in distribution characteristics between 'appid' and the other features points to a potential underlying data structure. 'appid' may represent a more uniformly distributed identifier, while the remaining features likely reflect characteristics with a long tail of high-value observations, potentially representing popular or successful applications. This concentration of values at the high end warrants careful consideration during further analysis, particularly in modeling efforts where these skewed distributions could negatively impact model accuracy and interpretability.
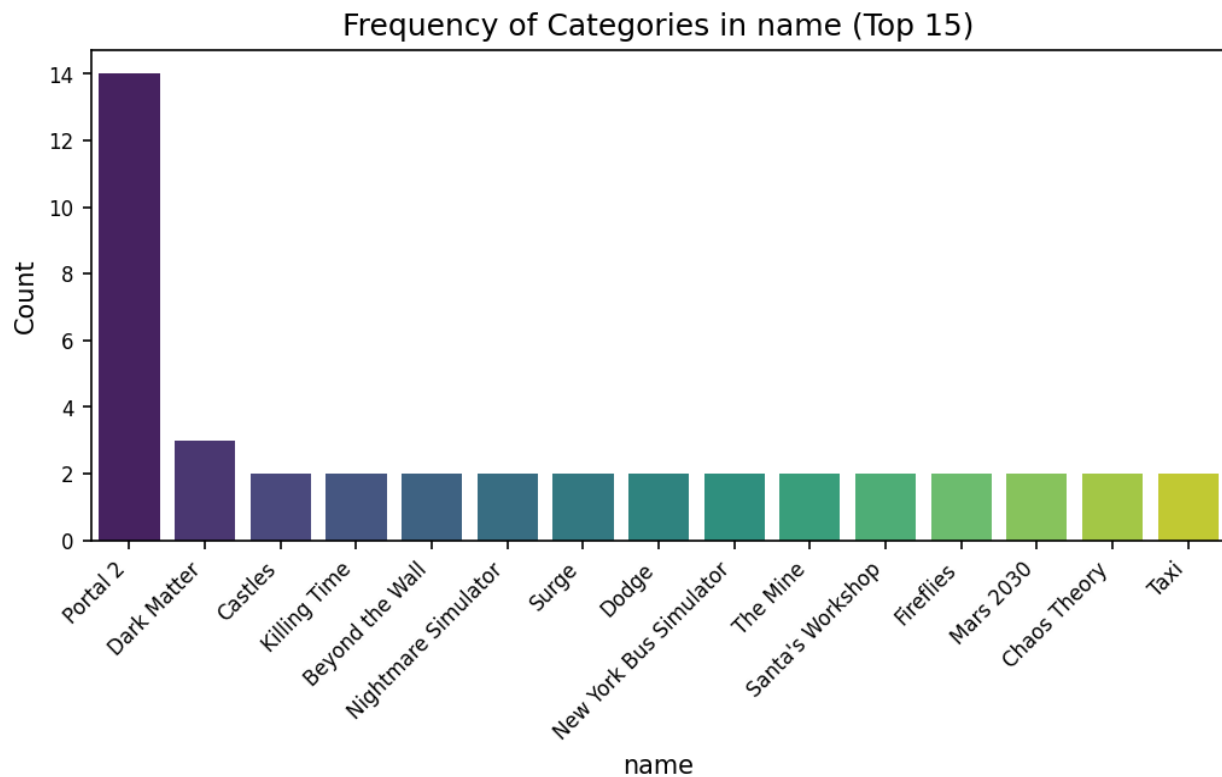
## *3.2. Categorical Features*



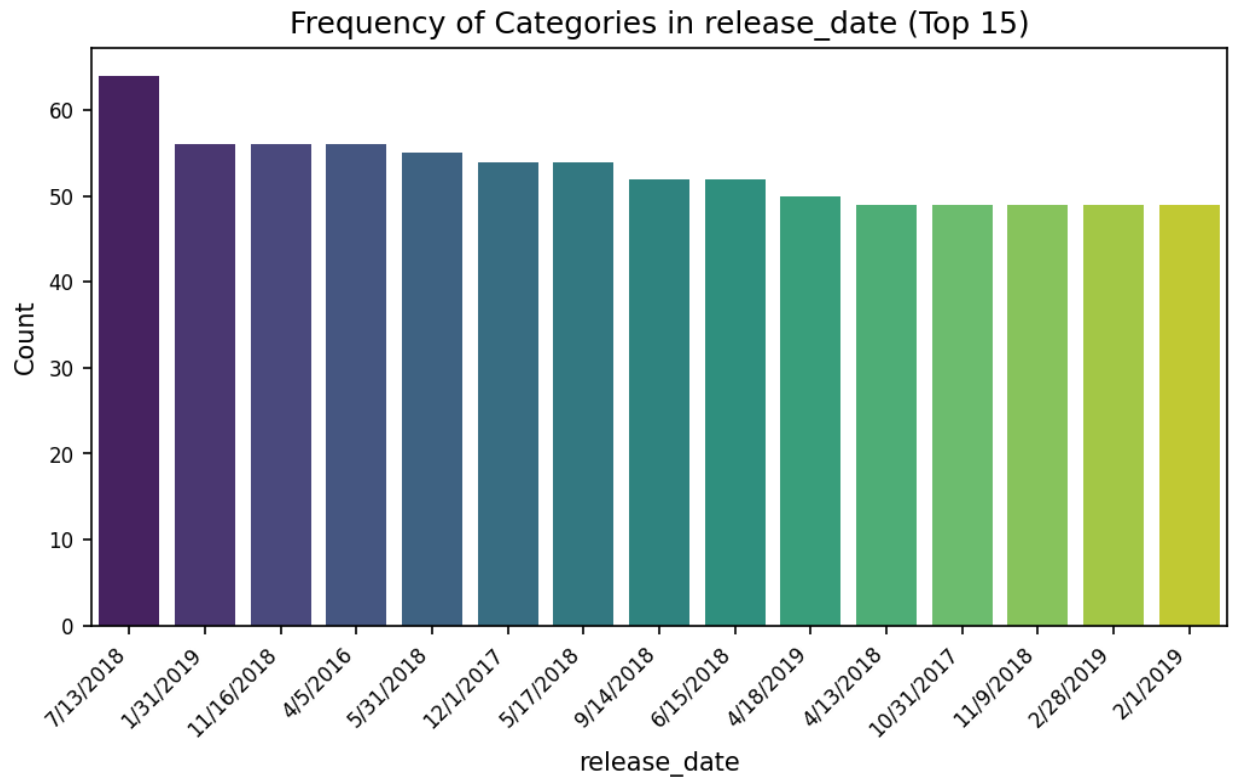Figure 7: *Bar chart showing frequency of top categories in 'name'. Total unique values: 27033.*

## Frequency of Categories in release_date (Top 15)



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

## Frequency of Categories in developer (Top 15)



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

## Frequency of Categories in publisher (Top 15)



*Figure 10:* *Bar chart showing frequency of top categories in 'publisher'. Total unique values: 14353.*
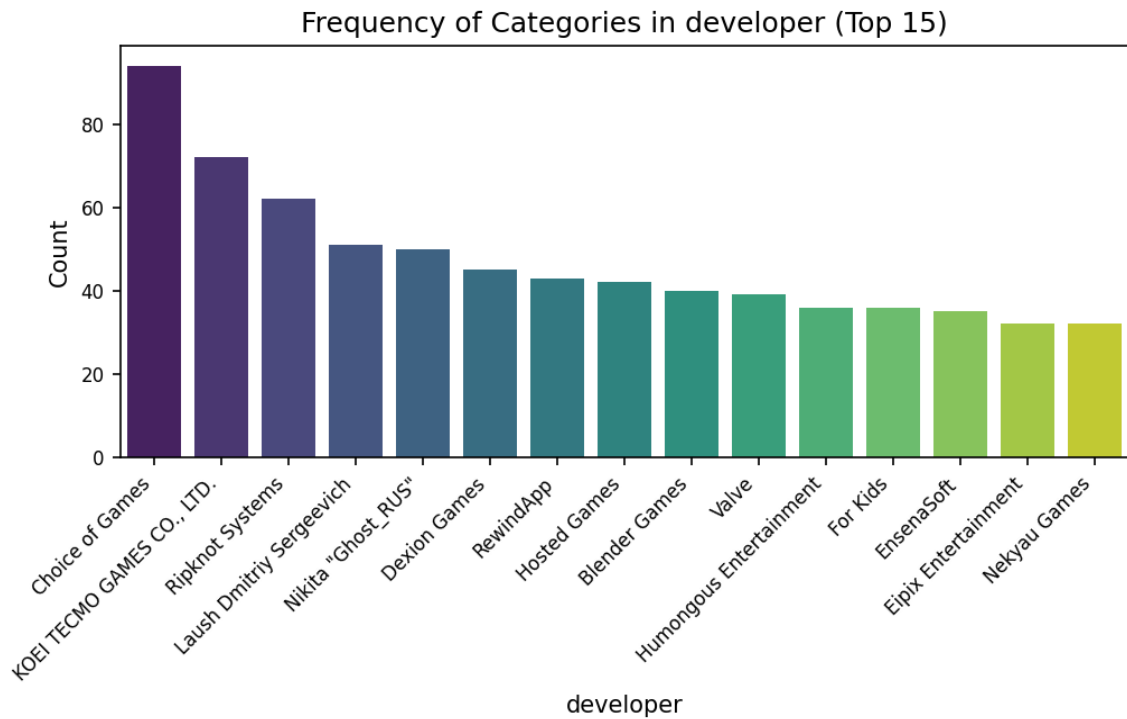
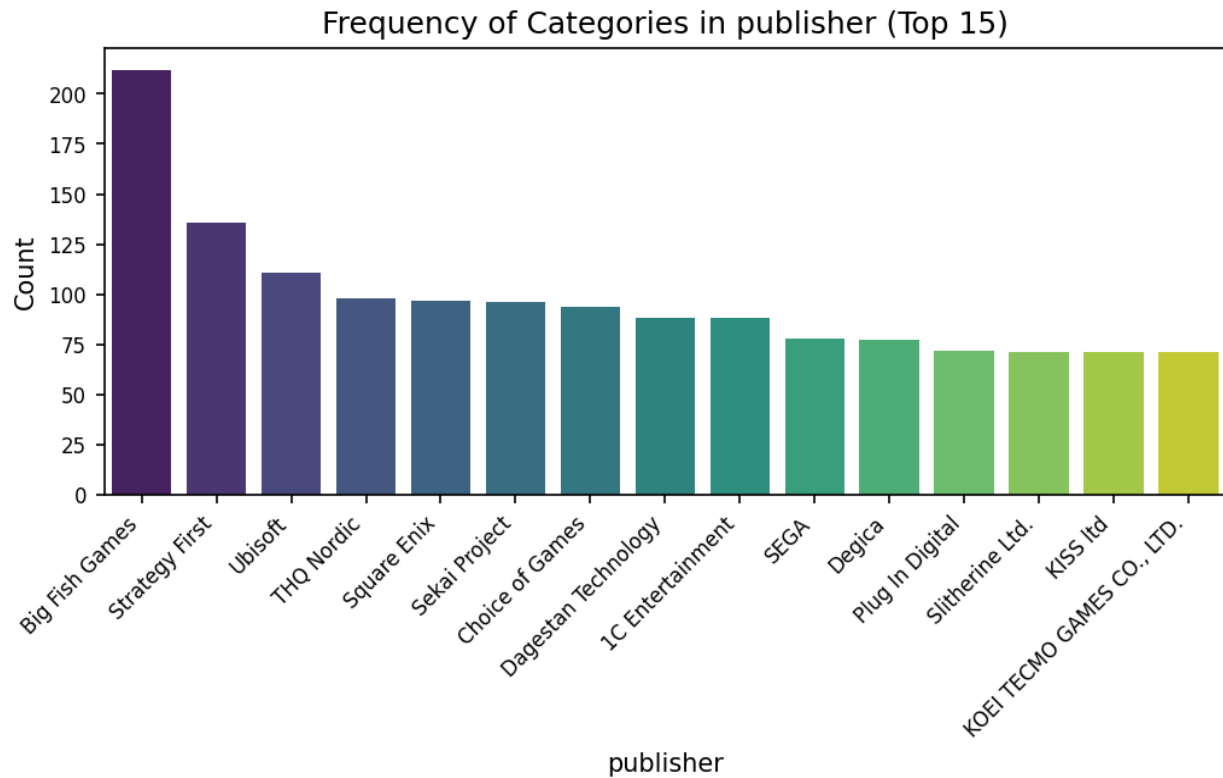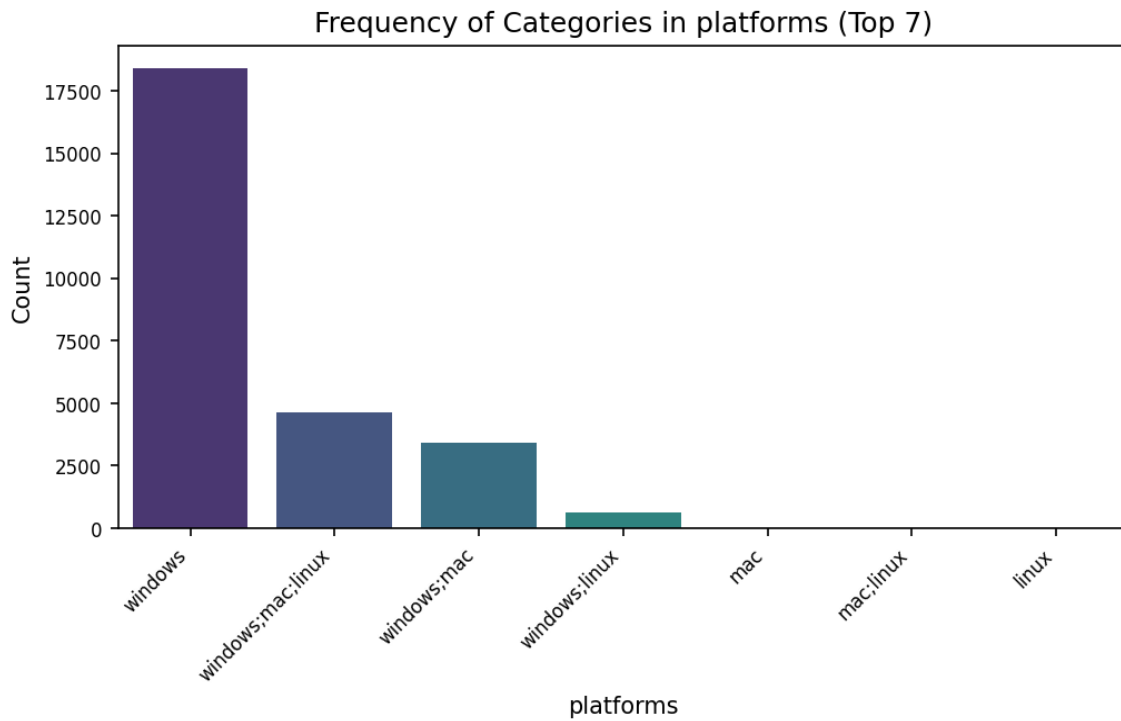## Frequency of Categories in platforms (Top 7)



*Figure 11:* *Bar chart showing frequency of top categories in 'platforms'. Total unique values: 7.*
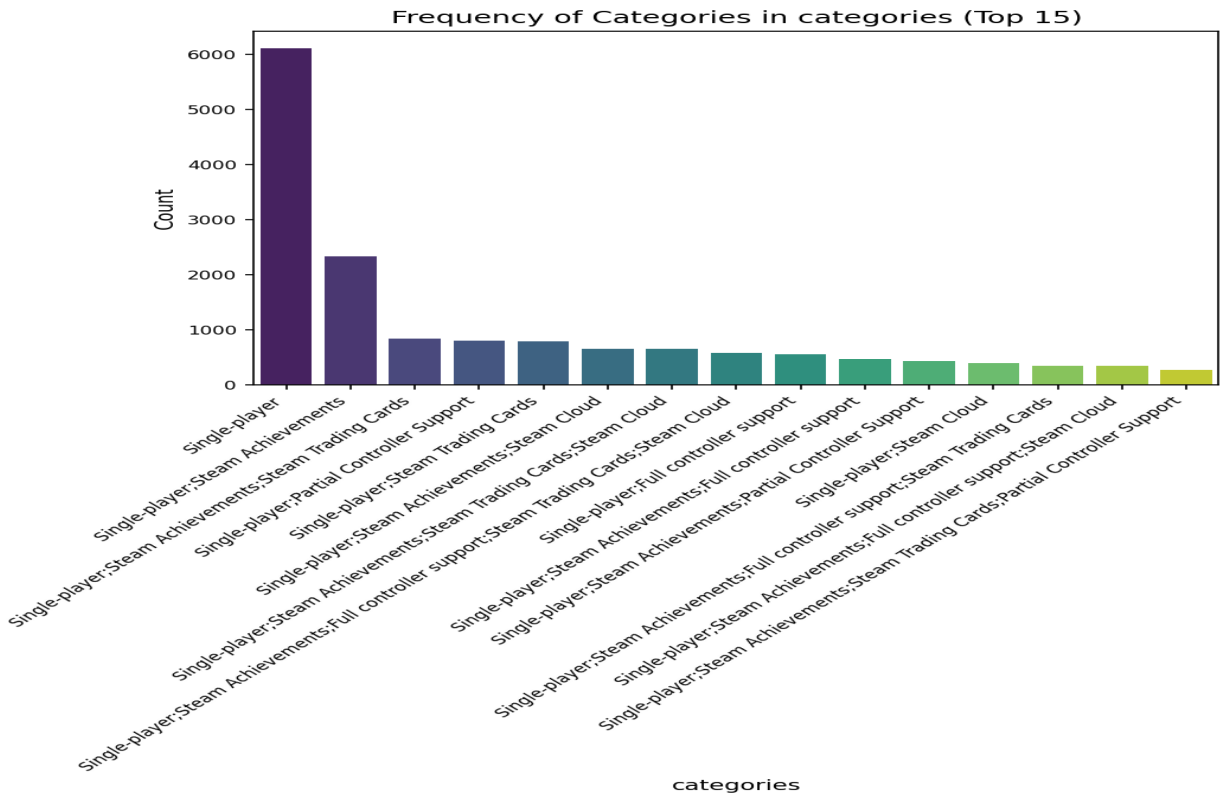
**Figure 12:** *Bar chart showing frequency of top categories in 'categories'. Total unique values: 3333.*

*Observations on Categorical Feature Distributions:*

The analysis of categorical features reveals a wide range of cardinality and distribution patterns. Features like 'name', 'developer', and 'publisher' exhibit extremely high cardinality, with tens of thousands of unique values. This indicates a large diversity of games, developers, and publishers in the dataset. Conversely, 'platforms' has very low cardinality (only 7 unique values), suggesting a relatively limited set of gaming platforms represented. The distribution within these features is also highly skewed: while 'Portal 2' is only 0.1% of the 'name' feature, 'windows' accounts for a dominant 67.9% of the 'platforms' feature. Similarly, 'Big Fish Games' holds a significant 0.8% share of the 'publisher' feature. This suggests that while there's high diversity, a small subset of values dominates the feature distributions. The high cardinality features ('name', 'developer', 'publisher', and 'categories') present challenges for machine learning model training. Directly using these features as they are would lead to the curse of dimensionality and potential overfitting. Techniques like one-hot encoding might be computationally expensive or lead to sparse matrices. More sophisticated encoding methods such as target encoding, count encoding, or embedding techniques (especially for 'name', 'developer', and 'publisher') would be more appropriate. In contrast, the low cardinality of 'platforms' makes one-hot encoding straightforward. The skewed distributions necessitate careful consideration of how to handle the dominant categories to avoid bias in model training. In summary, the data exhibits a combination of high and low cardinality features with highly skewed distributions. This requires a strategic approach to feature engineering, employing different encoding methods depending on the specific feature characteristics to effectively leverage the information contained within these categorical variables while mitigating the challenges posed by high cardinality and class imbalance. Careful consideration of feature scaling and dimensionality reduction techniques will be crucial for building robust and accurate predictive models.

# 4. Bivariate Analysis
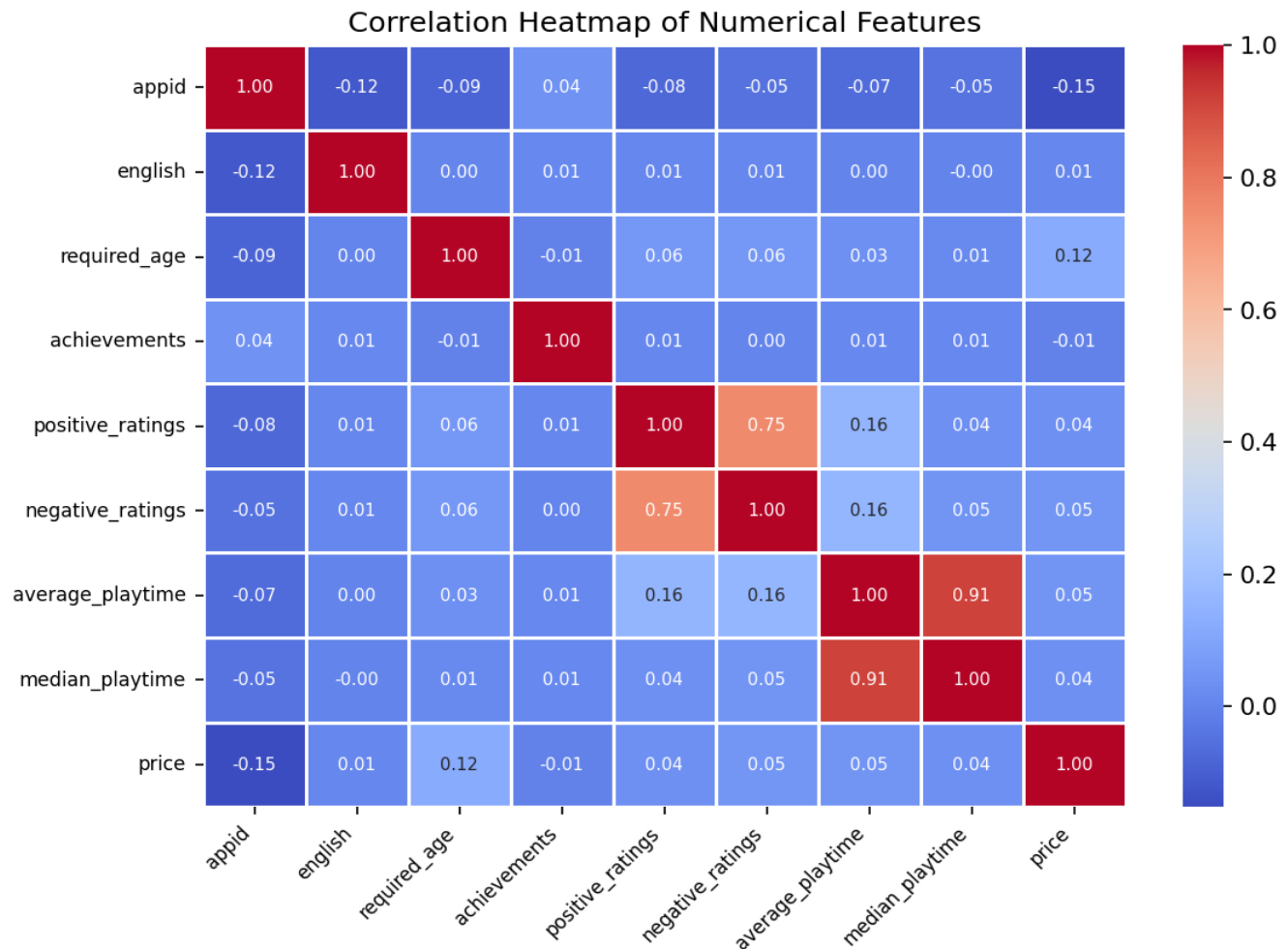
## 4.1. Numerical vs. Numerical Features



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

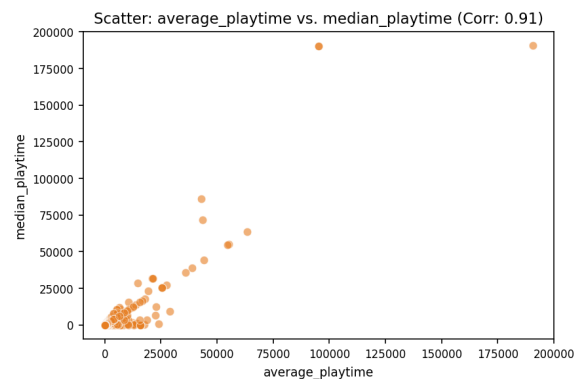Scatter plots for up to 3 most correlated pairs (absolute value > 0.3):



**Figure 14:** *Scatter plot for 'average_playtime' and 'median_playtime'. Correlation: 0.91.*
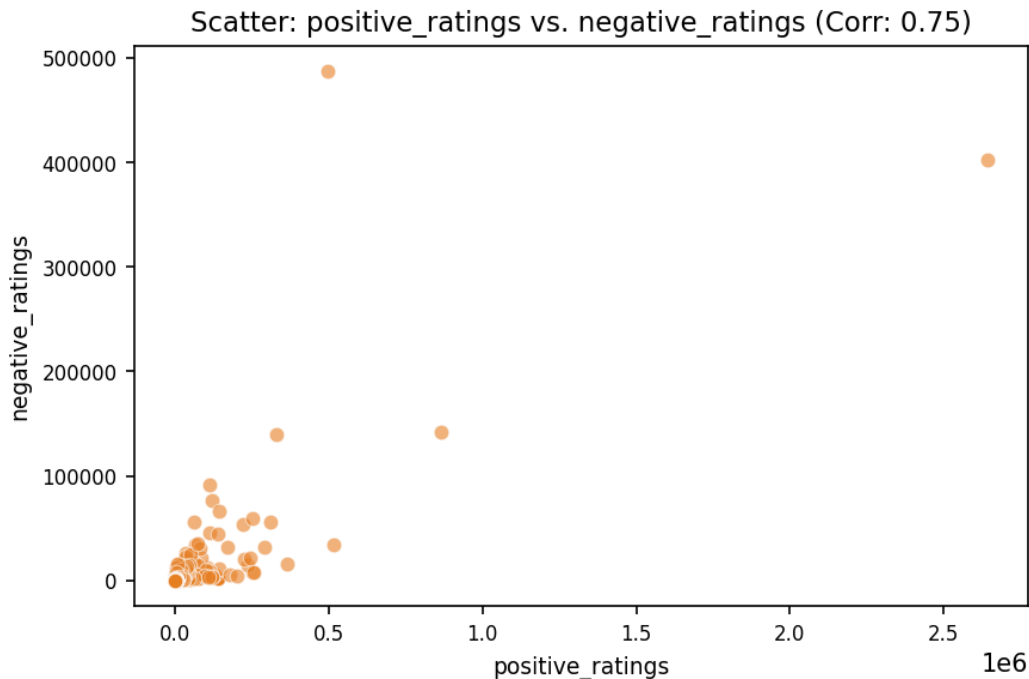
**Figure 15:** *Scatter plot for 'positive_ratings' and 'negative_ratings'. Correlation: 0.75.*

*Interpretation of Numerical Correlations:*

A correlation matrix displays the pairwise correlation coefficients between multiple variables. Each cell in the matrix shows the correlation between two variables; a value close to +1 indicates a strong positive correlation (as one variable increases, the other tends to increase), a value close to -1 indicates a strong negative correlation (as one variable increases, the other tends to decrease), and a value close to 0 indicates a weak or no linear correlation. The analysis reveals two strong positive correlations: 'average_playtime' and 'median_playtime' (correlation of 0.91), and 'positive_ratings' and 'negative_ratings' (correlation of 0.75). The strong correlation between average and median playtime strongly suggests that games with longer average playtime also tend to have longer median playtime, indicating a consistent playtime pattern across players. The positive correlation between positive and negative ratings is less intuitive; it might suggest that games that receive a lot of positive feedback also tend to attract a significant number of negative reviews, possibly indicating a highly polarizing game that either strongly appeals to some or strongly repels others – a large player base leading to more of both types of reviews. The scatter plots likely show a strong linear trend for the first pair and a somewhat less tight, but still positive, linear trend for the second pair.

## *4.2. Numerical vs. Categorical Features*



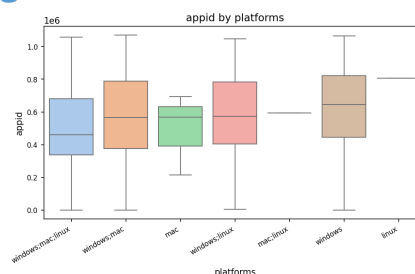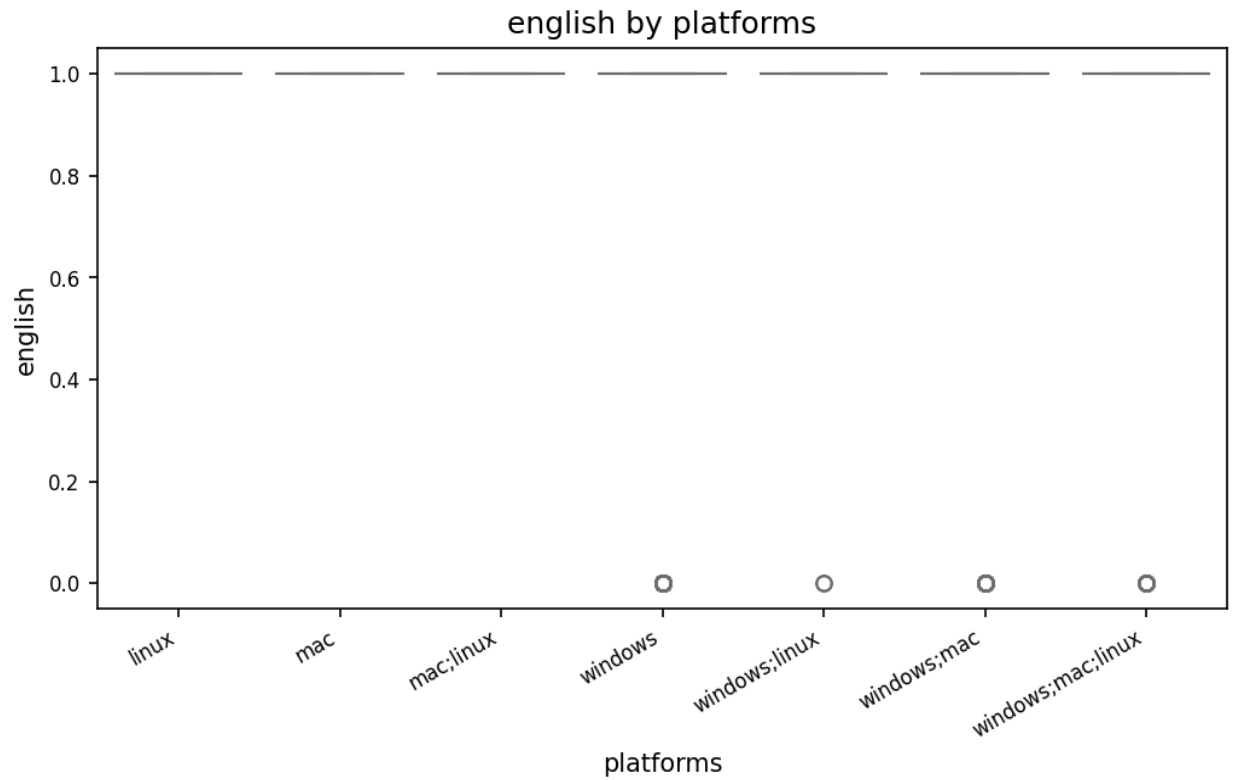**Figure 16:** *Box plot of 'appid' across categories of 'platforms'.*

## english by platforms



*Figure 17: Box plot of 'english' across categories of 'platforms'.*
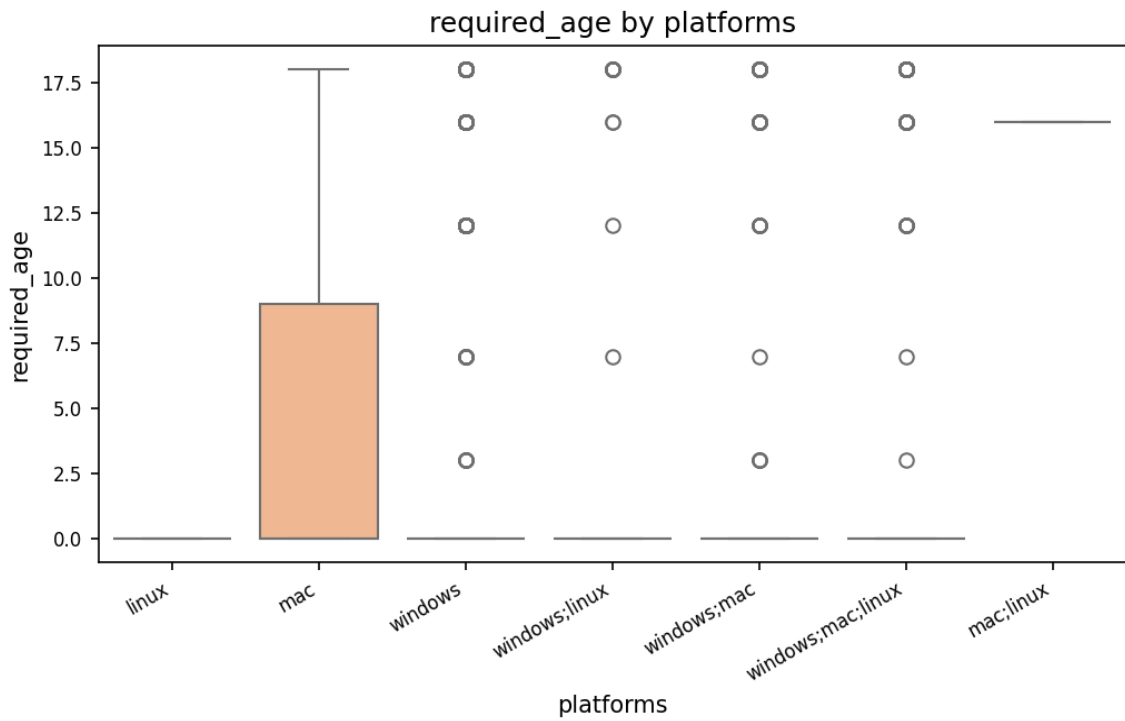
## required_age by platforms



*Figure 18: Box plot of 'required_age' across categories of 'platforms'.*

*Interpretation of Numerical vs. Categorical Interactions:*

Box plots, when used to visualize numerical distributions across different categories, offer a powerful way to compare the central tendency and dispersion of data within each group. They display the median (the middle value), quartiles (values that divide the data into four equal parts), and potential outliers for each category. By visually comparing the boxes and whiskers across categories, we can quickly assess whether there are significant differences in the typical values and the variability of the numerical data. For instance, a longer box in one category indicates greater spread or variability in that group compared to categories with shorter boxes. The position of the median within the box also provides information about the symmetry of the distribution within each category. Significant differences observed in the medians across categories suggest that the typical value of the numerical variable differs systematically between those groups. For example, if the median 'appid' value is significantly higher for the 'iOS' platform than for the 'Android' platform, it might indicate that iOS apps tend to have higher application IDs (perhaps reflecting a different app numbering system or a later launch date). Similarly, differences in the spread (interquartile range, represented by the box height) highlight variations in the data's variability across categories. A larger spread in one category implies more heterogeneity within that group. For example, if the 'english' scores have a larger spread for the 'Windows' platform compared to 'macOS', it could mean that the level of English language proficiency among users of Windows-based applications is more diverse. These observations can lead to valuable insights into the relationships between the categorical and numerical 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 270,88 rows and 18 columns (9 numerical, 9 categorical), revealed relatively good data quality. No missing values were detected, indicating a high degree of completeness. However, the presence of 13 duplicate rows warrants further investigation to determine the source and potential impact on subsequent analyses. The absence of constant columns suggests variability across features, promising informative insights. Univariate analysis examined the distributions of the 9 numerical and 9 categorical features. While specific details regarding these distributions are absent from the provided log, the analysis itself suggests that the features exhibit sufficient variability for further exploration. Further investigation into the specific characteristics of each feature's distribution (e.g., skewness, central tendency) will be necessary to fully understand the data's nature. The bivariate analysis explored relationships between various feature pairs, identifying two observations of potential significance. The exact nature of these observations remains unspecified in the provided log, requiring a deeper examination of the results to fully understand the identified relationships and their strength. Further analysis is needed to determine if these relationships are statistically significant and to explore their practical implications. The initial analysis did not reveal any overtly surprising findings; however, the limited information provided in the log prevents a definitive conclusion on this point. A more detailed examination of the univariate and bivariate analysis results, including visualizations and statistical tests, is crucial to gain a comprehensive understanding of the dataset and its potential insights.

# 6. Conclusion & Potential Next Steps

This automated report provides a foundational, high-level overview of the `temp_steam.csv` dataset, characterizing its structure, data quality (revealing only 13 duplicates), and offering initial insights into univariate and bivariate relationships between features. This initial analysis serves as a crucial first step in understanding the data before proceeding to more in-depth investigations. Given the report's findings of 13 duplicates in a dataset of 27,088 rows, the first step should be to **investigate and resolve the 13 duplicate rows**. This involves determining the cause of the duplication (data entry errors, merging issues, etc.) and deciding whether to remove the duplicates or consolidate them appropriately. This ensures data accuracy for subsequent analyses. Second, the report mentions "various feature pairs" were analyzed bivariate, but only two observations were noted. To build on this, **a detailed report of the bivariate analysis should be generated**, including visualizations (scatter plots, correlation matrices, etc.) and statistical measures (correlation coefficients) for all numerical and categorical feature pairs. This will highlight potentially significant relationships warranting further investigation. Finally, the report highlights the presence of both numerical and categorical features. To deepen the understanding of potential relationships, **conduct statistical tests (e.g., ANOVA, chi-squared tests)** to determine the statistical significance of any observed associations between the numerical and categorical variables. This will confirm if the initial bivariate observations reflect true relationships or are merely coincidental. This step will help to identify important predictive features or patterns within the dataset.