

Monthly Intership Program for Professionals@ THE CODERS READY

THE CODERS READY

Data Science and Business Analytic

Name:Ashima Sapehia

Trainee Intern @CodersReady

Task 3: Dataset of INDIAN PREMIER LEAGUE

Problem Statement: 1.Perform 'Exploratory Data Analysis' on dataset of 'INDIAN PLEAGUE LEAGUE'. 2.As a Sports analyst,try to find out the most successful teams,players and factor contributing to the wins or loss of team. 3.What other suggestions/insights you can derive through this EDA thus suggest teams or players a company should endorse for its products. 4.You can use Python tool to perform this analysis .

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

In [2]: matches_df = pd.read_csv("matches.csv")
deliveries_df = pd.read_csv("deliveries.csv")

In [3]: matches_df.head()
```

	id	season	city	date	team1	team2	toss_winner	toss_decision	result	dl_applied	winner	win_by_runs	win_by_wickets	player_of_match	venue	umpire1	umpire2	umpire3
0	1	2017	Hyderabad	05-04-17	Sunrisers Hyderabad	Royal Challengers Bangalore	Royal Challengers Bangalore	field	normal	0	Sunrisers Hyderabad	35	0	Yuvraj Singh	Rajiv Gandhi International Stadium, Uppal	AY Dandekar		NJ Lio
1	2	2017	Pune	06-04-17	Mumbai Indians	Rising Pune Supergiant	Rising Pune Supergiant	field	normal	0	Rising Pune Supergiant	0	7	SPD Smith	Maharashtra Cricket Association Stadium	A Nand Kishore		S R
2	3	2017	Rajkot	07-04-17	Gujarat Lions	Kolkata Knight Riders	Kolkata Knight Riders	field	normal	0	Kolkata Knight Riders	0	10	CA Lynn	Saurashtra Cricket Association Stadium	Nitin Menon		CK Nand
3	4	2017	Indore	08-04-17	Rising Pune Supergiant	Kings XI Punjab	Kings XI Punjab	field	normal	0	Kings XI Punjab	0	6	GJ Maxwell	Holkar Cricket Stadium	AK Chaudhary		Shamshud
4	5	2017	Bangalore	08-04-17	Royal Challengers Bangalore	Dehi Daredevils	Royal Challengers Bangalore	bat	normal	0	Royal Challengers Bangalore	15	0	KM Jadhav	M Chinnaswamy Stadium	NaN		N

```
In [4]: deliveries_df.head()
```

	match_id	inning	battling_team	bowling_team	over	ball	batsman	non_striker	bowler	is_super_over	...	bye_runs	legbye_runs	noball_runs	penalty_runs	batsman_runs	extra_runs	total_runs
0	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	1	DA Warner	S Dhawan	TS Mills	0	...	0	0	0	0	0	0	0
1	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	2	DA Warner	S Dhawan	TS Mills	0	...	0	0	0	0	0	0	0
2	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	3	DA Warner	S Dhawan	TS Mills	0	...	0	0	0	0	0	4	0
3	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	4	DA Warner	S Dhawan	TS Mills	0	...	0	0	0	0	0	0	0
4	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	5	DA Warner	S Dhawan	TS Mills	0	...	0	0	0	0	0	2	2

5 rows × 21 columns

DATA INFORMATION

```
In [5]: print(matches_df.info())
print(deliveries_df.info())

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 756 entries, 0 to 755
Data columns (total 18 columns):
# Column Non-Null Count Dtype
---
0 id 756 non-null int64
1 season 756 non-null int64
2 city 749 non-null object
3 date 756 non-null object
4 team1 756 non-null object
5 team2 756 non-null object
6 toss_winner 756 non-null object
7 toss_decision 756 non-null object
8 result 756 non-null object
9 dl_applied 756 non-null object
10 winner 752 non-null object
11 win_by_runs 756 non-null int64
12 win_by_wickets 756 non-null int64
13 player_of_match 752 non-null object
14 venue 756 non-null object
15 umpire1 754 non-null object
16 umpire2 754 non-null object
17 umpire3 119 non-null object
dtypes: int64(5), object(13)
memory usage: 186.4+ KB
None

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 17978 entries, 0 to 17977
Data columns (total 21 columns):
# Column Non-Null Count Dtype
---
0 match_id 17978 non-null int64
1 inning 17978 non-null int64
2 batting_team 17978 non-null object
3 bowling_team 17978 non-null object
4 over 17978 non-null int64
5 ball 17978 non-null int64
6 batsman 17978 non-null object
7 non_striker 17978 non-null object
8 bowler 17978 non-null object
9 is_super_over 17978 non-null int64
10 wide_runs 17978 non-null int64
11 bye_runs 17978 non-null int64
12 legbye_runs 17978 non-null int64
13 noball_runs 17978 non-null int64
14 penalty_runs 17978 non-null int64
15 batsman_runs 17978 non-null int64
16 extra_runs 17978 non-null int64
17 total_runs 17978 non-null int64
18 player_dismissed 8834 non-null object
19 dismissal_kind 8834 non-null object
20 fielder 6448 non-null object
dtypes: int64(13), object(8)
memory usage: 28.7+ MB
None

In [6]: matches_df["umpire3"].isnull().sum()

Out[6]:
637

In [7]: matches_df["umpire3"].tail(10)

Out[7]:
746 Nanda Kishore
747 KN Ananthapadmanabhan
748 Nitin Menon
749 Ullhas Gandhe
750 Bruce Oxenford
751 S Ravi
752 Ian Gould
753 NaN
754 Chettithody Shamsuddin
755 Nigel Long
Name: umpire3, dtype: object

In [8]: matches_df.describe()
```

	id	season	dl_applied	win_by_runs	win_by_wickets
count	756.000000	756.000000	756.000000	756.000000	756.000000
mean	1792.178571	2013.444444	0.025132	13.283069	3.350529
std	3464.478148	3.366895	0.156330	23.471144	3.387963
min	1.000000	2008.000000	0.000000	0.000000	0.000000
25%	189.750000	2011.000000	0.000000	0.000000	0.000000
50%	378.500000	2013.000000	0.000000	0.000000	0.000000
75%	567.250000	2016.000000	0.000000	19.000000	6.000000
max	11415.000000	2019.000000	1.000000	146.000000	10.000000

```
In [9]: matches_df['id'].max() # Matches we have got in the dataset

Out[9]:
11415

In [10]: matches_df['season'].unique()

Out[10]:
array([2017, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2018,
       2019], dtype=int64)

Team won by Maximum Runs

In [11]: matches_df.iloc[matches_df['win_by_runs'].idxmax()]

Out[11]:
id season 44
city 2017
date Delhi
date 06-05-17
team1 Mumbai Indians
team2 Kolkata
toss_winner Delhi Daredevils
toss_decision bat
result normal
dl_applied 0
winner Mumbai Indians
win_by_runs 146
win_by_wickets 0
player_of_match LMP Simmons
venue Feroz Shah Kotla
umpire1 Nitin Menon
umpire2 CK Nandan
umpire3 NaN
Name: 43, dtype: object

Team won by minimum runs

In [12]: matches_df.iloc[matches_df['win_by_runs'].idxmin()]['winner']

Out[12]:
'Mumbai Indians'

Team won by Maximum Wickets

In [13]: matches_df.iloc[matches_df['win_by_wickets'].idxmax()]['winner']

Out[13]:
'Kolkata Knight Riders'

Team won by minimum runs

In [14]: matches_df.iloc[matches_df[matches_df['win_by_runs'].ge(1)].win_by_runs.idxmin()]['winner']

Out[14]:
'Mumbai Indians'

Team won by Minimum Wickets

In [15]: matches_df.iloc[matches_df[matches_df['win_by_wickets'].ge(1)].win_by_wickets.idxmin()]

Out[15]:
id season 569
city Kolkata
date 09-05-16
team1 Kings XI Punjab
team2 Kolkata Knight Riders
toss_winner Kings XI Punjab
toss_decision bat
result normal
dl_applied 0
winner Kolkata Knight Riders
win_by_runs 0
win_by_wickets 1
player_of_match AD Russell
venue Eden Gardens
umpire1 AK Chaudhary
umpire2 HDPK Dharmasena
umpire3 NaN
Name: 559, dtype: object

In [16]: matches_df.iloc[matches_df[matches_df['win_by_wickets'].ge(1)].win_by_wickets.idxmin()]['winner']

Out[16]:
'Kolkata Knight Riders'

Season Which had most number of matches

In [17]: plt.figure(figsize=(12,6))
sns.countplot(x='season', data=matches_df)
plt.show()
```

Observation

1.Mumbai Indians is the team which won by maximum and minimum runs. 2.Kolkata Knight Riders is the team which won by maximum and minimum wickets.

```
In [30]: top_players = matches_df.match_value_counts()[0:10]
fig, ax = plt.subplots(figsize=(15,8))
ax.set_ylim([0,20])
ax.set_ylabel("Count")
ax.set_title("Top player of the match Winners")
top_players.plot.bar()
sns.barplot(x = top_players.index, y = top_players, orient='v', palette="Blues");
plt.show()
```

Top player of the match Winners

CH Gayle is the most Successful player in all match winners

```
In [35]: deliveries_df.plot(kind='area')
plt.show()
```

```
In [36]: matches_df.plot(kind='area')
plt.show()
```

Number of matches in each venue:

```
In [19]: plt.figure(figsize=(12,6))
sns.countplot(x='venue', data=matches_df)
plt.xticks(rotation='vertical')
plt.show()
```

Number of wins per team:

```
In [20]: plt.figure(figsize=(12,6))
sns.countplot(x='winner', data=matches_df)
plt.xticks(rotation=90)
plt.show()
```

Champions each season:

```
In [21]: temp_df = matches_df.drop_duplicates(subset=['season'], keep='last')[['season', 'winner']].reset_index(drop=True)
temp_df
```

	season	winner
0	2017	Mumbai Indians
1	2008	Rajasthan Royals
2	2009	Deccan Chargers
3	2010	Chennai Super Kings
4	2011	Kolkata Knight Riders
5	2012	Kolkata Knight Riders
6	2013	Mumbai Indians
7	2014	Kolkata Knight Riders
8	2015	Mumbai Indians
9	2016	Sunrisers Hyderabad
10	2018	Chennai Super Kings
11	2019	Mumbai Indians

Toss decision:

```
In [22]: temp_series = matches_df.toss_decision.value_counts()
labels = (np.array(temp_series.index))
sizes = (np.array(temp_series / temp_series.sum())*100)
colors = ['gold', 'lightskyblue']
plt.pie(sizes, labels=labels, colors=colors,
        autopct='%1.1f%%', shadow=True, startangle=90)
plt.title("Toss decision percentage")
plt.show()
```

```
In [23]: plt.figure(figsize=(12,6))
sns.countplot(x='season', hue='toss_decision', data=matches_df)
plt.xticks(rotation='vertical')
plt.show()
```