

IPL DATA ANALYSIS



GALGOTIAS
UNIVERSITY



IPL

AKASH: 23SCSE1011203

ONKAR DUREJA: 23SCSE1180036

SANDEEP YADAV: 23SCSE1010514

YASH KUMAR GUPTA: 23SCSE1180115

INTRODUCTION

The Indian Premier League (IPL) is one of the most popular and competitive cricket leagues globally. With rich historical data and numerous matches played over the years, it provides a great opportunity for data analysis. This project involves analyzing IPL data to uncover trends, performance metrics, and winning strategies using data science techniques.

This project explores IPL match and delivery data using Python tools. We perform exploratory data analysis (EDA), visualize key metrics, and identify patterns among teams and players. The goal is to derive meaningful insights that help understand the game better from a data perspective.

PROBLEM STATEMENT

How can historical IPL data be analyzed to:

- Understand team and player performances?
- Discover the impact of toss decisions and venues?
- Identify factors contributing to winning matches?

By answering these questions, we aim to turn raw data into actionable insights.



Why We choose this ?

Cricket is a passion for millions, and IPL data is widely available and rich in patterns. Being a cricket enthusiast and aspiring data analyst, we chose this project to combine both interests. It also serves as an ideal case for applying and improving our data analysis and visualization skills.

Project Scope

1. Analyze IPL matches from 2008 to recent seasons.
2. Focus on team performance, player stats, and match outcomes.
3. Visualize findings using charts and graphs.
4. Scope does not include prediction or live data analysis (but can be extended in the future).


```

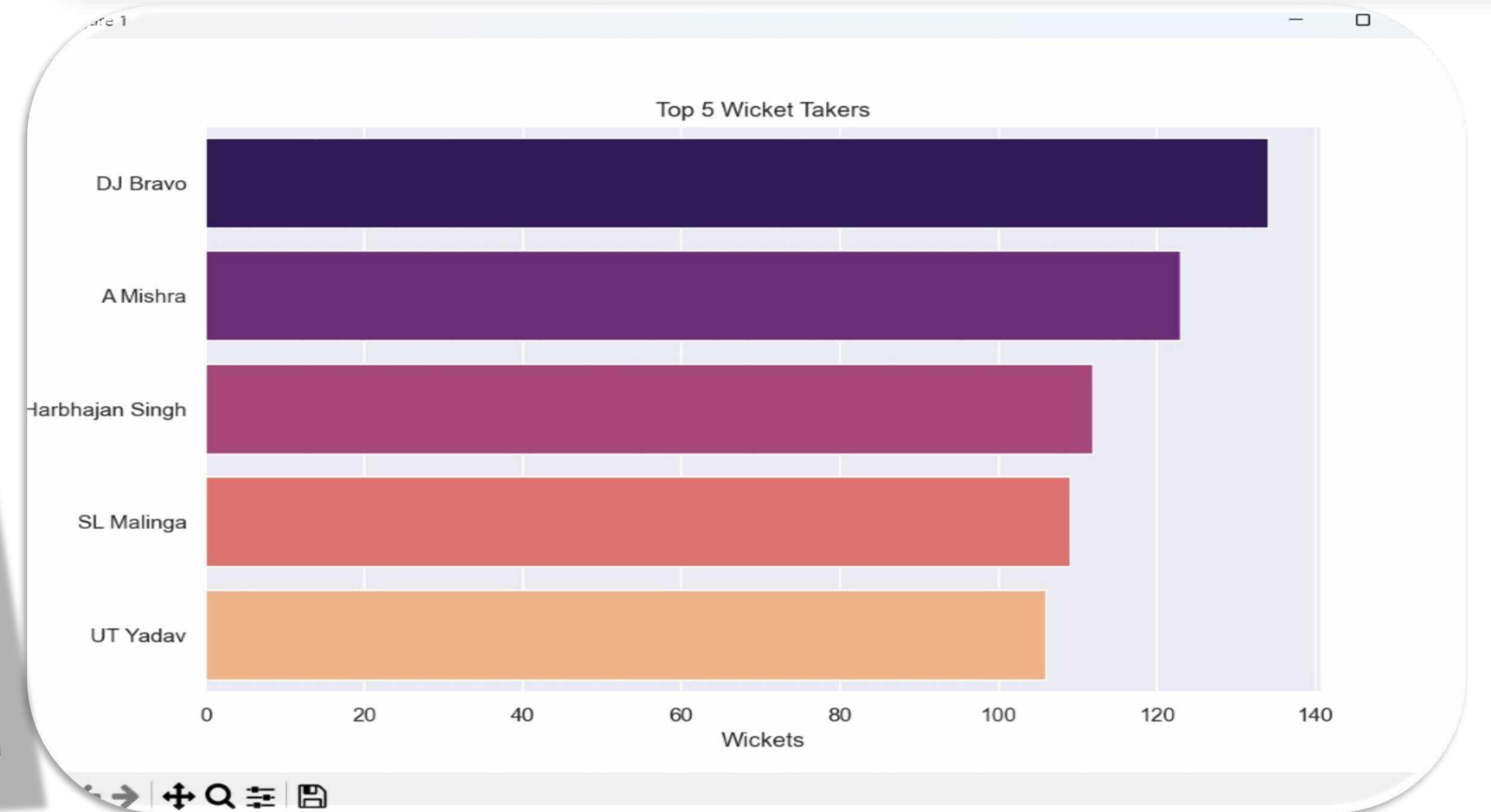
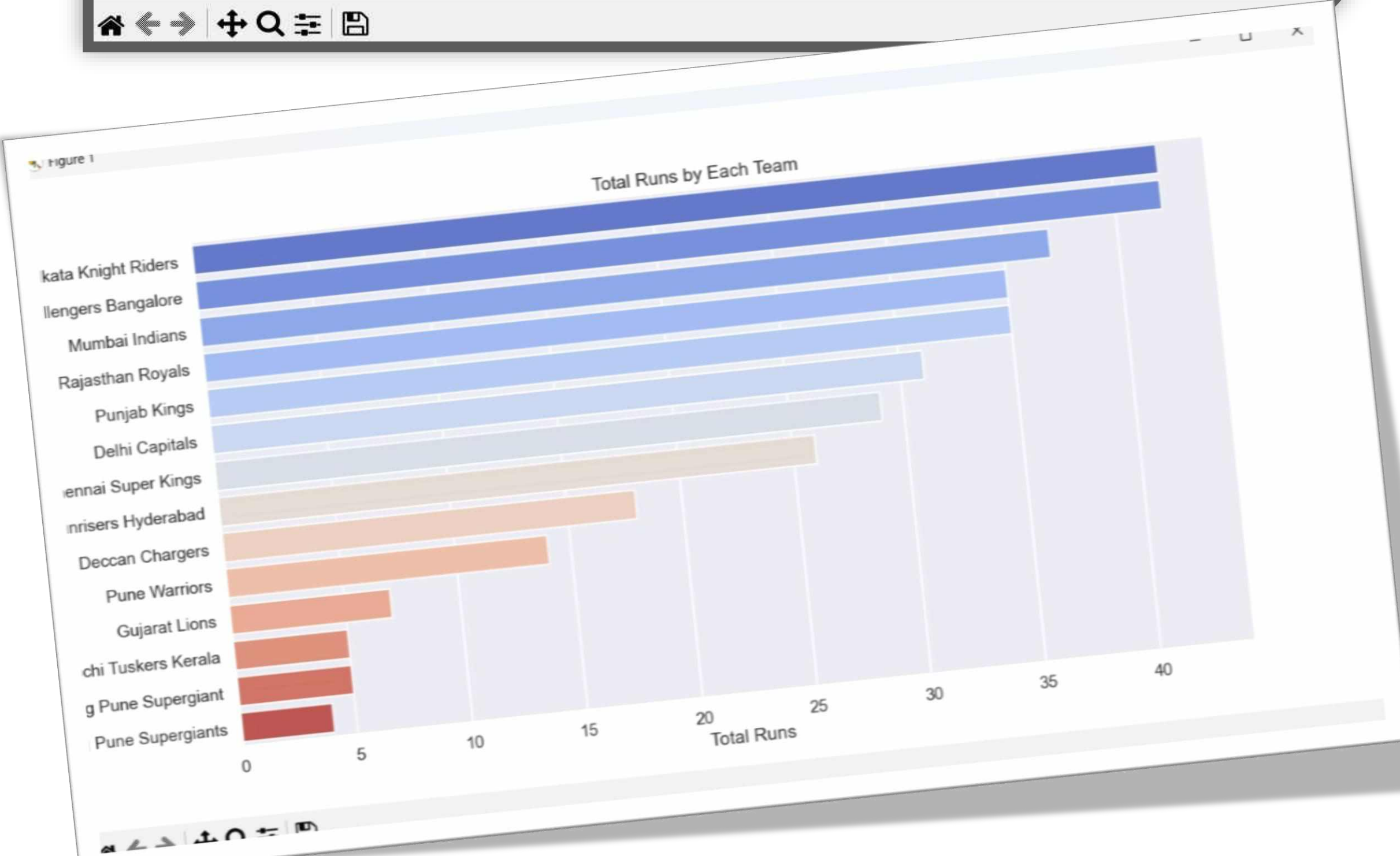
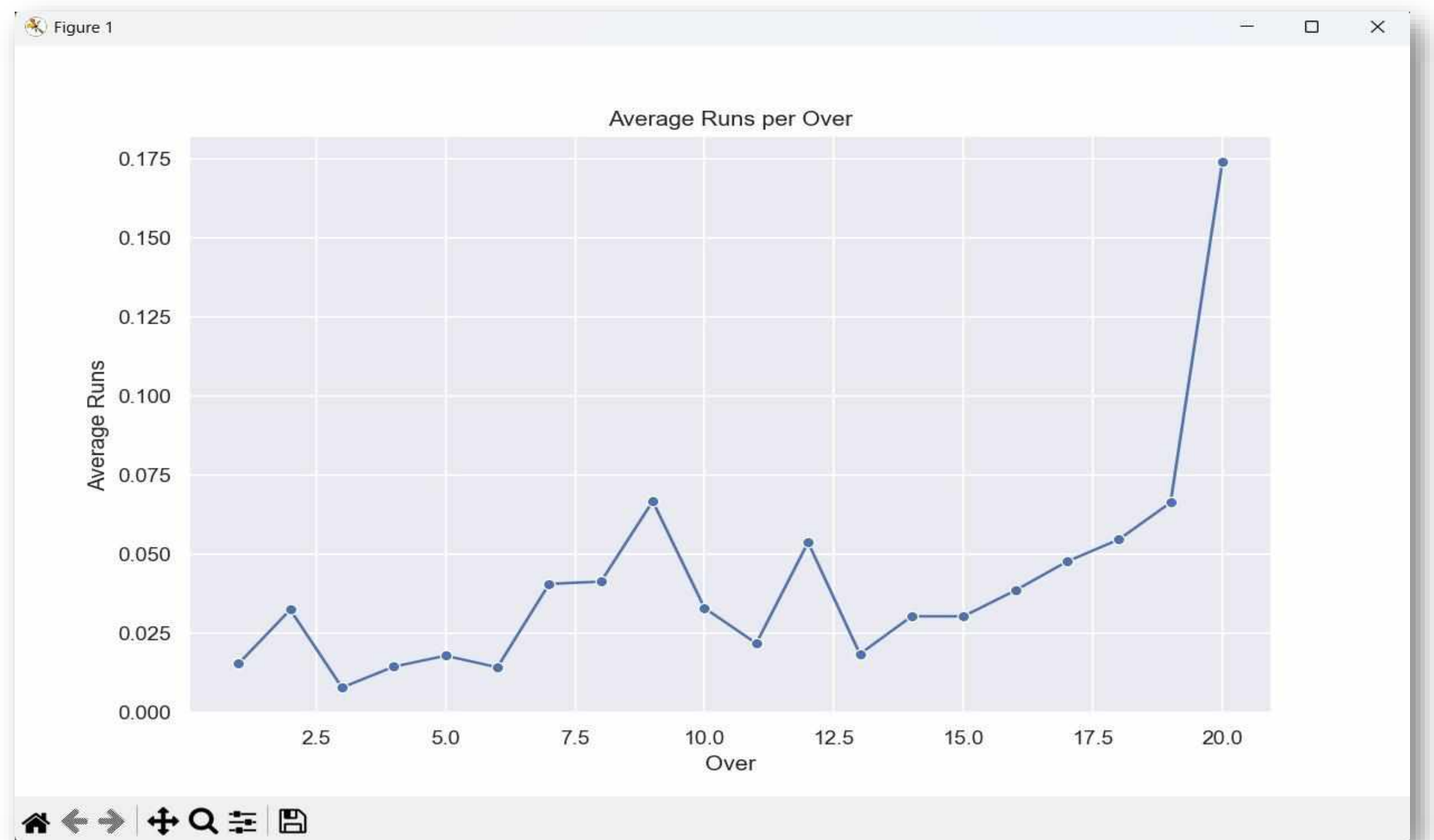
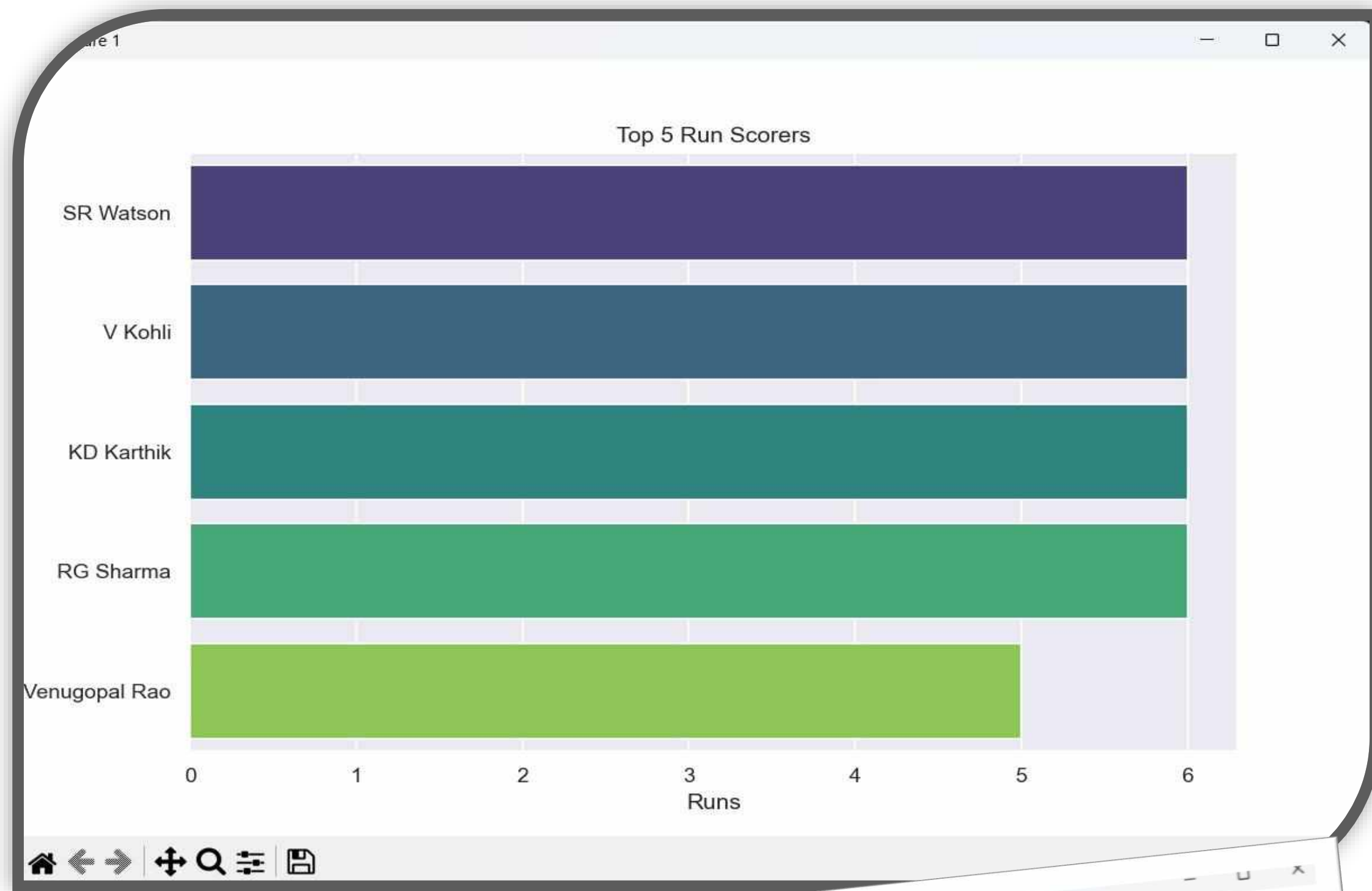
DA.py X
DA.py > ...
1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import plotly.express as px
7 import warnings
8 warnings.filterwarnings('ignore')
9
10 # Set plot styles
11 sns.set(style='darkgrid')
12 plt.rcParams['figure.figsize'] = (12, 6)
13 |
14 # Load datasets
15 matches = pd.read_csv('matches.csv')
16 deliveries = pd.read_csv('deliveries.csv')
17
18 # Display initial data overview
19 print("Matches Dataset:")
20 print(matches.head())
21 print("\nDeliveries Dataset:")
22 print(deliveries.head())
23
24 # -----
25 # 1. Data Cleaning and Handling Missing Values
26 # -----
27
28 # Check for missing values
29 print("\nMissing values in matches dataset:")
30 print(matches.isnull().sum())
31 print("\nMissing values in deliveries dataset:")
32 print(deliveries.isnull().sum())
33
34 # Fill missing 'player_of_match' with 'Unknown'
35 matches['player_of_match'].fillna('Unknown', inplace=True)
36
37 # Drop rows with missing values in deliveries
38 deliveries.dropna(inplace=True)
39
40 # -----
41 # 2. Feature Selection and Engineering
42 # -----
43

```

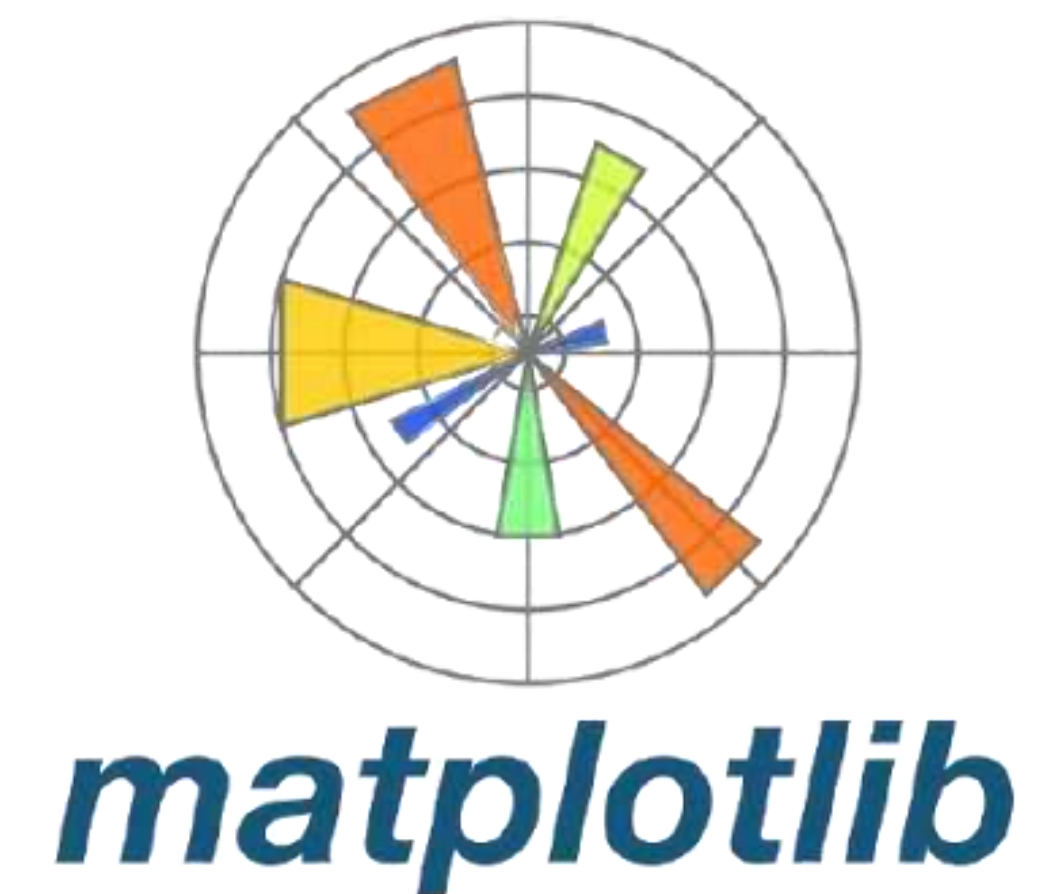
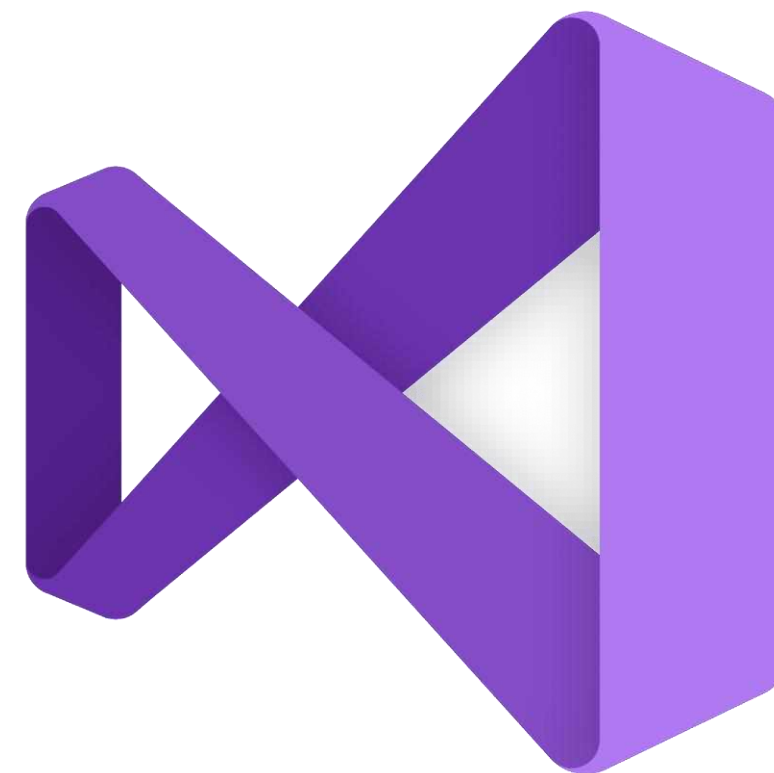
```

DA.py X
DA.py > ...
115
116 # Identify outliers in 'total_runs' using IQR
117 Q1 = deliveries['total_runs'].quantile(0.25)
118 Q3 = deliveries['total_runs'].quantile(0.75)
119 IQR = Q3 - Q1
120 outliers = deliveries[(deliveries['total_runs'] < Q1 - 1.5 * IQR) | (deliveries['total_runs'] > Q3 + 1.5 * IQR)]
121 print(f"\nNumber of outlier deliveries: {outliers.shape[0]}")
122
123 # -----
124 # 7. Initial Visual Representation of Key Findings
125 # -----
126
127 # Top 5 run scorers bar plot
128 plt.figure(figsize=(10,6))
129 sns.barplot(x=top_scorers.values, y=top_scorers.index, palette='viridis')
130 plt.title('Top 5 Run Scorers')
131 plt.xlabel('Runs')
132 plt.ylabel('Batsman')
133 plt.show()
134
135 # Runs per over line plot
136 plt.figure(figsize=(10,6))
137 sns.lineplot(x=runs_per_over.index, y=runs_per_over.values, marker='o')
138 plt.title('Average Runs per Over')
139 plt.xlabel('Over')
140 plt.ylabel('Average Runs')
141 plt.show()
142
DA.py X
DA.py > ...
44 # Convert 'date' column to datetime
45 matches['date'] = pd.to_datetime(matches['date'], dayfirst=True, errors='coerce')
46
47
48 # Create a new feature: 'match_year'
49 matches['match_year'] = matches['date'].dt.year
50
51 # Merge matches and deliveries data on 'match_id'
52 combined_df = deliveries.merge(matches, left_on='match_id', right_on='id')
53
54
55 # Calculate total runs per match
56 total_runs_per_match = combined_df.groupby('match_id')['total_runs'].sum().reset_index()
57 total_runs_per_match.rename(columns={'total_runs': 'total_runs_in_match'}, inplace=True)
58
59 # Merge total runs per match back to matches dataframe
60 matches = matches.merge(total_runs_per_match, left_on='id', right_on='match_id')
61
62
63 # -----
64 # 3. Ensuring Data Integrity and Consistency
65 # -----
66
67 # Standardize team names
68 team_replacements = {
69     'Delhi Daredevils': 'Delhi Capitals',
70     'Kings XI Punjab': 'Punjab Kings'
71 }
72 matches.replace({'team1': team_replacements, 'team2': team_replacements, 'winner': team_replacements}, inplace=True)
73 deliveries.replace({'batting_team': team_replacements, 'bowling_team': team_replacements}, inplace=True)
74
75 # -----
76 # 4. Summary Statistics and Insights
77 # -----
78
79 # Total matches played
80 total_matches = matches.shape[0]
81 print(f"\nTotal matches played: {total_matches}")
82
83 # Total runs scored
84 total_runs = deliveries['total_runs'].sum()
85 print(f"Total runs scored: {total_runs}")

```

TECHNOLOGY USED



KEY CONCEPT USED

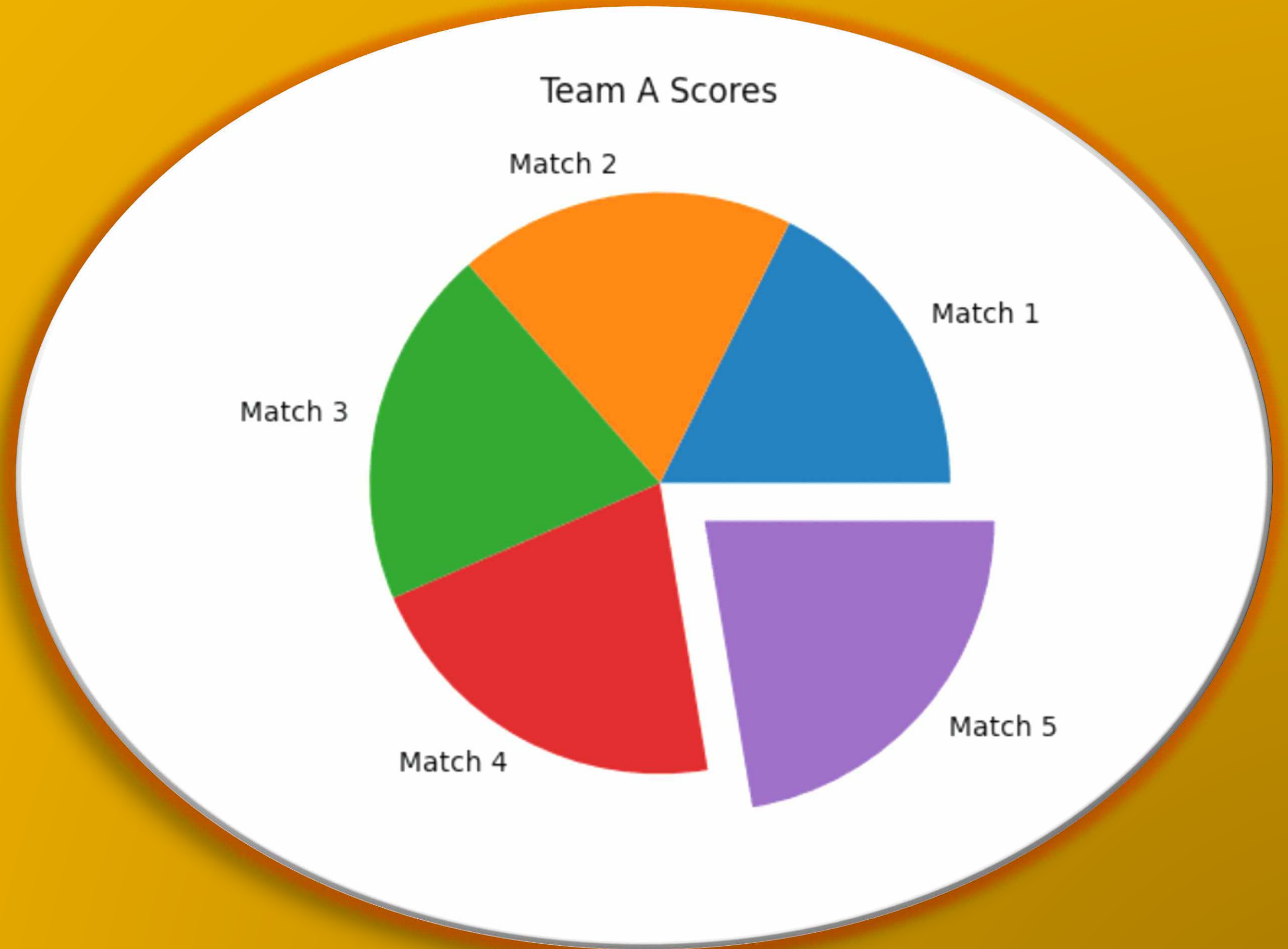
- I. Data Cleaning: Handling null values, formatting
- II. Exploratory Data Analysis (EDA): Summary statistics, data grouping
- III. Visualization: Bar plots, line graphs, heatmaps
- IV. Aggregation: GroupBy functions for team and player metrics
- V. Comparative Analysis: Toss vs result, batting first vs second, etc.

CHALLENGES FACED

- Inconsistent data formats in early seasons
- Handling large delivery-level dataset efficiently
- Interpreting cricket-specific statistics in data context
- Choosing the right visualizations to represent insights clearly

CONCLUSION

This project provided valuable insights into IPL matches and player performance. It reinforced the importance of data preprocessing and visualization in uncovering patterns. The project also strengthened our analytical thinking and Python skills, while opening possibilities for future work in predictive modeling.



THANK YOU