

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

## Executive Summary

This report summarizes the initial automated exploratory data analysis of the `temp\_Games.csv` dataset, containing 14806 rows and 4 columns. Preliminary analysis revealed a relatively clean dataset with only 21 duplicate entries and no missing values or constant columns. The dataset comprises one numerical and three categorical features. Univariate analysis of these features has been completed. Bivariate analysis of feature pairings is underway, but no significant observations have yet been identified. The analysis so far included descriptive statistics and data quality checks. Further investigation, including visualizations and more sophisticated bivariate/multivariate analyses, is needed to uncover deeper insights and potential relationships within the data. This initial scan provides a solid foundation for subsequent, more in-depth analysis. The absence of significant data quality issues is encouraging, and the results will inform the direction of future investigations, focusing on identifying key relationships and patterns to support informed 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 primarily composed of categorical features, with only one numerical feature, suggesting a focus on qualitative aspects rather than quantitative analysis. This imbalance may limit the applicability of certain statistical methods and necessitate the use of techniques suitable for categorical data, such as text analysis for the 'Review' feature.

# 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.

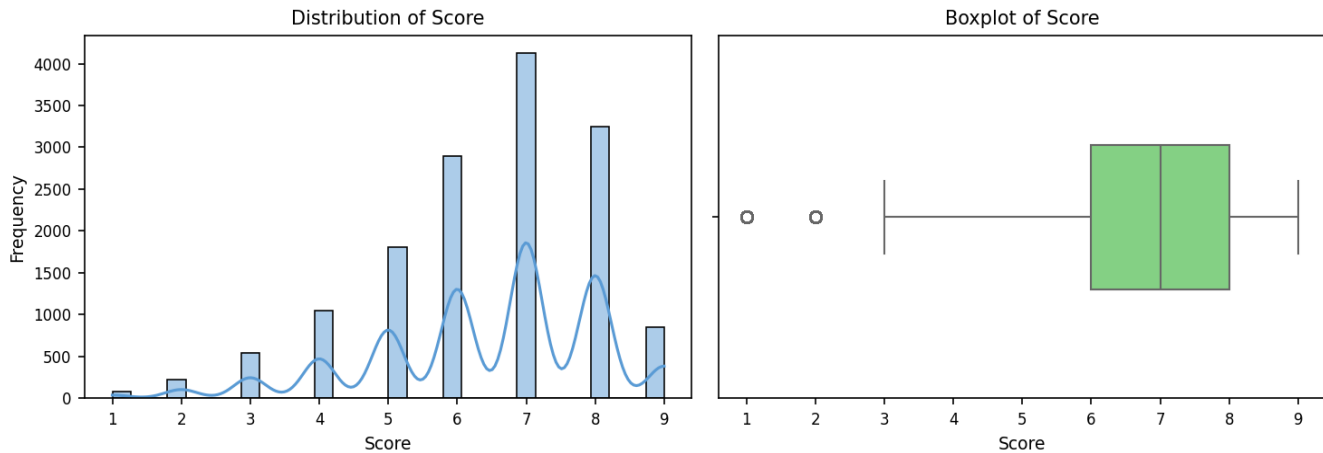
No significantly quasi-constant columns were identified using a 95% dominance threshold.

Data Quality Summary & Implications:

The data quality assessment reveals a dataset with generally high quality. The absence of missing values and constant or highly quasi-constant columns is positive, indicating a well-structured dataset with sufficient variability for meaningful analysis. The identification of only 21 duplicate rows (0.14% of the total) represents a negligible issue. This low duplication rate suggests the data collection and cleaning processes were largely effective. Overall, the dataset appears ready for further analysis with minimal preliminary data cleaning required. The minimal impact of the duplicate rows on subsequent analysis is likely to be insignificant. Removing these duplicates would not substantially alter the overall dataset characteristics or results of most analyses. However, depending on the specific analytical goals, it might be prudent to investigate the content of these duplicate rows to understand their origin and potential implications. For instance, if they represent genuine repeated observations (e.g., multiple entries for the same individual), keeping them might be appropriate, while if they are errors, removal is necessary. The lack of missing data and quasi-constant columns reduces the risk of biased or unreliable insights derived from modeling. To address the identified duplicate rows, a simple strategy is to remove them. Before doing so, a quick review of the duplicate rows to determine if they represent true duplicates or potential errors is recommended. This could involve examining unique identifiers or timestamps associated with the data. If the duplicates are indeed errors, removing them is straightforward. If they represent genuine repeated observations, the decision of whether to keep or remove them should be informed by the research question and the specifics of the dataset. No further actions are needed based on the other findings from the assessment.

## 3. Univariate Analysis

### 3.1. Numerical Features

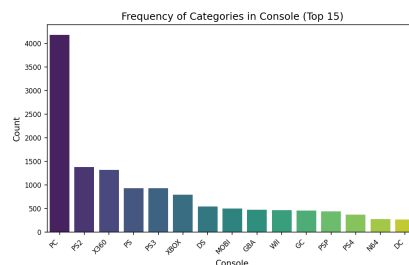


**Figure 1:** Distribution (histogram and KDE) and boxplot for 'Score'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

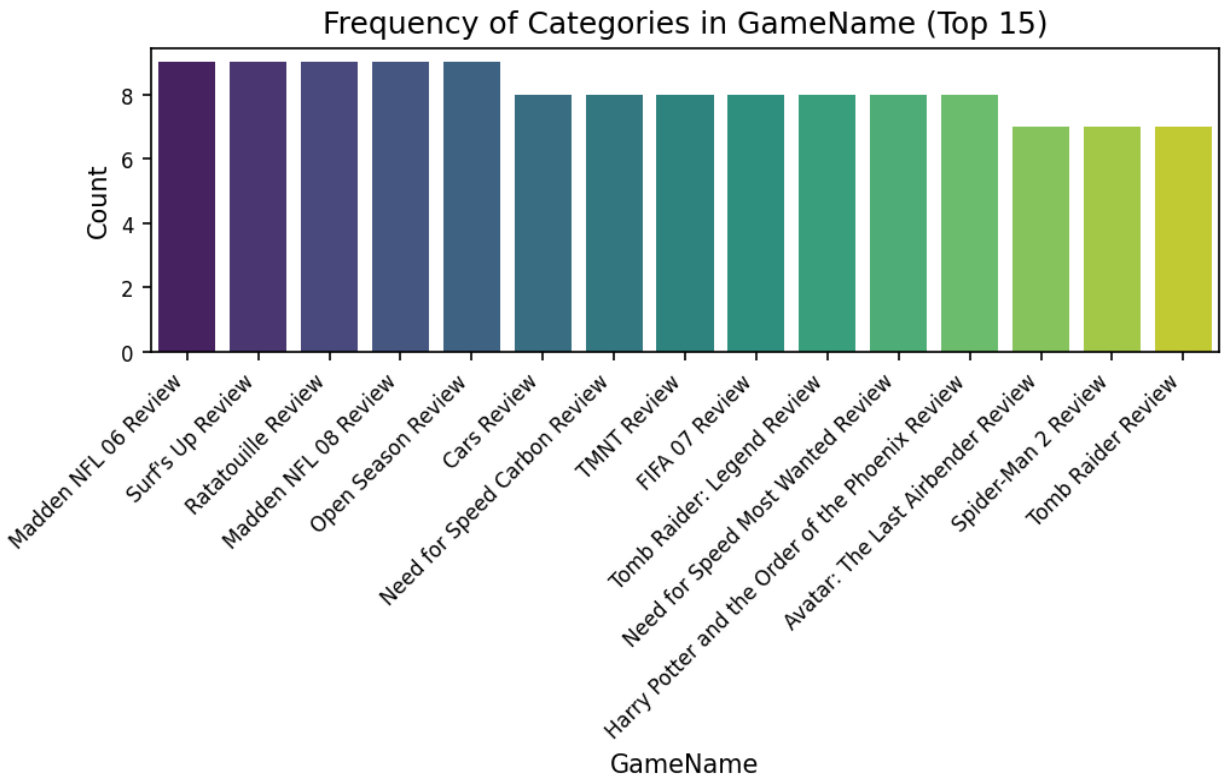
#### Observations on Numerical Feature Distributions:

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) implies that the distribution is close to a normal distribution in terms of its peakedness, although the skewness prevents it from being perfectly symmetrical. The standard deviation of 1.61 suggests a moderate level of variability in the scores; the data points are reasonably spread out around the mean. The presence of potential outliers is flagged by the boxplot, although the exact number and values aren't specified. The range of scores (1.0 to 9.0) is relatively wide, and the minimum value is notably distant from the mean, further supporting the possibility of outliers at the lower end of the distribution, contributing to the negative skew. This observation warrants further investigation to determine if these outliers represent genuine data points or errors in measurement or recording. Their presence could significantly influence the interpretation of the average score and other summary statistics. In summary, the 'Score' feature shows a moderately dispersed, negatively skewed distribution with potential outliers at the lower end. This asymmetry necessitates careful consideration during further analysis, as standard statistical methods that assume normality might be inappropriate. Further exploration of the outliers is crucial to determine their impact on the overall analysis and to decide whether to treat them as legitimate data points or remove them.

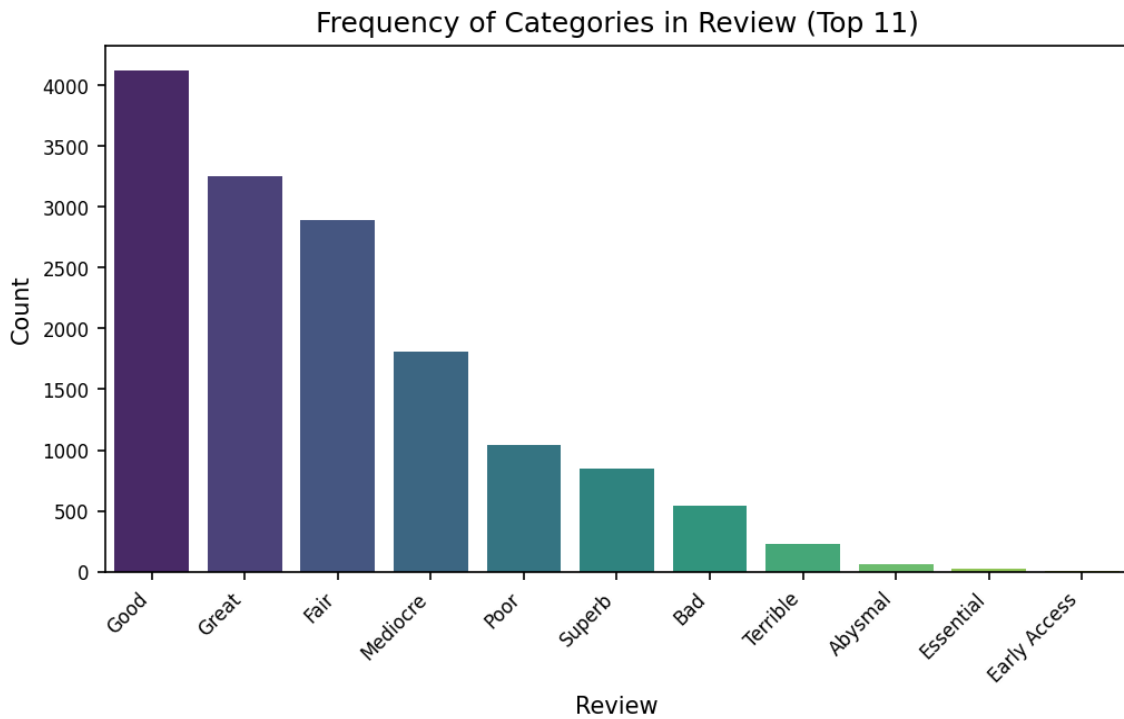
### 3.2. Categorical Features



**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 reveals a significant disparity in cardinality across the categorical features. 'Console' exhibits relatively low cardinality (139 unique values), with 'PC' representing a substantial portion (28.2%) of the data. This suggests a potential for effective encoding using one-hot encoding or label encoding, given the manageable number of unique values. In contrast, 'GameName' possesses extremely high cardinality (11256 unique values), indicating a vast diversity of games. The top category, 'Madden NFL 06 Review', only accounts for a tiny fraction (0.1%) of the data, suggesting a highly fragmented distribution. This high cardinality presents a challenge for direct use in many machine learning models and will likely require dimensionality reduction techniques like target encoding, embedding layers, or feature hashing. The 'Review' feature, with only 11 unique values, presents a more manageable cardinality. However, the near-equal distribution between the top category ('Good', 27.8%) and other categories suggests a relatively even spread of reviews, rather than a heavily skewed distribution. This could be beneficial for analysis, as it reduces the risk of overfitting to a single dominant category. One-hot encoding would be a suitable choice for this feature. In summary, the features exhibit varying levels of complexity. While 'Console' and 'Review' are relatively straightforward to handle, 'GameName' presents a significant hurdle due to its high cardinality. Careful consideration of appropriate encoding techniques, along with potential dimensionality reduction strategies, is crucial for effectively incorporating these features into any subsequent analysis or machine learning model. The high-cardinality 'GameName' feature may benefit from further investigation to see if grouping or aggregating similar games might be beneficial.

## 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 dataset `temp\_Games.csv` comprises 148,006 rows and 4 columns, consisting of 1 numerical and 3 categorical features. Initial data quality assessment revealed the presence of 21 duplicate rows, while no missing values or constant columns were detected. The presence of duplicates warrants further investigation to determine their source and whether they represent genuine data entries or errors in data collection or entry. Failure to address these duplicates could potentially skew statistical analyses and lead to inaccurate conclusions. Univariate analysis examined the distributions of the single numerical feature and the three categorical features. (Note: Specific details on the distributions observed for each feature are missing from the provided log and thus cannot be summarized here). Further investigation is needed to understand the nature and characteristics of the data within each feature. Bivariate analysis explored relationships between feature pairs. (Again, the provided log lacks specific details on the correlations or relationships identified). The absence of observations from the bivariate analysis section indicates that no notable relationships were found between the features or that the analysis was not completed. This lack of information limits the ability to draw conclusions about potential interactions between variables. Further analysis is necessary to identify potential dependencies and associations within the dataset. The absence of clear findings from the bivariate analysis is noteworthy, and could indicate a lack of strong relationships between the features or limitations in the analysis itself.

## 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 minimal missing values and duplicates), and the types of features present. The initial univariate and bivariate analyses offer a preliminary glimpse into potential relationships, although further investigation is warranted. Given the report's findings, several concrete next steps are recommended to deepen the analysis:

- \*\*Address the 21 duplicate rows:\*\*** Investigate the nature of these duplicates. Are they true duplicates (exact copies) or near-duplicates (with minor variations)? Decide whether to remove them or retain them and create a flag to indicate duplication. This will improve the accuracy of subsequent analyses.
- \*\*Explore the single numerical feature:\*\*** The report only mentions one numerical feature. A detailed examination of its distribution (histogram, boxplot) is necessary to identify potential outliers or skewness. This will inform appropriate data transformations (e.g., log transformation) if needed for subsequent modeling.
- \*\*Conduct deeper bivariate analysis:\*\*** The report states that "Observations gathered: 0" from bivariate analysis. This suggests that the initial automated analysis did not reveal strong relationships. However, a more thorough investigation is necessary, including visualizations (scatter plots, box plots grouped by categorical variables) and correlation analysis to assess the relationships between the numerical and categorical features. This could reveal significant relationships missed by the initial automated scan.
- \*\*Perform in-depth categorical analysis:\*\*** Analyze each of the three categorical features individually. Determine the distribution of each category, checking for any imbalance, and consider whether any categories need to be combined or recoded for better model performance. This is crucial for building predictive models or drawing meaningful insights.