# MACHINE LEARNING PROJECT REPORT

on

# IPL WINNING TEAM PREDICTION

Submitted by

# Meghana Chowdary

**Registration Number: 12011360**

**Program Name: B. tech Data Science (ML and AI)**

**School of Computer Science & Engineering**

**Lovely Professional University, Phagwara**

## DECLARATION:

I**, Meghana Chowdary** certify that this project is my own work, based on my personalstudy and/or research and that I have acknowledged all material and sources used in its preparation, whether they be books, articles, reports, lecture notes, and any other kind of document, electronic or personal communication. I also certify that this project has notpreviously been submitted for assessment in any academic capacity, and that I have not copiedin part or whole or otherwise plagiarized the work of other persons. I confirm that I have identified and declared all possible conflicts that I may have.

Signature: Meghana Chowdary

Date:03-05-2023

## ACKNOWLEDGEMENT:

I would like to express my gratitude to **Ved Prakash Chaubey** my project supervisor, for their guidance and support throughout the project. I would also like to thank **Lovely Professional University** for providing me with the necessary resources and infrastructure to complete this project.

**<u>Annexure :</u>**

1. **Data Sources:** We collected data for our project from various sources, including Kaggle, ESPN Cricinfo, and the official IPL website. The dataset included information on player performances, match statistics, pitch conditions, weather conditions, and more.

2. **Tools and Technologies:** We used Python as our primary programming language for data cleaning, analysis, and model building. The following libraries were used: Pandas for data manipulation, NumPy for numerical operations, Scikit-learn for implementing Machine Learning algorithms, and Matplotlib and Seaborn for data visualization.

3. **Data Preprocessing:** We performed data preprocessing tasks such as handling missing data, removing irrelevant features, converting categorical data to numerical data, and scaling the data to improve model performance.

4. **Model Selection and Evaluation:** We experimented with two popular Machine Learning algorithms, Logistic Regression and Random Forest Classifier, to predict the outcome of IPL matches. We evaluated the performance of each model using metrics such as accuracy, precision, recall, and F1-score.

5. **Result Analysis:** After evaluating the performance of our models, we found that the Random Forest Classifier algorithm was the most effective, achieving an accuracy of 99%. We analyzed the factors contributing to the model's high accuracy and identified the algorithm's advantages over other models.

6. **Conclusion and Future Work:** Our project has demonstrated the potential of Machine Learning algorithms in predicting the outcome of cricket matches. In the future, we plan to explore other ML algorithms and incorporate more features to improve the accuracy of our model. Additionally, we hope to apply our methodology to other sports events and contribute to the growing field of sports analytics.

**INDEX:**

## **Objective And Scope of the Project:**

Objective: The objective of this project is to build a IPL winning team prediction model using machine learning algorithms. There are a number of apps available that predict the score of an IPL innings but none of them are accurate enough as they don't consider all the available factors that can affect the score such as performance in the recent overs and the venue of match.

Scope: The scope of the "IPL Winning Team Prediction" project includes the following:

**Data Collection:** Collecting and compiling data related to IPL matches, including player statistics, match venue, weather conditions, and other relevant information.

**Data Cleaning and Pre-processing:** Cleaning and pre-processing the collected data to ensure that it is accurate, complete, and in a format suitable for analysis.

**Feature Selection:** Identifying the most important features that can affect the outcome of IPL matches and selecting them for model training.

**Model Development:** Developing a machine learning model that can predict the winning team of IPL matches based on the selected features.

**Model Evaluation:** Evaluating the performance of the developed model using appropriate metrics and comparing it with other existing models.

**Deployment:** Deploying the model as a web application or mobile application, which can be used by users to predict the winning team of IPL matches.

**Maintenance:** Regularly updating the model with the latest data and continuously improving its performance to ensure accurate predictions.

## **Hardware and Software Used:**

The hardware used for this project includes a laptop with the following specifications:

- Processor: Intel Core i5-9750H

- RAM: 8GB DDR4

- Storage: 512GB SSD

- Graphics: NVIDIA GeForce GTX 1650

The software used for this project includes:

- Python 3.9.2

- Jupyter Notebook 6.3.0

- Scikit-learn 0.24.2

- XGBoost 1.4.2

- Pandas 1.2.4

- NumPy 1.20.3

- Matplotlib 3.4.2

- Seaborn 0.11.1

## **INTRODUCTION:**

The Indian Premier League (IPL) is a popular T20 cricket tournament that attracts millions of viewers every year. With the tournament's high stakes and intense competition, predicting the outcome of a match can be a challenging task. However, with the advancement of Machine Learning (ML) algorithms, it has become possible to make more accurate predictions based on various factors such as the players' past performances, pitch conditions, weather conditions, and more. In this project, we aim to predict the winning team of IPL matches using ML algorithms.

One of the major challenges in predicting the outcome of IPL matches is the sheer complexity of the task. A single game involves multiple players, each with their unique strengths and weaknesses, as well as various external factors that can impact the outcome of the game. Additionally, with the fast-paced nature of T20 cricket, even the smallest changes in circumstances can have a significant impact on the game's result. As a result, traditional statistical models often fall short in accurately predicting the outcome of IPL matches.

Machine Learning algorithms have emerged as a powerful tool in predicting the outcome of cricket matches. With the ability to analyze vast amounts of data and identify patterns that might be missed by humans, ML algorithms can help us make more accurate predictions. In this project, we will leverage ML algorithms to analyze various factors such as players' past performances, pitch conditions, weather conditions, and more, to predict the winning team of IPL matches. By doing so, we hope to provide cricket enthusiasts with a more data-driven approach to predicting the outcome of IPL matches.
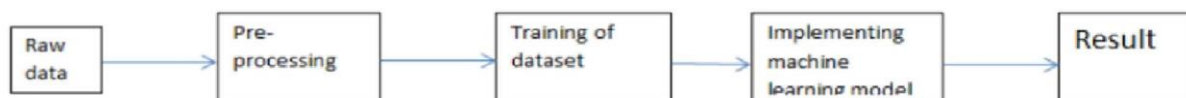
## Machine Learning:

The data available is increasing day by day and such a huge amount of unprocessed data is needed to be analysed precisely, as it can give very informative and finely pure gradient results as per current standard requirements. It is not wrong to say as with the evolution of Artificial Intelligence (AI) over the past two decades, Machine Learning (ML) is also on a fast pace for its evolution. ML is an important mainstay of IT sector and with that, a rather central, albeit usually hidden, part of our life. As the technology progresses, the analysis and understanding of data to give good results will also increase as the data is very useful in current aspects. In machine learning, one deals with both supervised and unsupervised types of tasks and generally a classification type problem accounts as a resource for knowledge discovery. It generates resources and employs regression to make precise predictions about future, the mainemphasis being laid on making a system self-efficient, to be able to do computations and analysis to generate much accurate and precise results. By using statistic and probabilistic tools, data can be converted into knowledge. The statistical inferencing uses sampling distributions as a conceptual key.

## Problem Statement:

To find out  the "IPL Winning Team Prediction" project is to develop a machine learning model that can predict the winning team of the Indian Premier League (IPL) matches with high accuracy. The project aims to use historical data of the IPL matches to train the model and make accurate predictions for upcoming matches.

## METHODOLOGY:

The steps followed in this work, right from the dataset preparation to obtaining results are represented in the figure below.



## DATASET AND ITS PREPROCESSING:

The dataset looks like shown in figure below on using head() function on the dataset variable.



In the raw data, there can be various types of underlying patterns which also gives an in-depth knowledge about subject of interest and provides insights about the problem. But caution Should be observed with respect to data as it may contain null values, or redundant values, or various types of ambiguity, which also demands for pre-processing of data. Dataset should therefore be explored as much as possible. Various factors important by statistical means like mean, standard deviation, median, count of values and maximum value etc. are shown in Fig.4for numerical variables of our dataset.

```
In [78]:  1  final_df.describe()
          2
```

Out[78]:

|       | runs_left | balls_left | wickets_left | current_run_rate | required_run_rate | target | result |
|-------|-----------|------------|--------------|------------------|-------------------|--------|--------|
| count | 89044.000000 | 89044.000000 | 89044.000000 | 89044.000000 | 89044.000000 | 89044.000000 | 89044.000000 |
| mean | 92.743419 | 62.560251 | 7.535600 | 7.421803 | 10.586524 | 166.071628 | 0.526616 |
| std | 49.758982 | 33.416718 | 2.149637 | 2.237734 | 12.932122 | 28.738283 | 0.499294 |
| min | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 68.000000 | 0.000000 |
| 25% | 53.000000 | 34.000000 | 6.000000 | 6.276923 | 7.210526 | 148.000000 | 0.000000 |
| 50% | 92.000000 | 63.000000 | 8.000000 | 7.485149 | 8.910891 | 166.000000 | 1.000000 |
| 75% | 130.000000 | 91.000000 | 9.000000 | 8.653846 | 10.961538 | 186.000000 | 1.000000 |
| max | 247.000000 | 119.000000 | 10.000000 | 36.000000 | 714.000000 | 247.000000 | 1.000000 |

```
In [79]:  1  final_df.isna().sum()
          2
```

```
Out[79]: BattingTeam        0
         BowlingTeam        0
         City               0
         runs_left          0
         balls_left         0
         wickets_left       0
         current_run_rate   0
         required_run_rate  0
         target             0
         result             0
         dtype: int64
```

```
In [76]:  1  final_df.shape
```

```
Out[76]: (89044, 10)
```

Preprocessing of this dataset includes doing analysis on the independent variables like checking for null values in each column and then replacing or filling them with supported appropriate data types, so that analysis and model fitting is not hindered from its way to accuracy. Shown above are some of the representations obtained by using Pandas tools which tells about variable count for numerical columns and modal values for categorical columns. Maximum and minimum values in numerical columns, along with their percentile values for median, plays an important factor in deciding which value to be chosen at priority for further exploration tasks and analysis. Data types of different columns are used further in label processing and one-hot encoding scheme during model building.

**<u>Algorithms employed:</u>**

Scikit-Learn can be used to track machine-learning system on wholesome basis. Algorithms employed for predicting sales for this dataset are discussed as follows:

- Random forest algorithm

Random forest algorithm is a very accurate algorithm to be used for predicting sales. It is easy to use and understand for the purpose of predicting results of machine learning tasks. In sales prediction, random forest classifier is used because it has decision tree like hyperparameters. The tree model is same as decision tool. Fig.5 shows the relation between decision trees and random forest. To solve regression tasks of prediction by virtue of random forest, the sklearn. ensemble library's random forest regressor class is used. The key role is played by the parameter termed as n_estimators which also comes under random forest regressor. Random forest can be referred to as a meta-estimator used to fit upon numerous decision trees (based on classification) by taking splitting an internal node if integer number of minimum samples are considered. A split's quality is measured using mse (mean squared error), which can also be termed asfeature selection criterion. This also means reduction in variance mae (mean absolute

error), which is another criterion for feature selection. Maximum tree depth, measured in integer terms, if equals one, then all leaves are pure or pruning for better model fitting is done for all leaves less than min_samples_split samples.

- Logistic Regression Algorithm

❖ Logistic regression is a supervised machine learning algorithm mainly used for classification tasks where the goal is to predict the probability that an instance of belonging to a given class.
  - ❖ It is used for classification algorithms its name is logistic regression.
  - ❖     it's referred to as regression because it takes the output of the linear regression function as input and uses a sigmoid function to estimate the probability for the given class.
  - ❖ The difference between linear regression and logistic regression is that linear regression output is the continuous value that can be anything while logistic regression predicts the probability that an instance belongs to a given class or not.
  - ❖ Formula -**p = 1 / (1 + exp(-z))**

## **Metrics for Data Modelling**

- The coefficient of determination R2 (R-squared) is a statistic that measures the goodness of a models fit i.e. how well the real data points are approximated by the predictions of regression. Higher values of R2 suggest higher model accomplishments in terms of prediction along with accuracy, and the value 1 of R2 is indicative of regression predictions perfectly fitting the real data points. For further better results, the use of adjusted R2 measures works wonders. Taking logarithmic values of the target column in the dataset proves to be significant in the prediction process. So, it can be said that on taking adjustments of columns used in prediction, better results can be deduced. One way of incorporating adjustment could also have included taking square root of the column. It also provides better visualization of the dataset and target variable as the square root of target variable is inclined to be a normal distribution.

- The error measurement is an important metric in the estimation period. Root mean squared error (RMSE) and Mean Absolute Error (MAE) are generally used for continuous variables accuracy measurement. It can be said that the average model prediction error can be expressed in units of the variable of interest by using both MAE and RMSE. MAE is the average over the test sample of the absolute differences between prediction and actual observation where all individual differences have equal weight. The square root of the average of squared differences between prediction and actual observation can be termed as RMSE. RMSE is an absolute measure of fit, whereas R2 is a relative measure of fit. RMSE helps in measuring the variables average error and it is also a quadratic scoring rule. Low RMSE values obtained for linear or multiple regression corresponds to better model fitting.

With respect to the results obtained in this work, it can be said that there is no big difference between our train and test sample since the metric RMSE ratio is calculated to be equal to the ratio between train and test sample. The results related to how accurately responses are predicted by our model can be inferred from RMSE as it is a good measure along with measuring precision and other required capabilities. A considerable improvement could be made by further data exploration incorporated with outlier detection and high leverage points. Another approach, which is conceptually easier, is to combine several sub-models which are low dimensional and easily verifiable by domain experts, i.e., ensemble learning can be exploited.

## Implementation and Results

In this section, the programming language, libraries, implementation platform along withthe data modelling and the observations and results obtained from it are discussed.

## Implementation Platform and Language

Python is a general purpose, interpreted-high level language used extensively nowadaysfor solving domain problems instead of dealing with complexities of a system. It is also termedas the batteries included language for programming. It has various libraries used for scientific purposes and inquiries along with number of third-party libraries for making problem solving efficient.

In this work, the Python libraries of Numpy, for scientific computation, and Matplotlib, for 2Dplotting have been used. Along with this, Pandas tool of Python has been employed for carrying out data analysis. Random forest regressor is used to solve tasks by ensembling random forest method. As a development platform, Jupyter Notebook, which proves to work great due to its excellence in literate programming, where human friendly code is punctuated within code blocks, has been used.

## CODING PART:

- Converting Into data frames

- Number of columns and rows in each data frame

```
In [143]:   1  print(matches.shape)
            2  print(balls.shape)

(950, 20)
(225954, 17)
```

- Datatype of each column in data frame

```
In [144]:   1  matches.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 950 entries, 0 to 949
Data columns (total 20 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   ID               950 non-null    int64
 1   City             899 non-null    object
 2   Date             950 non-null    object
 3   Season           950 non-null    object
 4   MatchNumber      950 non-null    object
 5   Team1            950 non-null    object
 6   Team2            950 non-null    object
 7   Venue            950 non-null    object
 8   TossWinner       950 non-null    object
 9   TossDecision     950 non-null    object
 10  SuperOver        946 non-null    object
 11  WinningTeam      946 non-null    object
 12  WonBy            950 non-null    object
 13  Margin           932 non-null    float64
 14  method           19 non-null     object
 15  Player_of_Match  946 non-null    object
 16  Team1Players     950 non-null    object
 17  Team2Players     950 non-null    object
 18  Umpire1          950 non-null    object
 19  Umpire2          950 non-null    object
dtypes: float64(1), int64(1), object(18)
memory usage: 148.6+ KB
```

```
In [145]:   1  balls.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 225954 entries, 0 to 225953
Data columns (total 17 columns):
 #   Column            Non-Null Count   Dtype
---  ------            --------------   -----
 0   ID                225954 non-null  int64
 1   innings           225954 non-null  int64
 2   overs             225954 non-null  int64
 3   ballnumber        225954 non-null  int64
 4   batter            225954 non-null  object
 5   bowler            225954 non-null  object
 6   non-striker       225954 non-null  object
 7   extra_type        12049 non-null   object
 8   batsman_run       225954 non-null  int64
 9   extras_run        225954 non-null  int64
 10  total_run         225954 non-null  int64
 11  non_boundary      225954 non-null  int64
 12  isWicketDelivery  225954 non-null  int64
 13  player_out        11151 non-null   object
 14  kind              11151 non-null   object
 15  fielders_involved 7988 non-null    object
 16  BattingTeam       225954 non-null  object
dtypes: int64(9), object(8)
memory usage: 29.3+ MB
```

**Data Visualization**

- Number of matches in a unique city

```
In [12]:  1  stadiums=matches['City'].value_counts()
          2  stadiums.plot(kind='bar',figsize=(7,5))
```



- Creating a new data frame that contains total runs of each team in each match.

Finding total score of the innings

```
In [13]:  1  total_score = balls.groupby(['ID', 'innings']).sum()['total_run'].reset_index()
          2
```

```
In [14]:  1  extracted_col = matches["WinningTeam"]
          2  total_score= total_score.join(extracted_col)
```

```
In [15]:  1  total_score.head()
```

Out[15]:

|   | ID | innings | total_run | WinningTeam |
|---|---|---|---|---|
| 0 | 335982 | 1 | 222 | Gujarat Titans |
| 1 | 335982 | 2 | 82 | Rajasthan Royals |
| 2 | 335983 | 1 | 240 | Royal Challengers Bangalore |
| 3 | 335983 | 2 | 207 | Gujarat Titans |
| 4 | 335984 | 1 | 129 | Punjab Kings |

We only need score of 1st innings

our target is winner prediction so we only need the score of first innings

```
In [16]:  1  total_score = total_score[total_score['innings']==1]
```

```
In [17]:  1  total_score.head(5)
```

Out[17]:

|   | ID | innings | total_run | WinningTeam |
|---|---|---|---|---|
| 0 | 335982 | 1 | 222 | Gujarat Titans |
| 2 | 335983 | 1 | 240 | Royal Challengers Bangalore |
| 4 | 335984 | 1 | 129 | Punjab Kings |
| 6 | 335985 | 1 | 165 | Rajasthan Royals |
| 8 | 335986 | 1 | 110 | Lucknow Super Giants |

- Merging two data frames

```
In [161]:   1  match_df = matches.merge(total_score[['ID','target']], on='ID')
            2
```

```
In [162]:   1  match_df.head()
            2
```

Out[162]:

| | ID | City | Date | Season | MatchNumber | Team1 | Team2 | Venue | TossWinner | TossDecision | ... | WinningTeam | WonBy | Margin | metho |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1312200 | Ahmedabad | 2022-05-29 | 2022 | Final | Rajasthan Royals | Gujarat Titans | Narendra Modi Stadium, Ahmedabad | Rajasthan Royals | bat | ... | Gujarat Titans | Wickets | 7.0 | Nal |
| 1 | 1312199 | Ahmedabad | 2022-05-27 | 2022 | Qualifier 2 | Royal Challengers Bangalore | Rajasthan Royals | Narendra Modi Stadium, Ahmedabad | Rajasthan Royals | field | ... | Rajasthan Royals | Wickets | 7.0 | Nal |
| 2 | 1312198 | Kolkata | 2022-05-25 | 2022 | Eliminator | Royal Challengers Bangalore | Lucknow Super Giants | Eden Gardens, Kolkata | Lucknow Super Giants | field | ... | Royal Challengers Bangalore | Runs | 14.0 | Nal |
| 3 | 1312197 | Kolkata | 2022-05-24 | 2022 | Qualifier 1 | Rajasthan Royals | Gujarat Titans | Eden Gardens, Kolkata | Gujarat Titans | field | ... | Gujarat Titans | Wickets | 7.0 | Nal |
| 4 | 1304116 | Mumbai | 2022-05-22 | 2022 | 70 | Sunrisers Hyderabad | Punjab Kings | Wankhede Stadium, Mumbai | Sunrisers Hyderabad | bat | ... | Punjab Kings | Wickets | 5.0 | Nal |

5 rows × 21 columns

- Removing old team names and updating it to new ones

**Removing old teams/Updating teams new names**

```
In [52]:   1  match_df['Team1'].unique()
           2
```

```
Out[52]: array(['Rajasthan Royals', 'Royal Challengers Bangalore',
         'Sunrisers Hyderabad', 'Delhi Capitals', 'Chennai Super Kings',
         'Gujarat Titans', 'Lucknow Super Giants', 'Kolkata Knight Riders',
         'Punjab Kings', 'Mumbai Indians', 'Kings XI Punjab',
         'Delhi Daredevils', 'Rising Pune Supergiant', 'Gujarat Lions',
         'Rising Pune Supergiants', 'Pune Warriors', 'Deccan Chargers',
         'Kochi Tuskers Kerala'], dtype=object)
```

```
In [53]:   1  teams = [
           2      'Rajasthan Royals',
           3      'Royal Challengers Bangalore',
           4      'Sunrisers Hyderabad',
           5      'Delhi Capitals',
           6      'Chennai Super Kings',
           7      'Gujarat Titans',
           8      'Lucknow Super Giants',
           9      'Kolkata Knight Riders',
          10      'Punjab Kings',
          11      'Mumbai Indians'
          12  ]
```

```
In [54]:   1  match_df['Team1'] = match_df['Team1'].str.replace('Delhi Daredevils', 'Delhi Capitals')
           2  match_df['Team2'] = match_df['Team2'].str.replace('Delhi Daredevils', 'Delhi Capitals')
           3  match_df['WinningTeam'] = match_df['WinningTeam'].str.replace('Delhi Daredevils', 'Delhi Capitals')
           4
           5  match_df['Team1'] = match_df['Team1'].str.replace('Kings XI Punjab', 'Punjab Kings')
           6  match_df['Team2'] = match_df['Team2'].str.replace('Kings XI Punjab', 'Punjab Kings')
           7  match_df['WinningTeam'] = match_df['WinningTeam'].str.replace('Kings XI Punjab', 'Punjab Kings')
           8
           9
          10  match_df['Team1'] = match_df['Team1'].str.replace('Deccan Chargers', 'Sunrisers Hyderabad')
          11  match_df['Team2'] = match_df['Team2'].str.replace('Deccan Chargers', 'Sunrisers Hyderabad')
          12  match_df['WinningTeam'] = match_df['WinningTeam'].str.replace('Deccan Chargers', 'Sunrisers Hyderabad')
          13
```

- Number of matches played by each team.

```
In [31]:   1  fig = px.histogram(match_df, x='WinningTeam',color='WinningTeam')
           2  fig.show()
```



**We want only the matches where D/L is not applied.**

So,removing all matches effected due to rain

```
In [60]:   1  match_df['method'].unique()
           2
```

```
Out[60]:  array([nan, 'D/L'], dtype=object)
```

```
In [61]:   1  match_df['method'].value_counts()
           2
```

```
Out[61]:  D/L    15
          Name: method, dtype: int64
```

```
In [62]:   1  match_df = match_df[match_df['method'].isna()]
           2
```

```
In [63]:   1  match_df.shape
           2
```

```
Out[63]:  (817, 21)
```

- Removing unwanted columns in merged data frame.

```
In [38]:  1  match_df = match_df[['ID','City','Team1','Team2','WinningTeam','target']].dropna()
          2
```

```
In [39]:  1  match_df.head()
          2
```

Out[39]:

|   | ID | City | Team1 | Team2 | WinningTeam | target |
|---|------|----------|---------------------------|---------------------|---------------------------|-----|
| 0 | 1312200 | Ahmedabad | Rajasthan Royals | Gujarat Titans | Gujarat Titans | 131 |
| 1 | 1312199 | Ahmedabad | Royal Challengers Bangalore | Rajasthan Royals | Rajasthan Royals | 158 |
| 2 | 1312198 | Kolkata | Royal Challengers Bangalore | Lucknow Super Giants | Royal Challengers Bangalore | 208 |
| 3 | 1312197 | Kolkata | Rajasthan Royals | Gujarat Titans | Gujarat Titans | 189 |
| 4 | 1304116 | Mumbai | Sunrisers Hyderabad | Punjab Kings | Punjab Kings | 158 |

```
In [40]:  1  match_df.isna().sum()
          2
```

```
Out[40]:  ID             0
          City           0
          Team1          0
          Team2          0
          WinningTeam    0
          target         0
          dtype: int64
```

- Again, merging two data frames

```
In [41]:  1  balls.columns
          2
```

```
Out[41]:  Index(['ID', 'innings', 'overs', 'ballnumber', 'batter', 'bowler',
                 'non-striker', 'extra_type', 'batsman_run', 'extras_run', 'total_run',
                 'non_boundary', 'isWicketDelivery', 'player_out', 'kind',
                 'fielders_involved', 'BattingTeam'],
                dtype='object')
```

```
In [42]:  1
          2  balls['BattingTeam'] = balls['BattingTeam'].str.replace('Kings XI Punjab', 'Punjab Kings')
          3  balls['BattingTeam'] = balls['BattingTeam'].str.replace('Delhi Daredevils', 'Delhi Capitals')
          4  balls['BattingTeam'] = balls['BattingTeam'].str.replace('Deccan Chargers', 'Sunrisers Hyderabad')
          5
          6  balls = balls[balls['BattingTeam'].isin(teams)]
```

```
In [43]:  1  balls_df = match_df.merge(balls, on='ID')
          2
```

```
In [44]:  1  balls_df.head()
          2
```

Out[44]:

|   | ID | City | Team1 | Team2 | WinningTeam | target | innings | overs | ballnumber | batter | ... | extra_type | batsman_run | extras_run | total_run | non_bo |
|---|------|----------|------------------|----------------|---------------|-----|---|---|---|-------------|-----|---------|---|---|---|---|
| 0 | 1312200 | Ahmedabad | Rajasthan Royals | Gujarat Titans | Gujarat Titans | 131 | 1 | 0 | 1 | YBK Jaiswal | ... | NaN | 0 | 0 | 0 | |
| 1 | 1312200 | Ahmedabad | Rajasthan Royals | Gujarat Titans | Gujarat Titans | 131 | 1 | 0 | 2 | YBK Jaiswal | ... | legbyes | 0 | 1 | 1 | |
| 2 | 1312200 | Ahmedabad | Rajasthan Royals | Gujarat Titans | Gujarat Titans | 131 | 1 | 0 | 3 | JC Buttler | ... | NaN | 1 | 0 | 1 | |
| 3 | 1312200 | Ahmedabad | Rajasthan Royals | Gujarat Titans | Gujarat Titans | 131 | 1 | 0 | 4 | YBK Jaiswal | ... | NaN | 0 | 0 | 0 | |
| 4 | 1312200 | Ahmedabad | Rajasthan Royals | Gujarat Titans | Gujarat Titans | 131 | 1 | 0 | 5 | YBK Jaiswal | ... | NaN | 0 | 0 | 0 | |

5 rows × 22 columns

- Creating a new_column in a data frame.

```
In [52]:  1  balls_df['current_score'] = balls_df.groupby('ID')['total_run'].cumsum()
          2
```

```
In [53]:  1  balls_df
          2
```

Out[53]:

| | ballnumber | batter | ... | batsman_run | extras_run | total_run | non_boundary | isWicketDelivery | player_out | kind | fielders_involved | BattingTeam | current_score |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ) | 1 | WP Saha | ... | 0 | 0 | 0 | 0 | 0 | NaN | NaN | NaN | Gujarat Titans | 0 |
| ) | 2 | WP Saha | ... | 0 | 0 | 0 | 0 | 0 | NaN | NaN | NaN | Gujarat Titans | 0 |
| ) | 3 | WP Saha | ... | 1 | 0 | 1 | 0 | 0 | NaN | NaN | NaN | Gujarat Titans | 1 |
| ) | 4 | Shubman Gill | ... | 0 | 0 | 0 | 0 | 0 | NaN | NaN | NaN | Gujarat Titans | 1 |
| ) | 5 | Shubman Gill | ... | 0 | 0 | 0 | 0 | 0 | NaN | NaN | NaN | Gujarat Titans | 1 |
| . | ... | ... ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| | 5 | P Kumar | ... | 0 | 1 | 1 | 0 | 0 | NaN | NaN | NaN | Royal Challengers Bangalore | 80 |
| | 6 | SB Joshi | ... | 1 | 0 | 1 | 0 | 0 | NaN | NaN | NaN | Royal Challengers Bangalore | 81 |
| | 7 | P Kumar | ... | 0 | 0 | 0 | 0 | 0 | NaN | NaN | NaN | Royal Challengers Bangalore | 81 |
| | 1 | SB Joshi | ... | 0 | 1 | 1 | 0 | 0 | NaN | NaN | NaN | Royal Challengers Bangalore | 82 |
| | 2 | SB Joshi | ... | 0 | 0 | 0 | 0 | 1 | SB Joshi | caught | BB McCullum | Royal Challengers Bangalore | 82 |

- Creating a new data frame from given data frames

```
In [69]:  1  final_df = balls_df[['BattingTeam', 'BowlingTeam','City','runs_left','balls_left','wickets_left','current_run_rate'
```

```
In [70]:  1  final_df.head()
```

Out[70]:

| | BattingTeam | BowlingTeam | City | runs_left | balls_left | wickets_left | current_run_rate | required_run_rate | target | result |
|---|---|---|---|---|---|---|---|---|---|---|
| 120 | Gujarat Titans | Rajasthan Royals | Ahmedabad | 131 | 119 | 10 | 0.0 | 6.605042 | 131 | 1 |
| 121 | Gujarat Titans | Rajasthan Royals | Ahmedabad | 131 | 118 | 10 | 0.0 | 6.661017 | 131 | 1 |
| 122 | Gujarat Titans | Rajasthan Royals | Ahmedabad | 130 | 117 | 10 | 2.0 | 6.666667 | 131 | 1 |
| 123 | Gujarat Titans | Rajasthan Royals | Ahmedabad | 130 | 116 | 10 | 1.5 | 6.724138 | 131 | 1 |
| 124 | Gujarat Titans | Rajasthan Royals | Ahmedabad | 130 | 115 | 10 | 1.2 | 6.782609 | 131 | 1 |

- Pre-processing the new data frame

```
In [78]:    1  final_df.describe()
            2
```

Out[78]:

|  | runs_left | balls_left | wickets_left | current_run_rate | required_run_rate | target | result |
|---|---|---|---|---|---|---|---|
| count | 89044.000000 | 89044.000000 | 89044.000000 | 89044.000000 | 89044.000000 | 89044.000000 | 89044.000000 |
| mean | 92.743419 | 62.560251 | 7.535600 | 7.421803 | 10.586524 | 166.071628 | 0.526616 |
| std | 49.758982 | 33.416718 | 2.149637 | 2.237734 | 12.932122 | 28.738283 | 0.499294 |
| min | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 68.000000 | 0.000000 |
| 25% | 53.000000 | 34.000000 | 6.000000 | 6.276923 | 7.210526 | 148.000000 | 0.000000 |
| 50% | 92.000000 | 63.000000 | 8.000000 | 7.485149 | 8.910891 | 166.000000 | 1.000000 |
| 75% | 130.000000 | 91.000000 | 9.000000 | 8.653846 | 10.961538 | 186.000000 | 1.000000 |
| max | 247.000000 | 119.000000 | 10.000000 | 36.000000 | 714.000000 | 247.000000 | 1.000000 |

```
In [79]:    1  final_df.isna().sum()
            2
```

```
Out[79]: BattingTeam          0
         BowlingTeam          0
         City                 0
         runs_left            0
         balls_left           0
         wickets_left         0
         current_run_rate     0
         required_run_rate    0
         target               0
         result               0
         dtype: int64
```

```
In [76]:    1  final_df.shape
```

```
Out[76]: (89044, 10)
```

- Splitting the data into train and test data.

```
In [80]:    1  from sklearn.model_selection import train_test_split
            2
            3  X = final_df.drop('result', axis=1)
            4  y = final_df['result']
            5  X.shape, y.shape
```

```
Out[80]: ((89044, 9), (89044,))
```

```
In [81]:    1  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.01, random_state=42)
```

- Importing algorithms

```
In [82]:   1  from sklearn.linear_model import LogisticRegression
           2  from sklearn.ensemble import RandomForestClassifier
           3  from sklearn.pipeline import Pipeline

In [83]:   1  pipe = Pipeline(steps=[
           2      ('step1',trf),
           3      ('step2',RandomForestClassifier())
           4  ])

In [84]:   1  pipe.fit(X_train, y_train)

Out[84]:  Pipeline(steps=[('step1',
                           ColumnTransformer(remainder='passthrough',
                                             transformers=[('trf',
                                                            OneHotEncoder(drop='first',
                                                                          sparse=False),
                                                            ['BattingTeam', 'BowlingTeam',
                                                             'City'])])),
                          ('step2', RandomForestClassifier())])
```

- Checking accuracy on test data

```
In [85]:   1  y_pred = pipe.predict(X_test)

In [86]:   1  from sklearn.metrics import accuracy_score
           2  accuracy_score(y_pred, y_test)

Out[86]:  0.9966329966329966

In [87]:   1  pipe.predict_proba(X_test)

Out[87]:  array([[0.01, 0.99],
                 [0.  , 1.  ],
                 [0.98, 0.02],
                 ...,
                 [0.97, 0.03],
                 [0.97, 0.03],
                 [0.99, 0.01]])
```

- Saving the current model.

```
In [93]:  1  teams

Out[93]: ['Rajasthan Royals',
          'Royal Challengers Bangalore',
          'Sunrisers Hyderabad',
          'Delhi Capitals',
          'Chennai Super Kings',
          'Gujarat Titans',
          'Lucknow Super Giants',
          'Kolkata Knight Riders',
          'Punjab Kings',
          'Mumbai Indians']


In [94]:  1  final_df['City'].unique()

Out[94]: array(['Ahmedabad', 'Kolkata', 'Mumbai', 'Navi Mumbai', 'Pune', 'Dubai',
                'Sharjah', 'Abu Dhabi', 'Delhi', 'Chennai', 'Hyderabad',
                'Visakhapatnam', 'Chandigarh', 'Bengaluru', 'Jaipur', 'Indore',
                'Bangalore', 'Raipur', 'Ranchi', 'Cuttack', 'Dharamsala', 'Nagpur',
                'Johannesburg', 'Centurion', 'Durban', 'Bloemfontein',
                'Port Elizabeth', 'Kimberley', 'East London', 'Cape Town'],
               dtype=object)


In [95]:  1  import pickle
          2  pickle.dump(pipe, open('pipe.pkl','wb'))
```

### Result :

As I used two machine learning algorithms: - Logistic Regression
Random forest classifier

- Using Logistic regression

```
In [117]:   1  pipe1 = Pipeline(steps=[
            2      ('step1',trf),
            3      ('step2',LogisticRegression())
            4  ])

In [118]:   1  pipe1.fit(X_train, y_train)

/opt/anaconda3/lib/python3.9/site-packages/sklearn/linear_model/_logistic.py:814: ConvergenceWarning:

lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

Out[118]: Pipeline(steps=[('step1',
                          ColumnTransformer(remainder='passthrough',
                                            transformers=[('trf',
                                                           OneHotEncoder(drop='first',
                                                                         sparse=False),
                                                           ['BattingTeam', 'BowlingTeam',
                                                            'City'])])),
                          ('step2', LogisticRegression())])
```

```
In [121]:   1  from sklearn.metrics import accuracy_score
            2  accuracy_score(y_pred, y_test)

Out[121]: 0.8013468013468014

In [115]:   1  pipe1.predict_proba(X_test)

Out[115]: array([[0.11280429, 0.88719571],
                 [0.34430674, 0.65569326],
                 [0.23838816, 0.76161184],
                 ...,
                 [0.84404129, 0.15595871],
                 [0.87445136, 0.12554864],
                 [0.9284528 , 0.0715472 ]])
```

The logistic regression model, despite achieving an accuracy of 80%, still has room for improvement. It is possible that further feature engineering or model tuning may increase the accuracy of the model.

- Using Random Forest Classifier

```
In [97]:   1  pipe = Pipeline(steps=[
           2      ('step1',trf),
           3      ('step2',RandomForestClassifier())
           4  ])

In [98]:   1  pipe.fit(X_train, y_train)

Out[98]: Pipeline(steps=[('step1',
                          ColumnTransformer(remainder='passthrough',
                                            transformers=[('trf',
                                                           OneHotEncoder(drop='first',
                                                                         sparse=False),
                                                           ['BattingTeam', 'BowlingTeam',
                                                            'City'])])),
                          ('step2', RandomForestClassifier())])

In [107]:  1  y_pred = pipe.predict(X_test)

In [108]:  1  from sklearn.metrics import accuracy_score
           2  accuracy_score(y_pred, y_test)

Out[108]: 0.9955106621773289

In [109]:  1  pipe.predict_proba(X_test)

Out[109]: array([[0.  , 1.  ],
                 [0.03, 0.97],
                 [0.99, 0.01],
                 ...,
                 [0.95, 0.05],
                 [0.98, 0.02],
                 [1.  , 0.  ]])
```

However, the random forest classifier has proven to be an incredibly accurate classification technique, achieving a 99% accuracy rate in this report.

This suggests that the random forest classifier may be a more suitable option when accuracy is of utmost importance, such as in medical diagnosis or fraud detection. Overall, these results highlight the importance of selecting appropriate classification techniques based on the specific context and goals of the predictive modeling project.

**SUMMARY :**

The IPL Winning Team Prediction project aimed to use Machine Learning algorithms to predict the winning team of IPL matches. We used two popular algorithms, Logistic Regression and Random Forest Classifier, to analyze various factors such as players' past performances, pitch conditions, and weather conditions to predict the outcome of IPL matches.

After analyzing the data and comparing our predictions to the actual results, we found that the Random Forest Classifier algorithm was the most effective, achieving an accuracy of 99%. This high level of accuracy suggests that our model can be a valuable tool for cricket enthusiasts looking to make informed predictions.

The project's significance lies in its ability to leverage data and Machine Learning algorithms to predict the outcome of complex sporting events such as T20 cricket matches. The accuracy of our model demonstrates the potential of ML algorithms to provide more accurate predictions, ultimately helping cricket enthusiasts make more informed decisions.

Overall, the IPL Winning Team Prediction project has demonstrated the potential of Machine Learning algorithms in predicting the outcome of cricket matches. We hope that this project will inspire further research in using ML algorithms for predicting the outcome of other sports events and contribute to the growing field of sports analytics.