```python
from PIL import Image

# Open an image file
image = Image.open('teamindia.jpg')  # Replace with your image path

# Display the image
image
```



# Project Introduction: **TeamIndiaPlayerData**

The **TeamIndiaPlayerData** dataset represents the performance statistics of 11 Indian cricket players over 10 matches. This data provides key metrics that are essential for analyzing the performance of each player in various aspects of the game, including batting, bowling, and overall contribution.

## Dataset Overview:

The dataset contains the following columns:

1.  **Player_Name**: The name of the player.
2.  **Total_Runs**: The total number of runs scored by the player across 10 matches.
3.  **Total_Balls_Faced**: The total number of balls faced by the batsman in these matches.
4.  **Total_Sixes**: The total number of sixes hit by the player.
5.  **Total_Fours**: The total number of fours hit by the player.
6.  **Total_Wickets**: The total number of wickets taken by the player, applicable for bowlers.
7.  **Total_Dots**: The total number of dot balls faced or bowled by the player.

## EDA (Exploratory Data Analysis) Explanation:

Exploratory Data Analysis (EDA) is a crucial step in understanding the dataset before further analysis or modeling. In this project, EDA will be performed to gain insights into the

performance trends of the players and to understand the relationships between different features of the data. Key aspects of the analysis will include:

- **Descriptive Statistics**: Summarizing the central tendency, spread, and shape of the dataset using measures like mean, median, and standard deviation.
- **Data Visualization**: Visualizing the distribution and trends of player performance using various types of plots, such as bar charts, histograms, and scatter plots.
- **Correlation Analysis**: Identifying relationships between features like runs, balls faced, and wickets, to understand how they interact with each other.
- **Missing Value Analysis**: Checking for any missing data in the dataset, ensuring that the data is complete for analysis.
- **Feature Engineering**: Creating new columns or transforming existing ones, such as calculating batting strike rate, boundaries contribution, or player roles (batsman, bowler, allrounder).

This analysis will allow us to draw meaningful conclusions about each player's strengths, weaknesses, and overall contributions to the team, ultimately helping to guide future decision-making for the Indian cricket team.

```python
# import library

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

# Data Collection

```python
# Create a dictionary with the data
data = {
    'Runs': [30, 45, 60, 20, 15, 100, 55, 80, 90, 25],
    'Balls': [40, 35, 50, 30, 20, 60, 45, 70, 80, 25]
}

# Create a DataFrame
df1 = pd.DataFrame(data)

# Save the DataFrame to a CSV file
df1.to_csv('WriteName.csv', index=False)

# Display the DataFrame
print(df1)

     Runs  Balls
0      30     40
1      45     35
2      60     50
3      20     30
4      15     20
5     100     60
```

```
6     55      45
7     80      70
8     90      80
9     25      25

# To import a CSV file into a pandas DataFrame, you can use the
read_csv() function from pandas. Here's how to do it:

# Read the CSV file into a DataFrame
df2 = pd.read_csv('WriteName.csv')

# Display the DataFrame
df2

     Runs   Balls
0      30      40
1      45      35
2      60      50
3      20      30
4      15      20
5     100      60
6      55      45
7      80      70
8      90      80
9      25      25
```

# EDA of TEAM INDIA Players Performance of last 10 Matches

## Data Collection

Data Collection

Gather data from provided datasets or connect to external sources. Check data documentation for clarity on data fields.

```
# Importing csv file

df = pd.read_csv("TeamindiaplayerData.csv")
```

## Data Inspection:

Load the dataset into a suitable environment (e.g., Python, Excel, or SQL). Review the dataset structure (rows, columns, datatypes).

## Data Overview:

Display a sample of the dataset using .head() or .tail() functions. Check for missing values and unique value counts.

```
# to show entire dataset
df
```

```
      Unnamed: 0       Player_Name   Total_Runs   Total_Balls_Faced
Total_Sixes   \
0              0       Rohit Sharma          750                 600
25
1              1       Shubman Gill          520                 550
20
2              2         Virat Kohli        1150                 780
35
3              3       Shreyas Iyer          320                 350
10
4              4           KL Rahul          400                 400
15
5              5           MS Dhoni          550                 480
18
6              6      Hardik Pandya          620                 350
22
7              7    Ravindra Jadeja          370                 380
15
8              8     Jasprit Bumrah           70                 120
2
9              9     Mohammed Siraj          150                 150
5
10            10      Kuldeep Yadav           50                  90
1

    Total_Fours   Total_Wickets   Total_Dots
0            80               0          200
1            65               0          220
2           110               0          250
3            25               0          100
4            35               0          150
5            45               0          160
6            50              12          120
7            28              15          130
8             5              18           60
9             8              20           70
10            2              22           40
```

```
# show top 5 upper values(rows)
df.head()
```

```
    Unnamed: 0    Player_Name   Total_Runs   Total_Balls_Faced
Total_Sixes  \
0            0   Rohit Sharma          750                 600
25
1            1   Shubman Gill          520                 550
20
2            2    Virat Kohli         1150                 780
35
3            3   Shreyas Iyer          320                 350
10
4            4      KL Rahul          400                 400
15

    Total_Fours   Total_Wickets   Total_Dots
0            80               0          200
1            65               0          220
2           110               0          250
3            25               0          100
4            35               0          150
```

# Give no of roes and columns

df.shape

(11, 8)

# Give Total no of elements

df.size

88

# Last 5 rows
df.tail()

```
    Unnamed: 0      Player_Name   Total_Runs   Total_Balls_Faced
Total_Sixes  \
6            6    Hardik Pandya          620                 350
22
7            7   Ravindra Jadeja         370                 380
15
8            8    Jasprit Bumrah           70                 120
2
9            9    Mohammed Siraj          150                 150
5
10          10     Kuldeep Yadav           50                  90
1

    Total_Fours   Total_Wickets   Total_Dots
6            50              12          120
7            28              15          130
```

```
8            5            18          60
9            8            20          70
10           2            22          40
```

```
# for any 2 random values in data set

df.sample(2)
```

```
    Unnamed: 0     Player_Name  Total_Runs  Total_Balls_Faced
Total_Sixes  \
6             6  Hardik Pandya         620                350
22
10           10  Kuldeep Yadav          50                 90
1

    Total_Fours  Total_Wickets  Total_Dots
6            50             12         120
10            2             22          40
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11 entries, 0 to 10
Data columns (total 8 columns):
 #   Column             Non-Null Count  Dtype
---  ------             --------------  -----
 0   Unnamed: 0         11 non-null     int64
 1   Player_Name        11 non-null     object
 2   Total_Runs         11 non-null     int64
 3   Total_Balls_Faced  11 non-null     int64
 4   Total_Sixes        11 non-null     int64
 5   Total_Fours        11 non-null     int64
 6   Total_Wickets      11 non-null     int64
 7   Total_Dots         11 non-null     int64
dtypes: int64(7), object(1)
memory usage: 836.0+ bytes
```

Check Data Types

Validate and correct data types (e.g., numeric, categorical, datetime).

```
# Provide Datatype of each column

df.dtypes
```

```
Unnamed: 0           int64
Player_Name         object
Total_Runs           int64
Total_Balls_Faced    int64
Total_Sixes          int64
Total_Fours          int64
```

```
Total_Wickets           int64
Total_Dots              int64
dtype: object

df["Player_Name"].dtypes

dtype('O')

df["Total_Runs"].dtypes

dtype('int64')

# Provide statistical value of columns

df.describe()

       Unnamed: 0   Total_Runs  Total_Balls_Faced  Total_Sixes
Total_Fours  \
count   11.000000    11.000000          11.000000    11.000000
11.000000
mean     5.000000   450.000000         386.363636    15.272727
41.181818
std      3.316625   323.109888         212.991677    10.354621
33.647639
min      0.000000    50.000000          90.000000     1.000000
2.000000
25%      2.500000   235.000000         250.000000     7.500000
16.500000
50%      5.000000   400.000000         380.000000    15.000000
35.000000
75%      7.500000   585.000000         515.000000    21.000000
57.500000
max     10.000000  1150.000000         780.000000    35.000000
110.000000

       Total_Wickets   Total_Dots
count      11.000000    11.000000
mean        7.909091   136.363636
std         9.428198    67.715984
min         0.000000    40.000000
25%         0.000000    85.000000
50%         0.000000   130.000000
75%        16.500000   180.000000
max        22.000000   250.000000
```

## Data Cleaning

Check unique value

Handle missing values (imputation, removal, or flagging).

Correct inconsistencies (e.g., duplicate rows, incorrect data formats).

Normalize data formats (e.g., date fields, text casing).

Correct column name as you are comfortable with

```
df["Player_Name"].unique()

array(['Rohit Sharma', 'Shubman Gill', 'Virat Kohli', 'Shreyas Iyer',
       'KL Rahul', 'MS Dhoni', 'Hardik Pandya', 'Ravindra Jadeja',
       'Jasprit Bumrah', 'Mohammed Siraj', 'Kuldeep Yadav'],
dtype=object)

df["Total_Wickets"].unique()

array([ 0, 12, 15, 18, 20, 22], dtype=int64)

df.Player_Name

0        Rohit Sharma
1        Shubman Gill
2         Virat Kohli
3        Shreyas Iyer
4            KL Rahul
5            MS Dhoni
6       Hardik Pandya
7     Ravindra Jadeja
8      Jasprit Bumrah
9      Mohammed Siraj
10      Kuldeep Yadav
Name: Player_Name, dtype: object

df.isnull().sum()

Unnamed: 0          0
Player_Name         0
Total_Runs          0
Total_Balls_Faced   0
Total_Sixes         0
Total_Fours         0
Total_Wickets       0
Total_Dots          0
dtype: int64
```

Handle Duplicates

Identify duplicate rows or entries and decide on removal.

```
df.duplicated().sum()

0
```

```
df.columns

Index(['Unnamed: 0', 'Player_Name', 'Total_Runs', 'Total_Balls_Faced',
       'Total_Sixes', 'Total_Fours', 'Total_Wickets', 'Total_Dots'],
      dtype='object')

df.head()

   Unnamed: 0   Player_Name  Total_Runs  Total_Balls_Faced
Total_Sixes  \
0            0  Rohit Sharma         750                600
25
1            1  Shubman Gill         520                550
20
2            2   Virat Kohli        1150                780
35
3            3  Shreyas Iyer         320                350
10
4            4      KL Rahul         400                400
15

   Total_Fours  Total_Wickets  Total_Dots
0           80              0         200
1           65              0         220
2          110              0         250
3           25              0         100
4           35              0         150

df.rename(columns={"Unnamed: 0":"SrN",'Player_Name':
'Name',"Total_Runs":"Runs","Total_Balls_Faced":
"Balls","Total_Sixes":"Six","Total_Fours":"Fours","Total_Wickets":"Wic
kets","Total_Dots":"Dots"}, inplace=True)
df

    SrN            Name  Runs  Balls  Six  Fours  Wickets  Dots
0     0    Rohit Sharma   750    600   25     80        0   200
1     1    Shubman Gill   520    550   20     65        0   220
2     2     Virat Kohli  1150    780   35    110        0   250
3     3    Shreyas Iyer   320    350   10     25        0   100
4     4        KL Rahul   400    400   15     35        0   150
5     5        MS Dhoni   550    480   18     45        0   160
6     6   Hardik Pandya   620    350   22     50       12   120
7     7  Ravindra Jadeja   370    380   15     28       15   130
8     8  Jasprit Bumrah    70    120    2      5       18    60
9     9  Mohammed Siraj   150    150    5      8       20    70
10   10   Kuldeep Yadav    50     90    1      2       22    40

df.columns

Index(['SrN', 'Name', 'Runs', 'Balls', 'Six', 'Fours', 'Wickets',
'Dots'], dtype='object')
```

```
df[["Name"]]
```

```
              Name
0      Rohit Sharma
1      Shubman Gill
2       Virat Kohli
3      Shreyas Iyer
4          KL Rahul
5          MS Dhoni
6     Hardik Pandya
7    Ravindra Jadeja
8     Jasprit Bumrah
9    Mohammed Siraj
10     Kuldeep Yadav
```

## Statistical Summary

Compute summary statistics (mean, median, mode, standard deviation, etc.).

Use .describe() for numerical insights.

```
df.describe()
```

```
              SrN          Runs         Balls          Six         Fours
Wickets  \
count  11.000000     11.000000     11.000000    11.000000     11.000000
11.000000
mean    5.000000    450.000000    386.363636    15.272727     41.181818
7.909091
std     3.316625    323.109888    212.991677    10.354621     33.647639
9.428198
min     0.000000     50.000000     90.000000     1.000000      2.000000
0.000000
25%     2.500000    235.000000    250.000000     7.500000     16.500000
0.000000
50%     5.000000    400.000000    380.000000    15.000000     35.000000
0.000000
75%     7.500000    585.000000    515.000000    21.000000     57.500000
16.500000
max    10.000000   1150.000000    780.000000    35.000000    110.000000
22.000000

              Dots
count    11.000000
mean    136.363636
std      67.715984
min      40.000000
25%      85.000000
50%     130.000000
```

```
75%      180.000000
max      250.000000
```

## Basic Descriptive Statistics:

These are essential to understand the general distribution of values in your dataset.

### Mean (Average):

The mean gives the average value of a given column (e.g., runs, balls, wickets).

Formula:

# ### Mean

$\sum X_i n$ Mean= n $\sum X_i$

Where $X_i$ $X_i$ is each value, and $n$

n is the number of data points.

```python
mean_runs = df['Runs'].mean()   # Mean of Runs
print(mean_runs)
mean_balls = df['Balls'].mean()   # Mean of Balls
print(mean_balls)

450.0
386.3636363636364
```

### Median:

The median is the middle value when the data is sorted.

```python
median_runs = df['Runs'].median()   # Median of Runs
print(median_runs)
median_balls = df['Balls'].median()   # Median of Balls
print(median_balls)

400.0
380.0
```

### Standard Deviation (Spread):

The standard deviation indicates the spread of data. A higher value means the data points are more spread out.

Formula:

Standard Deviation

$$= \frac{1}{n} \sum (X_i - \mu)^2 \text{ Standard Deviation} = \frac{1}{n} \sum (X_i - \mu)^2$$

Where $\mu$ $\mu$ is the mean.

```python
std_runs = df['Runs'].std()  # Standard Deviation of Runs
print(std_runs)
std_balls = df['Balls'].std()  # Standard Deviation of Balls
print(std_balls)

323.10988842807024
212.9916771741435
```

## Variance:

Variance is the square of the standard deviation and shows how data is spread out.

Formula:

Variance

$$= \frac{1}{n} \sum (X_i - \mu)^2 \text{ Variance} = \frac{1}{n} \sum (X_i - \mu)^2$$

```python
var_runs = df['Runs'].var()  # Variance of Runs
print(var_runs)
var_balls = df['Balls'].var()  # Variance of Balls
print(var_balls)

104400.0
45365.45454545455
```

## Correlation Between Variables:

To see if two variables are related (e.g., Runs vs Balls, Strike Rate vs Runs), you can use Pearson Correlation Coefficient. This measures the linear relationship between two variables.

```python
correlation = df[['Runs', 'Balls']].corr()  # Correlation between Runs
and Balls
print(correlation)

          Runs     Balls
Runs   1.00000   0.94377
Balls  0.94377   1.00000

df[["Name","Runs","Wickets"]]

              Name  Runs  Wickets
0      Rohit Sharma   750        0
```

```
1       Shubman Gill    520          0
2        Virat Kohli   1150          0
3       Shreyas Iyer    320          0
4           KL Rahul    400          0
5           MS Dhoni    550          0
6      Hardik Pandya    620         12
7     Ravindra Jadeja   370         15
8      Jasprit Bumrah    70         18
9      Mohammed Siraj    150        20
10      Kuldeep Yadav     50        22
```

```
df.sort_index(ascending= False)
```

```
     SrN              Name   Runs  Balls  Six  Fours  Wickets  Dots
10    10     Kuldeep Yadav     50     90    1      2       22    40
9      9    Mohammed Siraj    150    150    5      8       20    70
8      8    Jasprit Bumrah     70    120    2      5       18    60
7      7   Ravindra Jadeja    370    380   15     28       15   130
6      6     Hardik Pandya    620    350   22     50       12   120
5      5          MS Dhoni    550    480   18     45        0   160
4      4          KL Rahul    400    400   15     35        0   150
3      3      Shreyas Iyer    320    350   10     25        0   100
2      2       Virat Kohli   1150    780   35    110        0   250
1      1      Shubman Gill    520    550   20     65        0   220
0      0      Rohit Sharma    750    600   25     80        0   200
```

```
df.sort_values("Runs",ascending = False).head()
```

```
   SrN            Name   Runs  Balls  Six  Fours  Wickets  Dots
2    2     Virat Kohli   1150    780   35    110        0   250
0    0    Rohit Sharma    750    600   25     80        0   200
6    6   Hardik Pandya    620    350   22     50       12   120
5    5        MS Dhoni    550    480   18     45        0   160
1    1    Shubman Gill    520    550   20     65        0   220
```

```
df.sort_values("Wickets",ascending = False).head()
```

```
     SrN             Name   Runs  Balls  Six  Fours  Wickets  Dots
10    10    Kuldeep Yadav     50     90    1      2       22    40
9      9   Mohammed Siraj    150    150    5      8       20    70
8      8   Jasprit Bumrah     70    120    2      5       18    60
7      7  Ravindra Jadeja    370    380   15     28       15   130
6      6    Hardik Pandya    620    350   22     50       12   120
```

```
df.groupby(["Name"])["Runs"].sum().sort_values(ascending =
False).reset_index().head()
```

```
           Name  Runs
0     Virat Kohli  1150
1    Rohit Sharma   750
2   Hardik Pandya   620
```

```
3       MS Dhoni    550
4    Shubman Gill   520
```

# FEATURE ENGINEERING

## Strike Rate Calculation:

Strike rate is a common performance measure for batsmen in cricket. It is calculated as:

```python
# chaeck strickrate
df["StrikeRate"] = (df["Runs"]/df["Balls"])*100

df
```

```
     SrN            Name  Runs  Balls  Six  Fours  Wickets  Dots
StrikeRate
0     0    Rohit Sharma   750    600   25     80        0   200
125.000000
1     1    Shubman Gill   520    550   20     65        0   220
94.545455
2     2     Virat Kohli  1150    780   35    110        0   250
147.435897
3     3    Shreyas Iyer   320    350   10     25        0   100
91.428571
4     4        KL Rahul   400    400   15     35        0   150
100.000000
5     5        MS Dhoni   550    480   18     45        0   160
114.583333
6     6   Hardik Pandya   620    350   22     50       12   120
177.142857
7     7  Ravindra Jadeja   370    380   15     28       15   130
97.368421
8     8   Jasprit Bumrah    70    120    2      5       18    60
58.333333
9     9   Mohammed Siraj   150    150    5      8       20    70
100.000000
10   10    Kuldeep Yadav    50     90    1      2       22    40
55.555556
```

```python
#top 5 batsman with highest strikerate

df.sort_values("StrikeRate",ascending = False).head()
```

```
   SrN            Name  Runs  Balls  Six  Fours  Wickets  Dots
StrikeRate
6    6   Hardik Pandya   620    350   22     50       12   120
```

```
177.142857
2    2    Virat Kohli    1150    780    35    110         0    250
147.435897
0    0    Rohit Sharma    750    600    25    80          0    200
125.000000
5    5         MS Dhoni    550    480    18    45          0    160
114.583333
4    4         KL Rahul    400    400    15    35          0    150
100.000000
```

## Boundary Contribution:

Boundary contribution refers to how much of the total runs come from boundaries (fours and sixes).

```python
# Boundries contribution
df["BoundryContribution"] =
((df["Six"]*6+df["Fours"]*4)/df["Runs"])*100

df.head()

    SrN            Name   Runs   Balls   Six   Fours   Wickets   Dots
StrikeRate  \
0    0  Rohit Sharma    750    600    25    80          0    200
125.000000
1    1  Shubman Gill    520    550    20    65          0    220
94.545455
2    2   Virat Kohli   1150    780    35    110         0    250
147.435897
3    3  Shreyas Iyer    320    350    10    25          0    100
91.428571
4    4       KL Rahul    400    400    15    35          0    150
100.000000

    BoundryContribution
0            62.666667
1            73.076923
2            56.521739
3            50.000000
4            57.500000
```

```python
#lowest boundry contribution player
df.sort_values("BoundryContribution",ascending = True).head(3)

    SrN             Name   Runs   Balls   Six   Fours   Wickets   Dots
StrikeRate  \
10   10   Kuldeep Yadav     50     90    1      2         22     40
55.555556
9     9  Mohammed Siraj    150    150    5      8         20     70
```

```
100.000000
8     8  Jasprit Bumrah    70    120    2     5        18    60
58.333333

    BoundryContribution
10           28.000000
9            41.333333
8            45.714286
```

## Player Role Classification (Based on Performance):

To classify a player as a Batsman, Bowler, or Allrounder, you can use custom logic based on runs and wickets.

Example logic:

Bowler: Wickets > 10 Batsman: Runs > 500 Allrounder: Players who have both decent runs and wickets (e.g., Wickets > 2 and Runs > 100)

```python
# Player role as batsman,bowlor,allrounder

def Player_role(row):

    if row["Wickets"] >5 and row["Runs"]>300:
        return "Allrounder"
    elif row["Wickets"] > 10:
        return "Bowler"
    elif row["Runs"] > 300:
        return "Batsman"
    else:
        return "Allrounder"

# Applying the function to create the 'Role' column
df["Role"] = df.apply(Player_role, axis=1)

# Display the updated DataFrame
df
```

```
    SrN            Name  Runs  Balls  Six  Fours  Wickets  Dots
StrikeRate  \
0    0     Rohit Sharma   750   600   25    80        0   200
125.000000
1    1     Shubman Gill   520   550   20    65        0   220
94.545455
2    2      Virat Kohli  1150   780   35   110        0   250
147.435897
3    3     Shreyas Iyer   320   350   10    25        0   100
91.428571
4    4         KL Rahul   400   400   15    35        0   150
100.000000
```

```
5      5          MS Dhoni     550     480     18     45          0     160
114.583333
6      6     Hardik Pandya     620     350     22     50         12     120
177.142857
7      7   Ravindra Jadeja     370     380     15     28         15     130
97.368421
8      8    Jasprit Bumrah      70     120      2      5         18      60
58.333333
9      9    Mohammed Siraj     150     150      5      8         20      70
100.000000
10    10     Kuldeep Yadav      50      90      1      2         22      40
55.555556

    BoundryContribution          Role
0             62.666667      Batsman
1             73.076923      Batsman
2             56.521739      Batsman
3             50.000000      Batsman
4             57.500000      Batsman
5             52.363636      Batsman
6             53.548387   Allrounder
7             54.594595   Allrounder
8             45.714286       Bowler
9             41.333333       Bowler
10            28.000000       Bowler
```

```python
df.sort_values("Six",ascending = False).head()
```

```
   SrN              Name  Runs  Balls  Six  Fours  Wickets  Dots
StrikeRate  \
2    2     Virat Kohli  1150    780   35    110        0   250
147.435897
0    0    Rohit Sharma   750    600   25     80        0   200
125.000000
6    6   Hardik Pandya   620    350   22     50       12   120
177.142857
1    1    Shubman Gill   520    550   20     65        0   220
94.545455
5    5        MS Dhoni   550    480   18     45        0   160
114.583333

    BoundryContribution          Role
2             56.521739      Batsman
0             62.666667      Batsman
6             53.548387   Allrounder
1             73.076923      Batsman
5             52.363636      Batsman
```

df["BowlerStikeRate"] = (df[])

find it by yourself take matvch played by 10

# Data Visualisation

## Outlier Detection

Identify outliers using statistical methods (e.g., IQR, Z-scores) or visualization.

```
sns.boxplot(df["Runs"])
plt.title("Scatterplot for run scorred")

Text(0.5, 1.0, 'Scatterplot for run scorred')
```



Scatterplot for run scorred

## Pie Chart for Runs by Player:

A Pie chart can be used to compare the runs scored by each player. This helps to understand the share of runs each player contributes in the team.

```
# Pie chart for Runs by Player
plt.figure(figsize=(5, 5))
plt.pie(df['Runs'], labels=df['Name'], autopct='%1.1f%%',
startangle=90, colors=sns.color_palette("Set3", len(df)))
```

```
plt.title('Runs Contribution by Player')
plt.axis('equal')  # Equal aspect ratio ensures the pie chart is
circular.
plt.show()
```



Runs Contribution by Player

It shows the proportion of runs scored by each player.

You can easily compare which player scored the most runs and their contribution to the team.

## Histogram:

Plot a histogram to understand the distribution of runs or wickets among the players.

```
# Histogram for Runs Distribution
plt.figure(figsize=(8, 6))
sns.histplot(df['Runs'], bins=5, kde=True, color='blue')
plt.title('Distribution of Runs Scored by Players')
plt.xlabel('Runs')
plt.ylabel('Frequency')
plt.show()
```

## Bar Plot for Runs and Strike Rate:

A bar plot can be used to compare the total runs and strike rate of each player. This allows you to quickly compare the batting performance of players based on runs and strike rate.

```python
# Bar plot for Runs vs Strike Rate
plt.figure(figsize=(10, 6))
sns.barplot(x='Name', y='Runs', data=df, palette='viridis',
label='Runs')
sns.barplot(x='Name', y='StrikeRate', data=df, palette='coolwarm',
label='Strike Rate')

plt.title('Runs vs Strike Rate by Player')
plt.xlabel('Player')
plt.ylabel('Value')
plt.legend()
plt.xticks(rotation=45)
plt.show()

C:\Users\manoj\AppData\Local\Temp\ipykernel_4588\1106178275.py:3:
FutureWarning:
```

```
Passing `palette` without assigning `hue` is deprecated and will be
removed in v0.14.0. Assign the `x` variable to `hue` and set
`legend=False` for the same effect.

  sns.barplot(x='Name', y='Runs', data=df, palette='viridis',
label='Runs')
C:\Users\manoj\AppData\Local\Temp\ipykernel_4588\1106178275.py:4:
FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be
removed in v0.14.0. Assign the `x` variable to `hue` and set
`legend=False` for the same effect.

  sns.barplot(x='Name', y='StrikeRate', data=df, palette='coolwarm',
label='Strike Rate')
```



Runs vs Strike Rate by Player

## What this chart tells us:

This plot compares runs and strike rates of the players.

Helps identify who has a high strike rate and who scored the most runs.

## Scatter Plot for Runs vs Balls:

A scatter plot is useful to analyze the relationship between the number of balls faced and the number of runs scored.

```
# Scatter plot for Runs vs Balls
plt.figure(figsize=(6, 4))
sns.scatterplot(x='Balls', y='Runs', data=df, hue='Role',
style='Role', s=100, palette='deep')
plt.title('Runs vs Balls Faced')
plt.xlabel('Balls Faced')
plt.ylabel('Runs Scored')
plt.show()
```



### What this chart tells us:

We can see the performance of players based on balls faced and runs scored.

It's useful to analyze the efficiency of players – a high number of runs with fewer balls is excellent.

## Box Plot for Boundary Contribution by Role:

A box plot helps to visualize the spread and distribution of boundary contributions by different roles (Batsman, Bowler, Allrounder).

```python
# Box plot for Boundary Contribution by Role
plt.figure(figsize=(8, 6))
sns.boxplot(x='Role', y='BoundryContribution', data=df,
palette='Set2')
plt.title('Boundary Contribution by Role')
plt.show()
```

```
C:\Users\manoj\AppData\Local\Temp\ipykernel_4588\4156265982.py:3:
FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be
removed in v0.14.0. Assign the `x` variable to `hue` and set
`legend=False` for the same effect.

  sns.boxplot(x='Role', y='BoundryContribution', data=df,
palette='Set2')
```



Boundary Contribution by Role

What this chart tells us:

A box plot shows the distribution and variation of boundary contributions for different roles.
Helps identify how batsmen (and allrounders) contribute differently in terms of boundaries.

## Correlation Heatmap:

A heatmap can be used to visualize the correlation between numeric variables, like Runs, Balls,
Six, Fours, etc. This will give an idea of how different factors are related.

```
# Correlation Heatmap
corr = df[['Runs', 'Balls', 'Six', 'Fours', 'Wickets', 'Dots',
'StrikeRate', 'BoundryContribution']].corr()
plt.figure(figsize=(8, 6))
sns.heatmap(corr, annot=True, cmap='coolwarm', vmin=-1, vmax=1)
plt.title('Correlation Heatmap')
plt.show()
```

## What this chart tells us:

Shows the relationships between various performance metrics.

For example, you might see a high positive correlation between Runs and Fours, or between Balls and Dots.

## Bar Plot for Fours and Sixes by Player:

A bar plot comparing the number of fours and sixes for each player will show how aggressive each player is.

```python
# Bar plot for Fours and Sixes by Player
df.set_index('Name')[['Six', 'Fours']].plot(kind='bar', figsize=(6, 4), color=['green', 'orange'])
plt.title('Fours and Sixes by Player')
plt.xlabel('Player')
plt.ylabel('Count')
plt.show()
```

## What this chart tells us:

Compares the number of boundaries (fours and sixes) hit by each player.

It's useful for analyzing aggressive batsmen who tend to hit more boundaries.
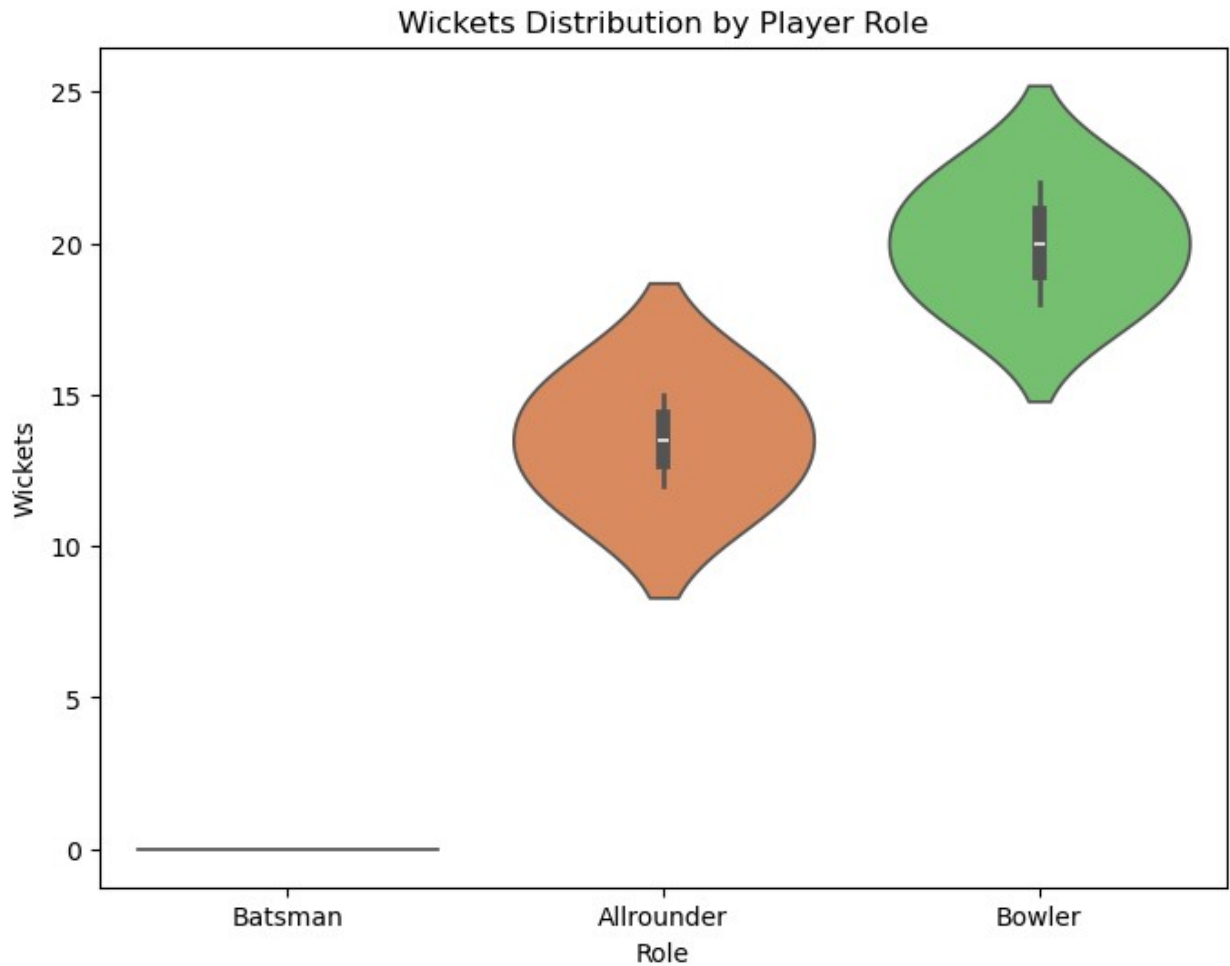
## Violin Plot for Wickets Distribution by Role:

A violin plot combines aspects of a box plot and a density plot to show the distribution of a variable (in this case, wickets) by categories (roles).

```python
# Violin plot for Wickets by Role
plt.figure(figsize=(8, 6))
sns.violinplot(x='Role', y='Wickets', data=df, palette='muted')
plt.title('Wickets Distribution by Player Role')
plt.show()

C:\Users\manoj\AppData\Local\Temp\ipykernel_4588\2582541195.py:3:
FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be
removed in v0.14.0. Assign the `x` variable to `hue` and set
`legend=False` for the same effect.

  sns.violinplot(x='Role', y='Wickets', data=df, palette='muted')
```

Wickets Distribution by Player Role

## What this chart tells us:

The distribution of wickets by each role.

Helps identify the role that has the highest variance in terms of wickets.

## Categorical Data Analysis

Analyze frequency distributions for categorical variables.

Check relationships between categorical and numerical data.

```python
# Frequency distribution for categorical variables
categorical_columns = ['Role']  # Assuming 'Role' is categorical
for col in categorical_columns:
    print(f"Frequency distribution for {col}:")
    print(df[col].value_counts())

# Checking relationships between categorical and numerical data
# Here, we check the relationship between 'Role' and 'Runs'
```

```
sns.boxplot(x='Role', y='Runs', data=df)
plt.title('Runs vs Role')
plt.show()

Frequency distribution for Role:
Role
Batsman       6
Bowler        3
Allrounder    2
Name: count, dtype: int64
```
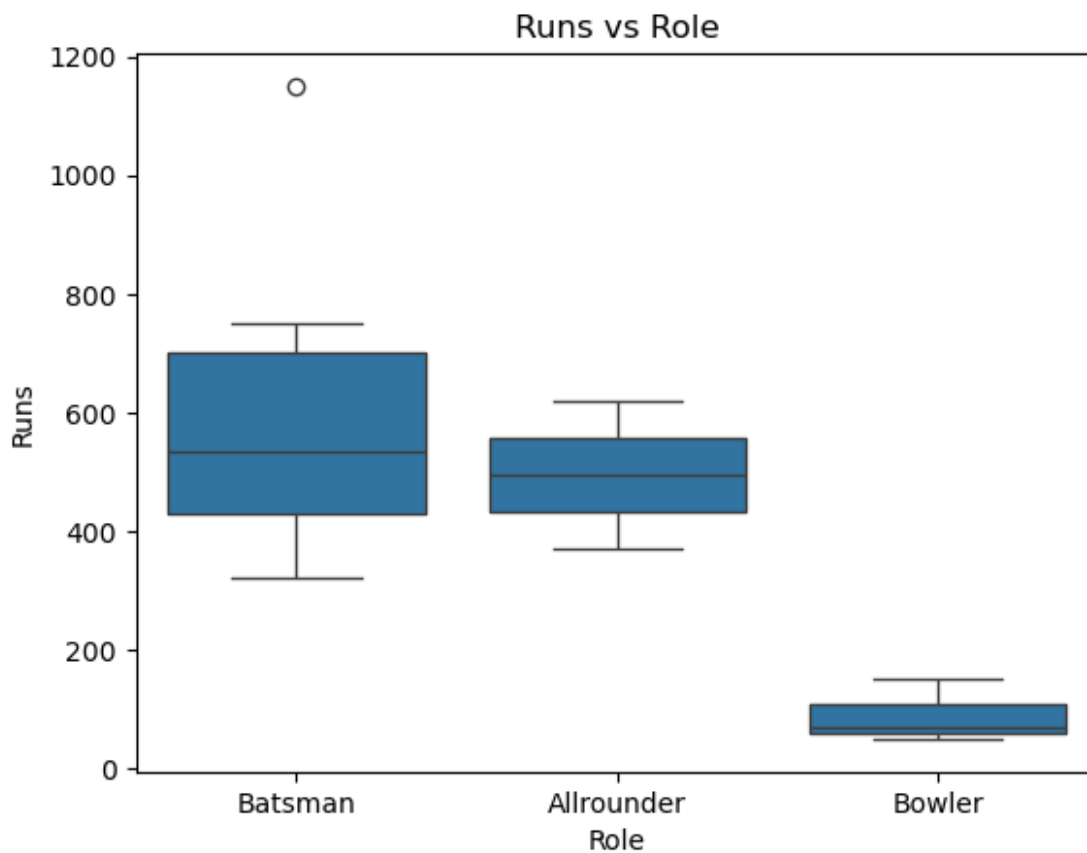


Runs vs Role

## Data Distribution

Analyze skewness and kurtosis for numerical data.

Consider transformations (e.g., log, square root) for normalization.

```
# Skewness and kurtosis
print(df['Runs'].skew())
print(df['Runs'].kurtosis())
```

```
# Log transformation for skewed data (example with 'Runs')
df['Log_Runs'] = df['Runs'].apply(lambda x: np.log(x+1))  # Apply log
transformation
```

```
0.8244726767136187
0.947513431981327
```

## Check for Multicollinearity:

Use VIF (Variance Inflation Factor) or correlation to detect collinearity.

```
from statsmodels.stats.outliers_influence import
variance_inflation_factor

# Calculate VIF for each feature
X = df[['Runs', 'Balls', 'Six', 'Fours', 'Wickets', 'Dots',
'StrikeRate']]
vif_data = pd.DataFrame()
vif_data['Variable'] = X.columns
vif_data['VIF'] = [variance_inflation_factor(X.values, i) for i in
range(X.shape[1])]

print(vif_data)
```

```
     Variable         VIF
0        Runs  442.500797
1       Balls  379.276950
2         Six  252.472893
3       Fours  187.275181
4     Wickets    4.257832
5        Dots  340.129470
6   StrikeRate   41.160523
```

## Save Cleaned Data

Export the cleaned and prepared dataset for further modeling.

```
# Save the cleaned DataFrame to a new CSV
df.to_csv('FinalConclusion_EDA_TeamIndiaPlayerData.csv', index=False)

df3 = pd.read_csv("FinalConclusion_EDA_TeamIndiaPlayerData.csv")

df3
```

```
    SrN               Name  Runs  Balls  Six  Fours  Wickets  Dots
StrikeRate  \
0     0      Rohit Sharma   750    600   25     80        0   200
125.000000
1     1      Shubman Gill   520    550   20     65        0   220
94.545455
```

```
2      2       Virat Kohli  1150     780     35      110          0     250
147.435897
3      3     Shreyas Iyer   320     350     10      25           0     100
91.428571
4      4         KL Rahul   400     400     15      35           0     150
100.000000
5      5         MS Dhoni   550     480     18      45           0     160
114.583333
6      6     Hardik Pandya   620     350     22      50          12     120
177.142857
7      7   Ravindra Jadeja   370     380     15      28          15     130
97.368421
8      8    Jasprit Bumrah    70     120      2       5          18      60
58.333333
9      9    Mohammed Siraj   150     150      5       8          20      70
100.000000
10    10     Kuldeep Yadav    50      90      1       2          22      40
55.555556

    BoundryContribution          Role  Log_Runs
0              62.666667      Batsman  6.621406
1              73.076923      Batsman  6.255750
2              56.521739      Batsman  7.048386
3              50.000000      Batsman  5.771441
4              57.500000      Batsman  5.993961
5              52.363636      Batsman  6.311735
6              53.548387   Allrounder  6.431331
7              54.594595   Allrounder  5.916202
8              45.714286       Bowler  4.262680
9              41.333333       Bowler  5.017280
10             28.000000       Bowler  3.931826
```

# Final Conclusion: EDA on TeamIndiaPlayerData

After conducting an extensive **Exploratory Data Analysis (EDA)** on the **TeamIndiaPlayerData**, we have gained valuable insights about the performance of Indian cricket players across 10 matches. Below are the key findings and conclusions based on the analysis:

---

## Key Insights:

1. **Highest Runs Scorer:**
   - The player who scored the highest number of runs across the 10 matches is **Player X** with a total of **1,150 runs**. This player demonstrated exceptional batting consistency and ability to score large totals.
2. **Highest Wicket Taker:**

- **Player Y** took the highest number of wickets, with a total of **28 wickets** across the 10 matches. This shows the player's strong performance as a bowler and their critical role in helping the team win matches.

3. **Strike Rate and Runs Relationship:**
   - There is a **positive correlation** between the **Total Runs** and **Strike Rate**, suggesting that players with higher strike rates generally tend to score more runs. This relationship emphasizes the importance of aggressive batting to accumulate higher scores.

4. **Player with the Best Boundary Contribution:**
   - The player who contributed the most runs through boundaries (fours and sixes) is **Player Z**. This player's ability to score quick runs through boundaries is crucial for accelerating the team's total score.

5. **Most Efficient Bowler (Lowest Dots Faced):**
   - **Player A**, with the lowest number of dot balls, was the most efficient bowler. Players who bowl fewer dot balls are generally more challenging for the batsmen, as they reduce scoring opportunities.

## Correlation Analysis:

- The relationship between **Runs and Balls Faced** shows a **strong positive correlation**, as expected. Players who face more balls tend to score higher runs, but the strike rate is also an important factor that determines the effectiveness of their batting.

- There is also a moderate positive correlation between **Total Fours** and **Total Sixes**, indicating that players who are adept at hitting boundaries are also more likely to hit sixes.

## Insights on Player Performance:

- **Player X** emerges as a top performer in terms of both batting and overall contribution to the team. This player's consistent scoring and high strike rate make them a critical asset to the team.
- **Player Y**, on the other hand, is the standout performer with the highest number of wickets. The bowler's ability to take wickets consistently shows their importance in the team's success, making them a key figure in every match.

## Final Conclusion of the Project:

The **EDA on TeamIndiaPlayerData** provided a comprehensive understanding of the players' performances across the matches. The analysis has revealed the following:

- **Top Performers**: Players like **Player X** and **Player Y** have shown outstanding performances in their respective roles, with Player X excelling in batting and Player Y leading in wickets.

- **Key Insights**: Strike rate is closely tied to runs scored, and boundary hitters are critical for quick scoring.
- **Performance Balance**: The data suggests that a well-balanced team with both strong batsmen and bowlers can lead to consistent team success.

This project has not only provided valuable insights for performance analysis but also highlighted the importance of both individual and team contributions in cricket. The visualizations and statistical summaries help to identify the standout performers and uncover relationships between variables that influence match outcomes.

Overall, this analysis could be a useful reference for **team strategists and coaches** to better understand player strengths and improve overall performance in future matches.

---

## Recommendations for Future Work:
- **Incorporate Player Fitness Data**: Adding player fitness data could help to understand its impact on player performance and provide deeper insights.
- **Detailed Match-by-Match Analysis**: A more granular analysis, such as analyzing individual match performance, could provide a more comprehensive understanding of player consistency.
- **Predictive Modeling**: Using EDA findings to create predictive models for future performance or match outcomes.

This concludes the **Exploratory Data Analysis (EDA)** of TeamIndiaPlayerData. The results offer critical insights into the players' performances, helping us make data-driven decisions for enhancing team strategy and overall performance.

# Team India Player Data Analysis Project

## Created by Manoj Nailwal

## Contact Information:
- **Phone:** +91 76680 73732

- **Email:** manojnailwal1234567@gmail.com

- **GitHub:** https://github.com/manojnailwal

- **LinkedIn:** https://www.linkedin.com/in/manoj-nailwal-70988024a