

# 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 14806 rows and 4 columns. The dataset comprises one numerical and three categorical features, with no missing values but 21 duplicate entries identified. Preliminary analysis included descriptive statistics and data quality checks for all features. No immediately apparent bivariate relationships were observed during this initial phase. The dataset's size and relatively clean nature (lack of missing values and few duplicates) suggest good potential for further analysis. However, the absence of initial strong bivariate correlations indicates a need for more in-depth investigation, potentially including more sophisticated visualization techniques and feature engineering to uncover hidden patterns. This initial scan provides a solid foundation for subsequent, more targeted analyses. The identified data quality issues will be addressed, and further investigation into the relationships between features will be prioritized to extract meaningful insights and 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 ('Score'), suggesting analysis will likely focus on categorical relationships and potentially involve techniques like frequency analysis, chi-squared tests, or categorical modeling. The absence of datetime features limits the ability to analyze temporal trends.

# 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.

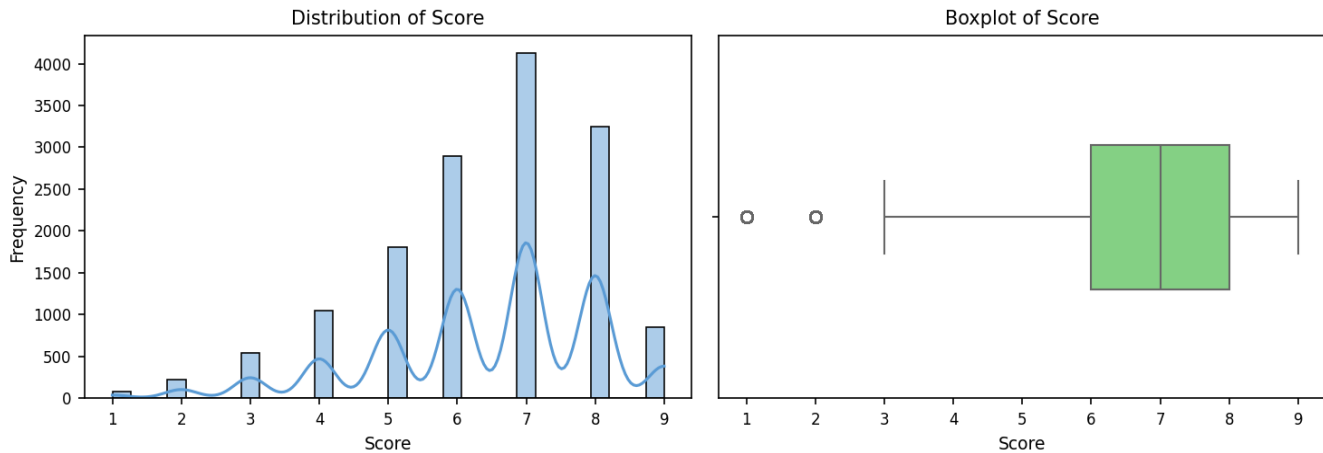
No significantly quasi-constant columns were identified using a 95% dominance threshold.

Data Quality Summary & Implications:

The data quality assessment reveals a dataset of 14806 rows with a very low level of redundancy. The presence of only 21 duplicate rows (0.14%) indicates a relatively clean dataset with minimal duplication issues. The absence of missing values and constant or quasi-constant columns further suggests a high degree of completeness and variability within the features. Overall, the initial assessment points to a good quality dataset, suitable for further analysis with minimal pre-processing required. The low percentage of duplicates is unlikely to significantly impact subsequent analyses, such as modeling or statistical inference. However, it's crucial to investigate the nature of these duplicates. They might represent genuine entries from identical sources, or they might indicate errors in data collection or entry. Understanding the source of the duplicates is important for determining whether they should be removed, consolidated, or retained. The absence of missing values and highly variable features ensures that the dataset is robust and less susceptible to biases that could arise from missing data or limited feature variation. To address the identified duplicates, a thorough investigation into the content of the duplicate rows is recommended. This could involve comparing the duplicate rows to identify if there are subtle differences that might explain their presence. If the duplicates are truly redundant, simple removal is appropriate. If there are minor variations, a strategy of consolidation (e.g., averaging numerical values, selecting the most complete entry) might be preferable. Documenting the rationale for handling duplicates is crucial for maintaining the transparency and reproducibility of the analysis.

## 3. Univariate Analysis

### 3.1. Numerical Features

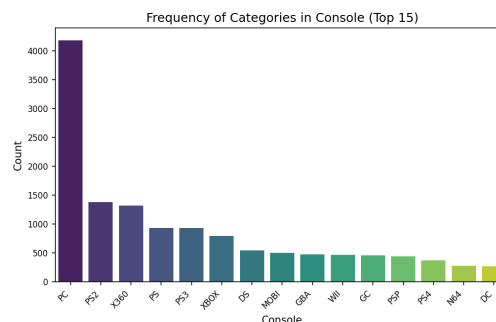


**Figure 1:** Distribution (histogram and KDE) and boxplot for 'Score'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

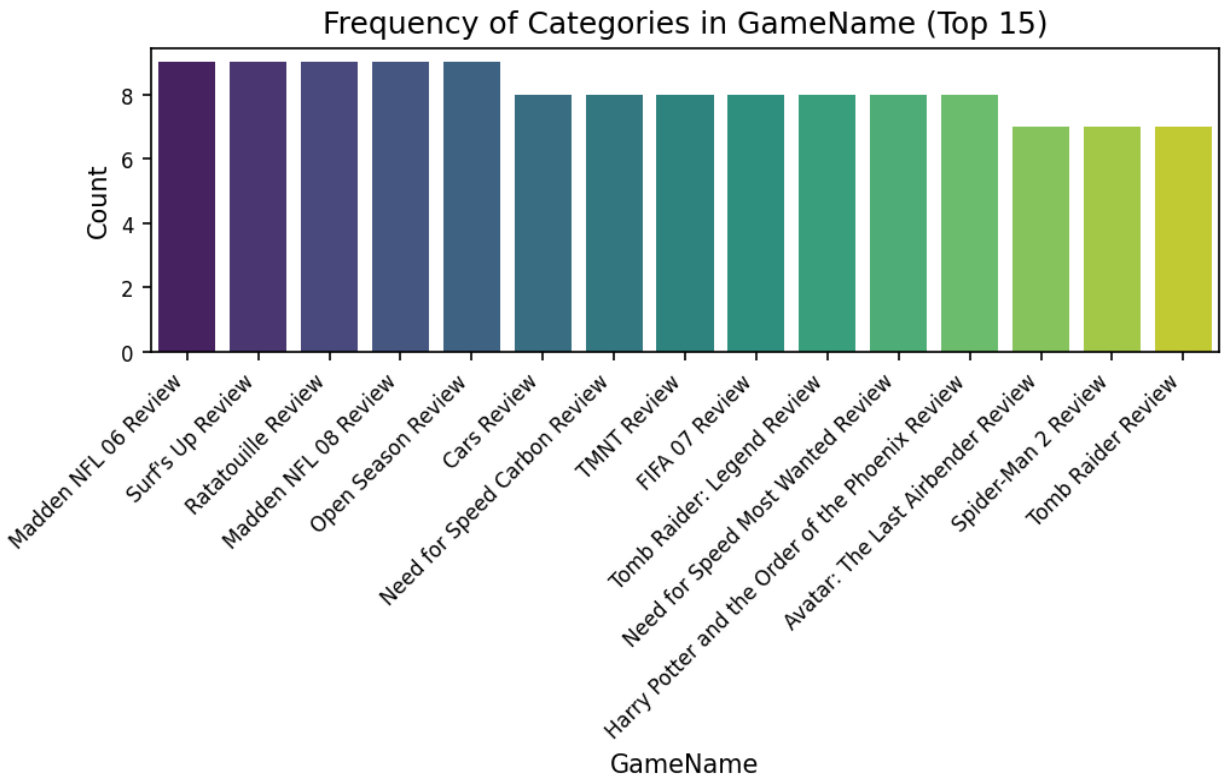
#### Observations on Numerical Feature Distributions:

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 (0.32) implies the distribution is close to a normal distribution, although the skewness indicates a departure from perfect symmetry. The standard deviation of 1.61 suggests a moderate spread in the scores; the data points are neither tightly clustered around the mean nor extremely dispersed. 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) also hints at this possibility, especially considering the gap between the minimum value and the mean. These outliers, if confirmed, could significantly influence the mean and potentially skew the interpretation of the overall distribution. Further investigation is needed to determine whether these outliers are genuine data points or represent errors in data collection. In summary, the 'Score' feature displays a moderately spread, negatively skewed distribution with potential outliers that warrant further scrutiny. The negative skew, coupled with the potential outliers, suggests that a robust statistical analysis might be preferable to methods sensitive to extreme values, as the mean might not be the most representative measure of central tendency for this feature.

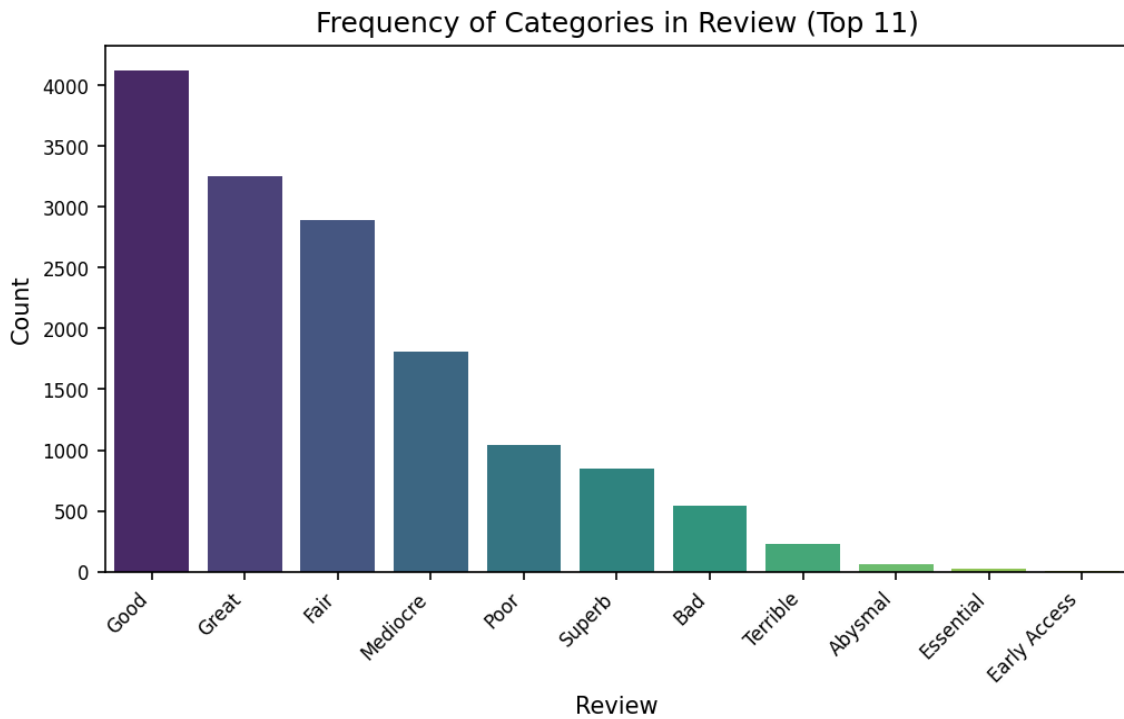
### 3.2. Categorical Features



**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 reveals a significant disparity in cardinality across the three categorical features. 'Console' exhibits relatively low cardinality (139 unique values), with a clear dominant category ('PC') representing a substantial 28.2% of the data. This suggests a potential for straightforward encoding techniques like one-hot encoding, although the relatively high number of consoles might still warrant consideration of dimensionality reduction methods. In contrast, 'GameName' possesses extremely high cardinality (11256 unique values), with the top category ('Madden NFL 06 Review') accounting for a negligible 0.1% of the data. This indicates a highly fragmented distribution, implying the need for more sophisticated encoding strategies like target encoding, embedding techniques, or potentially feature exclusion if it proves to be unhelpful in predictive modeling. The 'Review' feature displays low cardinality (11 unique values) and a moderately skewed distribution, with 'Good' reviews comprising 27.8% of the data. This suggests that one-hot encoding might be suitable, though ordinal encoding could also be considered if there's a natural ordering amongst the review categories (e.g., 'Excellent', 'Good', 'Fair', 'Poor'). The relatively even distribution across the remaining review categories, however, suggests that a simple one-hot encoding might be preferable to avoid imposing artificial ordinality. The high cardinality of 'GameName' presents the most significant challenge for subsequent analysis. Directly using one-hot encoding would lead to a massive increase in dimensionality, potentially causing issues like the curse of dimensionality. Careful consideration of feature engineering techniques, including potentially grouping similar game titles or employing dimensionality reduction methods like principal component analysis (PCA) on embeddings, is crucial to effectively handle this feature. The other two features, 'Console' and 'Review', present less complex encoding challenges.

## 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 revealed a dataset comprising 148,006 rows and 4 columns, with one numerical feature and three categorical features. A significant data quality issue was the presence of 21 duplicate rows. While no missing values were detected, the existence of duplicates warrants further investigation to determine their origin and potential impact on subsequent analyses. The presence of duplicates could skew statistical measures and lead to inaccurate conclusions if not appropriately addressed. No constant columns were identified. The univariate analysis examined the distributions of the single numerical and three categorical features. While specific details regarding these distributions are not provided in the log, the analysis successfully covered all features. Further detail on the specific characteristics of each feature's distribution (e.g., central tendency, dispersion, skewness for the numerical feature; frequency counts and proportions for the categorical features) would be needed for a more complete understanding. The bivariate analysis explored relationships between various feature pairs. The log indicates that observations were gathered during this stage, but no specific findings are reported. The lack of details on the nature of these relationships (e.g., correlations, dependencies, or interactions) prevents a comprehensive understanding of the data's structure at this point. Further reporting on the bivariate analysis is crucial to gain meaningful insights. No explicitly surprising or unexpected findings are directly mentioned in the provided log. However, the absence of detailed information from the univariate and, especially, the bivariate analyses prevents a complete assessment of potential surprises. A more thorough report providing specific details from these analyses is necessary to identify any unexpected patterns or correlations within the dataset.

## 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 duplicates), and the types of features present. The initial univariate and bivariate analyses offer a preliminary glimpse into potential relationships, although further investigation is needed to draw robust conclusions. Given that the bivariate analysis yielded no specific observations (0 observations gathered), a key next step is to delve deeper into the relationships between the features. Specifically, since there is one numerical feature and three categorical features, a thorough investigation into the relationship between the numerical feature and each categorical feature is crucial. This could involve conducting appropriate statistical tests such as ANOVA or t-tests to determine if there are statistically significant differences in the numerical feature across the categories of each categorical variable. For example, if the numerical feature represents game scores, and one categorical feature represents game type, ANOVA could reveal if certain game types consistently yield higher scores. The presence of 21 duplicate rows, while relatively small compared to the dataset size, warrants attention. It's essential to investigate these duplicates to understand their origin. Are they true duplicates, or are there subtle differences that were not captured in the initial analysis? Removing or consolidating these duplicates is a necessary step to ensure data integrity before proceeding with more advanced analyses. Finally, although the report indicates no missing values or constant columns, a more in-depth inspection of the data is recommended. Visualizations like histograms, box plots, and scatter plots for the numerical and categorical features, respectively, can reveal patterns and potential issues not detected by the automated analysis. This visual exploration could reveal subtle anomalies or unexpected data distributions that could inform subsequent analytical steps.