

▼ Descriptive Statistics

Another important aspect of Descriptive Analysis is visualization of data to understand it better and d

```

1  # Importing Libraries
2  import math
3
4  import numpy as np
5  import pandas as pd
6
7  import matplotlib.pyplot as plt
8
9  import seaborn as sns
10
11 import scipy.stats as stats
12
13 import warnings
14 warnings.filterwarnings('ignore')

1  deliveries = pd.read_csv('/content/deliveries.csv')
2  matches = pd.read_csv('/content/matches.csv')
3  matches.head()

```



	id	season	city	date	team1	team2	toss_winner	toss_decision	resu
0	1	2017	Hyderabad	2017-04-05	Sunrisers Hyderabad	Royal Challengers Bangalore	Royal Challengers Bangalore	field	norm
1	2	2017	Pune	2017-04-06	Mumbai Indians	Rising Pune Supergiant	Rising Pune Supergiant	field	norm
2	3	2017	Rajkot	2017-04-07	Gujarat Lions	Kolkata Knight Riders	Kolkata Knight Riders	field	norm
3	4	2017	Indore	2017-04-08	Rising Pune Supergiant	Kings XI Punjab	Kings XI Punjab	field	norm

▼ Bar Plot

A bar chart is a graph with rectangular bars. The graph usually compares different categories. Although (bars standing up) or horizontally (bars laying flat from left to right), the most usual type of bar graph

The horizontal (x) axis represents the categories; The vertical (y) axis represents a value for those cat

A bar graph is useful for looking at a set of data and making comparisons. For example, it's easier to :
chunk of your budget by glancing at the above chart rather than looking at a string of numbers.

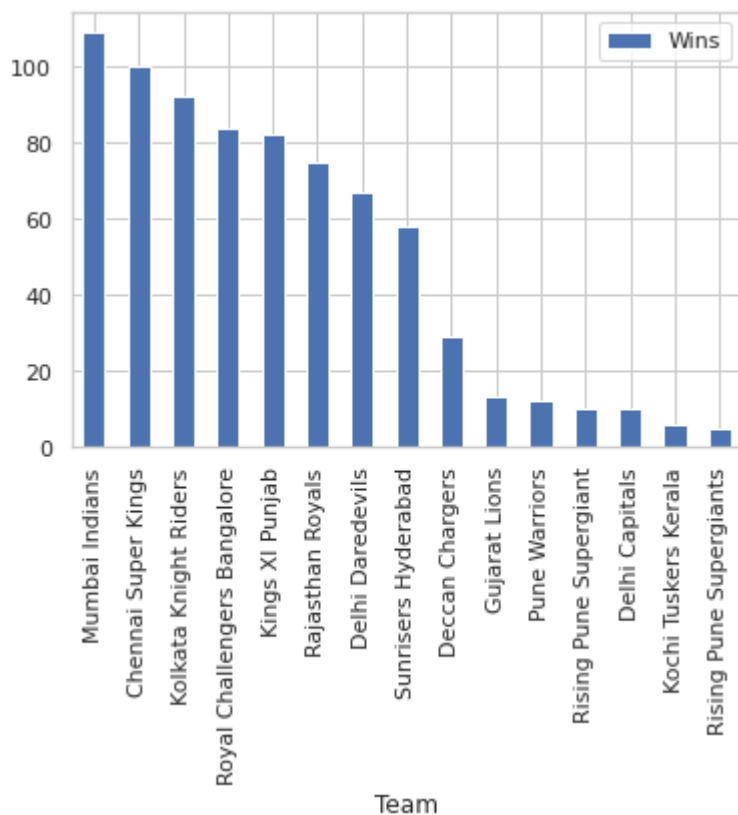
Bar charts can also represent more complex categories with stacked bar charts or grouped bar charts

Although they bar charts and histograms look the same, but they have one important difference: **they**
used to plot discrete data whereas the continous data is plotted using Histograms.

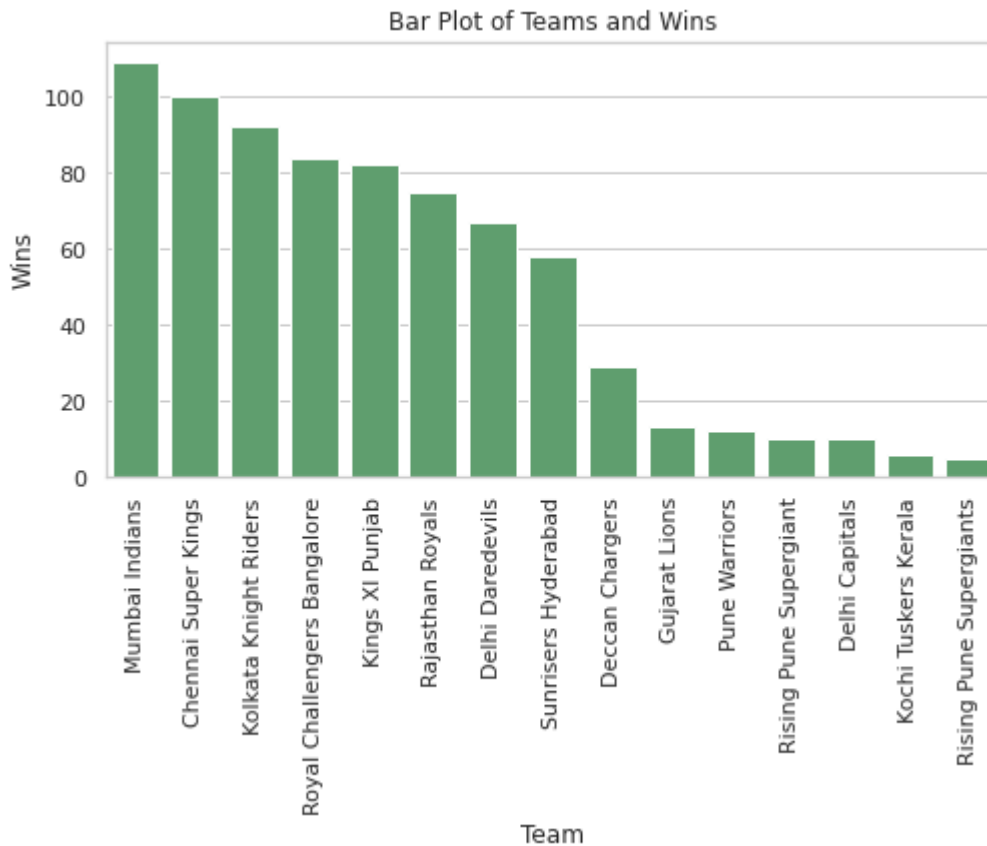
```
1 Wins = matches.winner.value_counts().rename_axis('Team').to_frame('Wins').reset_index()
```

```
1 # using Default Pandas Function
2 Wins.plot.bar(x = 'Team', rot = 90)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f86b615bb70>
```



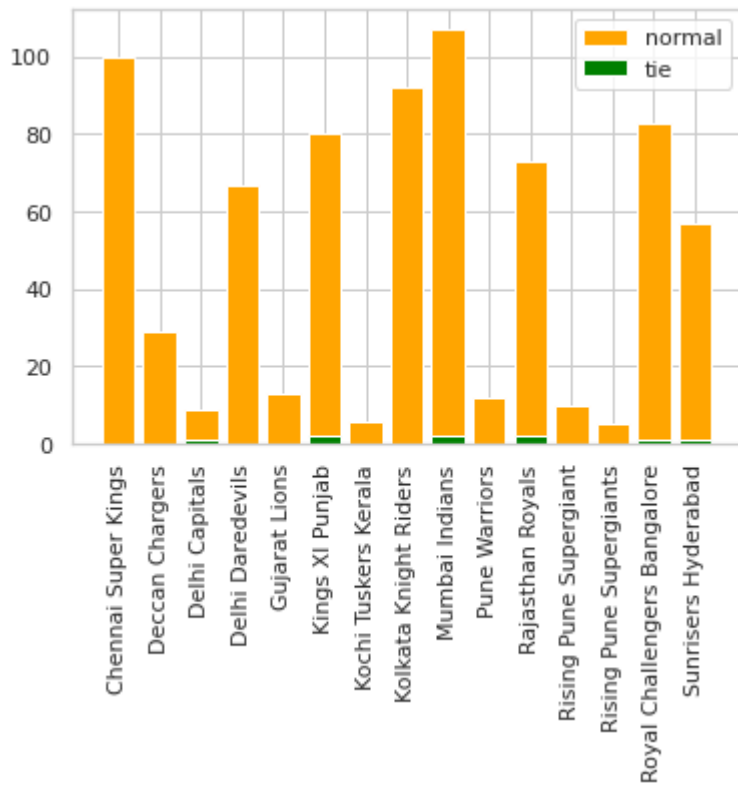
```
1 # Bar Plot to visualize the number of wins by a team
2
3 sns.set(style="whitegrid")
4
5 # Initialize the matplotlib figure
6 f, ax = plt.subplots(figsize=(8, 4))
7
8 # Plot the total wins by each team
9 sns.barplot(x="Team", y="Wins", data=Wins, color="g")
10 ax.set_xticklabels(Wins.Team, rotation = 90)
11 plt.title('Bar Plot of Teams and Wins')
12 plt.show()
```



```
1 xtab = pd.crosstab(matches['winner'],matches['result']).reset_index()
```

```
1 # Stacked Chart
2 ax = plt.subplot(111)
3 ax.bar(xtab.winner,xtab.normal, color = 'orange')
4 ax.bar(xtab.winner,xtab.tie, color = 'green')
5 ax.set_xticklabels(xtab.winner,rotation = 90)
6 ax.legend(['normal','tie'])
7 plt.show()
```

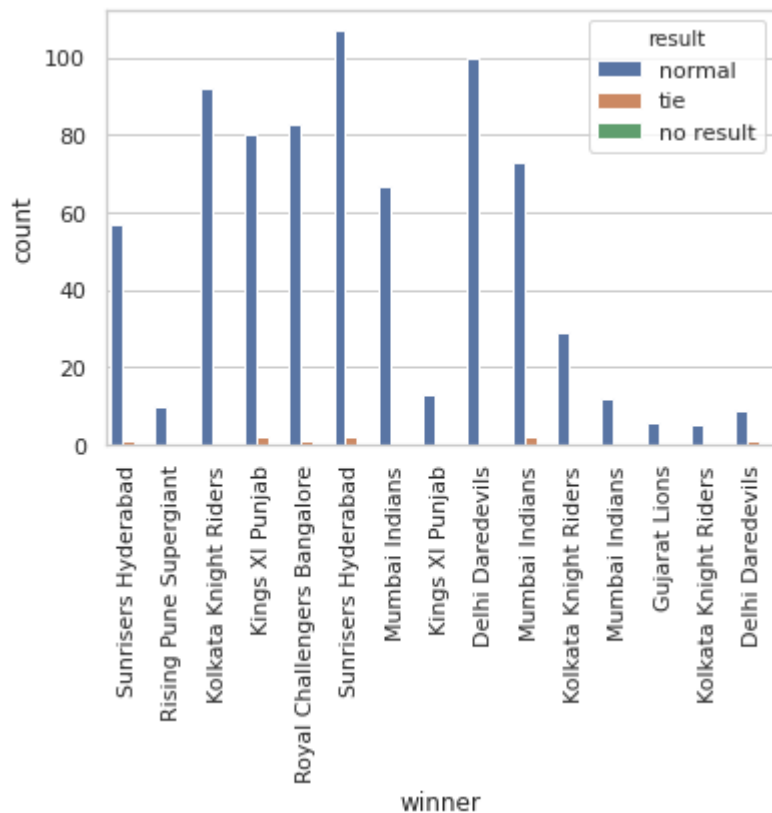




```

1 # Categorical Count Plot
2 fi = sns.countplot(matches.winner, hue = matches.result)
3 fi.set_xticklabels(matches.winner, rotation = 90)
4 plt.show()

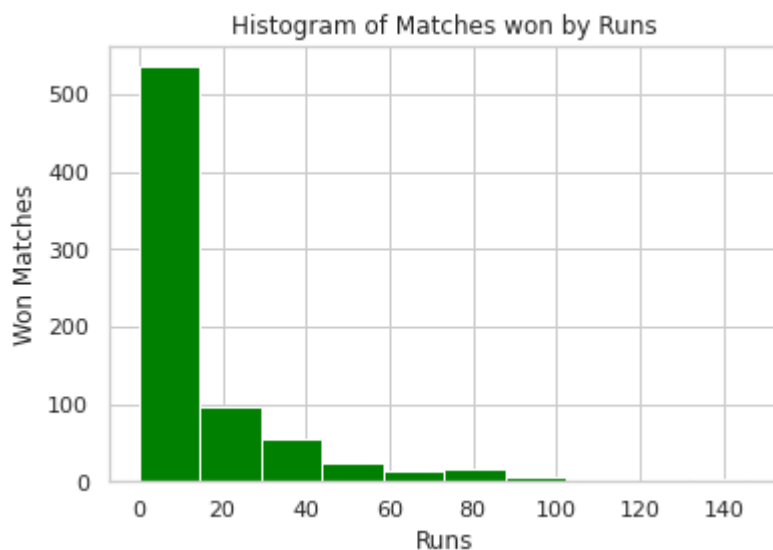
```



▼ Histograms

Histograms are similar to bar charts; they are a way to display counts of data. A bar graph charts actual values, while a histogram charts the number of items in that category. A histogram displays the same categorical data as a bar chart, but the height of the bar indicates the number of items in that category. A histogram divides the variable into bins, counts the data points in each bin, and shows the bins on the x-axis.

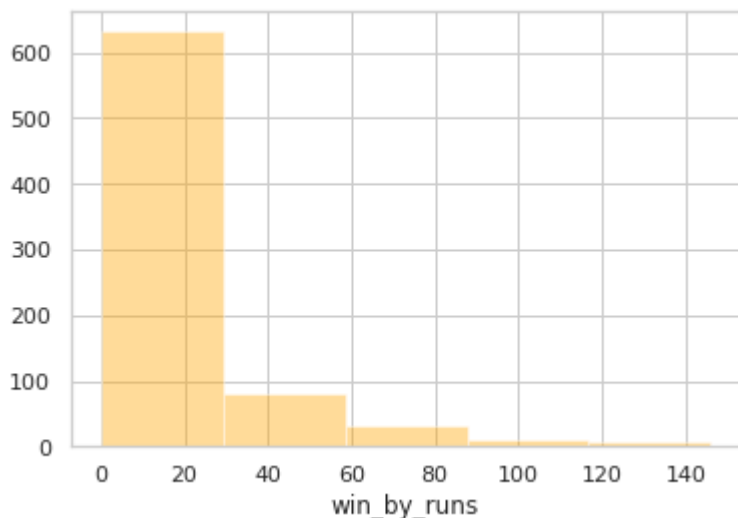
```
1 # Matplotlib histogram for matches won by runs
2 plt.hist(matches.win_by_runs, color = 'green')
3
4 # Add labels
5 plt.title('Histogram of Matches won by Runs')
6 plt.xlabel('Runs')
7 plt.ylabel('Count of Won Matches')
8 plt.show()
```



```
1 # seaborn histogram
2 sns.distplot(matches.win_by_runs, hist=True, kde=False, bins=int(140/28), color = 'orange')
```



<matplotlib.axes._subplots.AxesSubplot at 0x7f86b3bdf98>

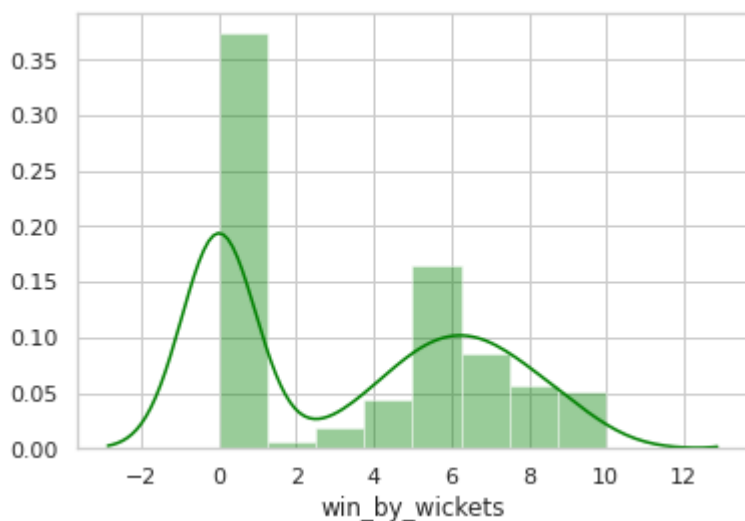


▼ Density Plot

A density plot is a smoothed, continuous version of a histogram estimated from the data. The most common method for this is **Kernel Density Estimation** (KDE). KDE is a non-parametric way to estimate the probability density function of a continuous variable. A continuous curve (the kernel) is drawn at every individual data point and all of these curves are then summed to produce a density estimation. The kernel most often used is a Gaussian (which produces a Gaussian bell curve).

```
1 sns.distplot(matches.win_by_wickets, hist=True, kde=True, color = 'green') # Seaborn
```

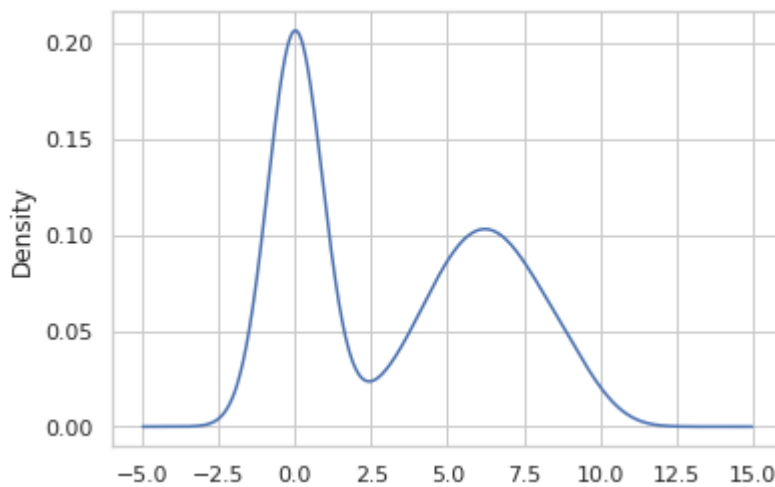
↳ <matplotlib.axes._subplots.AxesSubplot at 0x7f86b3cde160>



```
1 matches.win_by_wickets.plot(kind='density') # Matplotlib
```

↳

<matplotlib.axes._subplots.AxesSubplot at 0x7f86b3ce7b70>

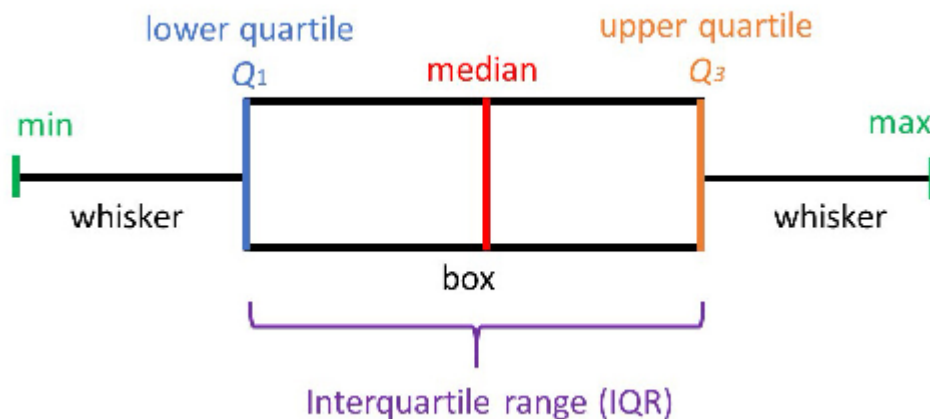


▼ **Box-Whisker Plot**

Box and whisker plots (or box plots) represents five-number summary of the dataset. The five-number

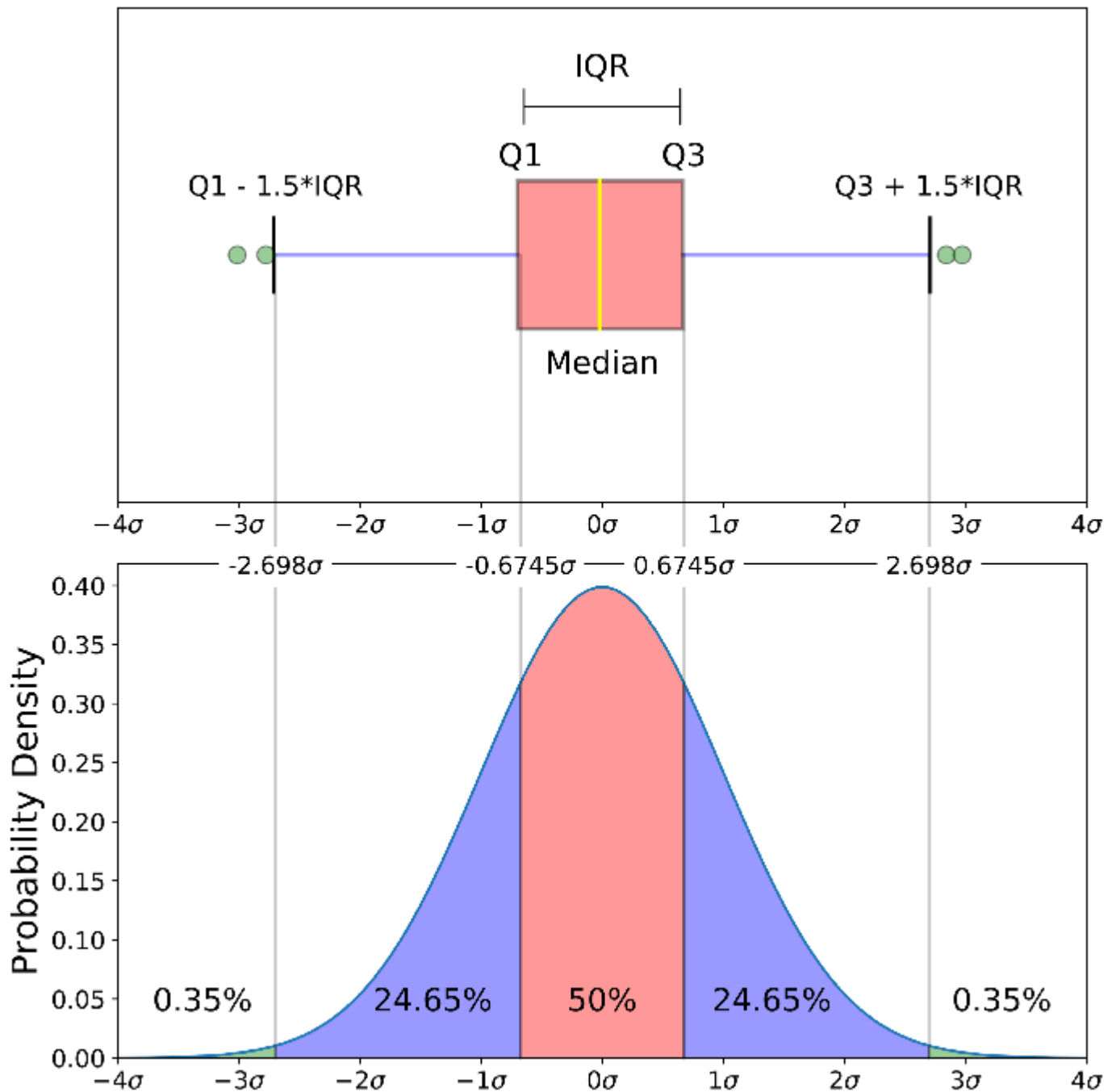
- Minimum
- First quartile (25th percentile)
- Median (50th percentile)
- Third quartile (75th percentile)
- Maximum

The following is a representation of box-and-whisker plot:



Data sets can sometimes contain outliers that are suspected to be anomalies (perhaps because of data flukes). If outliers are present, the whisker on the appropriate side is drawn to $1.5 \times \text{IQR}$ rather than the full range. Small circles or unfilled dots are drawn on the chart to indicate where suspected outliers lie. Filled circles indicate confirmed outliers.

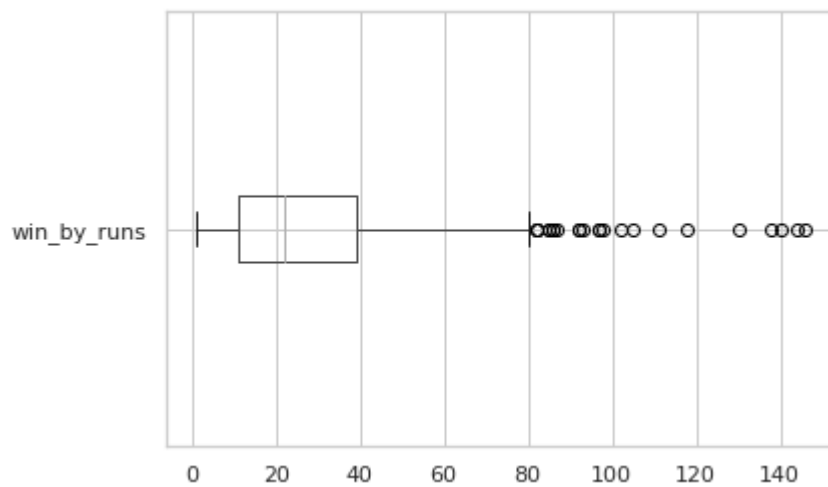
Box Plot of Normal Distribution:



```
1 matches.loc[matches['win_by_runs']>0, 'win_by_runs'].to_frame().boxplot(vert=False)
```

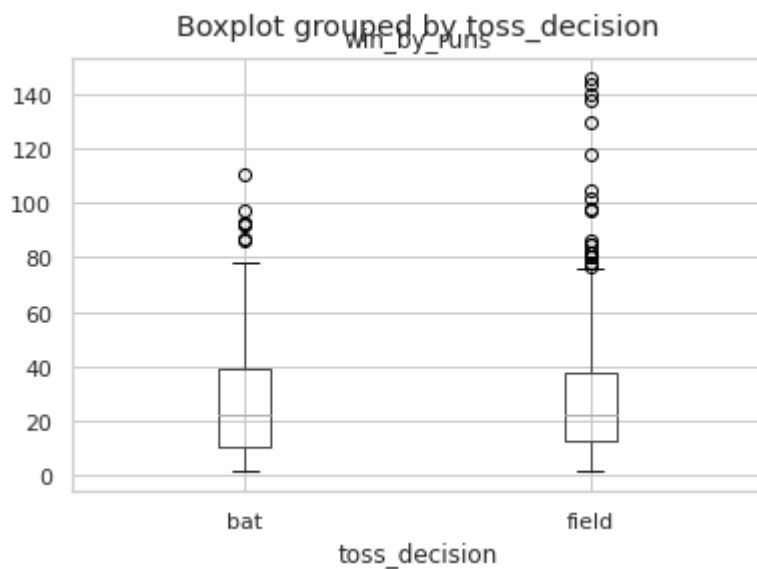


<matplotlib.axes._subplots.AxesSubplot at 0x7f86b3605e48>



```
1 runs = matches.loc[matches['win_by_runs']>0,:]
2 runs.boxplot(column= 'win_by_runs', by = 'toss_decision')
```

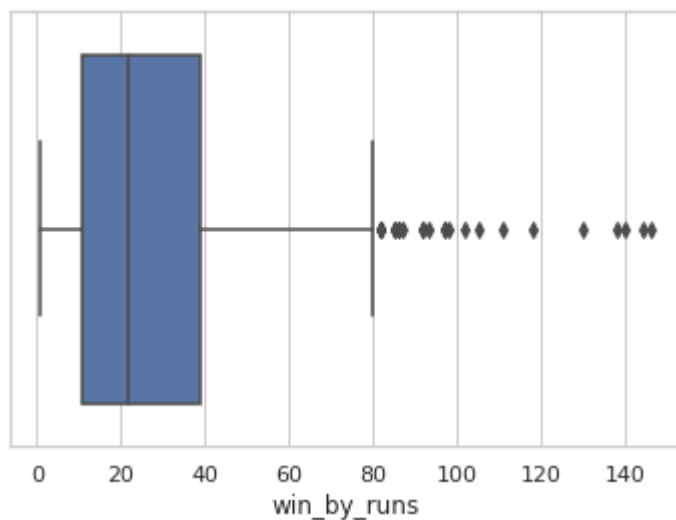
↳ <matplotlib.axes._subplots.AxesSubplot at 0x7f86b3580d30>



```
1 sns.boxplot(x = 'win_by_runs', data = runs)
```

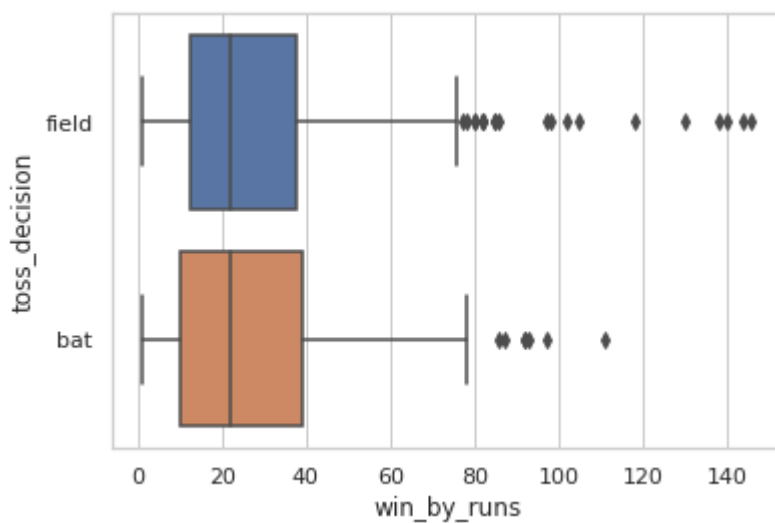
↳

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f86b35987f0>
```



```
1 sns.boxplot(x = 'win_by_runs', y='toss_decision', data = runs)
```

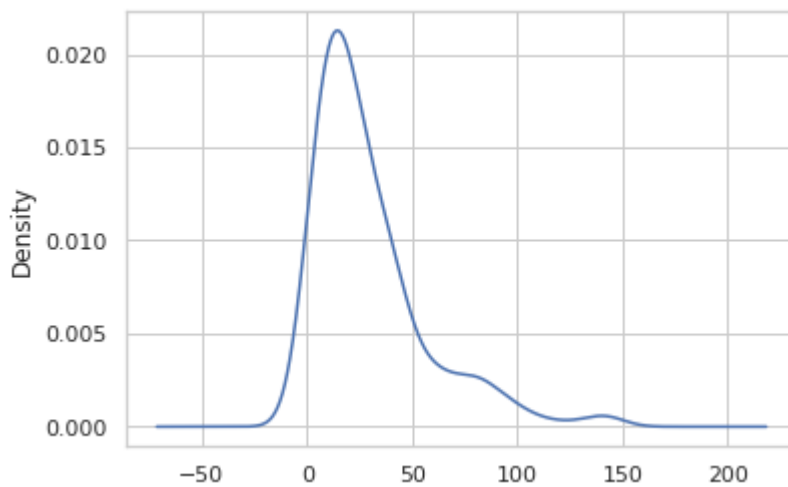
```
<matplotlib.axes._subplots.AxesSubplot at 0x7f86b3300e48>
```



```
1 runs.win_by_runs.plot(kind= 'density') # right-skewed
```

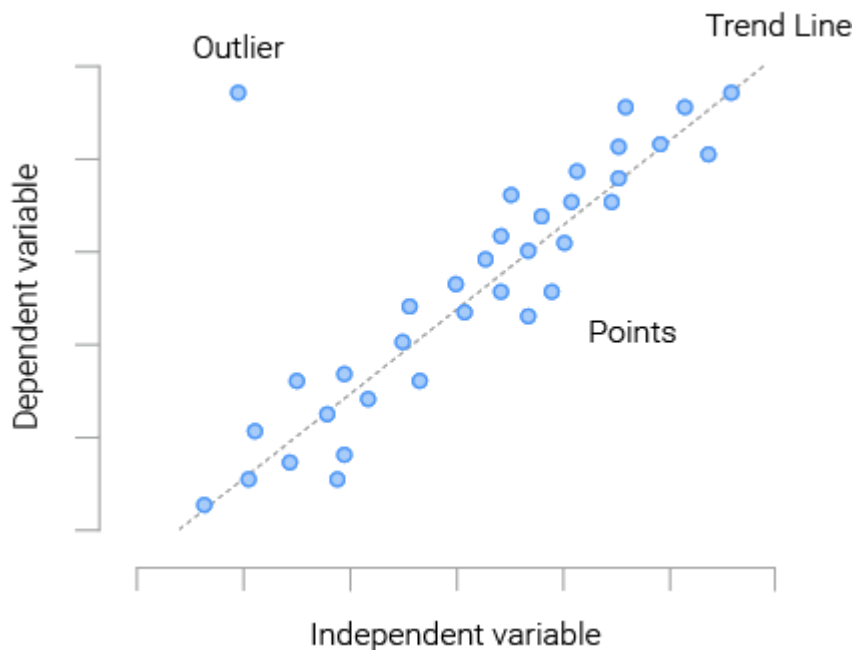
```
<matplotlib.axes._subplots.AxesSubplot at 0x7f86b3300e48>
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f86b320d470>



▼ Scatter Plot

Scatter plot Scatter plots help us see relationship between variables or how much one variable is affected by the other. From the plot, whether there is a positive or negative or no association between the variables. Following is a scatter plot of Variable A and Variable B.



We can see that there is a positive association between Variable A and Variable B. i.e. Variable A increases as Variable B increases. Some datasets may contain a few data points that don't fit the pattern. They're called outliers. A scatter chart can suggest a linear relationship between 2 variables (i.e. a straight line).

```
1 deliveries.head()
```



	match_id	inning	batting_team	bowling_team	over	ball	batsman	non_striker	bowler
0	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	1	DA Warner	S Dhawan	TS Mill
1	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	2	DA Warner	S Dhawan	TS Mill
2	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	3	DA Warner	S Dhawan	TS Mill
3	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	4	DA Warner	S Dhawan	TS Mill
4	1	1	Sunrisers Hyderabad	Royal Challengers Bangalore	1	5	DA Warner	S Dhawan	TS Mill

```

1 highest_score_batsman = deliveries.groupby(['match_id', 'batting_team', 'batsman'])['bat
2
3 total_runs = deliveries.groupby(['match_id', 'batting_team'])['total_runs'].sum()
4
5 match_total_highest_combined = pd.concat([highest_score_batsman, total_runs], axis=1)
6 match_total_highest_combined.head()

```



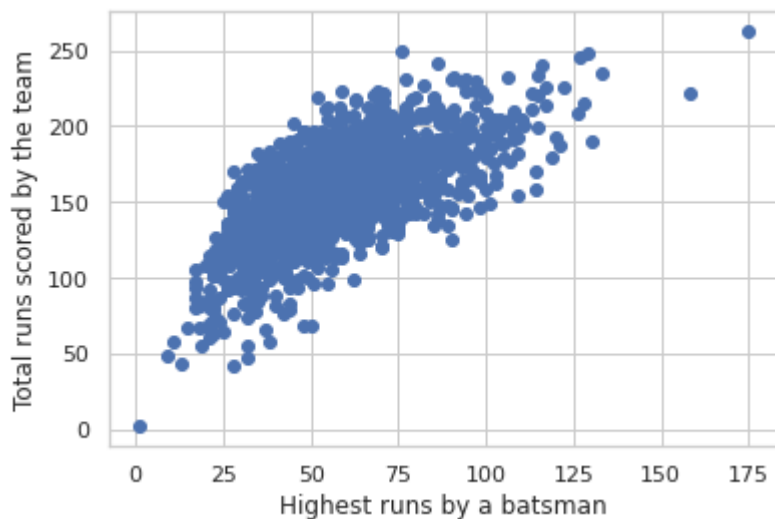
		batsman_runs	total_runs
match_id	batting_team		
1	Royal Challengers Bangalore	32	172
	Sunrisers Hyderabad	62	207
2	Mumbai Indians	38	184
	Rising Pune Supergiant	84	187
3	Gujarat Lions	68	183

```

1 plt.scatter(match_total_highest_combined.batsman_runs, match_total_highest_combined.tota
2 plt.xlabel("Highest runs by a batsman")
3 plt.ylabel("Total runs scored by the team")
4 plt.show()

```

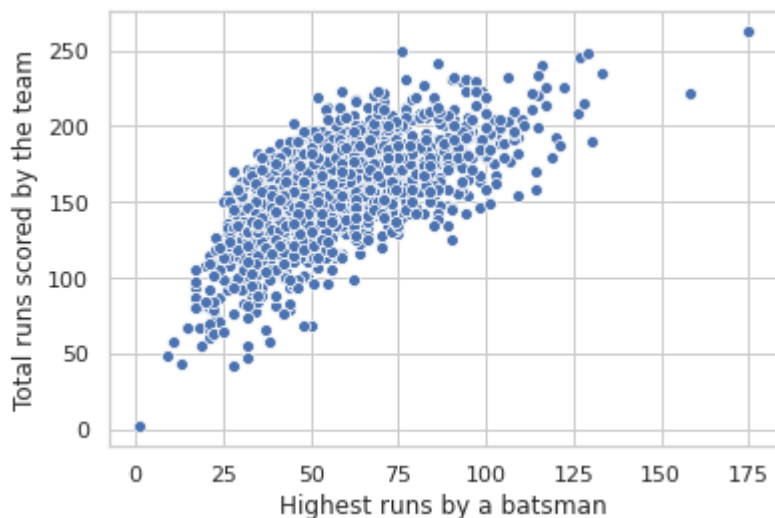




```

1 sns.scatterplot(match_total_highest_combined.batsman_runs, match_total_highest_combined.
2 plt.xlabel("Highest runs by a batsman")
3 plt.ylabel("Total runs scored by the team")
4 plt.show()

```



▼ Pair Plot

The pair plot technique allows us to visualize distributions of individual numerical features, as well as numerical features. Pair plot is a handy technique because it is very easy to implement and allows us between each of our features by pairing them together.

Seaborn pair plot command will not run if there are NaN values in the data frame.

Output of Pair Plot has two basic types of graphs: histograms for each numerical feature showing the the numerical values plotted against each other.

The histograms are displayed in a diagonal fashion starting in the upper left cell, or the (0, 0)th cell, of right cell, or the (n-1, n-1)th cell, each histogram can be looked by the corresponding labels on either t

The scatter plots can be located in all other cells on the grid. We can identify the individual features p corresponding labels on either the x-axis or the y-axis.

```
1 sns.pairplot(matches.iloc[:,1:]) # not including id column here
```

```
<seaborn.axisgrid.PairGrid at 0x7f86b1162d30>
```

