

# 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. 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, revealing no immediately apparent patterns in the bivariate analysis. The dataset's relatively large size (14806 rows) suggests sufficient data for meaningful analysis. The absence of missing values and identification of duplicates are significant quality indicators, informing subsequent data cleaning strategies. Further investigation is required to uncover relationships between features and to generate insights. This initial scan provides a foundational understanding of the dataset's structure and quality. Further analysis, focusing on bivariate and multivariate relationships, including appropriate visualizations, will be crucial for extracting actionable insights and informing 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 for analysis. This limited numerical data may restrict the application of certain analytical techniques, potentially necessitating the use of methods suitable for categorical data or requiring further feature engineering to create more numerical variables.

# 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.

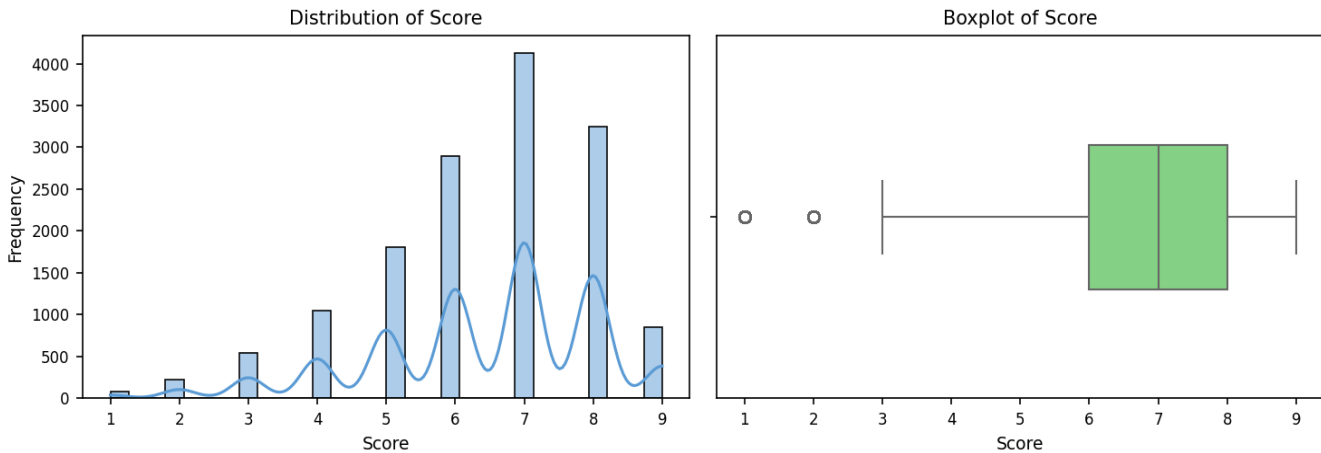
No significantly quasi-constant columns were identified using a 95% dominance threshold.

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 variation in features for meaningful analysis. The presence of only 21 duplicate rows (0.14% of the total) represents a negligible level of redundancy. This low percentage suggests that the duplication is unlikely to significantly skew results or impact the reliability of analyses conducted on this dataset. Overall, the data appears to be clean and ready for further processing and analysis. The minimal number of duplicate rows poses little risk to the integrity of subsequent analyses. Their impact on modeling would be minimal, especially given their low percentage. Inferences drawn from this data are likely to be reliable. However, it's important to understand the nature of these duplicates. If they represent genuine entries (e.g., identical twins in a medical study), then they should be retained. If they are errors in data entry or data ingestion, then they should be removed to ensure data accuracy. The absence of missing values simplifies the preprocessing step, avoiding the need for imputation or removal of incomplete records. To address the identified duplicate rows, a thorough investigation into their origin is recommended. This may involve examining the source data and data entry processes. If the duplicates are judged to be erroneous, they should be removed. Simple deduplication techniques, such as identifying and removing rows with identical values across all columns, would suffice given the small number of duplicates. If the duplicates are legitimate, they can be retained. The overall conclusion is that the data quality is excellent and poses minimal challenges for further analysis, requiring only a minor deduplication step.

## 3. Univariate Analysis

### 3.1. Numerical Features

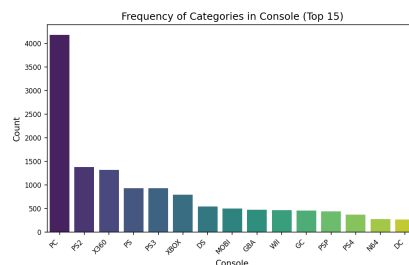


**Figure 1:** Distribution (histogram and KDE) and boxplot for 'Score'. The histogram shows shape, central tendency, and spread. The boxplot highlights median, quartiles, and potential outliers.

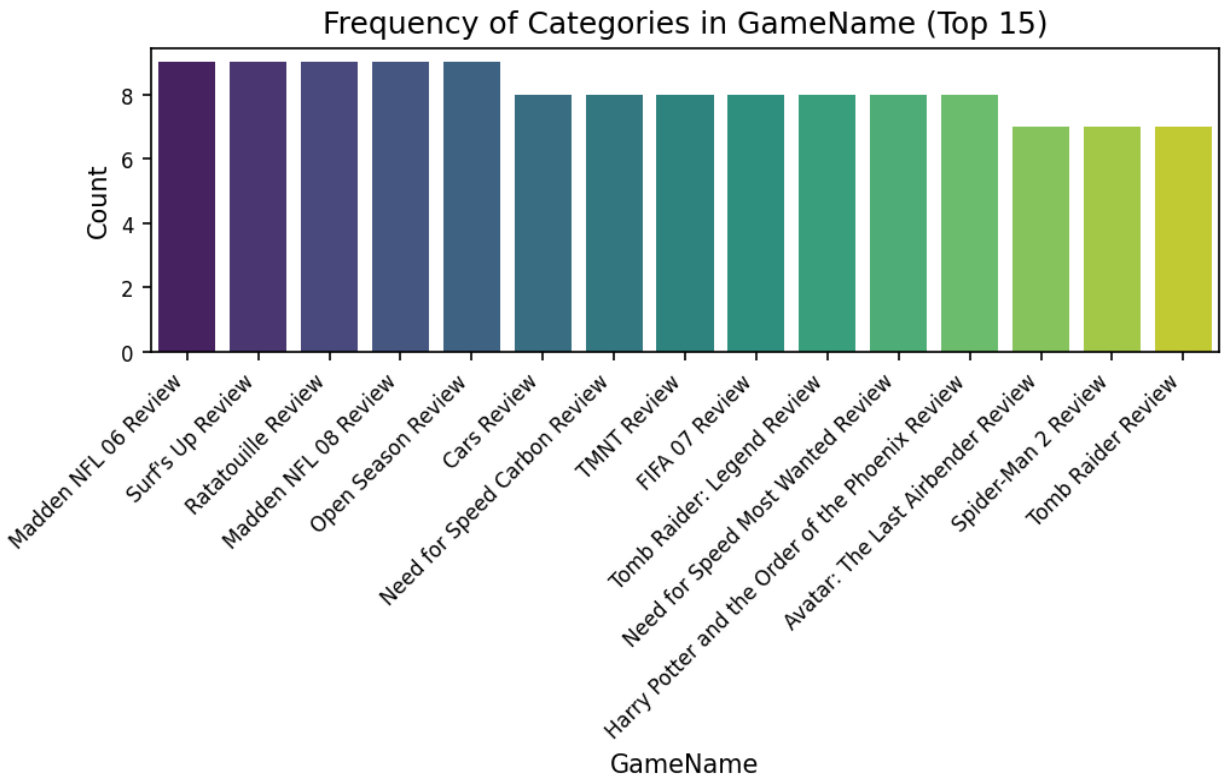
#### Observations on Numerical Feature Distributions:

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, although the negative skew is a significant departure from perfect symmetry. The standard deviation of 1.61 indicates a moderate spread in the scores, meaning there's a noticeable degree of variability among the observations. The presence of potential outliers is flagged by the boxplot analysis and the relatively large difference between the minimum (1.0) and maximum (9.0) values compared to the mean and median. These extreme values at both ends of the distribution, coupled with the negative skew, warrant further investigation. Are these outliers genuine data points or are they errors? Understanding their origin is crucial for determining whether to retain or remove them from the dataset, as their presence could significantly influence subsequent analyses and model building. The relatively small range (8) suggests that even with outliers, the majority of scores are clustered within a reasonable range. In summary, the 'Score' feature shows a moderately spread, negatively skewed distribution with potential outliers that require closer examination. The combination of skewness and potential outliers highlights the importance of using robust statistical methods in subsequent analyses which are less sensitive to extreme values. Further investigation into the nature of these outliers is necessary to determine their impact on the overall interpretation of the data.

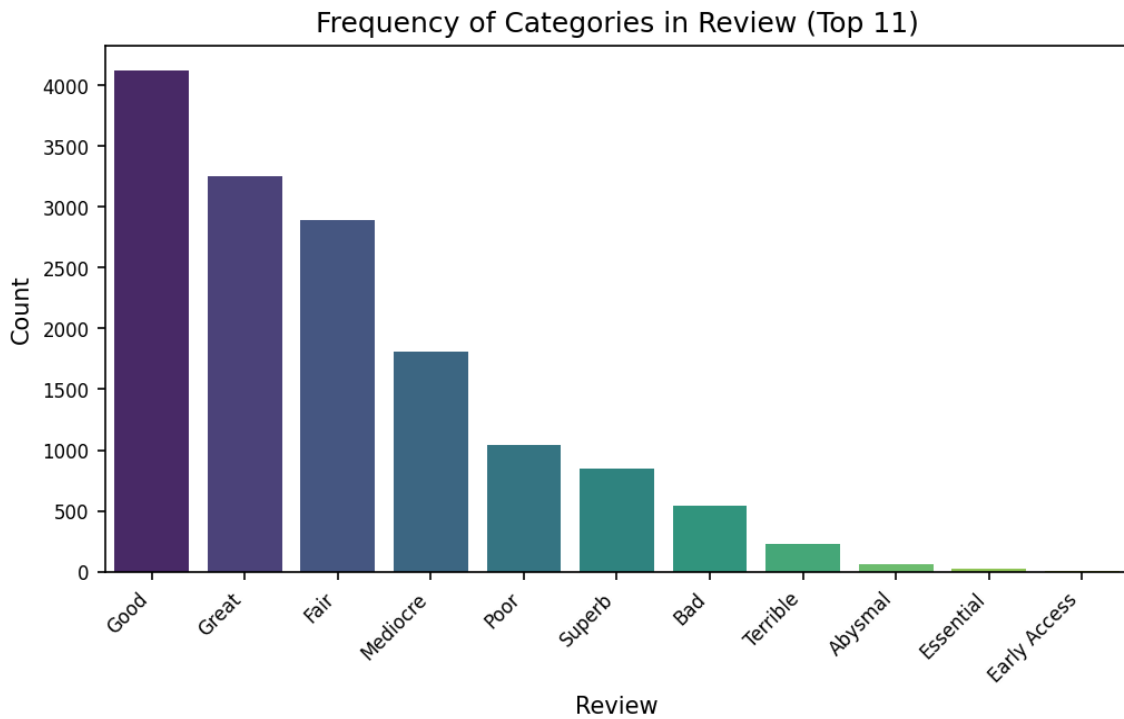
### 3.2. Categorical Features



**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. In contrast, 'GameName' possesses extremely high cardinality (11256 unique values), with the top category ('Madden NFL 06 Review') accounting for only 0.1% of the data. This indicates a highly skewed distribution and necessitates careful consideration of encoding strategies to avoid the curse of dimensionality. Techniques like target encoding, embedding layers (in neural networks), or potentially feature hashing might be more appropriate. The 'Review' feature demonstrates intermediate cardinality (11 unique values) and a moderately skewed distribution, with 'Good' reviews comprising 27.8% of the data. While one-hot encoding remains a viable option, the relatively small number of unique values makes this feature less prone to dimensionality issues compared to 'GameName'. The presence of a dominant category in 'Review' suggests potential for ordinal encoding if the order of the review categories has inherent meaning (e.g., 'Excellent' > 'Good' > 'Average'). In summary, the diverse cardinality and distributions of these features highlight the need for tailored encoding approaches. Low-cardinality features like 'Console' can be easily handled with standard techniques, while the ultra-high cardinality of 'GameName' requires more sophisticated dimensionality reduction or embedding methods. The moderate cardinality and distribution of 'Review' allows for flexibility in encoding choice, depending on the desired level of detail and potential for ordinal interpretation.

## 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, comprising 148,06 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 analyses depending on their nature and the analytical methods employed. Further investigation into the source and characteristics of these duplicates is recommended to assess their impact on subsequent analyses. The absence of constant columns indicates that all features contribute some level of variability to the dataset. Univariate analysis explored the distributions of the single numerical and three categorical features. While specific details of these distributions are not provided in the log, the analysis itself suggests that each feature offers unique insights into the dataset's characteristics. Further information regarding the specific distributions (e.g., skewness, central tendency) of these features is needed for a more complete understanding. The bivariate analysis examined relationships between various feature pairs, though the log offers no specific observations from this phase. The absence of concrete findings here highlights the need for a more detailed report that includes the specific relationships explored and the results obtained. This lack of detail prevents any conclusions regarding correlations or significant interactions between features. The log only indicates that some bivariate analysis was performed, but the results are not reported. The overall analysis reveals a dataset that is largely clean but requires further investigation into the identified duplicate rows and a comprehensive report of the univariate and bivariate analysis results to gain a complete understanding of the data's characteristics and relationships between its features. The current log provides a high-level overview but lacks the granular detail needed to draw substantive conclusions.

## 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 duplicate entries identified), and the initial characteristics of its numerical and categorical features. The lack of missing values and constant columns suggests a relatively clean dataset ready for further exploration. Given the report's findings, several concrete next steps are recommended:

- Investigate the 21 duplicate rows:** Determine the nature of these duplicates. Are they true duplicates (identical across all columns) or near-duplicates (with minor variations)? Understanding the source of these duplicates is crucial for data cleaning and ensuring data integrity. If they are true duplicates, they can be removed. If near-duplicates, careful investigation is needed to determine which entries to keep.
- Perform a deeper univariate analysis:** While a basic univariate analysis was conducted, more in-depth exploration of the single numerical feature is needed. This should include calculating descriptive statistics (mean, median, standard deviation, quartiles, etc.), visualizing its distribution (histogram, box plot), and identifying potential outliers. Similarly, a more detailed look at the distributions of the three categorical features is needed, focusing on frequency counts and visualizations to understand the class imbalance, if any.
- Conduct a thorough bivariate analysis:** The report notes that bivariate analysis was performed but yielded no observations. This warrants a more systematic approach. Specifically, explore the relationships between the numerical feature and each of the categorical features using appropriate visualization techniques (box plots, violin plots) and statistical tests (e.g., ANOVA, t-tests) to determine if there are statistically significant differences in the numerical feature across different categories. Further, explore the relationships between the categorical features themselves using contingency tables and appropriate visualizations (e.g., heatmaps).
- Explore potential feature engineering:** Based on the findings from the deeper univariate and bivariate analyses, consider whether new features can be engineered from the existing ones to improve model performance in subsequent analyses. This might involve creating interaction terms between categorical features or transforming the numerical feature if its distribution is non-normal.