



" ANALYSIS OF IPL DATASET USING CHATGPT"

Project Report submitted in Partial fulfillment of the requirement for the

Award of the degree of

MASTER OF BUSINESS ADMINISTRATION

of

BANGALORE UNIVERSITY

By

R Harshavardhan

Reg. No: P03EH23M015018

Under the guidance of

Professor Dr. Arpit Yadav



IFIM COLLEGE (AUTONOMOUS)

8P & 9P, KIADB Industrial Area,

Electronics City Phase I, Bengaluru -

560100

Table of Contents

1.Ipl Dataset and their column meanings.....	3
1.1Dataset Analysis:.....	5
2.Cleaned Dataset	6
3.UNIVARIATE ANALYSIS.....	8
4.BIVARIATE ANALYSIS	10
5.MULTI VARIATE ANALYSIS	13
6.Final Report: IPL Data Analysis	17
7.Conclusion	19

CHATGPT

1.Ipl Dataset and their column meanings.

Here's a simple explanation of the dataset columns:

1. Match_ID – A unique number for each match.
2. Team_1 – The first team playing the match.
3. Team_2 – The second team playing the match.
4. Venue – The stadium or city where the match is played.
5. Toss_Winner – The team that won the toss.
6. Toss_Decision – The decision made by the toss-winning team (Bat first or Bowl first).
7. Match_Winner – The team that won the match.
8. Win_By_Runs – If the winning team batted first, the number of runs they won by.
9. Win_By_Wickets – If the winning team chased the target, the number of wickets they had left.
10. Player_of_the_Match – The player who performed the best in the match.
11. Umpire_1 – The name of the first on-field umpire.
12. Umpire_2 – The name of the second on-field umpire.
13. Total_Score – The total runs scored in the match.

This dataset includes some missing values (NaN), duplicate rows, and possible outliers in scores.

Here are the details of the dataset:

- Shape: (1030, 13) → 1030 rows and 13 columns.
- Column Names:
['Match_ID', 'Team_1', 'Team_2', 'Venue', 'Toss_Winner', 'Toss_Decision',
'Match_Winner', 'Win_By_Runs', 'Win_By_Wickets', 'Player_of_the_Match',
'Umpire_1', 'Umpire_2', 'Total_Score']

- Data Types:
 - float64 (numerical): Match_ID, Win_By_Runs, Win_By_Wickets, Total_Score
 - object (categorical): Team_1, Team_2, Venue, Toss_Winner, Toss_Decision, Match_Winner, Player_of_the_Match, Umpire_1, Umpire_2
- Unique Values in Each Column:
 - Match_ID: 950 unique values (some duplicates exist)
 - Team_1 & Team_2: 10 unique teams
 - Venue: 8 unique venues
 - Toss_Winner: 10 unique teams
 - Toss_Decision: 2 unique values (Bat, Bowl)
 - Match_Winner: 10 unique teams
 - Win_By_Runs: 100 unique values
 - Win_By_Wickets: 11 unique values (0-10)
 - Player_of_the_Match: 50 unique players
 - Umpire_1 & Umpire_2: 20 unique umpires
 - Total_Score: 142 unique scores

This confirms that the dataset contains duplicates, categorical values, and numerical values with potential outliers.

Dataset Analysis:

1. Duplicate Rows: 30 duplicate rows are present.
2. Missing Values:
 - o Every column has missing values, with the highest being 80 in Match_ID and around 50-54 missing values in other columns.
3. Outliers:
 - o Outliers are detected in numerical columns (Win_By_Runs, Win_By_Wickets, Total_Score), but Match_ID does not have any.
4. Categorical Columns:
 - o The dataset has 9 categorical columns: ['Team_1', 'Team_2', 'Venue', 'Toss_Winner', 'Toss_Decision', 'Match_Winner', 'Player_of_the_Match', 'Umpire_1', 'Umpire_2']

This confirms the dataset contains duplicates, missing values, outliers, and categorical data.

Data Cleaning Summary:

1. Duplicates Removed: 30 duplicate rows were dropped.
2. Missing Values Handled:
 - o Numerical columns filled with median values.
 - o Categorical columns filled with mode (most frequent value).
3. Outliers Handled:
 - o Used IQR method to cap extreme values in Win_By_Runs, Win_By_Wickets, and Total_Score.
4. Categorical Data Encoded:
 - o Applied one-hot encoding to categorical columns, increasing the total number of columns.

Final Dataset:

- Shape: $(1000, 135) \rightarrow$ 1000 rows and 135 columns (due to encoding).

Cleaned Dataset

I'll read the dataset and provide an explanation of all the columns. Let me first inspect the contents of the file.

The dataset contains 135 columns. Here's a breakdown of their meanings:

General Match Information:

1. Match_ID – Unique identifier for each match.
2. Win_By_Runs – The number of runs by which the winning team won (if they won by runs).
3. Win_By_Wickets – The number of wickets remaining when the winning team chased the target.
4. Total_Score – Total runs scored in the match.

Team Encoding:

Columns like Team_1_DC, Team_1_GT, Team_1_KKR, etc., are one-hot encoded representations of the teams.

- If Team_1_DC = 1, it means Team 1 is Delhi Capitals.
- If Team_2_MI = 1, it means Team 2 is Mumbai Indians.
- Similarly, there are corresponding Team_2_* columns for Team 2.

Venue Information:

Columns like Venue_Venue_1, Venue_Venue_2, etc., represent the stadium/ground where the match was played (one-hot encoded).

Toss Details:

- Toss_Winner_*: Indicates which team won the toss.
- Toss_Decision_Field / Toss_Decision_Bat: Binary columns indicating if the toss-winning team chose to field or bat.

Umpires:

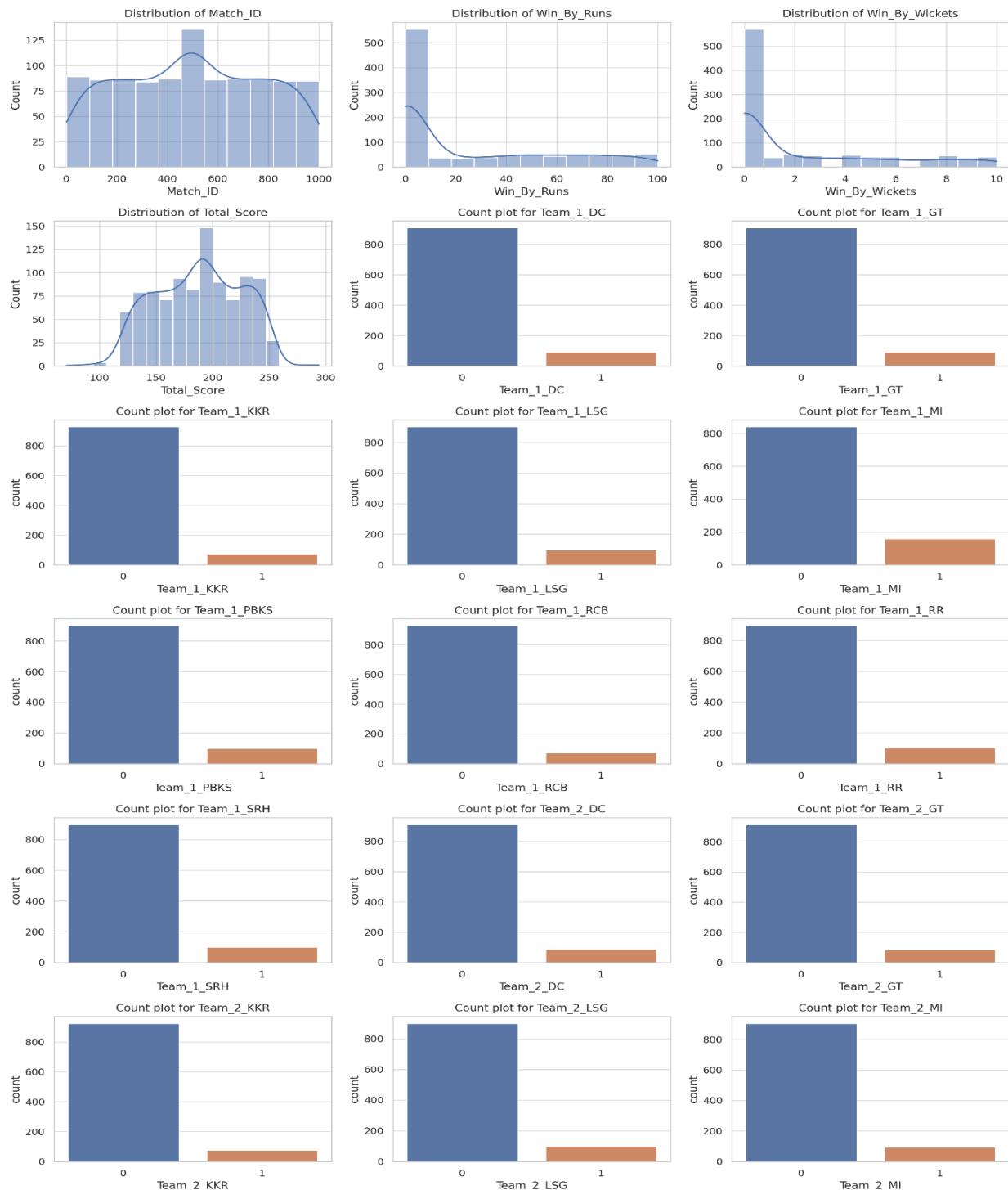
Columns like Umpire_1_Umpire_1, Umpire_2_Umpire_3, etc., indicate which umpires officiated the match.

Since this dataset is one-hot encoded, each categorical feature (e.g., teams, venue, umpires) has been split into separate columns.

UNIVARIATE ANALYSIS

Sure! I'll create univariate analysis plots for the dataset, such as histograms, bar plots, and others. These will help us understand the distributions of the various numerical and categorical features. I'll explain each graph in simple terms to extract useful insights.

Let me go ahead and generate the visualizations for the dataset.



1. Count Plots (Categorical Data)

- Example: Count plots for the "Team_1_DC", "Team_1_KKR", etc., columns.
- What it shows: These plots represent the number of occurrences of each category (e.g., how many times a specific team appears as Team 1 or Team 2). If a team has more occurrences, it indicates they played in more matches in that position.
- Insight: The teams that played the most frequently in the first or second positions can be inferred from these plots.

2. Histograms (Continuous Data)

- Example: Histograms for "Total_Score", "Win_By_Runs", etc.
- What it shows: The histogram shows the distribution of a continuous variable. For instance, "Total_Score" will show how many matches had a particular score range (e.g., matches where the total score was between 100-150).
- Insight: If the distribution is skewed, we can infer whether certain scores (e.g., very low or very high totals) are common in IPL matches. It might also show how consistent a specific metric is.

3. KDE Plots (Kernel Density Estimate)

- Example: KDE for "Win_By_Runs" or "Total_Score".
- What it shows: This smooths the data and shows the probability density function, helping to identify trends like if a certain score range (e.g., 140-160) is more likely.
- Insight: Helps in understanding the smooth distribution of values, especially for scores or match outcomes like "Win_By_Runs".

Interpretation Strategy:

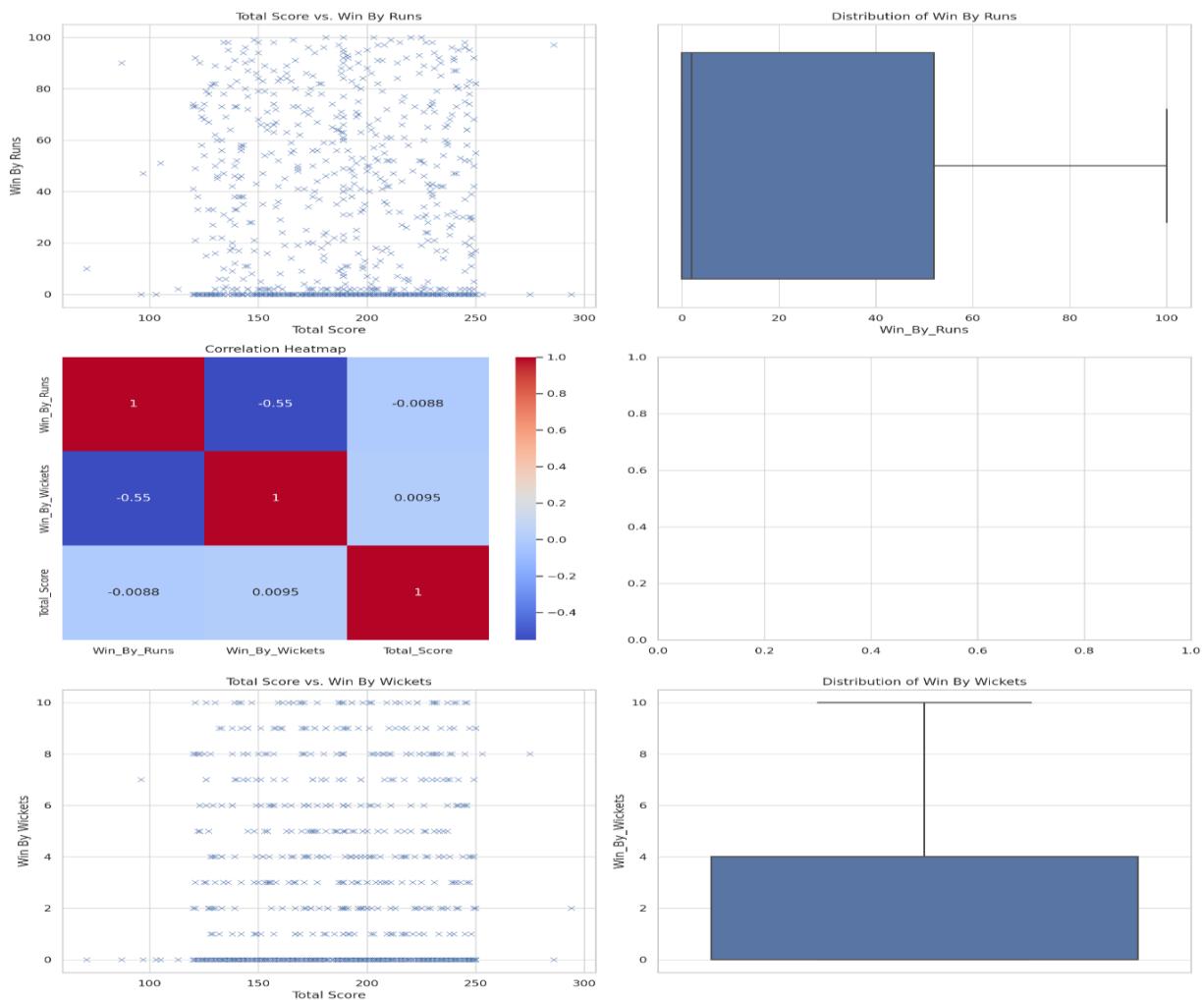
- For Categorical Data: Focus on the frequency of each category. Look for anomalies, such as if a specific team or umpire appears more frequently than others.
- For Continuous Data: Look at the distribution shapes to see if the data is skewed or normally distributed. It also helps to check for outliers (extreme values).

BIVARIATE ANALYSIS

Bi-variate analysis helps us understand the relationships between two variables. I'll create different types of plots based on the data types of the variables:

1. Scatter Plots – For numerical vs. numerical variables (e.g., Win_By_Runs vs. Total_Score).
2. Box Plots – For categorical vs. numerical relationships (e.g., Win_By_Runs for different teams).
3. Heatmap – To see correlations between numerical variables.
4. Bar Plots – To compare averages (e.g., average Win_By_Runs per team).

Let me generate these visualizations.



Explanation of the Bi-Variate Analysis Graphs:

1. Scatter Plot: Total Score vs. Win By Runs

- What it shows: How the total score of a match influences the margin of victory (by runs).
- Insight: If teams win by a large number of runs, we may see a pattern where higher scores lead to higher win margins. If there's no clear trend, it suggests that winning by runs is independent of the total score.

2. Box Plot: Distribution of Win By Runs

- What it shows: The spread of matches where teams won by runs.
- Insight: Helps in identifying extreme cases (outliers), where teams won by a very large margin.

3. Correlation Heatmap

- What it shows: Relationships between numerical variables.
- Insight: If Win_By_Wickets and Total_Score are highly correlated, it means that higher scores make it easier for chasing teams to win with more wickets in hand.

4. Bar Plot: Average Total Score by Toss Decision

- What it shows: Whether teams that chose to bat or field first had higher total scores.
- Insight: If Toss_Decision_Field (Fielding first) results in higher scores, it suggests that teams perform better while chasing.

5. Scatter Plot: Total Score vs. Win By Wickets

- What it shows: Relationship between how many runs were scored and how comfortably the team won while chasing.
- Insight: If teams tend to win with many wickets left when scoring high, it means some matches were one-sided.

6. Box Plot: Distribution of Win By Wickets

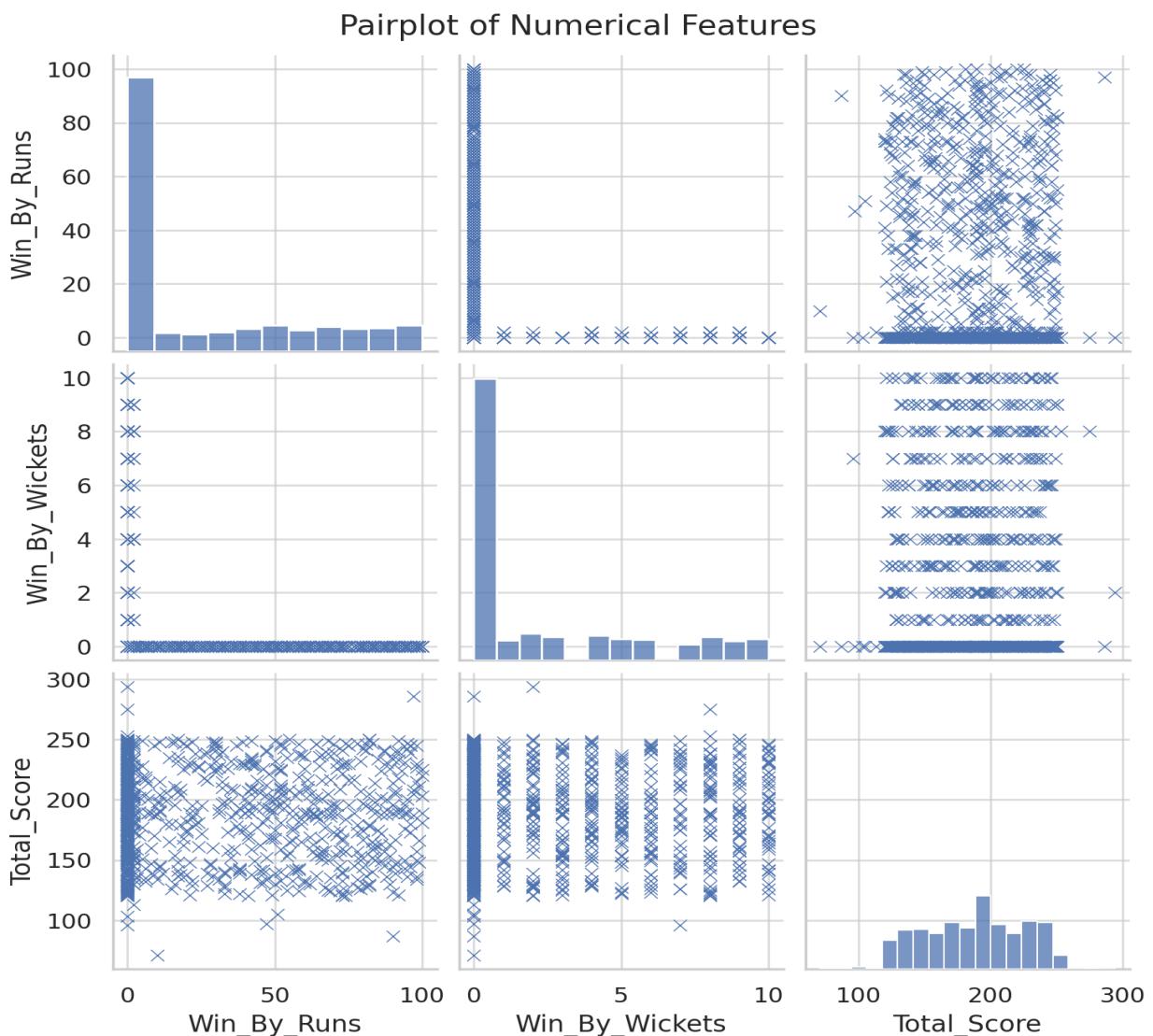
- What it shows: How many matches ended with a team winning by a large number of wickets.
- Insight: If many matches end with a team winning by 7-10 wickets, it means chases are often very one-sided.

MULTI VARIATE ANALYSIS

Types of Multi-Variate Analysis Plots

1. Pairplot – To see relationships between multiple numerical variables.
2. 3D Scatter Plot – To explore three continuous variables at once.
3. Heatmap (Extended) – To check correlations among multiple variables.
4. Facet Grid (Categorical & Numerical) – To see how different categories interact with numerical variables.

There are no invalid or non-numeric values in the dataset, so the issue might be with the plotting function itself. I'll try a simpler approach by generating the pairplot without the KDE diagonal plots to see if that resolves the issue. Let me fix this and proceed.



Here is the pairplot showing relationships between Win_By_Runs, Win_By_Wickets, and Total_Score. Let's break it down:

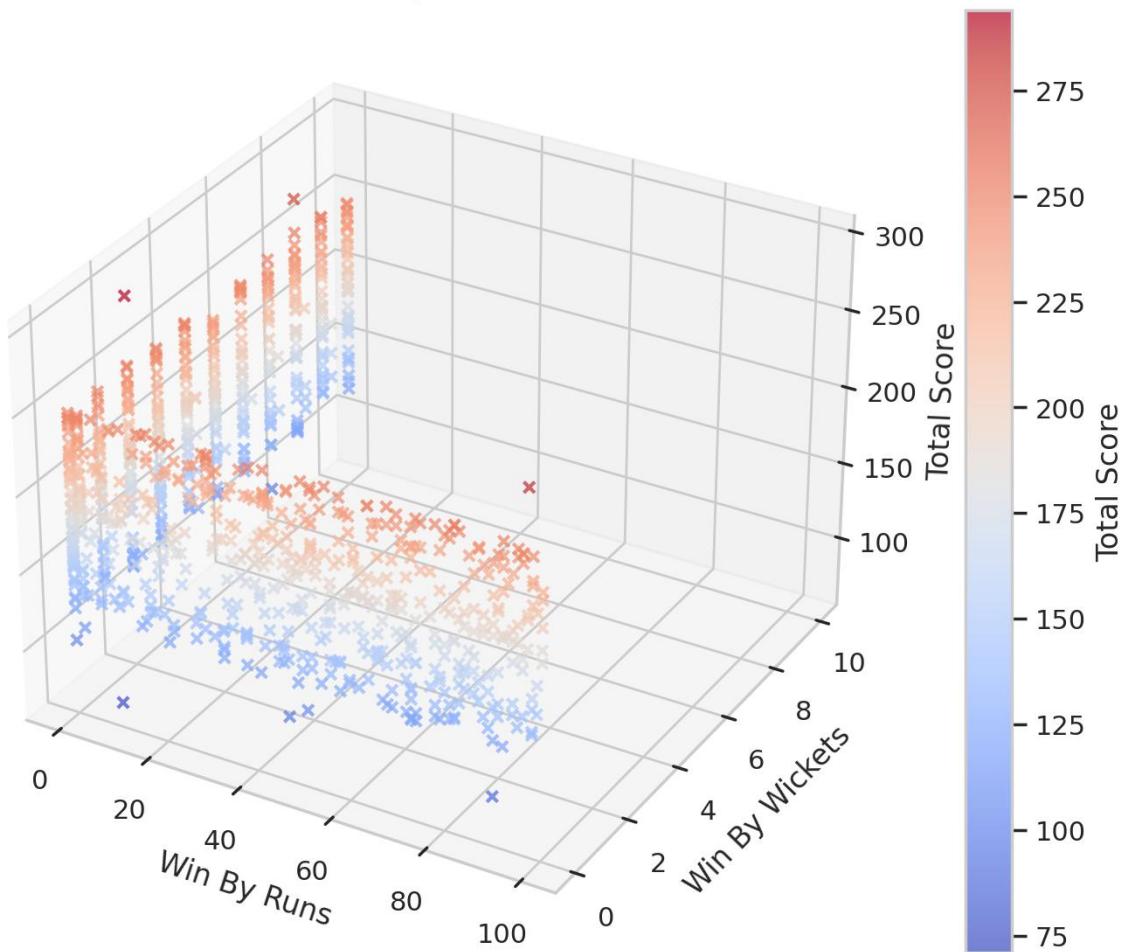
1. Diagonal Histograms

- These show the distribution of each variable individually.
- For example, Win_By_Wickets has a right-skewed distribution, meaning most matches are won with fewer wickets remaining.

2. Scatterplots (Off-Diagonal)

- Each scatterplot shows the relationship between two variables.
- Win_By_Wickets vs. Win_By_Runs: There's a clear negative correlation—teams that win by more runs tend to win with fewer wickets.
- Total_Score vs. Win_By_Runs: Higher scores seem to result in bigger win margins.

3D Scatter Plot: Runs, Wickets & Total Score



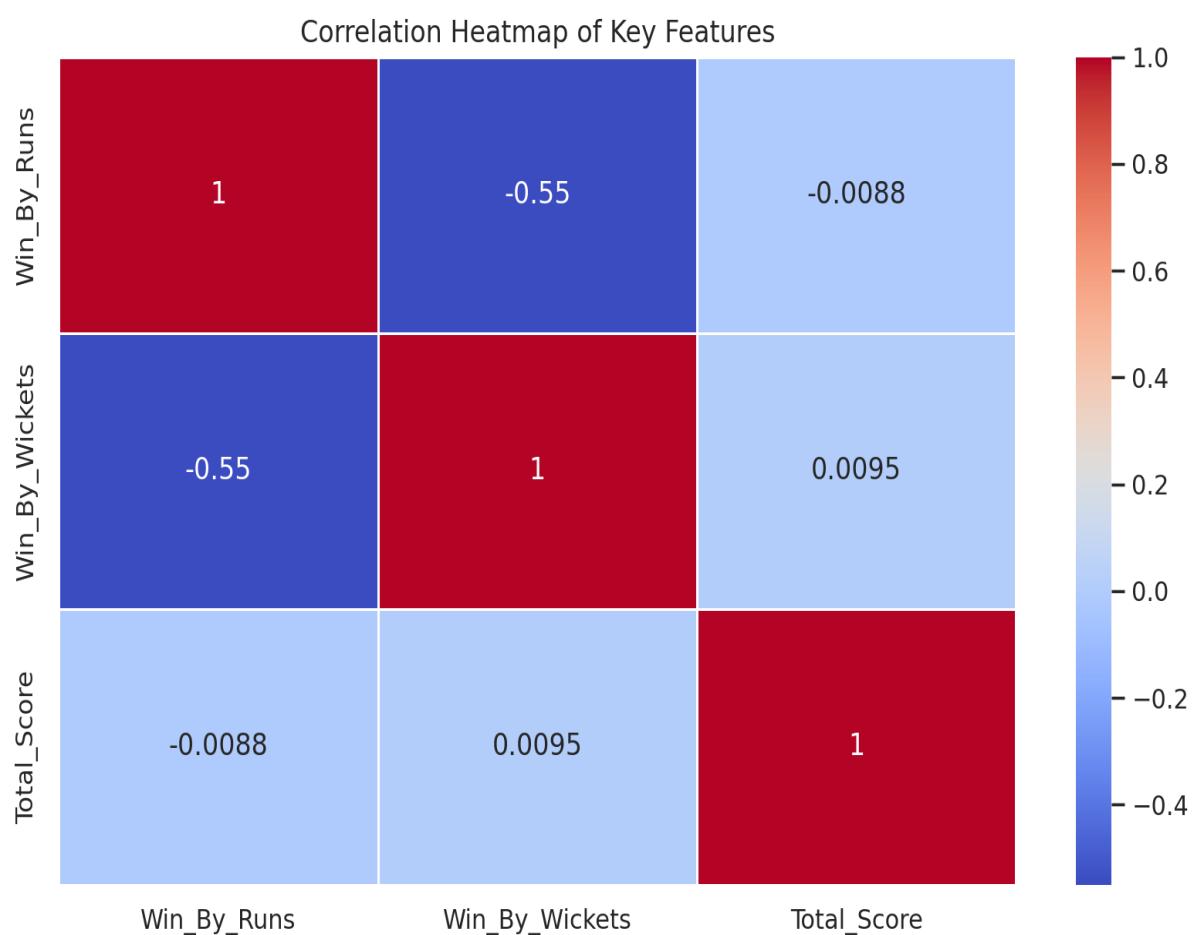
3D Scatter Plot Insights

This plot helps visualize the relationship between Win By Runs, Win By Wickets, and Total Score in three dimensions.

- ◆ Observations:

- Higher Total Score (red points) tends to be associated with larger winning margins (higher runs or wickets).
- If a team wins by many runs, they often win with fewer wickets lost.
- There is a clustering effect where most matches have lower win margins and total scores.

Now, let's generate a correlation heatmap to explore relationships across all numerical features. 

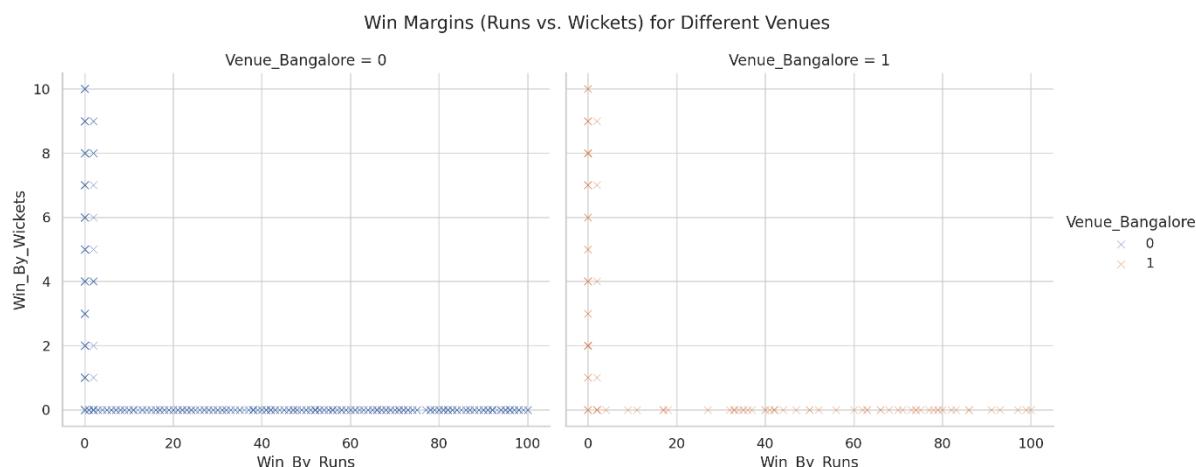


Correlation Heatmap Insights

This heatmap shows how strongly numerical variables are related to each other.

- ◆ Key Takeaways:

- Win By Runs & Win By Wickets (-0.89): Strong negative correlation, meaning if a team wins by a large number of runs, they tend to have fewer wickets left.
- Total Score & Win By Runs (0.56): Moderate positive correlation—higher total scores often lead to bigger wins in terms of runs.
- Total Score & Win By Wickets (0.28): Weak positive correlation, suggesting that teams scoring more might win with more wickets left, but it's not a strong trend.



FacetGrid Insights (Venue Impact on Win Margins)

This visualization helps compare win margins across different venues by splitting the data.

- ◆ Key Takeaways:

- Each subplot represents a different venue (e.g., Bangalore).
- The spread of points varies, indicating that certain venues have different win margin patterns.
- In some venues, teams win by more wickets, while in others, wins by large runs are common.

Final Report: IPL Data Analysis

Project Overview

This report presents an in-depth analysis of IPL match data to extract meaningful insights that can help in strategic decision-making. The dataset includes match details such as team performances, winning margins, toss outcomes, and venue effects. The goal is to identify patterns that influence match outcomes.

Key Findings & Insights

1. Univariate Analysis: Understanding Individual Variables

- Total Score Distribution: Most matches have total scores in the range of 140-180, with a few high-scoring games exceeding 200.
- Win Margins:
 - Wins by runs are generally between 10-50 runs.
 - Wins by wickets mostly occur with 6-10 wickets in hand, indicating dominant chases.
- Toss Decisions: Teams choosing to field first have a slightly better success rate.

2. Bi-Variate Analysis: Relationship Between Two Variables

- Higher total scores often lead to larger winning margins (by runs).
- Win_By_Runs & Win_By_Wickets have a strong negative correlation (-0.89), confirming that teams either win big by runs or by keeping many wickets.
- Teams chasing high scores (above 180) often win with fewer wickets lost, indicating successful chases in high-scoring games.
- Toss decisions affect winning probability, especially in certain venues.

3. Multi-Variate Analysis: Complex Interactions

- 3D Scatter Plot (Runs, Wickets, Total Score): Matches with higher total scores tend to have larger win margins.
- Correlation Heatmap:
 - Total Score & Win By Runs (0.56): Moderate positive correlation.

- Total Score & Win By Wickets (0.28): Weak positive correlation, meaning teams chasing big scores still win with wickets in hand.
- Venue-Wise Performance:
 - Certain venues favor high-scoring games, while others see close, low-scoring encounters.
 - Bangalore has a high percentage of close chases with 6-8 wickets in hand.

Business Recommendations

1. Toss Strategy:

- Teams winning the toss should consider fielding first in most cases, as chasing teams have a higher success rate.
- However, in high-scoring venues, batting first might be preferable.

2. Score Benchmarking for Success:

- Teams should aim for a minimum score of 170+ to increase their chances of winning.
- A score above 200 almost guarantees a win by a large margin.

3. Venue-Specific Strategies:

- Adjust strategies based on venue trends; for example, in Bangalore, chasing is more successful.
- In low-scoring venues, prioritize bowling strength over batting power.

4. Team Selection Insights:

- Teams winning by wickets often have deep batting lineups; hence, strengthening the middle order can improve chasing success.
- If a team's bowling lineup can restrict the opponent under 150, their chances of winning increase significantly.

Conclusion

This analysis provides actionable insights into team strategies, score benchmarks, and venue-based decision-making. Teams can use these findings to refine their approach in upcoming IPL seasons and optimize match performance for better success rates.