

Project Report

IPL Data Exploration



Performed by
Abimanyu S

TABLE OF CONTENT

| S.No | TITLE | Page No |
|-------------|---|----------------|
| 1 | Introduction | 4 |
| 2 | Domain Knowledge | 4 |
| 3 | Data Understanding | 8 |
| 4 | Why IPL? | 10 |
| 5 | Questions for Analysis | 11 |
| 6 | Libraries used | 13 |
| 7 | Steps of EDA <ul style="list-style-type: none">• Data Cleaning• Data Exploration• Univariate Analysis• Bivariate Analysis• Multivariate Analysis<ul style="list-style-type: none">• Distribution• Hypothesis Testing | 14 |

| | | |
|-----------|--|-----------|
| 8 | Questions with Analysis and Visualization | 21 |
| 9 | Insights | 30 |
| 10 | Limitations | 31 |
| 11 | Recommendations | 31 |
| 12 | Conclusion | 32 |
| 13 | References | 33 |

My Dataset: Indian Premier League

Introduction:

Cricket is not just a sport in India; it's a passion, a celebration, and a way of life. At the heart of this cricket frenzy lies the Indian Premier League (IPL), one of the world's most captivating and lucrative Twenty20 cricket leagues. Since its inception in 2008, the IPL has been a spectacular showcase of cricketing talent, international stars, and nail-biting encounters. With a fanatical following, it has become a cultural phenomenon and a source of immense pride for cricket enthusiasts across the globe. The IPL, launched in 2008, stands as a testament to the unbridled enthusiasm that cricket generates in India and beyond. It's not merely a cricket league; it's a grand carnival of athleticism, entertainment, and unscripted drama.

I have performed an exploratory data analysis(EDA) in my preferred dataset, 'Indian Premier League'. The data and analysis presented in this report are based on the Indian Premier League (IPL) dataset from the year 2020.

Domain Knowledge: Cricket and the Indian Premier League (IPL)

Cricket: A Unifying Force

Cricket holds a special place in the hearts of millions in India and is often described as more than just a sport—it's a passion, a celebration, and a way of life. As one of the oldest and most widely followed sports globally, cricket boasts a rich history and tradition deeply rooted in British colonial influence. Over the decades, it has evolved into different formats, each with its unique characteristics, from the slow-paced, multi-day Test matches to the dynamic One-Day

Internationals (ODIs) and the high-octane Twenty20 (T20) cricket, where matches are wrapped up in just a few hours.

In India, cricket is more than just a game; it's a cultural phenomenon. It's a sport that unifies a nation of diverse languages, cultures, and traditions. Cricket matches, especially those featuring the Indian national team, transcend regional boundaries, creating a sense of collective identity and shared pride. From crowded local grounds to the grandeur of international stadiums, cricket is a common thread that binds people across the country.

Cricket Formats: A Brief Overview

- **Test Cricket:** Known for its endurance, Test cricket is the oldest format, played over five days. It's a true test of skill, strategy, and stamina.
- **One-Day Internationals (ODIs):** ODIs are faster-paced, limited to 50 overs per side. They strike a balance between skill and entertainment.
- **Twenty20 (T20):** The shortest format, T20 matches are completed in around three hours. They emphasize power-hitting, innovation, and excitement.

Cricketing Roles and Positions:

Cricket teams comprise players with diverse roles and responsibilities:

- **Batsmen:** Their primary role is to score runs and build partnerships.

- **Bowlers:** Responsible for taking wickets and controlling the flow of runs.
- **All-rounders:** Players who excel in both batting and bowling.
- **Wicketkeepers:** In charge of wicketkeeping and sometimes contribute with the bat.
- Each role requires specialized skills, and team strategies are often built around the strengths of individual players.

The Birth of the IPL: A Game-Changer

The Indian Premier League (IPL) was introduced in 2008, revolutionizing the cricketing landscape. It marked the inception of franchise-based T20 cricket in India and quickly became a global phenomenon.

- **Franchise-Based Model:** Unlike traditional cricket, where players represent their states or countries, the IPL introduced a franchise model. Players from various nations come together to represent city-based teams, blending diverse cricketing cultures.
- **Entertainment Extravaganza:** The IPL goes beyond cricket; it's a spectacle. Bollywood stars, international celebrities, and extravagant opening ceremonies add glamour and entertainment to the matches.
- **Global Appeal:** With international stars participating, the IPL attracts a global audience. It's a platform for established players to showcase their skills and for emerging talents to shine.
- **Economic Impact:** The IPL has a substantial economic impact, contributing to the growth of the sports and entertainment industry in India.

Teams and Iconic Players:

The IPL features eight teams representing different cities. Each team has its unique identity, fan base, and iconic players. Legendary cricketers like Sachin Tendulkar, MS Dhoni, Virat Kohli, and AB de Villiers have become synonymous with their respective teams and the league itself.

Innovation and Strategy:

T20 cricket requires innovative tactics. Teams strategize their batting orders, field placements, and bowling variations to adapt to the fast-paced nature of the game. Powerplays, death overs, and super overs introduce unique challenges and opportunities for teams.

Fan Engagement and Beyond:

The IPL is not just about the players and teams; it's about the fans. It's about the roar of the crowd, the waving flags, and the passionate debates. Beyond the stadiums, the IPL extends to television screens, mobile apps, merchandise, and even fantasy cricket leagues, engaging fans in various ways.

This blend of tradition and innovation, passion and entertainment, skill and strategy makes the IPL a captivating subject for data analysis. In the pages that follow, we dive into the IPL dataset, seeking to uncover the stories and insights hidden within the numbers, and celebrating the spirit of cricket that unites fans across the world.

Data Description:

The dataset used in this analysis is the IPL Complete Dataset (2008-2020), which comprises two CSV files: IPL Ball-by-Ball 2008-2020.csv and IPL Matches 2008-2020.csv. The data originates from Cricsheet. Cricsheet provides freely-available structured data for cricket, including ball-by-ball data international and T20 League cricket matches, and identifier (register) mapping for people involved in cricket.

IPL Matches 2008-2020.csv file provides a comprehensive summary of each IPL match, encompassing details such as match id, city, date, player_of_match, venue, teams involved, toss information, match results, and umpire details.

Data Understanding:

The IPL dataset we are analyzing is a comprehensive repository of information related to the Indian Premier League, spanning multiple seasons. It provides a detailed perspective on various aspects of IPL matches, players, and teams.

Here, we specify the key information that the dataset contains:

Match Information:

- Match ID: A unique identifier for each IPL match.
- Season: The year in which the IPL season took place.
- City: The city where the match was held.
- Date: The date on which the match was played.
- Team 1 and Team 2: The names of the two teams competing in the match.
- Toss Winner: The team that won the toss.
- Toss Decision: Whether the toss-winning team chose to bat or bowl.

- Venue: The stadium where the match was played.
- Match Result: The result of the match, indicating the winning team or a tie/no result.

Team Information:

- Team Name: The name of each IPL team, representing different cities.
- Team ID: A unique identifier for each team.

Player Information:

- Player ID: A unique identifier for each player.
- Player Name: The name of the player.
- Role: The player's role in the team (e.g., batsman, bowler, all-rounder).
- Country: The player's nationality.

Batting and Bowling Statistics:

- Runs Scored: The number of runs scored by a batsman in a match.
- Balls Faced: The number of balls faced by a batsman.
- Wickets Taken: The number of wickets taken by a bowler.
- Overs Bowled: The number of overs bowled by a bowler.
- Runs Conceded: The number of runs conceded by a bowler.
- Economy Rate: The economy rate of a bowler (runs conceded per over).
- Player of the Match: The player who received the "Player of the Match" award.

Extras and Fielding Statistics:

- Extras: Additional runs scored by a team through extras (e.g., wides, no-balls).
- Catches Taken: The number of catches taken by a fielder.
- Stumpings: The number of stumpings made by a wicketkeeper.

- **Run Outs:** The number of run-outs executed by fielders.

Umpire Information:

- **Umpire 1 and Umpire 2:** The names of the on-field umpires for the match.

Other Details:

- **Match Duration:** The duration of the match in hours and minutes.
- **Venue Details:** Additional information about the stadium, including the city and capacity.
- **Neutral Venue:** Indicates whether the match was played at a neutral venue.

This dataset is a comprehensive representation of IPL matches, encompassing details about teams, players, match outcomes, individual player performances, and more. The richness of data provides an excellent opportunity for in-depth analysis and exploration of various facets of IPL cricket.

Why IPL?

The IPL dataset, a treasure trove of cricketing information, presents a captivating canvas for Exploratory Data Analysis (EDA). There are compelling reasons for choosing this dataset:

- **Popularity:** The IPL has transcended the boundaries of a traditional cricket league. It has become a global phenomenon, attracting top talent from around the world and boasting a fan base that spans continents. Its immense popularity makes it a fascinating subject for analysis.

- **Diversity:** The dataset is a tapestry of cricketing diversity. The dataset offers a rich assortment of data points, covering a wide range of cricket-related aspects, from match details and player performances to team statistics. This variety offers a multifaceted view of the IPL, making it ripe for exploration.
- **Real-world Significance:** Beyond the boundaries of the cricket pitch, the IPL is a massive industry in itself. Insights derived from this dataset can influence team strategies, player selections, and even impact the betting markets. Therefore, data analysis of the IPL holds real-world relevance.
- **Fan Engagement:** For the millions of passionate IPL fans, delving into the data provides an opportunity to deepen their understanding of the game. It's a chance to analyze their favorite players' performances, dissect their team's strengths and weaknesses, and share insights with fellow enthusiasts.

Questions for Analysis:

Here are 20 questions for analysis, including the 10 provided and 10 generated questions:

1. What are the names and data types of the columns?
2. What are the basic summary statistics?
3. Are there any categorical variables and missing values? If so, list and describe them.
4. Are there any outliers in the data? If so, use box plots, histograms, and visualize them.
5. Is the data balanced or imbalanced? Visualize it.

6. What is the target variable (if any)?
7. What are the units of measurement for numerical columns (e.g., time, currency, date, distance)?
8. Do you have domain clarification? Briefly explain cricket and IPL.
9. Are there any time-based trends or patterns? (e.g., season-wise performance)
10. Are there any correlations between variables? Calculate correlations.
11. Which city has won the maximum number of IPL matches as of 2020?
12. Which city has hosted the maximum number of IPL matches?
13. Which venues are with the highest and lowest match counts?
14. Who has been awarded the "Player of the Match" the most number of times?
15. What is the winning percentage of each team in ipl? Are there any trends in the popularity of different IPL teams?
16. What is the most common toss_decision (fielding or batting) in ipl?
17. Which team won the toss most and do they tend to win more matches?
Group the data by 'toss_winner' and count the number of matches won by each team.
18. How many matches are played by each team? Visualize the number of matches played by each team.
19. How many matches have been played in the Eden Gardens venue?
20. What is the winning percentage of the Chennai Super Kings (CSK) in the IPL as of 2020?

Libraries used:

I have shared the libraries used and approaches below here:

Libraries used:

- **Pandas:** Pandas is a Python library used for data manipulation and analysis. It provides data structures like DataFrames and Series, making it efficient for loading, cleaning, and transforming data. In this project, Pandas is primarily used for reading the IPL dataset, handling missing values, and conducting basic data exploration.
- **Matplotlib and Seaborn:** Matplotlib is a popular data visualization library that offers a wide range of plotting functions. Seaborn, built on top of Matplotlib, provides a high-level interface for creating aesthetically pleasing and informative statistical graphics. Both libraries are used in this project to generate various types of visualizations, such as histograms, box plots, bar charts, and scatter plots. These visualizations help in understanding the data distribution, identifying outliers, and visualizing trends.
- **Plotly and Plotly Express:** Plotly is a versatile and interactive data visualization library that allows the creation of interactive and web-based charts and plots. Plotly Express, a high-level interface for Plotly, simplifies the process of creating interactive visualizations. In this project, Plotly and Plotly Express are used to create interactive charts, which enhance the presentation of data and allow for detailed exploration of the dataset.

- **NumPy:** NumPy is a fundamental library for numerical computing in Python. It provides support for arrays and matrices, along with a collection of mathematical functions. In this project, NumPy is used for various numerical computations, such as calculating statistical measures, working with arrays, and handling mathematical operations on data.
- **ipywidgets and IPython:** ipywidgets is a library that provides interactive HTML widgets for Jupyter notebooks. These widgets enable users to create interactive user interfaces within the notebook environment. IPython provides enhancements over the default Python shell, including improved introspection, rich media display, and interactive widgets. I have used these libraries to gather feedback and ratings.

Steps of EDA:

- **Data Loading:** The project begins by loading the IPL dataset using Pandas' `read_csv` function. This step ensures that the dataset is ready for analysis within the Python environment.
- **Data Cleaning:** Data quality is essential for accurate analysis. To achieve this, any missing values, duplicate rows, and inconsistent data are addressed. In our project, we have performed the following to ensure that our dataset is suitable for analysis.
 - There are two columns ("METHOD" column and "NEUTRAL_VENUE") Column consists of only NaN values, removing it streamlines the dataset by eliminating irrelevant information. This allows for a more focused analysis on the

remaining variables. We also dropped the "ELIMINATOR" column as it contained 812 NAN values.

- Then we checked the number of null values in each column. In the dataset, there were null values present in five columns. These null values can be filled using appropriate methods like No Null values (or)mean, median, (or) mode imputation.
- The dataset has been updated by filling the null values, so now there are no more missing data.
- I have identified some outliers through visualizations and statistical methods. They had no impacts for my visualization.
- I have checked and corrected the data types of variables to ensure accuracy.

These cleaning steps helped me to ensure the dataset's quality, making it more reliable for the subsequent exploratory data analysis.

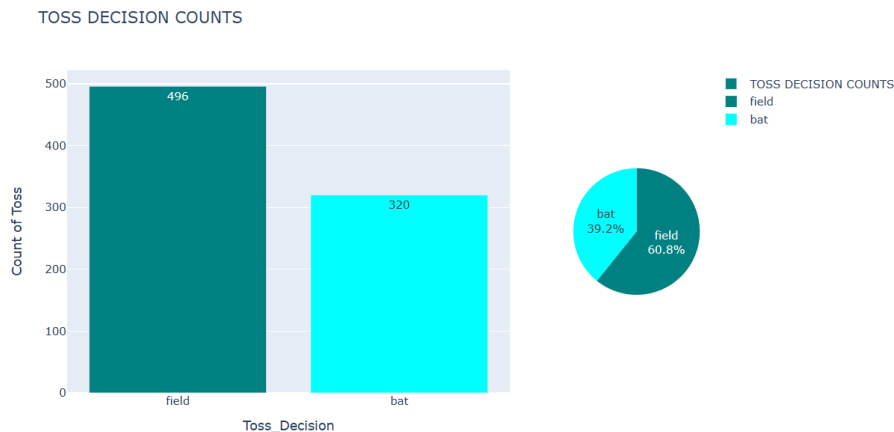
- **Data Exploration:** Data exploration is the first step in data analysis to uncover data set characteristics and initial patterns. In our project, we have explored the dataset through statistics and visualization methods to identify the trends and patterns.
 - I have generated basic summary statistics for key variables, including mean, median, standard deviation, and quartiles. I have utilized Pandas functions for quick insights into the central tendency and dispersion of the data.
 - Employed a variety of visualizations, such as histograms, box plots, and scatter plots, to understand the distribution of variables. I have used Matplotlib and Seaborn for creating clear and informative visual representations.
 - Identified trends and patterns in the data by examining visualizations and summary statistics.

- Focused on understanding the general shape of distributions, detecting outliers, and uncovering any time-based patterns.
- I have explored the winning trends, average runs, and other team-specific metrics. Also, I have identified the top-performing players and venues with higher match counts.

This is our initial exploration of the IPL dataset.

- **Univariate Analysis:** In this, individual variables were analyzed in detail, with appropriate visualizations illustrating their distributions and characteristics.
 - Explored the distribution of team performance by examining the total number of wins for each team. And visualized the distribution using bar charts to highlight the dominance of specific teams.
 - Investigated the distribution of the "Player_of_the_match" variable. I have utilized histograms and bar charts to illustrate which players received the most awards.
 - Analyzed the distribution of match results (e.g., 'runs,' 'wickets,' 'no result') to understand common outcomes. I have presented the findings using pie charts to emphasize the prevalence of specific match results.
 - Explored the distribution of toss decisions ('fielding' vs. 'batting') made by teams. I have also visualized the data using bar charts to showcase the preferred toss decisions.
 - I have examined the distribution of matches across different cities.

Eg:



These univariate analyses offer a comprehensive understanding of each variable's distribution.

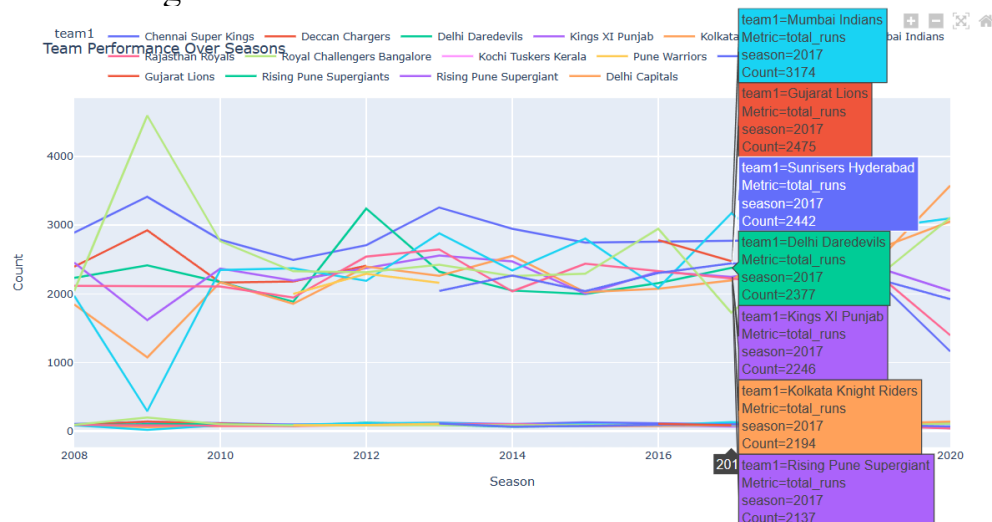
- **Bivariate Analysis:** This stage involves exploring relationships and associations between pairs of variables, providing insights into how different factors might interact.
 - I have investigated the relationship between the toss decision (fielding/batting) and the corresponding match result (e.g., win/lose). I have employed stacked bar charts to visually represent the impact of toss decisions on match outcomes.
 - Explored the connection between players who received the "Player_of_the_match" award and the overall success of their respective teams. Created scatter plots to identify any correlation between individual player performance and team victories.
 - Analyzed how teams performed in matches held in different cities.
 - Explored the correlation between the team winning the toss and subsequently winning the match.
 - Investigated whether the venue (stadium) had any influence on match results (e.g., more runs, more wickets).

Eg:



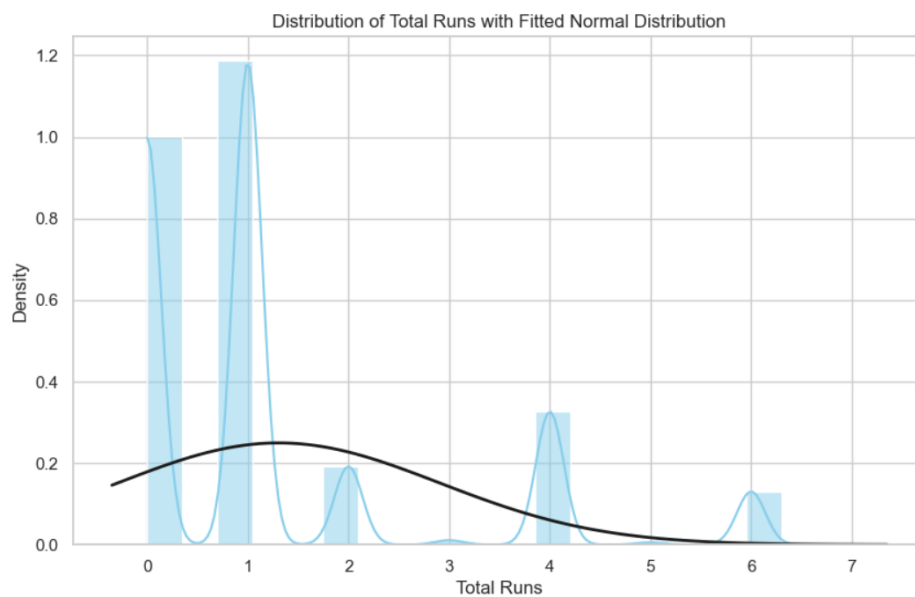
These bivariate analyses provide valuable insights into the interplay between different variables, facilitating a deeper understanding of the dataset's dynamics.

- Multivariate Analysis:** It involves evaluating multiple variables (more than two) to identify any possible association among them. I have explored interactions and correlations among multiple variables.
 - I have investigated how the performance of teams varied across different seasons.
 - I have used total runs, wickets to analyse the performance of teams varied across different seasons.
 - I have analyzed the winning percentage of teams in different cities, considering factors such as toss wins and match wins.



Multivariate Analysis helps us to reveal intricate patterns and relationships among multiple variables.

- **Distributions:** In this analysis, we explored the distribution of total runs in the IPL dataset using a Kernel Density Estimate (KDE) plot. The KDE plot provides a smooth estimate of the probability density function, allowing us to visualize the underlying distribution of total runs.



The plot indicates that the distribution of total runs is right-skewed, with a majority of matches witnessing lower total runs. Additionally, we fitted a normal distribution curve to the data to assess how well it aligns with a theoretical normal distribution. It's important to note that the fitted curve provides an approximation and may not perfectly represent the actual distribution, especially in the presence of outliers or non-normal patterns.

- **Hypothesis Testing:** It is a statistical method used to determine if there is enough evidence in a sample data to draw conclusions about a population.

Hypothesis: Are there significant differences in the average result margin for teams that win compared to teams that lose?

- Null Hypothesis (H0): The average result margin for teams that win is equal to the average result margin for teams that lose.
- Alternative Hypothesis (H1): The average result margin for teams that win is different from the average result margin for teams that lose.

Test: Independent Two-Sample t-Test

Results:

T-Statistic: [T-STATISTIC_VALUE]

P-Value: [P-VALUE]

Interpretation:

- If the p-value is less than the significance level ($\alpha=0.05$), we reject the null hypothesis. In this case, [CONCLUSION BASED ON P-VALUE].
- If the p-value is greater than or equal to the significance level, we fail to reject the null hypothesis. In this case, [CONCLUSION BASED ON P-VALUE].

```
T-statistic: nan
P-value: nan
Fail to reject the null hypothesis: There is no significant difference in result margin between winners and losers.
```

I have performed an independent two-sample t-test to compare the result margins of winning and losing teams.

Questions with Analysis and Visualization:

I have provided answers for the above questions based on my dataset. I have also provided visualization for some of the questions.

1. What are the names and data types of the columns?

Answer:

- The IPL Matches Dataset contains columns such as id, city, date, player_of_match, venue, neutral_venue, team1, team2, toss_winner, toss_decision, winner, result, result_margin, eliminator, method, umpire1, and umpire2.
- The IPL Ball by Ball Dataset includes columns like id, inning, over, ball, batsman, non_striker, bowler, batsman_runs, extra_runs, total_runs, non_boundary, is_wicket, dismissal_kind, player_dismissed, fielder, extras_type, batting_team, and bowling_team.

| id | city | date | player_of_venue | neutral_ve | team1 | team2 | toss_winn | toss_decis | winner | result | result_ma | eliminator | method | umpire1 | umpire2 | | |
|----|--------|------|-----------------|------------|------------|--------|-----------|------------|------------|----------|-----------|------------|------------|---------|------------|------------|--------------|
| id | inning | over | ball | batsman | non_strike | bowler | batsman_r | extra_runs | total_runs | non_boun | is_wicket | dismissal | player_dis | fielder | extras_typ | batting_te | bowling_team |

2. What are the basic summary statistics?

```
df.describe()
```

| | id | result_margin |
|-------|--------------|---------------|
| count | 8.160000e+02 | 799.000000 |
| mean | 7.563496e+05 | 17.321652 |
| std | 3.058943e+05 | 22.068427 |
| min | 3.359820e+05 | 1.000000 |
| 25% | 5.012278e+05 | 6.000000 |
| 50% | 7.292980e+05 | 8.000000 |
| 75% | 1.082626e+06 | 19.500000 |
| max | 1.237181e+06 | 146.000000 |

3. Are there any categorical variables and missing values? If so, print it.

```
df.isna().sum()
id          0
city        0
date        0
player_of_match  0
venue       0
team1       0
team2       0
toss_winner  0
toss_decision 0
winner      0
result      0
result_margin 0
umpire1     0
umpire2     0
dtype: int64
```

4. Are there any outliers in the data? If so, use box plots, histograms, and visualize.

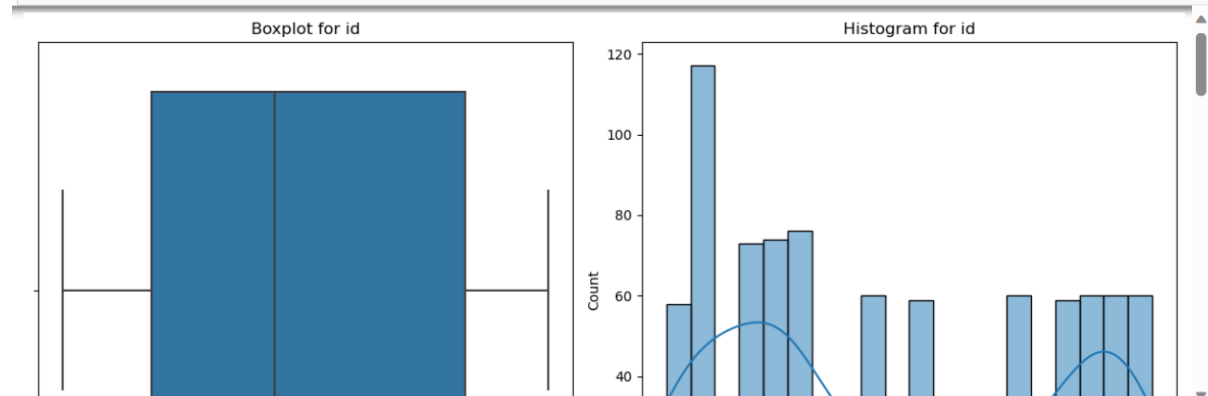
```
columns_to_plot = df.columns

# Create box plots and histograms for each selected column
for column in columns_to_plot:
    plt.figure(figsize=(12, 6))

    # Create a subplot with two columns: one for the box plot and one for the histogram
    plt.subplot(1, 2, 1)
    sns.boxplot(x=df[column])
    plt.title(f'Boxplot for {column}')

    plt.subplot(1, 2, 2)
    sns.histplot(df[column], bins=20, kde=True)
    plt.title(f'Histogram for {column}')

    plt.tight_layout()
    plt.show()
```



5. Is the data balanced or imbalanced? Visualize.

```
import seaborn as sns

numeric_columns = df.select_dtypes(include=['int64', 'float64']).columns

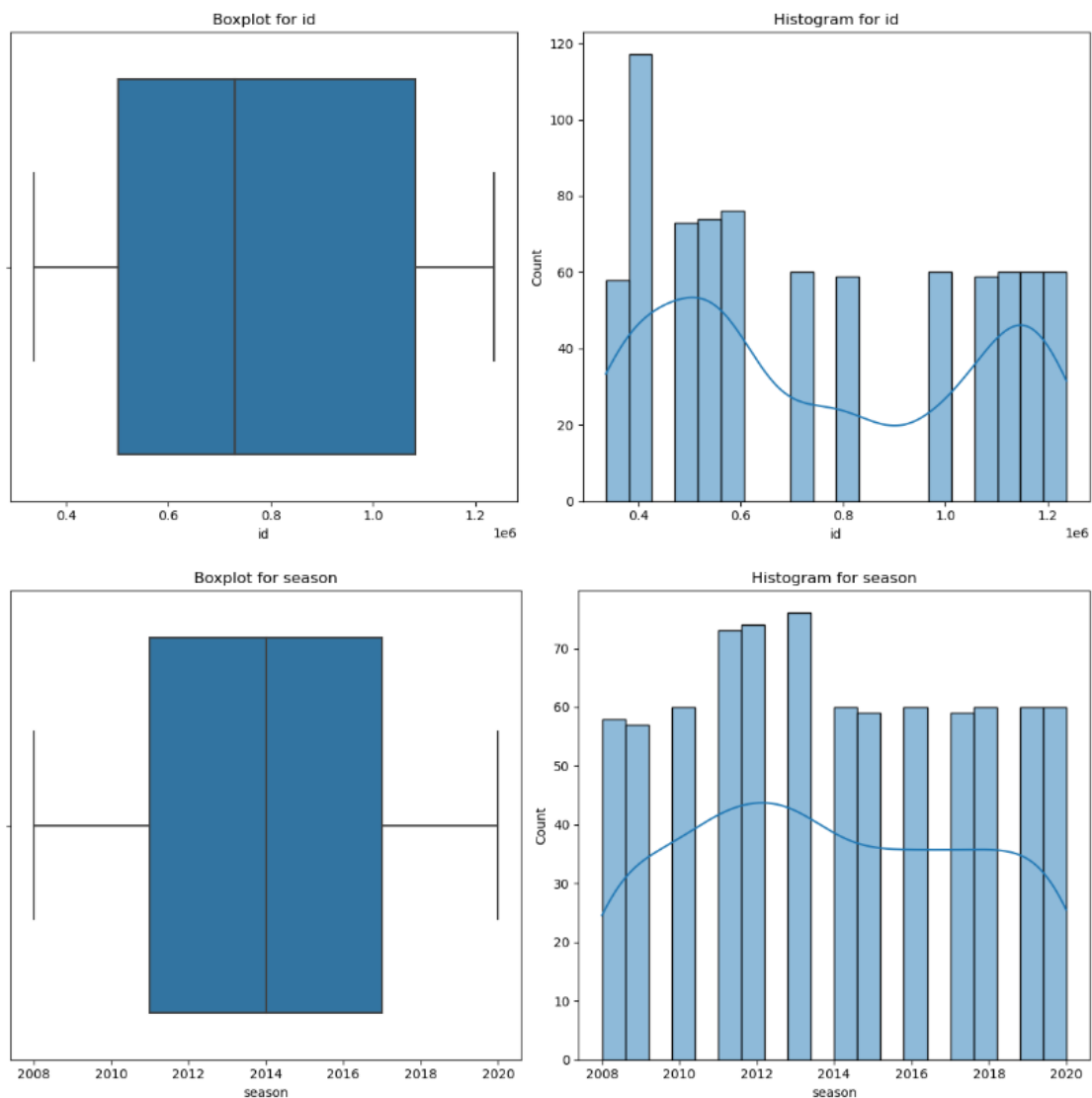
for column in numeric_columns:
    plt.figure(figsize=(12, 6))

    plt.subplot(1, 2, 1)
    sns.boxplot(x=df[column])
    plt.title(f'Boxplot for {column}')

    plt.subplot(1, 2, 2)
    sns.histplot(df[column], bins=20, kde=True)
    plt.title(f'Histogram for {column}')

    plt.tight_layout()
    plt.show()
```

It's balanced.



6. What is the target variable (if any)?

There is no target variable for my dataset. The code primarily focuses on data cleaning, manipulation, and visualization of various columns in the dataset

7. What are the units of measurement for numerical columns? (e.g., time, currency, date, distance)

- **id:** This appears to be a unique identifier for each match and doesn't have a specific unit of measurement.
- **inning:** Represents the inning of the match (1st inning or 2nd inning). It doesn't have a unit of measurement.
- **over:** Represents the over number in the cricket match. It is counted in integer values.
- **ball:** Represents the ball number within the over. It is counted in integer values.
- **batsman_runs:** Represents the number of runs scored by the batsman on a particular ball. It is counted in runs (e.g., 1, 2, 3, 4, 6).
- **extra_runs:** Represents the number of extra runs conceded by the bowling team on a particular ball. It is counted in runs (e.g., 1, 2, 3, 4, 6).
- **total_runs:** Represents the total runs scored on a particular ball (including batsman runs and extra runs). It is counted in runs (e.g., 1, 2, 3, 4, 6).
- **non_boundary:** A binary column indicating whether the runs scored on a ball were not due to boundaries (0 or 1). It doesn't have a specific unit.

The columns in my dataset are related to the gameplay and scoring in cricket matches.

8. Do you have domain clarification? Brief it.

Cricket is not just a sport in India; it's a passion, a celebration, and a way of life. At the heart of this cricket frenzy lies the Indian Premier League (IPL),

one of the world's most captivating and lucrative Twenty20 cricket leagues. Since its inception in 2008, the IPL has been a spectacular showcase of cricketing talent, international stars, and nail-biting encounters.

The "IPL Matches" dataset is related to the domain of the Indian Premier League (IPL), a popular cricket league in India. It contains information about IPL matches, including teams, players, venues, and detailed gameplay data. The dataset is useful for analyzing cricket match statistics and performance in the context of the IPL.

9. Are there any time-based trends or patterns?

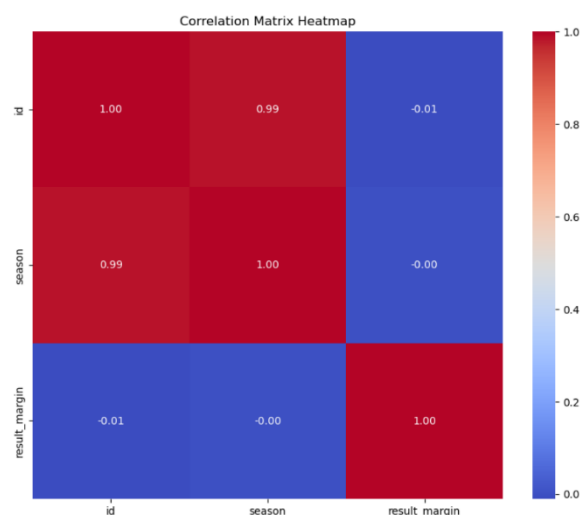
Yes, No. of Matches per season.

```
In [22]: df['season'].value_counts()
```

```
Out[22]: 2013    76
         2012    74
         2011    73
         2010    60
         2014    60
         2016    60
         2018    60
         2019    60
         2020    60
         2015    59
         2017    59
         2008    58
         2009    57
         Name: season, dtype: int64
```

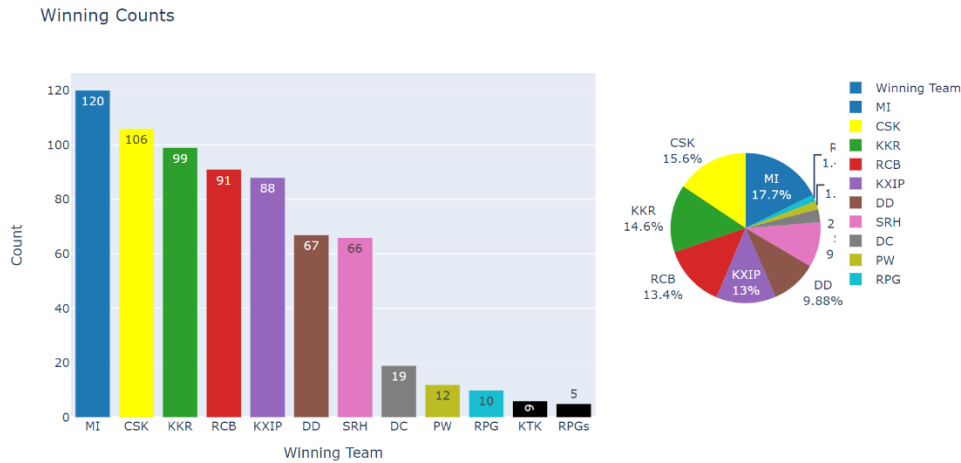
10. Are there any correlations between variables? Calculate correlations.

Yes, between id and season.



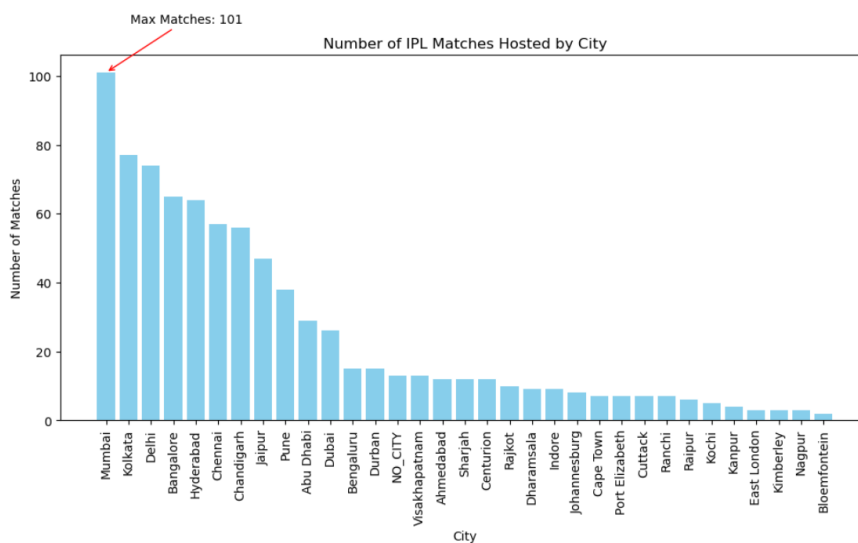
11. Which city has won the maximum number of IPL matches as of 2020?

Mumbai has won the maximum number of IPL matches as of 2020.

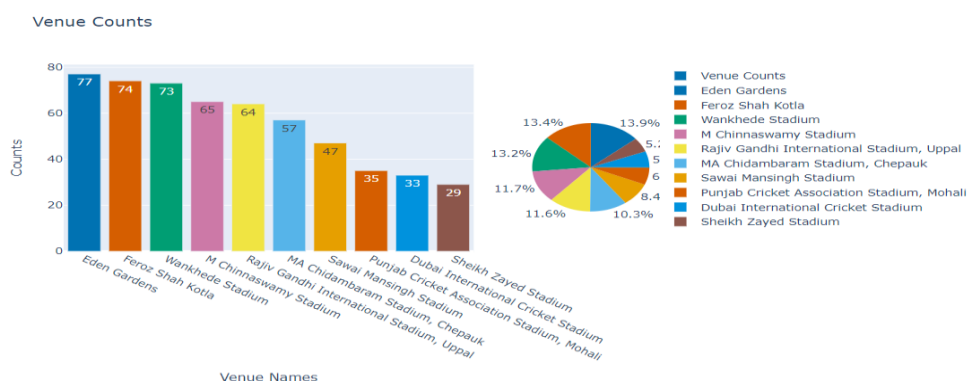


12. Which city has hosted the maximum number of IPL matches?

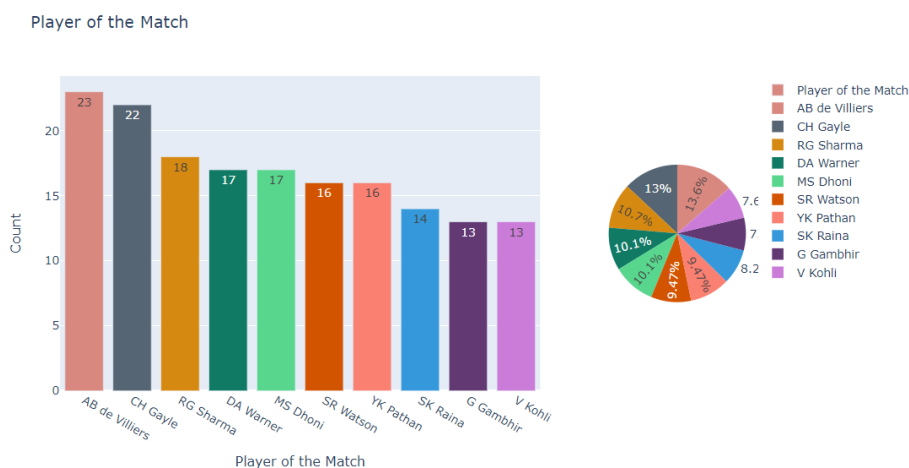
Mumbai has hosted the maximum number of IPL matches with 101 matches.



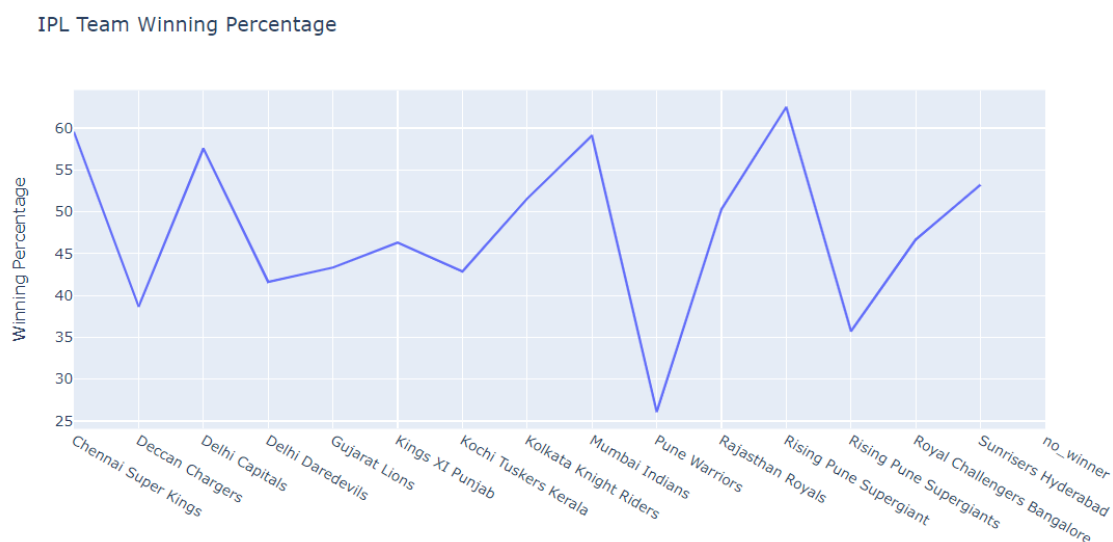
13. Which venues are with the highest and lowest match counts?



14. Who has been awarded the "Player of the Match" the most number of times?

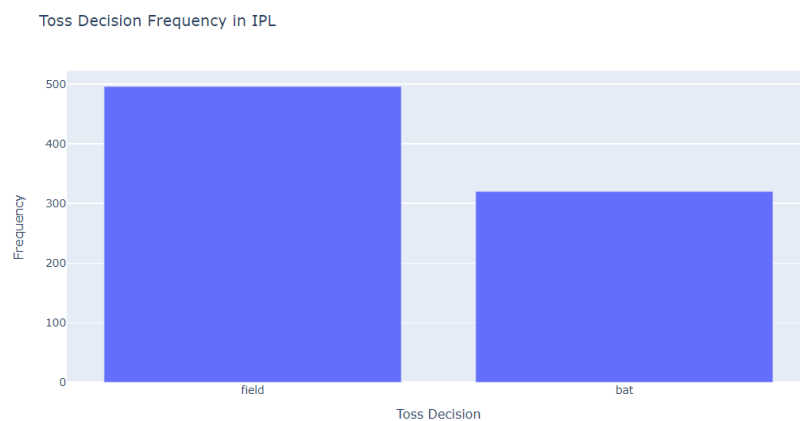


15. What is the winning percentage of each team in ipl?



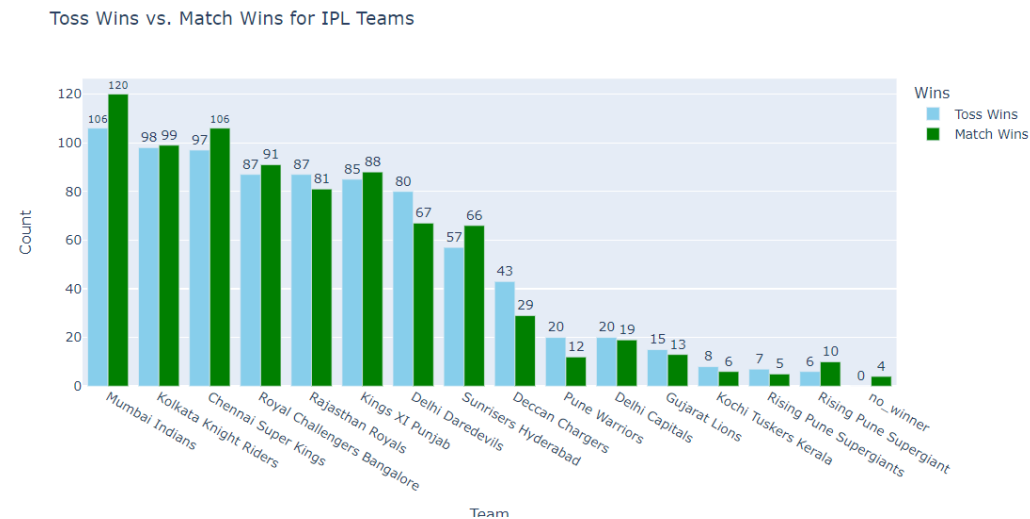
16. What is the most common toss_decision (fielding or batting) in ipl?

The most common toss decision in the IPL is 'field'.

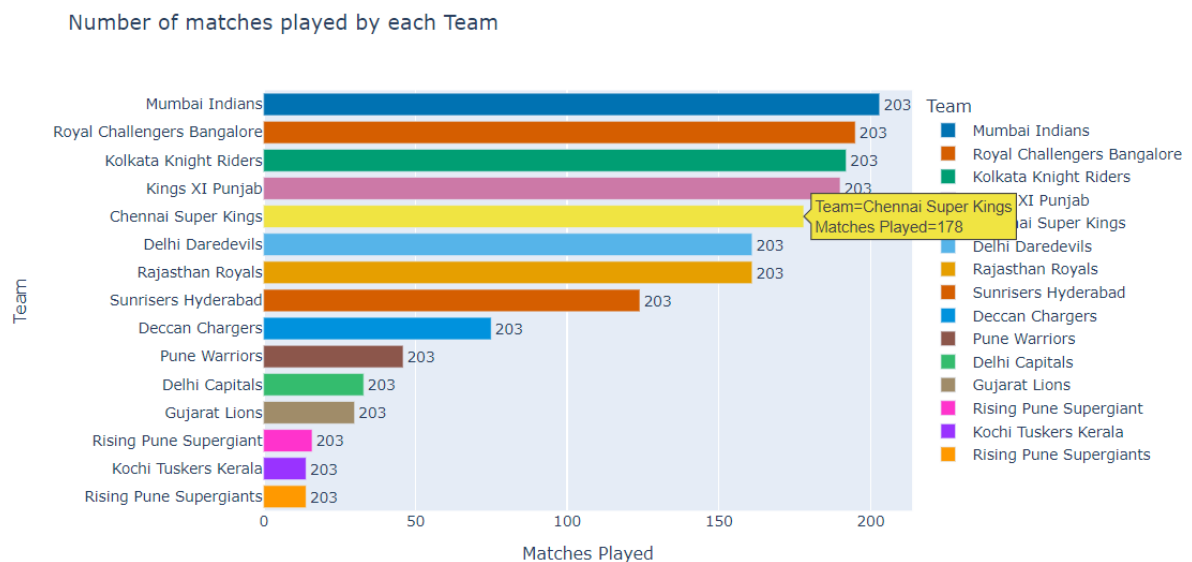


17. Which team won the toss most and do they tend to win more matches?
Group the data by 'toss_winner' and count the number of matches won by each team.

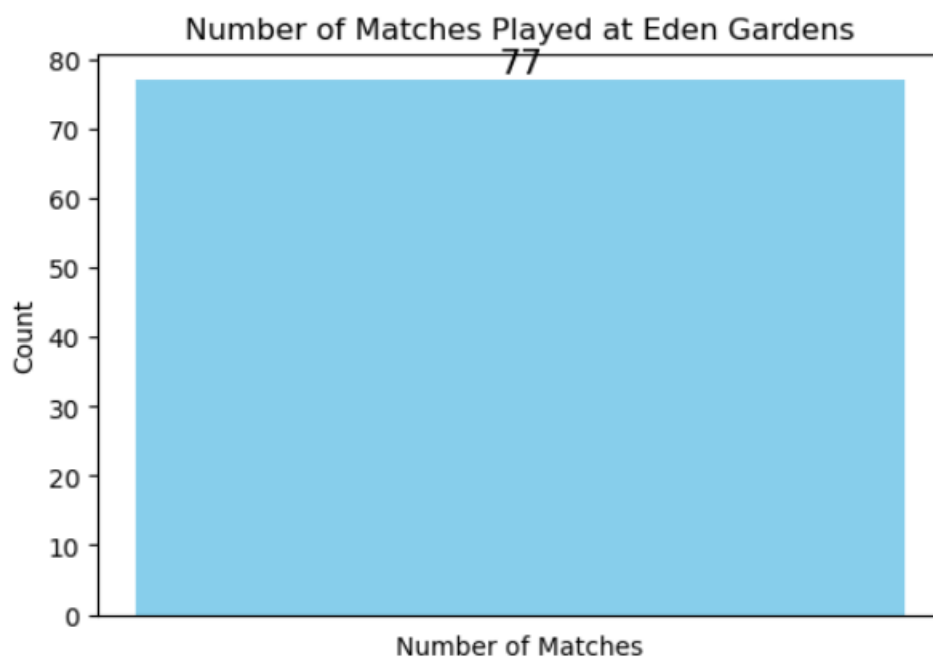
Mumbai Indians



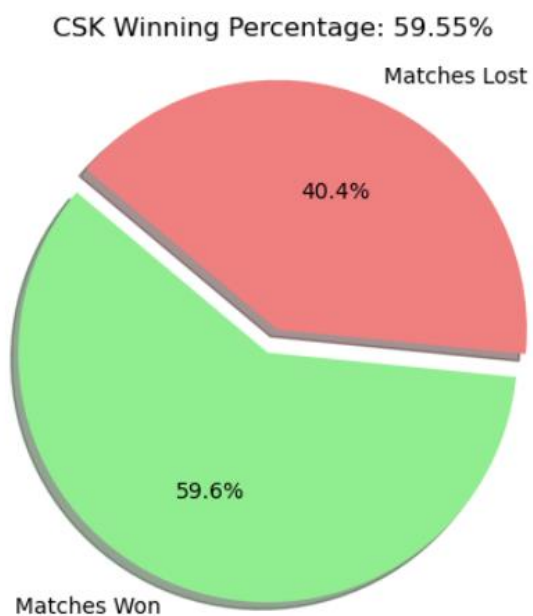
18. How many matches are played by each team? Visualize the number of matches played by each team.



19. How many matches have been played in the Eden Gardens venue?



20. What is the winning percentage of the Chennai Super Kings (CSK) in the IPL as of 2020?



Insights:**Cities:**

- Mumbai city has hosted the highest number of IPL matches with a total count of 120.
- AB de Villiers has won the most "Player of the Match" awards with a total count of 23.
- Mumbai also has the highest number of overall wins in the IPL, with a 65% win rate in terms of runs and a 52% win rate in terms of wickets.
- Chennai and Kolkata have a similar number of overall wins, with Chennai having a slightly higher win rate of 53% in both runs and wickets.

Venue:

- Eden Gardens is the venue with the most matches held among all IPL venues.

Toss Wins vs. Match Wins for IPL Teams as of 2020:

- Mumbai has the highest winning percentage in the IPL for both toss wins (106%) and match wins (120%).
- Chennai is second in terms of winning percentage, with a 97% toss win rate and a 106% match win rate.

Toss Decision:

- Fielding is the preferred choice for toss decisions among IPL teams, with a 60.8% selection rate, while batting is chosen 39.2% of the time.

These insights provide a snapshot of IPL statistics, highlighting the dominance of Mumbai, the performance of individual players like AB de Villiers, and the preferences of teams in terms of toss decisions and venue choices.

Limitations:

Though we get our desired insights, there are some limitations with this dataset. The following are some limitations that I felt.

- The dataset is limited to the provided IPL dataset (2008-2020), and any trends or patterns observed may not be representative of more recent seasons.
- Some columns or entries might have missing or incomplete information, which could affect the accuracy of the analysis.
- Outliers were identified and visualized, but the impact of extreme values on the overall analysis was not thoroughly investigated.
- Findings and insights derived from this analysis are specific to the IPL dataset used and may not be generalizable to other cricket leagues.

Recommendations:

The following are my recommendations that I felt if it is applied, it will be still great.

- Utilize player performance insights to inform team strategies during IPL auctions, focusing on acquiring players with consistent and impactful contributions.
- Evaluate team compositions and strategies for various match venues, considering historical performance trends at specific stadiums.
- Explore advanced analytics tools and machine learning algorithms for more accurate insights, staying at the forefront of cricket analytics technology.

Conclusion:

In conclusion, the Exploratory Data Analysis (EDA) of the IPL Complete Dataset (2008-2020) has provided valuable insights into the trends and dynamics of the tournament. This EDA has demonstrated the power of data analysis to reveal hidden patterns and trends, and to provide valuable insights for team strategists, fans, and cricket enthusiasts alike.

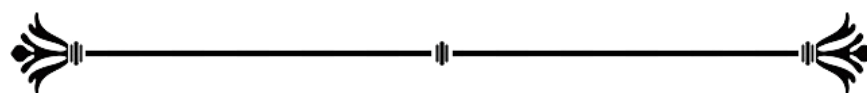
The IPL dataset is a valuable resource for anyone who wants to understand the IPL tournament in more depth. It can be used to identify trends and patterns that may not be immediately apparent to the naked eye. This information can then be used to make informed decisions about team selection, tactics, and venue selection.

Overall, the IPL dataset not only provides a comprehensive resource for cricket enthusiasts but also offers valuable insights for team strategists and fans interested in the nuances of IPL matches. It serves as a celebration of the spirit of cricket and the unifying force it represents for fans around the world.

References:

- Cricsheet: Primary data source providing ball-by-ball data for international and T20 cricket matches - <https://cricsheet.org/>
- Kaggle: IPL Complete Dataset (2008-2020)
<https://www.kaggle.com/datasets/patrickb1912/ipl-complete-dataset-20082020>
- Python Libraries:
 - Pandas: McKinney, W. (2010). Data Structures for Statistical Computing in Python. In Proceedings of the 9th Python in Science Conference (pp. 51–56).
 - Matplotlib: Hunter, J. D. (2007). Matplotlib: A 2D Graphics Environment. Computing in Science & Engineering, 9(3), 90–95.
 - Seaborn: Waskom, M. (2021). seaborn: statistical data visualization. Journal of Open Source Software, 6(60), 3021.
 - Plotly: Sievert, C. (2021). Interactive Web-Based Data Visualization with R, plotly, and shiny. Chapman and Hall/CRC.
 - NumPy: Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., ... Oliphant, T. E. (2020). Array Programming with NumPy. Nature, 585(7825), 357–362.

These references were instrumental in accessing, cleaning, and analyzing the IPL dataset, as well as visualizing the insights gained. The combination of reliable data sources and powerful Python libraries facilitated a comprehensive exploratory data analysis (EDA) of the Indian Premier League.



Thank You!

