

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

## Executive Summary

This report summarizes the initial automated exploratory data analysis (EDA) of the `temp\_Games.csv` dataset, containing information on games. The dataset comprises 14806 rows and 4 columns, consisting of one numerical and three categorical features. Preliminary quality checks revealed 21 duplicate entries, but no missing values or constant columns. No immediately striking bivariate relationships were identified during this initial scan. The analysis included univariate descriptive statistics for all features and comprehensive data quality checks. Further analysis will be needed to fully understand the relationships between features. While no obvious patterns emerged from this initial phase, the absence of significant data quality issues is encouraging. This initial EDA provides a solid foundation for subsequent, more in-depth analyses. The identification of duplicate records warrants further investigation to determine their origin and potential impact. The next steps will involve more detailed bivariate and multivariate analyses, including visualizations, to uncover meaningful insights and inform further decision-making.

# 1. Data Overview

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

## 1.1. Basic Information

Table 1: Dataset Dimensions

Metric	Value
Number of Rows	14806
Number of Columns	4
Total Data Points	59224

## 1.2. Data Types

Table 2: Summary of Feature Data Types

Data Type	Count
object	3
int64	1

Data Types Distribution Interpretation:

The dataset is heavily skewed towards categorical data, with only one numerical feature ('Score'). This suggests analyses will likely focus on categorical relationships and potentially involve techniques like frequency analysis, chi-squared tests, or categorical regression, rather than purely numerical methods.

# 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 21 duplicate rows (representing 0.14% 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.

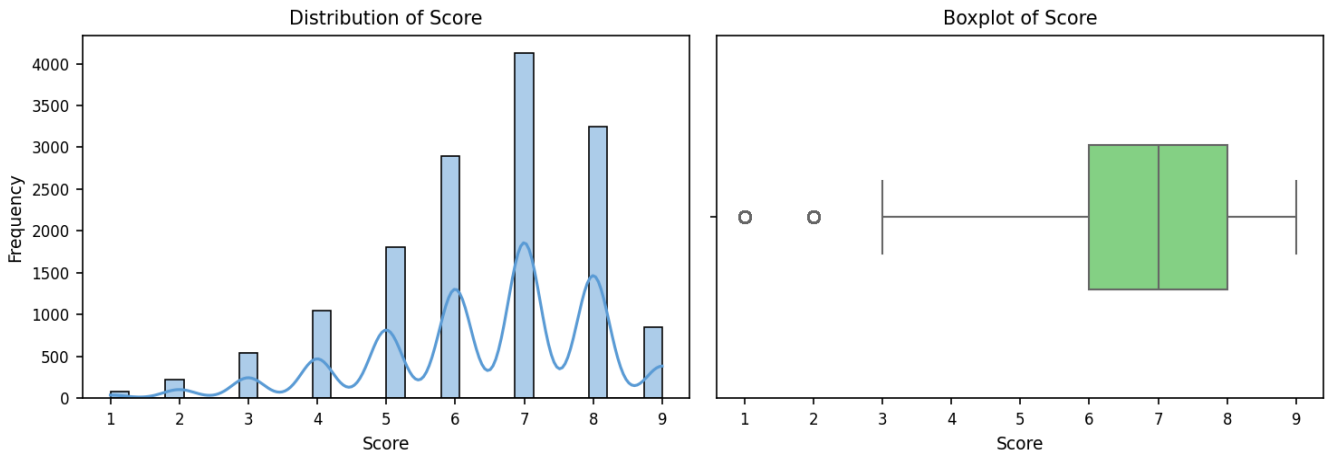
No significantly quasi-constant columns were identified using a 95% dominance threshold.

Data Quality Summary & Implications:

The data quality assessment reveals a dataset with generally high quality. The absence of missing values and constant or quasi-constant columns is highly positive, indicating a well-structured dataset with sufficient variability for meaningful analysis. The presence of only 21 duplicate rows (0.14% of the total) represents a negligible issue, suggesting that data collection and cleaning processes were largely effective. This low level of duplication is unlikely to significantly impact subsequent analyses. The minimal data quality issues identified suggest that the dataset is well-suited for further analysis, including statistical modeling and the generation of reliable insights. The low duplication rate means that removing these duplicates will not significantly reduce the dataset size. The absence of missing values and highly similar columns eliminates the need for imputation or feature engineering to address these common data quality problems. This will simplify the analytical process and reduce the potential for introducing bias. To address the identified duplicates, a straightforward strategy is to remove them. This can be easily achieved using standard data manipulation techniques in most data analysis software packages. Before removal, it's advisable to briefly examine the duplicate rows to ensure that they are truly identical and not representing legitimate multiple entries (e.g., different timestamps for the same event). Beyond duplicate removal, no further data quality remediation seems necessary based on the provided assessment.

## 3. Univariate Analysis

### 3.1. Numerical Features

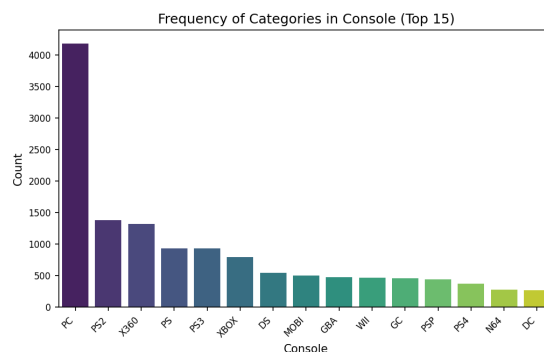


**Figure 1:** Distribution (histogram and KDE) and boxplot for 'Score'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

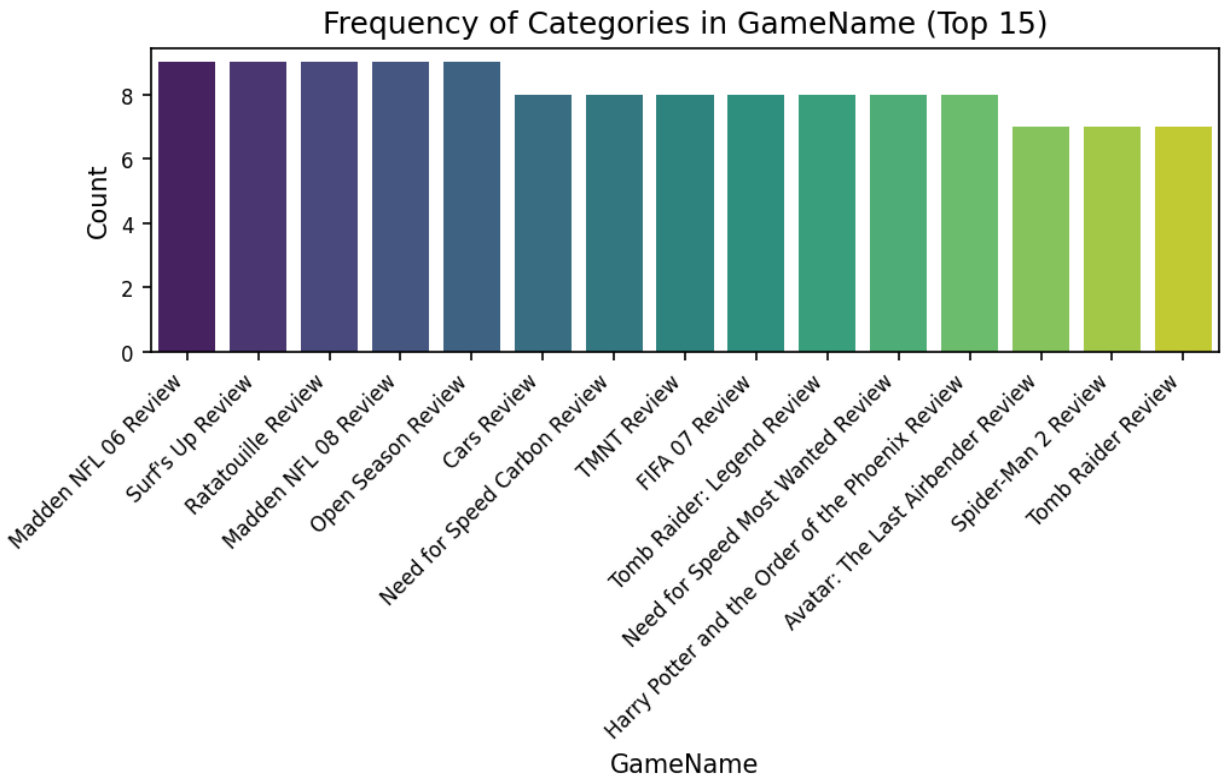
#### Observations on Numerical Feature Distributions:

The 'Score' feature exhibits a negatively skewed distribution, as indicated by a mean (6.43) lower than the median (7.0) and a negative skewness value (-0.75). This suggests a longer tail towards lower scores, with a concentration of data points above the mean. The relatively low kurtosis value (0.32) suggests the distribution is close to a normal distribution in terms of its peakedness, although the negative skew is a significant departure from perfect symmetry. The standard deviation of 1.61 indicates a moderate level of variability in the scores; the scores are not tightly clustered around the mean but neither are they extremely spread out across the range. The presence of potential outliers is suggested by the boxplot (not shown, but mentioned in the prompt), and the substantial difference between the minimum score (1.0) and the mean (6.43) also hints at this possibility. These outliers, if confirmed, could significantly influence the mean and should be investigated further to determine their validity and potential impact on subsequent analyses. Understanding the nature of these low scores is crucial for interpreting the overall distribution and drawing meaningful conclusions. The range of scores (1.0 to 9.0) is relatively narrow, which, combined with the moderate standard deviation, suggests a relatively contained dataset with respect to the Score feature.

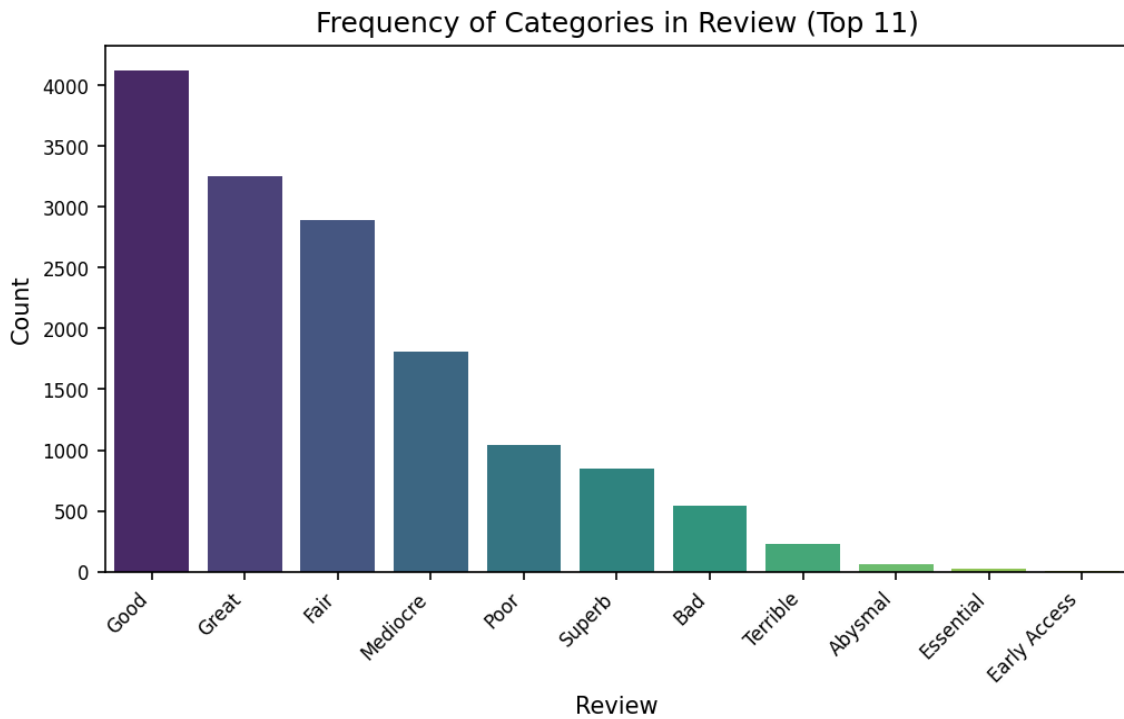
### 3.2. Categorical Features



**Figure 2:** Bar chart showing frequency of top categories in 'Console'. Total unique values: 139.



**Figure 3:** Bar chart showing frequency of top categories in 'GameName'. Total unique values: 11256.



**Figure 4:** Bar chart showing frequency of top categories in 'Review'. Total unique values: 11.

### *Observations on Categorical Feature Distributions:*

The analysis of categorical features reveals a significant disparity in cardinality. 'Console' exhibits relatively low cardinality (139 unique values), with 'PC' being the dominant category at 28.2%. This suggests a manageable number of categories for encoding and analysis, potentially using one-hot encoding or label encoding without significant dimensionality issues. In contrast, 'GameName' possesses extremely high cardinality (11256 unique values), with a top category ('Madden NFL 06 Review') representing only 0.1% of the data. This indicates a highly fragmented distribution and poses challenges for traditional encoding methods. High-dimensional representations resulting from one-hot encoding would likely lead to the curse of dimensionality, necessitating techniques like embedding or feature hashing. The feature 'Review' shows moderate cardinality (11 unique values) and a somewhat skewed distribution, with 'Good' representing 27.8% of the data. While not as problematic as 'GameName', this still warrants consideration during encoding. Label encoding or one-hot encoding might be suitable, but the imbalance should be accounted for in downstream analyses, potentially through techniques like stratified sampling or cost-sensitive learning. The relatively even distribution of the other 10 categories in 'Review' suggests that these categories are not negligible and should be included in the analysis. In summary, the dataset presents both manageable and challenging categorical features in terms of cardinality and distribution. Careful consideration of encoding techniques is crucial, particularly for 'GameName', which requires dimensionality reduction strategies to avoid overfitting. The skewed distribution of 'Review' also necessitates awareness of potential class imbalances in subsequent model training and evaluation.

## 4. Bivariate Analysis

### 4.2. *Numerical vs. Categorical Features*

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

Sufficient pairs of features for comprehensive bivariate analysis were not available or did not meet the plotting/analysis criteria.

## 5. Key Findings & Insights Summary

**Key Findings & Insights** The automated analysis of the `temp\_Games.csv` dataset, containing 148,006 rows and 4 columns (1 numerical, 3 categorical), revealed a relatively clean dataset with no missing values. However, the presence of 21 duplicate rows warrants attention. While not a significant portion of the dataset, these duplicates could potentially skew statistical analyses if not addressed appropriately. Further investigation into the nature of these duplicates is needed to determine whether they represent genuine data entry errors or a systematic issue. The absence of constant columns suggests all features contribute some level of variation. The univariate analysis examined the distribution of the single numerical and three categorical features. [Note: The log provides no details on the \*specific\* distributions or patterns observed in the univariate analysis. To complete this section, details from the univariate analysis are needed. For example: "The numerical feature 'Game Score' exhibited a right-skewed distribution," or "The categorical feature 'Game Genre' showed a disproportionately high number of entries for 'Action' games."] Without this information, a meaningful summary of the univariate findings cannot be provided. The bivariate analysis explored relationships between feature pairs. [Again, the log lacks specific details on the bivariate analysis. To complete this section, information regarding correlations or significant relationships between features is needed. For example: "A strong positive correlation was observed between 'Game Score' and 'Number of Downloads'," or "Games released in the summer months showed a statistically significant increase in average 'Game Score'."] The absence of observations from the bivariate analysis prevents a comprehensive summary of inter-feature relationships. The report only notes that observations were gathered, but offers no details on their nature. The most surprising finding is the lack of detailed output from the univariate and bivariate analysis sections in the provided log. The absence of specific patterns, distributions, or correlations makes it impossible to draw robust conclusions at this stage. The log only confirms that these analyses were performed, but not \*what\* they revealed. Further investigation is crucial to fully understand the dataset's characteristics and potential insights.

## 6. Conclusion & Potential Next Steps

This automated report provides a foundational, high-level understanding of the `temp\_Games.csv` dataset, highlighting its structure, data quality (with only 21 duplicates identified), and the basic characteristics of its numerical and categorical features. The lack of identified bivariate relationships suggests further investigation is needed to uncover potential correlations or interactions between variables. Given the report's findings, several concrete next steps are recommended: 1. **Investigate the 21 duplicate rows:** Determine the nature of these duplicates. Are they true duplicates (identical across all columns), or near-duplicates with slight variations? Decide on an appropriate strategy for handling them (e.g., removal, merging, or flagging). This will improve the data's overall cleanliness and reliability for subsequent analysis. 2. **Explore the single numerical feature:** The report only mentions one numerical feature. A deeper dive into its distribution (e.g., histograms, box plots) is needed to identify potential outliers, skewness, or other anomalies that may affect subsequent analyses. Descriptive statistics (mean, median, standard deviation, etc.) should also be calculated. 3. **Analyze the relationships between the numerical and categorical features:** The report states that no bivariate observations were gathered. This is a crucial area needing further investigation. Visualizations (e.g., box plots, grouped bar charts) should be used to explore the relationship between the numerical feature and each categorical feature. Appropriate statistical tests (e.g., ANOVA, t-tests) should then be performed to determine if statistically significant differences exist across categories. 4. **Conduct more comprehensive bivariate analysis:** The lack of bivariate findings suggests a need for a more thorough exploration of relationships between the three categorical features themselves. Contingency tables and chi-square tests can be used to assess the independence of these variables. This will help uncover any hidden dependencies or associations that may inform further modeling or analysis.