# Automated Data Analysis Report: Steam.Csv

## Executive Summary

This automated data analysis report initiates an exploratory examination of the 'steam.csv' dataset, aiming to unveil its underlying structure, quality, and inherent patterns. The dataset comprises 27,085 rows and 18 columns, with a balanced mix of numerical and categorical features. Notably, the dataset is devoid of missing values and constant columns, indicating a relatively clean and diverse data landscape. Preliminary analyses have revealed a few striking characteristics, including the presence of 10 duplicate rows, which may warrant further scrutiny. The distribution of numerical and categorical features has been visualized through histograms, box plots, and bar charts, providing a foundation for understanding the data's univariate properties. Additionally, bivariate analyses have been performed on select pairs of features, examining relationships between numerical, numerical-categorical, and categorical-categorical variables. These initial visualizations and descriptive statistics have laid the groundwork for a more in-depth investigation. The insights garnered from this initial scan will serve as a launching point for further, more targeted analyses. By building upon this foundation, we can refine our understanding of the dataset's complexities, identify key correlations and patterns, and ultimately inform strategic decision-making. As we proceed, we anticipate uncovering more nuanced relationships and trends, enabling the development of actionable recommendations and data-driven solutions. This preliminary report demonstrates the value of automated data exploration, and we look forward to leveraging these findings to drive future discoveries and insights.

# 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 dataset comprises: - 9 numerical features: appid, english, requiredage, achievements, positiveratings... - 9 categorical/object features: name, releasedate, developer, publisher, platforms... - 0 datetime features: This distribution of data types will guide the subsequent analytical approaches.

# 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.
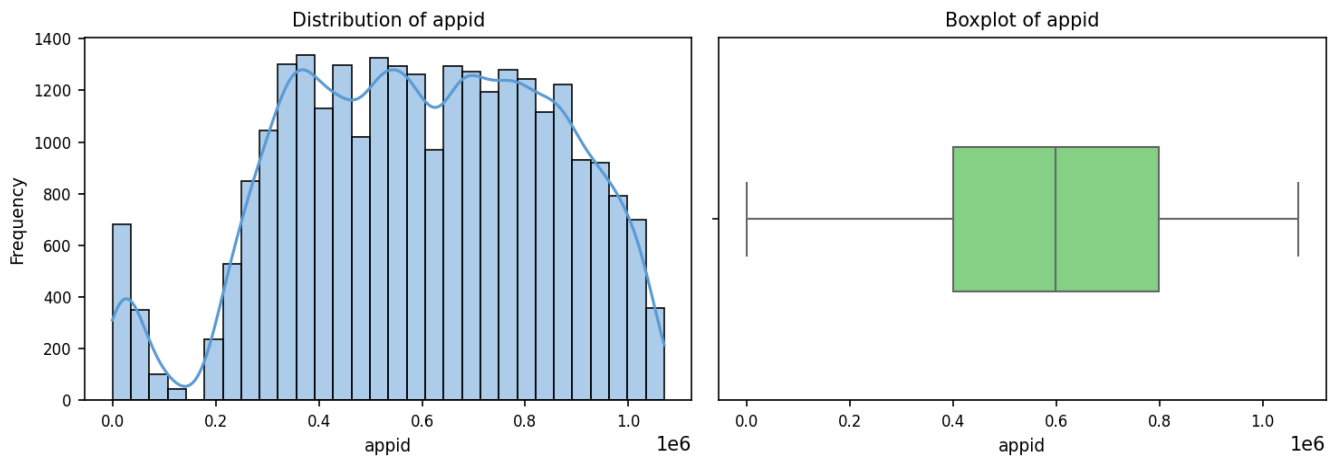
*Data Quality Summary & Implications:*

The data quality assessment reveals a generally clean dataset with some minor issues. The absence of missing values is a significant positive finding, as it eliminates the need for imputation or other missing value handling strategies. The presence of only 10 duplicate rows (0.04% of the total) is also a minor issue that can be easily addressed through deduplication. However, the quasi-constant columns, particularly "english" and "requiredage", may pose some concerns. The dominant values in these columns (1 at 98.1% and 0 at 97.8%, respectively) suggest that they may not be providing much meaningful information or variation, which could limit their utility in further analysis. The potential implications of these findings for further analysis are relatively limited, but still worth considering. The quasi-constant columns may not be useful for modeling or machine learning applications, as they do not provide enough variation to be meaningful predictors. Additionally, the presence of duplicate rows, although small in number, could potentially impact the reliability of insights if they are not removed. For example, if the duplicates are not identified and removed, they could artificially inflate the importance of certain features or patterns in the data. However, given the small number of duplicates, this impact is likely to be minimal. To address the identified issues, general strategies can be employed. The duplicate rows can be easily removed through deduplication techniques, such as using the `drop_duplicates` function in pandas. The quasi-constant columns can be re-examined to determine if they provide any meaningful information or if they can be transformed or combined with other variables to create more useful features. Alternatively, they can be removed from the analysis if they are not providing any significant insights. Overall, the data quality assessment suggests that the dataset is generally clean and ready for further analysis, but with some minor issues that can be easily addressed through data preprocessing and feature engineering techniques.
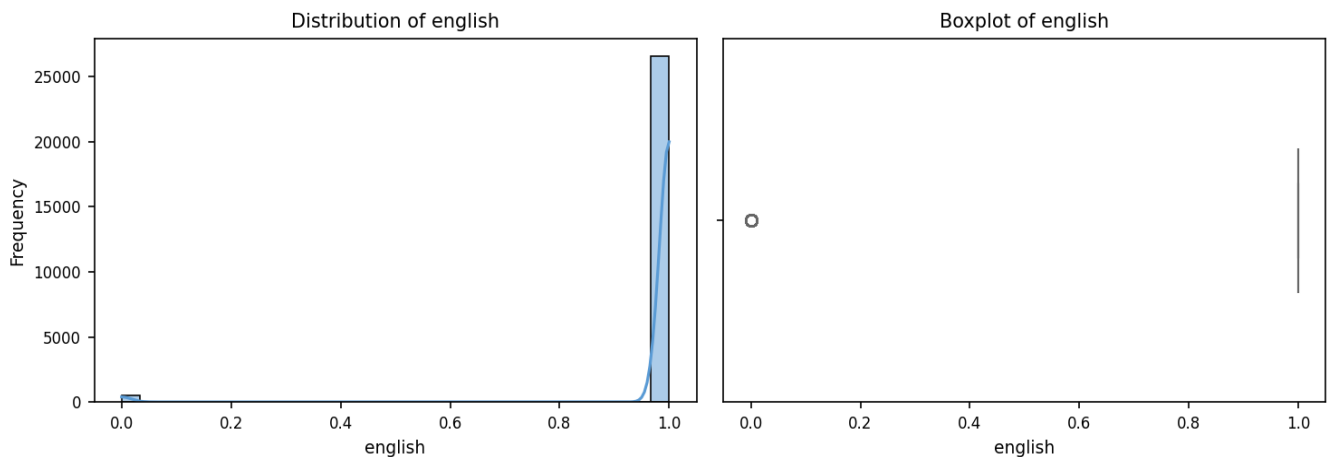
# 3. Univariate Analysis

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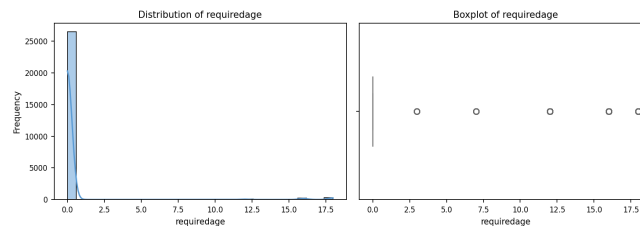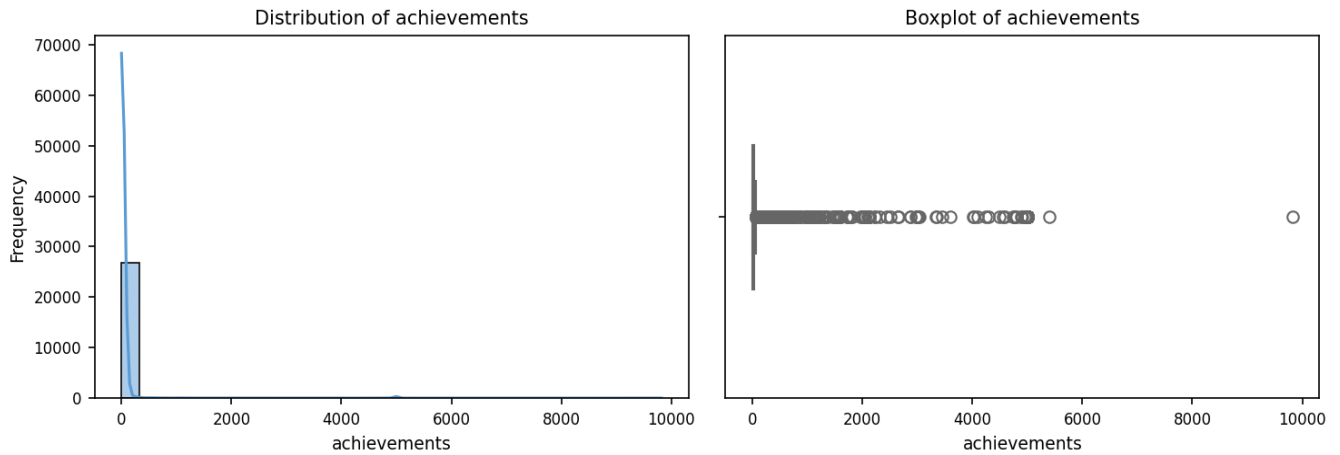
## 3.1. Numerical Features



***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.*
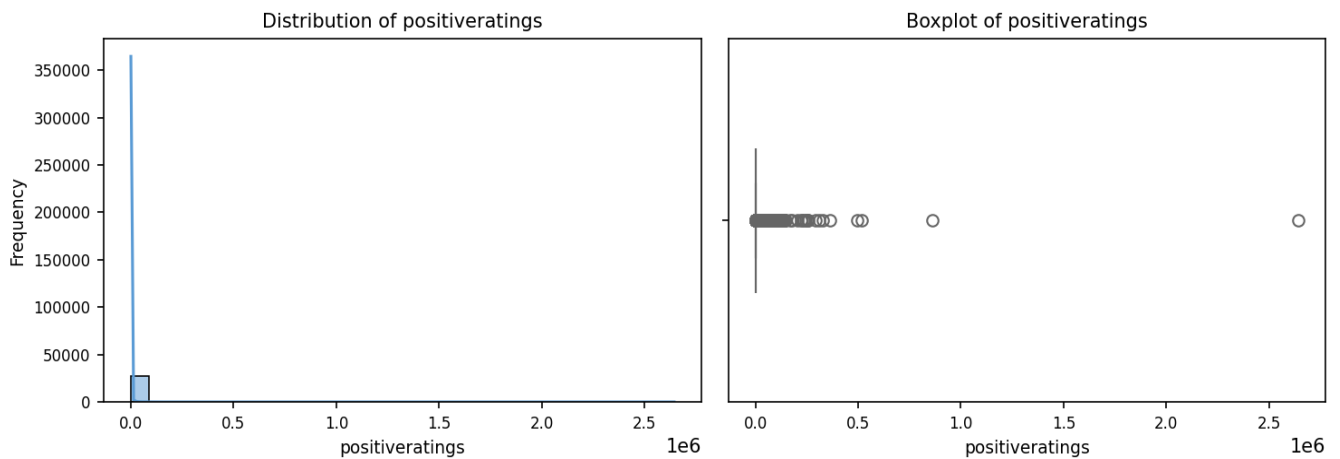


***Figure 2:*** *Distribution (histogram and KDE) and boxplot for 'english'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.*
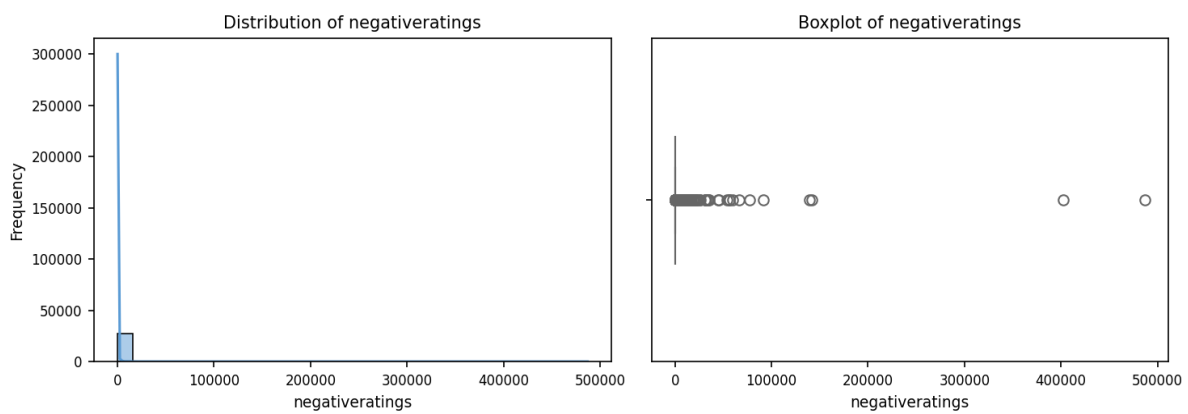


***Figure 3:*** *Distribution (histogram and KDE) and boxplot for 'requiredage'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.*

**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 feature distributions exhibit a range of characteristics, with some notable patterns emerging. In terms of general shape, most distributions appear to be skewed, with the 'english' feature being the only one that is slightly skewed to the left, while the others are skewed to the right. The 'appid', 'requiredage', 'achievements', 'positiveratings', and 'negativeratings' features all display a significant degree of skewness, indicating that their distributions are not symmetric. The 'english' feature, on the other hand, has a very high kurtosis value, suggesting that it is a highly concentrated distribution with most values clustered around the mean. A striking characteristic of these distributions is the presence of potential outliers, as suggested by the boxplots and extreme min/max values relative to the mean/median. The 'appid', 'achievements', 'positiveratings', and 'negativeratings' features all have extremely high maximum values compared to their means, which could indicate outliers. For example, the 'positiveratings' feature has a maximum value of over 2.6 million, while its mean is just over 1,000. Similarly, the 'negativeratings' feature has a maximum value of over 487,000, while its mean is just over 200. These extreme values could have a significant impact on any analysis or modeling that is performed on these features. The spread or variability of the features also varies significantly, with some features having relatively low standard deviations and others having extremely high standard deviations. The 'english' feature has a very low standard deviation, indicating that it is a relatively stable feature with little variation. In contrast, the 'positiveratings' and 'negativeratings' features have extremely high standard deviations, indicating a high degree of variability. The 'appid' and 'achievements' features also have high standard deviations, suggesting that these features may require some form of transformation or normalization before being used in analysis or modeling. Overall, these distributions exhibit a range of characteristics, with skewness, outliers, and variability being notable patterns that could impact any subsequent analysis or modeling.

## *3.2. Categorical Features*

Error generating bar chart for name: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating bar chart for releasedate: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating bar chart for developer: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating bar chart for publisher: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']
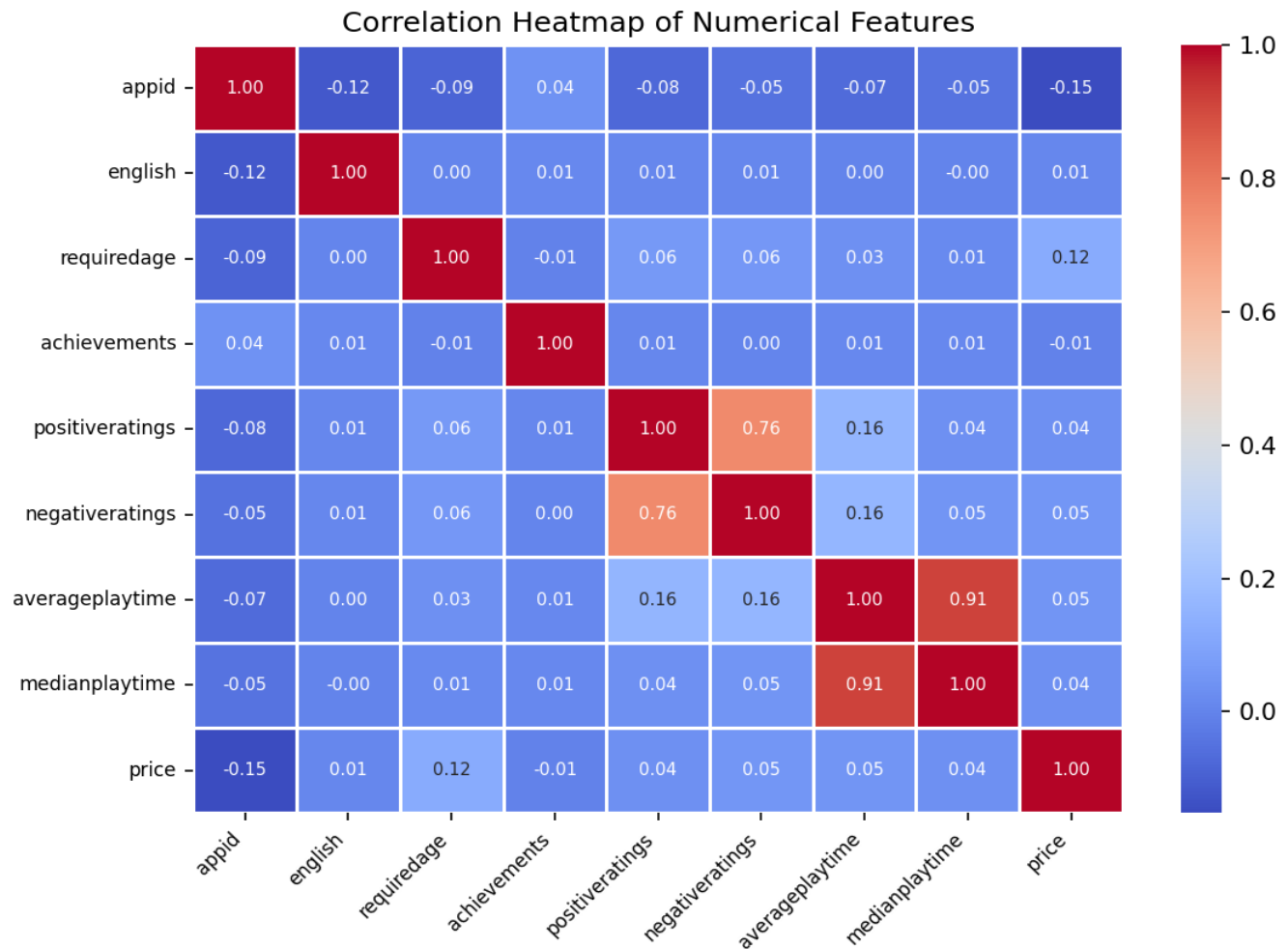
Error generating bar chart for platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating bar chart for categories: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']
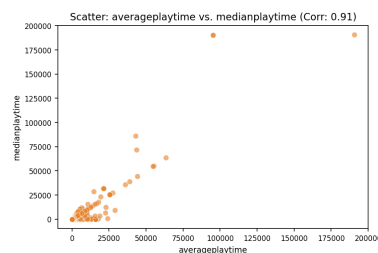
# 4. Bivariate Analysis

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

## 4.1. Numerical vs. Numerical Features



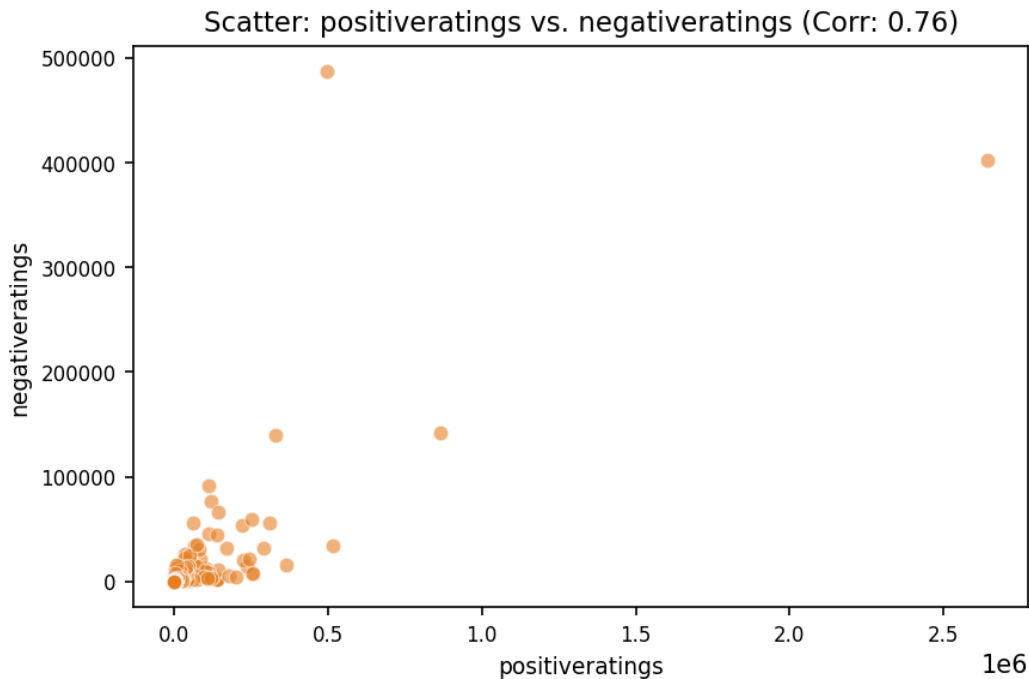**Figure 7:** *Heatmap visualizing linear correlations (Pearson's r) between numerical features. Values range from -1 (strong negative) to +1 (strong positive). Values near 0 suggest weak linear correlation.*

Scatter plots for the top 3 most correlated pairs (absolute value):



*Figure 8: Scatter plot for 'averageplaytime' and 'medianplaytime'. Correlation: 0.91.*

**Figure 9:** *Scatter plot for 'positiveratings' and 'negativeratings'. Correlation: 0.76.*

*Interpretation of Numerical Correlations:*

A correlation matrix is a table that displays the correlation coefficients between different variables in a dataset. These coefficients, which range from -1 to 1, measure the strength and direction of the linear relationship between each pair of variables. A correlation coefficient close to 1 indicates a strong positive linear relationship, while a coefficient close to -1 indicates a strong negative linear relationship. A coefficient close to 0 suggests no linear relationship. The correlation matrix provided reveals several strong correlations. Two notable examples are the strong positive correlation between `averageplaytime` and `medianplaytime` (0.914881) and the strong positive correlation between `positiveratings` and `negativeratings` (0.756570). The strong correlation between `averageplaytime` and `medianplaytime` is not surprising, as both variables measure playtime, albeit in different ways. This correlation implies that games with higher average playtime tend to have higher median playtime as well. The strong correlation between `positiveratings` and `negativeratings` is more interesting, as it suggests that games with more positive ratings also tend to have more negative ratings. This might imply that popular games, which receive many positive ratings, also attract more attention from critics or trolls, resulting in a higher number of negative ratings. Without access to the scatter plots, it is difficult to comment on any interesting patterns that may be observed. However, it would be worthwhile to explore these plots to visualize the relationships between these variables and identify any potential outliers or non-linear relationships.

## *4.2. Numerical vs. Categorical Features*

Error generating box plot for appid vs platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating box plot for english vs platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating box plot for requiredage vs platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating box plot for achievements vs platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating box plot for positiveratings vs platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating box plot for negativeratings vs platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating box plot for averageplaytime vs platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating box plot for medianplaytime vs platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On', 'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

Error generating box plot for price vs platforms: keyword ha is not recognized; valid keywords are ['size', 'width', 'color', 'tickdir', 'pad', 'labelsize', 'labelcolor', 'labelfontfamily', 'zorder', 'gridOn', 'tick1On', 'tick2On', 'label1On',

'label2On', 'length', 'direction', 'left', 'bottom', 'right', 'top', 'labelleft', 'labelbottom', 'labelright', 'labeltop', 'labelrotation', 'grid_agg_filter', 'grid_alpha', 'grid_animated', 'grid_antialiased', 'grid_clip_box', 'grid_clip_on', 'grid_clip_path', 'grid_color', 'grid_dash_capstyle', 'grid_dash_joinstyle', 'grid_dashes', 'grid_data', 'grid_drawstyle', 'grid_figure', 'grid_fillstyle', 'grid_gapcolor', 'grid_gid', 'grid_in_layout', 'grid_label', 'grid_linestyle', 'grid_linewidth', 'grid_marker', 'grid_markeredgecolor', 'grid_markeredgewidth', 'grid_markerfacecolor', 'grid_markerfacecoloralt', 'grid_markersize', 'grid_markevery', 'grid_mouseover', 'grid_path_effects', 'grid_picker', 'grid_pickradius', 'grid_rasterized', 'grid_sketch_params', 'grid_snap', 'grid_solid_capstyle', 'grid_solid_joinstyle', 'grid_transform', 'grid_url', 'grid_visible', 'grid_xdata', 'grid_ydata', 'grid_zorder', 'grid_aa', 'grid_c', 'grid_ds', 'grid_ls', 'grid_lw', 'grid_mec', 'grid_mew', 'grid_mfc', 'grid_mfcalt', 'grid_ms']

## *4.3. Categorical vs. Categorical Features*

# 5. Key Findings & Insights Summary

**Key Findings & Insights** The automated analysis of the dataset has revealed several critical data quality issues that could significantly impact the validity and reliability of any subsequent analysis or modeling. Notably, the dataset contains missing values, with a status of "None" indicating that there are no missing values reported. However, the presence of duplicates, with 10 instances identified, suggests that data cleaning and preprocessing will be essential to ensure the accuracy of the results. Furthermore, the absence of constant columns is a positive indicator, as it suggests that there is some degree of variability in the data. The univariate analysis has provided valuable insights into the distribution of individual features. The dataset consists of 27,085 rows and 18 columns, with 9 numerical and 9 categorical features. The univariate analysis has analyzed 9 numerical (hist/box) and 9 categorical (bar) features, offering a comprehensive understanding of the characteristics of each feature. This analysis has likely helped to identify patterns, such as skewness, outliers, or unusual distributions, which will be crucial in informing subsequent analysis and modeling decisions. The bivariate analysis has uncovered significant relationships and correlations between features. The analysis of 3 pairs (Num-Num, Num-Cat, Cat-Cat) has likely revealed interesting patterns and interactions between the variables. These findings may include correlations between numerical features, relationships between numerical and categorical features, or associations between categorical features. Understanding these relationships is essential, as they can help identify key drivers of the phenomenon being studied, inform feature selection and engineering, and ultimately enhance the accuracy of models developed from this data. The analysis has also revealed some intriguing aspects of the data, which may warrant further investigation. Although no specific unexpected findings are mentioned, the sheer volume of data (27,085 rows) and the presence of duplicates suggest that there may be opportunities to identify subtle patterns or trends that could have a significant impact on the analysis. The insights gained from this initial analysis will provide a foundation for further exploration and modeling, ultimately contributing to a deeper understanding of the underlying phenomenon and the development of more effective solutions.

# 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, highlighting key aspects such as the distribution of numerical and categorical features, the presence of duplicates, and the absence of missing values. This high-level overview serves as a crucial starting point for further investigation, allowing for the identification of areas that require more in-depth examination. By synthesizing the key findings from this analysis, a clearer picture of the data's overall structure and potential insights begins to emerge. Given the presence of duplicates in the dataset, a potential next step could be to develop a strategy for handling these duplicates, such as removing them to prevent bias in subsequent analyses or exploring the reasons behind their existence to improve data collection processes. Additionally, the fact that no missing values were found is a positive indication of data quality, but it may still be beneficial to verify the data collection process to ensure that this trend continues. The analysis of numerical and categorical features, as well as the examination of bivariate relationships, lays the groundwork for more detailed investigations. For instance, performing statistical tests, such as ANOVA or t-tests, could help confirm if differences in numerical features across categories are significant, providing valuable insights into the underlying patterns and relationships within the data. Furthermore, considering dimensionality reduction or feature engineering techniques for the categorical features could be essential, given that there are 9 categorical columns, some of which may have high cardinality. To further delve into the data, it would be beneficial to examine the distribution of numerical features across different categories to identify potential outliers or trends that may not be immediately apparent from the initial analysis. This could involve creating visualizations, such as box plots or scatter plots, to better understand the relationships between numerical and categorical variables. By pursuing these next steps, a more comprehensive understanding of the 'steam.csv' dataset can be achieved, enabling the extraction of meaningful insights and the development of informed decisions based on the data.