

# Introduction

This automated data analysis report is based on the 'steam.csv' dataset, which contains 27,085 rows and 18 columns of data related to steam games. The dataset provides a comprehensive overview of various aspects of steam games, including their characteristics, ratings, and user engagement metrics. With a wide range of columns, including appid, name, releasedate, english, developer, publisher, platforms, requiredage, categories, genres, steamspytags, achievements, positiveratings, negativeratings, averageplaytime, medianplaytime, owners, and price, this dataset offers a unique opportunity to explore the steam gaming ecosystem in depth. The primary purpose of this report is to provide an initial overview of the dataset, highlighting its key characteristics, identifying potential data quality issues, and offering some basic insights into the data. By examining the distribution of values, correlation between columns, and summary statistics, we aim to gain a deeper understanding of the dataset and its potential applications. This report will also serve as a foundation for further analysis, highlighting areas that may require additional attention or specialized treatment. Through this report, we hope to uncover trends, patterns, and relationships within the data that can inform decision-making or guide future research. As we delve into the analysis, we will pay close attention to potential data quality issues, such as missing or duplicate values, outliers, and inconsistencies in formatting. By addressing these issues upfront, we can ensure that our analysis is built on a solid foundation and that our findings are reliable and accurate. With a thorough understanding of the dataset and its limitations, we can begin to uncover insights that can inform our understanding of the steam gaming ecosystem, from the characteristics of popular games to the behaviors of their players. This report will provide a comprehensive starting point for exploring the 'steam.csv' dataset and unlocking its full potential.

## Descriptive Statistics

This section provides a statistical summary of the dataset.

Summary for Numerical Features

Feature	count	mean	std	min	25%	50%	75%	max
appid	27085.0	595983.61	251108.67	10.0	400910.0	598990.0	798710.0	1069460.0
english	27085.0	0.98	0.14	0.0	1.0	1.0	1.0	1.0
requiredage	27085.0	0.35	2.41	0.0	0.0	0.0	0.0	18.0
achievements	27085.0	45.25	352.61	0.0	0.0	7.0	23.0	9821.0
positiveratings	27085.0	1051.22	19167.37	0.0	6.0	24.0	126.0	2644404.0
negativeratings	27085.0	211.65	4284.27	0.0	2.0	9.0	42.0	487076.0
averageplaytime	27085.0	150.16	1826.79	0.0	0.0	0.0	0.0	190625.0
medianplaytime	27085.0	146.19	2353.46	0.0	0.0	0.0	0.0	190625.0
price	27085.0	6.08	7.87	0.0	1.69	3.99	7.19	421.99

Summary for Categorical Features

Feature	count	unique	top	freq
name	27085	27033	Portal 2	11
releasedate	27085	2619	7/13/2018	64
developer	27085	17112	Choice of Games	94
publisher	27085	14353	Big Fish Games	212

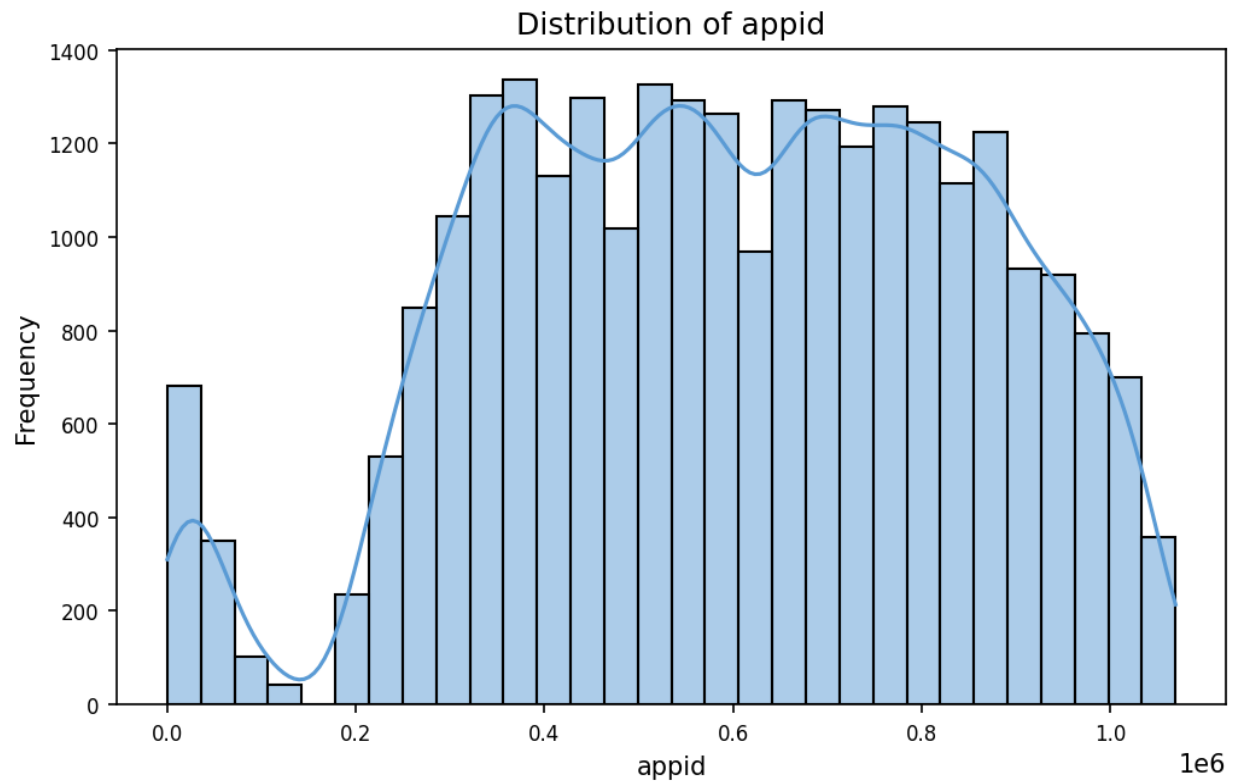
Feature	count	unique	top	freq
platforms	27085	7	windows	18398
categories	27085	3333	Single-player	6110
genres	27085	1552	Action;Indie	1852
steamspytags	27085	6423	Action;Indie;Casual	845
owners	27085	13	0-20000	18596

Interpretation: The provided tables offer a summary of a dataset, detailing the characteristics of both numerical and categorical features. The numerical features table presents an overview of the distribution of various metrics, including counts, means, standard deviations, and percentiles for each feature. This information helps understand the central tendency and variability of each numerical attribute in the dataset. Conversely, the categorical features table lists the frequency and uniqueness of categories within each feature, giving insight into the diversity and commonality of different attributes. Upon reviewing the numerical features, several observations stand out. For instance, the "achievements," "positiveratings," and "negativeratings" features exhibit wide ranges between their minimum and maximum values, with standard deviations significantly larger than their means. This suggests the presence of potential outliers, as the bulk of the data may be concentrated near the lower end of the scale, with a few extremely high values skewing the mean. Additionally, features like "averageplaytime" and "medianplaytime" have minimum, 25%, and 50% percentiles at 0, indicating that many games may have zero or very low playtime, possibly due to low popularity or recent release. The "price" feature seems relatively stable, with a small standard deviation and a range that, while wide, is still contained within a reasonable scale for game prices. The categorical features table reveals a different set of insights. The "name," "developer," "publisher," and "steamspytags" features have a very high number of unique values, suggesting a great diversity in game titles, developers, publishers, and tags. On the other hand, the "platforms" feature is dominated by "windows," with over 68% of the data (18398 out of 27085), indicating a strong focus on Windows as a gaming platform. Similarly, the "owners" feature is mostly concentrated in the "0-20000" category, which might reflect the challenging nature of achieving high ownership numbers for games. These observations highlight the complexity and diversity of the gaming dataset, with both numerical and categorical features showing a range of distributions and concentrations that could be useful for further analysis or modeling.

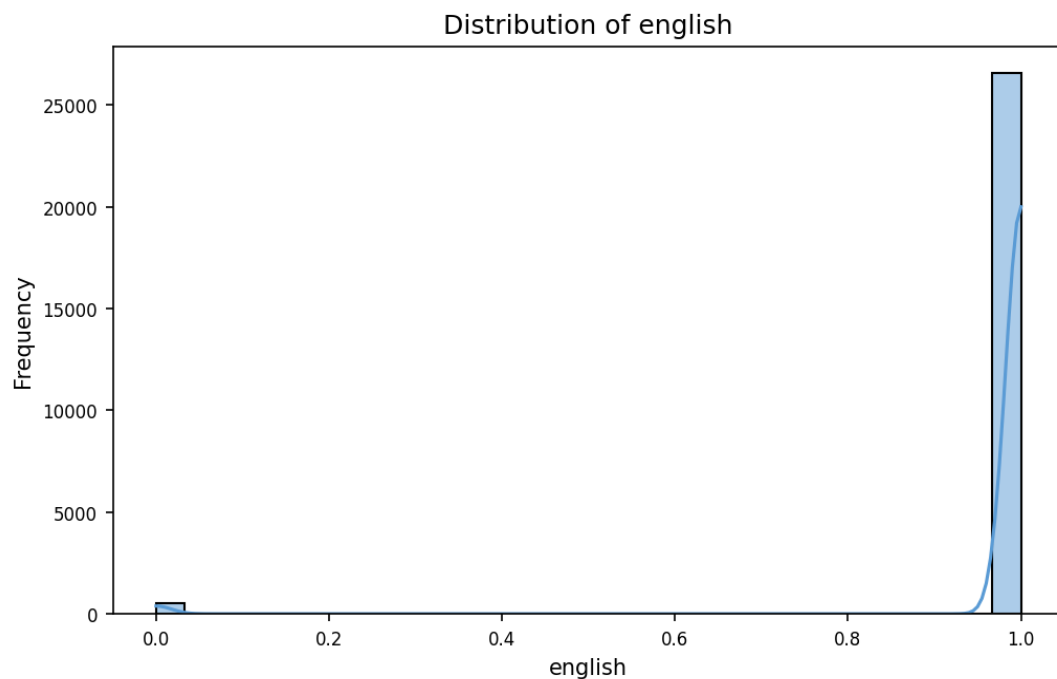
# Data Visualizations

Visual exploration can reveal patterns, trends, and outliers that are not obvious from summary statistics alone.

## Distribution of Numerical Features



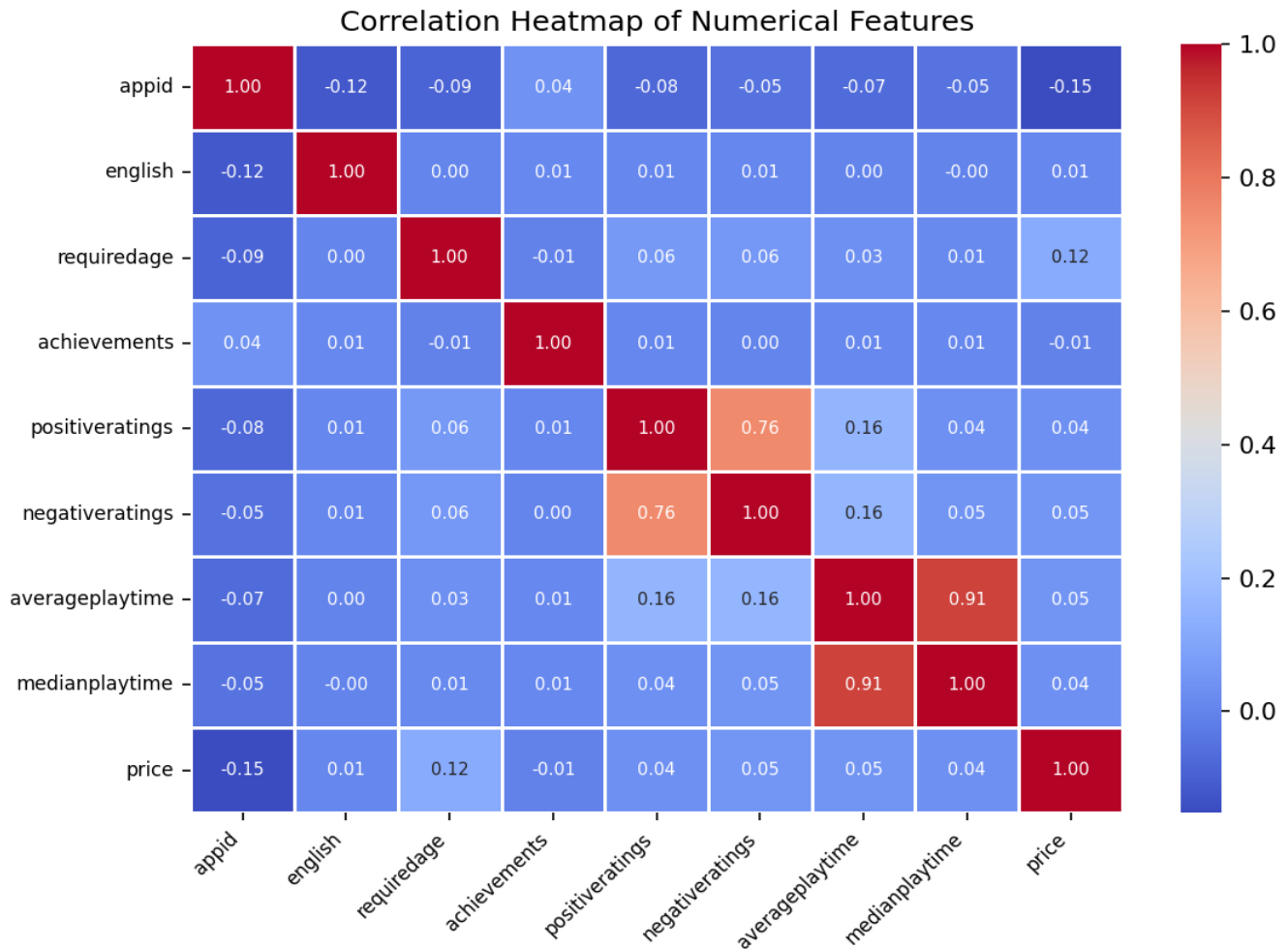
*Histogram showing the distribution of 'appid'. This helps understand its central tendency, spread, and skewness.*



*Histogram showing the distribution of 'english'. This helps understand its central tendency, spread, and skewness.*

## ***Distribution of Categorical Features***

## Correlation Analysis



Heatmap visualizing linear correlations between numerical features. Values range from -1 (strong negative) to +1 (strong positive).

Interpretation: A correlation matrix is a table that displays the correlation coefficients between different variables in a dataset. These coefficients measure the strength and direction of the linear relationship between each pair of variables, ranging from -1 (perfect negative correlation) to 1 (perfect positive correlation), with 0 indicating no correlation. The matrix provides a concise way to visualize and identify patterns, relationships, and potential dependencies between variables. Upon examining the given correlation matrix, one of the strongest correlations observed is between `positiveratings` and `negativeratings`, with a coefficient of 0.756570. This strong positive correlation implies that as the number of positive ratings increases, the number of negative ratings also tends to increase. This might seem counterintuitive, but it could suggest that games with more overall engagement (i.e., more ratings, both positive and negative) are more likely to have a larger number of both positive and negative reviews. Another notable correlation is between `averageplaytime` and `medianplaytime`, with a coefficient of 0.914881, indicating a very strong positive relationship between these two variables. This is not surprising, as both variables measure similar aspects of gameplay duration, and it is expected that they would be closely related. Overall, these correlations provide insights into the relationships between different variables in the dataset, which can be useful for further analysis and modeling.

# Data Quality Assessment

## ***Missing Values***

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

## ***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.

Overall Data Quality Remarks: Based on the key findings from the data quality assessment, the overall data quality appears to be relatively good. With a total of 27,085 rows, the dataset is substantial, providing a solid foundation for analysis. The absence of missing values is a significant positive aspect, as it eliminates the need for imputation or other methods to handle missing data, which can be time-consuming and potentially introduce bias. This suggests that the data collection process was thorough and well-managed. However, the presence of 10 duplicate rows may indicate some issues with data management or processing. While the number of duplicates is relatively small compared to the total number of rows (approximately 0.037%), it is still important to investigate and address this issue to ensure the accuracy and reliability of the data. The potential implications of these findings for further analysis are that the duplicates may lead to overrepresentation of certain data points, potentially skewing results. Removing or handling these duplicates will be necessary to prevent any potential biases or inaccuracies in subsequent analyses. Additionally, the fact that there are no missing values may also imply that the data collection process was robust, but it may also raise questions about whether the data collection process was too restrictive, potentially excluding important information. Overall, while the data quality is generally good, careful attention to the duplicate rows and potential implications for analysis is necessary.

## Conclusion and Potential Next Steps

The initial automated analysis of the dataset has provided a foundational understanding of the data, revealing key characteristics, distributions, and relationships between the 18 variables. This report serves as a crucial starting point for further investigation, offering a glimpse into the underlying structure and patterns within the data. Building upon this foundation, several potential next steps can be explored to gain a deeper understanding of the data. Conducting hypothesis tests to validate observed relationships can help confirm the significance of correlations and trends identified in the initial analysis. Additionally, performing feature engineering to create new insightful variables can uncover hidden patterns and relationships that may not be immediately apparent from the existing features. Further analysis can also involve segmenting the data based on key categorical features, such as genres or platforms, to facilitate comparative analysis and identify distinct trends or characteristics within specific subsets of the data. By pursuing these avenues of investigation, a more comprehensive and nuanced understanding of the data can be developed, ultimately informing decision-making and driving actionable insights.