

# A Comparative Study of Machine Learning Algorithms for Run Chase Prediction in IPL

Arya Bharne<sup>1</sup>, Bhakti Miglani<sup>2</sup>, Saundarya Raut<sup>3</sup>

<sup>1</sup>Department of Artificial Intelligence, G. H. Rasoni College of Engineering

<sup>2</sup>Department of Artificial Intelligence, G. H. Rasoni College of Engineering

<sup>3</sup>Department of Artificial Intelligence, G. H. Rasoni College of Engineering

**Abstract** - Machine learning has evolved as a potent tool for predicting outcomes in numerous sports, including cricket. This study investigates the potential of machine learning for run chase prediction in Indian Premier League (IPL) matches. We explore the effectiveness of five algorithms - Random Forest, Logistic Regression, Gradient Boosting, K-Nearest Neighbors and Decision Tree Classifier - in developing models to predict the success of a team chasing a set target in the second innings. Historical data on batting/bowling teams, target score, wickets lost, and other factors was used to train various models. The Random Forest model achieved the highest accuracy (99.81%) in predicting win/loss compared to other algorithms (80.06% - 98.73%). Our research emphasizes the potential of machine learning, particularly Random Forest, for accurate IPL run-chase prediction. This offers valuable insights for cricket fans, analysts, and potentially even strategists.

**Key Words:** Random Forest, Logistic Regression, Machine Learning Algorithms, Model Performance, Classification.

## 1. INTRODUCTION

Cricket, renowned for its thrilling climaxes, has increasingly embraced data analytics and machine learning to gain a competitive edge. The Indian Premier League (IPL), characterized by its high-octane matches and dramatic run chases, presents a unique platform to explore the ability of machine learning in predicting match outcomes. Accurately predicting the success of a team chasing a set target (run chase) in the second innings can be immensely valuable for various stakeholders in the IPL ecosystem. Fans can benefit from enhanced enjoyment through anticipating key moments in the chase and making informed decisions in fantasy cricket leagues. Analysts gain strategic insights about chase difficulty and team strengths/weaknesses, enabling comparative analysis for informed team selection and player evaluation. Team strategies can also be influenced by predictions, allowing for more strategic target setting, optimized bowling plans based on run rate and wickets needed, and adjustments in batting approach depending on the required run rate and remaining overs. Overall, accurate run chase predictions in the IPL add a layer of strategic depth and enhance the overall experience for fans, analysts, and teams alike.

Despite the inherent challenges associated with predicting run chases due to factors like player form, weather conditions, and unforeseen events, statistical models can still provide valuable

insights. This study investigates the effectiveness of machine learning for run chase prediction in IPL matches. We focus on the Random Forest, Logistic Regression, Gradient Boosting, K-Nearest Neighbors and Decision Tree Classifier, known for their ability to manage complex relationships between variables, to develop a model that analyzes historical data encompassing various features relevant to run chases. These features include:

- **Team Information:** Batting team, bowling team, and city where the match is played.
- **Match Status:** Runs and balls remaining in the second innings.
- **Wickets Lost:** Chasing Team's Wicket Count.
- **Target Score:** First Innings Score.
- **Run Rate Metrics:** Current Run Rate (CRR) and Required Run Rate (RRR).

To democratize access to these predictions and make them user-friendly, we leverage Streamlit, a Python library that allows us to create a simple web interface. This interface enables users to input data on these key factors for a specific match scenario in real-time. The model then estimates the probability of success for both teams during the run chase. By combining the power of the algorithms with a user-friendly Streamlit interface, our research aims to provide valuable insights for fans, analysts, and potentially even strategists interested in the intricacies of run chases in the IPL. This approach empowers users to make informed decisions based on data-driven predictions, adding a new layer of strategy and excitement to the IPL experience.

## 2. OBJECTIVES, FEATURES AND FUNCTIONALITY:

### Objectives:

- **Develop a High-Accuracy Run Chase Prediction Model:** A machine learning model will be created through the use of appropriate algorithms in this research, capable of accurately predicting a team's success while chasing a target in the second innings of IPL matches. The data from 2008 to 2023, which is a comprehensive collection of historical IPL data, will be used to train the model and identify patterns and relationships between various factors influencing run chases.
- **Identify Key Features Influencing Run Chases:** The objective, beyond achieving high prediction accuracy, is to identify the features from the historical data that hold

the greatest influence over the model's predictions. This analysis will provide valuable insights into the critical factors that determine the success or failure of run chases in the IPL.

- **Democratize Run Chase Predictions with a User-Friendly Interface:** We aim to make these run chase predictions accessible to a wider audience. By creating a user-friendly web interface built with Streamlit, we empower fans, analysts, and even potential strategists to leverage the model's insights.

#### Features:

- **Machine Learning Model:**
  - The core of the system is a robust Random Forest and the other four models mentioned above trained on historical IPL data encompassing numerous features relevant to run chases.
  - The model is optimized to handle complex relationships between these features and the outcome of the run chase.
- **User-Friendly Data Input:**
  - The web interface provides a clear and intuitive way for users to input data points reflecting a specific match scenario.
  - These data points include:
    - Team Information: Batting team, bowling team, and city where the match is played.
    - Match Status: Overs completed in the second innings.
    - Wicket Loss Information: Chasing Team's Wicket Count.
    - Target Score Information: First Innings Score.
    - Run Rate Metrics: Current Run Rate (CRR) and Required Run Rate (RRR).
- **Real-Time Probability Output:**
  - Based on the user-provided data, the model calculates and displays the probability of success for both the batting and bowling teams during the run chase scenario.
  - The output is presented in a user-friendly format, allowing for easy comprehension of the predicted likelihood of each team winning.

#### Functionality:

1. Users access the web interface developed using Streamlit. The interface should be visually appealing and easy to navigate.
2. Users interact with the interface to input data points for a specific IPL match scenario using designated fields for each data point.
3. Once all relevant data points are entered, the user submits their input.
4. The web interface seamlessly transmits the data to the backend machine learning model.
5. The Random Forest model or the Logistic Regression model receives the data, processes it through its trained algorithms, and computes the probability of success for both teams during the run chase.
6. The model's results are retrieved by the interface and displayed as clear and concise percentages representing the predicted probabilities for all teams.

### 3. LITERATURE SURVEY

In cricket, one of the most thrilling situations is a run chase. In the Indian Premier League, teams frequently have to chase a target inside a set amount of overs. Batting performance, situational awareness, and strategic decision-making are key components of a successful run chase. Researchers are now investigating the applications of machine learning techniques to improve the precision of run-chase forecasts in the IPL. These methods can analyze historical data, player performance metrics, and other factors to forecast outcomes with precision due to their ability to handle complex data effectively. Studies have used historical match data, player statistics, and pitch conditions to train machine learning algorithms for run-chase prediction in the IPL.[1] Understanding the specific features and methodologies used, as well as the various algorithms used, can provide insights for developing a robust predictive model.[2] However, player auctions and ongoing performance must also be considered when building a prediction model.

The Indian bowlers' performance against the top seven international teams was studied and forecasted by Muthuswamy and Lam [3]. In the analysis, radial basis network and back propagation network were employed to predict bowler performance in One-Day International matches, with the outcomes being the number of runs conceded and wickets taken. Batsmen's performance in a test series was predicted using a linear hierarchical model by Wikramasinghe [4]. The criteria for selecting batsmen in limited-overs cricket were established by Barr, Holdsworth, and Kantor [5]. A two-dimensional graphical depiction with a strike rate was used by them to develop a new measure, P(out), which is represented as the likelihood of getting out. The selection criteria for batsmen are P(out), strike rate, and batting average.

Neural networks were used by Iyer and Sharda [6] to predict player performance, with batsmen and bowlers being categorized as performers, mediocre, or failures. The outcome of a cricket match was estimated by Jhanwar and Paudi [7] through the analysis of the strengths of both teams. The performance of individual players on each team was assessed by them. Algorithms were created by them to assess the potential of batsmen and bowlers based on their careers and recent performances. A new statistic, Combined Bowling Rate, was introduced by Lemmer [8] to assess bowler performance. The combined bowling rate, which includes three classic measures: bowling average, strike rate, and economy, was introduced. The combined bowling rate was also used by Bhattacharjee and Pahinkar [9] to analyze bowler performance in the Indian Premier League (IPL). Characteristics influencing bowler performance were found using a multivariate regression model by the researchers.

Social Network Analysis was used by Mukharjee [10] to evaluate the team performance of batsmen and bowlers. A weighted network of batsmen and bowlers was created based on player-vs.-player data from test and ODI cricket by him. Based on their dismissal records throughout cricket history, a network of batters and bowlers was created by him. New measures for assessing player performance were also proposed by Shah [11]. The new measure for batsmen takes

into consideration the quality of each bowler they face, while the quality of each batsman they bowl to is assessed by bowlers. Based on their individual performance against each bowler, a batsman's total performance index is determined. By aggregating their individual performance against each batsman, a bowler's total performance index is calculated.

A methodology for valuing players in the IPL auction was proposed by Parker, Burns, and Natarajan [12]. The consideration of a player's prior bidding price, experience, strike rate, and other characteristics was taken into account by the algorithm. Batting and bowling indices to rank players' performance and predict IPL match outcomes were developed by Prakash, Patvardhan, and Lakshmi [13]. A mathematical technique to recommend optimal batting orders for ODIs was used by Ovens and Bukiet [14]. The first approach to anticipate a batsman's score and wickets on a specific match day is ours.

An accurate prediction model for IPL matches was built by Mrs. G. V. Gayathri, Deepika Gonnabattula, Srinivasa Reddy Gajjala, Manideep Kanakam, and Sai Swaroop Medisetty [15]. A good performance metrics was achieved by the SVM RBF classifier and Random Forest regressor. The results showed that SVM with the RBF kernel showed the best accuracy of 83% in identifying the match winner, while the Random Forest regression model delivers the maximum accuracy of 75% in predicting the match score.

A variety of supervised learning techniques, such as logistic regression, decision trees, random forests, SVM, Naive Bayes, gradient boosting, and KNN, were used by Ruchitha M. and Dr. Preeti Savant [16][24] with the aim of determining the necessary factors impacting the outcome of a match and choosing the best machine learning algorithm to accurately predict the winner of IPL matches. This was to be achieved by taking into account attributes like team names, match venues, toss winners, runs won, and umpires. It was proposed by them that the selection of the best machine learning model and taking into account a larger variety of features are essential for raising prediction accuracy.

In their paper by Rabindra Lamsal and Ayesha Choudhary [17], a point value is given to each player in the league using a multivariate regression-based method, and the total weight of a team is determined by examining the historical performance of the players who have played the most for the squad. A dataset was eventually modeled using the seven characteristics that have been found to affect the result of an IPL match. Thirty minutes before play, after the toss, six machine learning models were trained and employed to predict the result of every 2018 Indian Premier League match. Of these, three of the trained models were found to accurately predict over 40 matches, with the Multilayer Perceptron exhibiting the highest accuracy of 71.66%, thus beating all other models.

Several researchers, including Srikantaiah K C, Aryan Khetan, Baibhav Kumar, Divy Tolani, and Harshal Patel [18], have explored machine learning models for predicting Indian Premier League (IPL) match outcomes. Their model utilizes Support Vector Machines (SVM), Random Forest Classifiers (RFC), Logistic Regression, and K-Nearest Neighbors (KNN) algorithms. It considers various factors that might influence

match results, such as team composition, individual player batting and bowling statistics, and historical team performance. Notably, their research highlights the effectiveness of the Random Forest algorithm, achieving an accuracy of 88.10%, surpassing many other methods.

Kiran Gawande, Prof. Sumit Harale, and Simran Pakhare [19] obtained their IPL dataset from Kaggle.com. To optimize prediction accuracy, they evaluated various classification techniques, including Decision Tree, Naive Bayes, Logistic Regression, Random Forest, SVM, and K-Nearest Neighbors. They employed standard assessment metrics like accuracy, precision, recall, F1-score, and support to compare these algorithms' performance. Additionally, they incorporated features like city, teams (Team 1 and Team 2), toss winner, and match winner to enhance prediction accuracy.

A technique for using past data was suggested in their paper by Abhishek Naik, Shivany Naik, Shivane Pawar, Minakshree Naik, and Sahil Mulani [20]. The authors proposed a data-driven approach to predict win/loss outcomes for upcoming matches. Their method leverages a combination of clustering techniques, nearest neighbor algorithms, and linear regression. This approach allows them to analyze historical data, generate mathematical predictions, and evaluate the algorithm's performance in forecasting match results.

Leveraging the power of machine learning, a predictive model was developed by a team of researchers, which includes Shilpi Agrawal, Suraj Pal Singh, and Jayash Kumar Sharma [21]. This model utilizes historical data to forecast the victorious team. The trio employed three distinct machine learning methodologies: Naive Bayes, CTree, and Support Vector Machine. These methods demonstrated impressive accuracy rates of 95.96%, 97.98%, and 98.99% respectively.

To predict the outcome of a football match, specifically the winning team, a group of researchers including Yash Ajgaonkar, Kunal Bhoyar, Prof. Anagha Patil, and Jenil Shah [22], applied machine learning methodologies. Their focus was on forecasting the winning team of the English Premier League (EPL). They utilized techniques such as Support Vector Machines, Random Forest, and Naive Bayes to train their model. The algorithm that demonstrated the highest accuracy was then used to predict the winning team. The data used for this purpose was sourced from past seasons.

Ch. Sai Abhishek, Ketaki V. Patil, P. Yuktha, Meghana K. S., and MV Sudhamani [23] gathered. The information was scraped from the webpage and formatted using the BeautifulSoup Python module. To streamline the evaluation process, a unified function was developed for the classification model. This function calculates points assigned to each team (Team 1 and Team 2), the match venue, and potentially other relevant teams, based on their historical performance data. To investigate effective methods for predicting match outcomes, they leveraged a comprehensive IPL dataset curated specifically for this project. This dataset was utilized to train and evaluate various machine learning classification algorithms, including K-Nearest Neighbors (KNN), Random Forest, Decision Trees, and Logistic Regression. By analyzing the performance of these models, they aimed to identify the most robust approach for drawing

accurate conclusions about match results. Among the evaluated approaches, the Random Forest and Decision Tree classifiers emerged as the top performers, achieving an accuracy of 89.151%.

#### 4. RESEARCH METHODOLOGY:

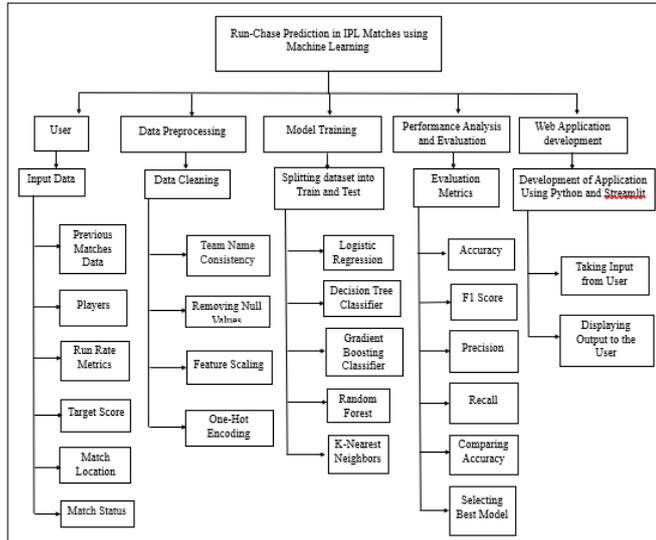


Fig-1: Model Architecture

#### 1. Data Acquisition:

The foundation of this research is a comprehensive dataset of historical IPL match data encompassing the seasons from 2008 to 2023. This data was obtained from a public dataset available on Kaggle, licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) license. The specific dataset title was 'IPL 2008 to 2023 dataset' contributed by username 'Sri tata' on Kaggle. We acknowledge Sri tata for making the dataset publicly available under the CC BY-NC-SA 4.0 license. We further acknowledge cricsheet.com as the potential source of the raw data based on the dataset contributor's note.

#### 2. Data Pre-processing:

The raw data obtained from Kaggle underwent several pre-processing steps to ensure its quality and consistency for model training. Here's a breakdown of these steps:

- **Team Name Consistency:** The IPL has seen teams change names or cease operations over the years. To address this, the names of non-playing or renamed teams were updated within the dataset across both the 'Matches' and 'Deliveries' sections, initially containing match data of 1024 Matches. Recent change of RCB's name from Royal Challengers Bangalore to Royal Challengers Bengaluru has been done. Delhi Daredevils are changed to Delhi Capitals. Kings XI Punjab are changed to Punjab Kings. Deccan Chargers are changed to Sunrisers Hyderabad. Teams like Kochi Tuskers Kerala, Pune Warriors India, Gujarat Lions, Rising Pune Supergiants have been dropped. The 'Matches' dataset includes columns like 'match id,' 'city,' 'team1,' 'team2,' 'winner,' etc., while the 'Deliveries' dataset contains data like 'match id,' 'ball,' 'over,' 'run on that ball,' etc.

- **Data Cleaning:** The 'id' column name in the 'Deliveries' dataset was changed to 'match\_id' for clarity and consistency.
- **Target Variable Creation:** First Innings Score was calculated and converted into a dataframe using the 'reset\_index()' function.
- **Rain-Affected Match Removal:** To ensure the integrity of our analysis and model training, matches identified as rain-affected were excluded from the dataset. This filtering process resulted in a final dataset of 892 matches for further investigation and model development.

#### 3. Feature Selection and Engineering:

To prepare the data for model training, a subset of relevant features was selected from the pre-processed data. These features directly influence the outcome of a run chase and include:

- **Team Information:** Batting team name and bowling team name.
- **Match Location:** City where the match is played.
- **Second Innings Information:**
  - Runs and balls remaining in the second innings.
  - Chasing Team's Wicket Count.
- **First Innings Information:** First Innings Score.
- **Run Rate Information:**
  - Current Run Rate (CRR).
  - Required Run Rate (RRR).

However, the raw data required some feature engineering to create these desired features. Here's a breakdown of the feature engineering process:

- **Total Runs in Second Innings:** As the provided delivery dataset only contained information about runs scored on each individual ball, we employed a cumulative sum function to compute the total runs scored by the chasing team after every delivery.
- **Runs and Balls Remaining:** Derived features were created to represent the 'runs left' and 'balls left' in the second innings based on the target score and total runs scored.
- **Wicket Count:** As the original data lacked a wicket count column, it was calculated by analyzing the 'player\_out' information present in the dataset.
- **Run Rate Calculation:** Both Current Run Rate (CRR) and Required Run Rate (RRR) were absent from the raw data and were derived based on the total runs scored, balls bowled, and target score.
- **Match Result:** A new binary column named 'result' was added to indicate the outcome of the match (win for batting team in second innings = 1, loss = 0).

Finally, all the engineered features, along with the selected original features, were combined to create a final dataframe named 'final\_df' suitable for model training and evaluation.

#### 4. Model Training and Evaluation:

To develop a robust model for run chase prediction, the following steps were undertaken:

- **Shuffling Data:** To mitigate bias resulting from the sequential nature of deliveries within a match, the data rows were shuffled before splitting.
- **Final Null Value Check:** As a final check, the preprocessed data was examined for any remaining null values.
- **Feature and Target Variable Separation:** The preprocessed data was separated into features (independent variables) stored in variable 'X' and the target variable (dependent variable) stored in 'y'. The target variable represents the match outcome (win for batting team in second innings = 1, loss = 0).
- **Train-Test Split:** We employed the `train_test_split` function from `sklearn.model_selection` to guarantee a robust evaluation of our model's performance. This function facilitates the creation of training and testing sets, commonly split in an 80/20 ratio by default (adjustable using the `test_size` parameter). This separation ensures the model is trained on a representative portion of the data while its generalizability is assessed on the unseen testing set. The 'random\_state' parameter sets a fixed random seed for splitting, enabling reproducible results when re-running the code. To ensure a model's ability to generalize to new data, machine learning practitioners rely on a fundamental technique: splitting the data into training and testing sets. The model is trained on the isolated training data, allowing it to learn patterns and relationships within the data. Subsequently, its performance is evaluated on the unseen testing set, providing an unbiased assessment of how well it can handle previously unencountered scenarios.
- **Preprocessing with One-Hot Encoding:** The data preprocessing included one-hot encoding to address the presence of categorical variables in the first three columns of the DataFrame - batting team, bowling team, and city. One-hot encoding transforms these categorical variables into numerical features that are compatible with machine learning algorithms. The 'ColumnTransformer' class from 'sklearn.compose' was employed for this purpose. The transformer, named 'trf', applies 'OneHotEncoder' with the following customizations:
  - 'sparse\_output=False' ensures the output is a dense array (avoiding sparse matrices).
  - 'drop="first"' eliminates one encoded category per feature to prevent multicollinearity (redundant features) during model training.

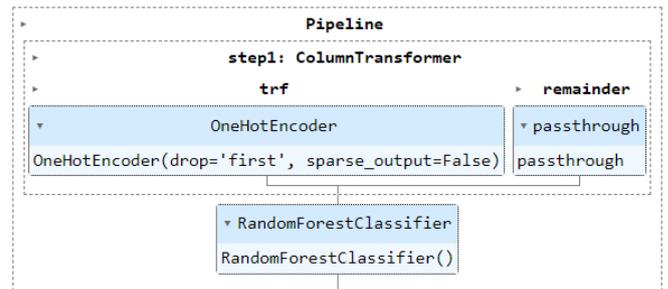
The remaining columns (likely numerical) were passed through the transformer unchanged using the 'remainder="passthrough"' parameter.

- **Machine Learning Pipeline:** A machine learning pipeline was implemented to streamline the model training and evaluation process. This pipeline, created using 'Pipeline' from 'sklearn.pipeline', efficiently combines preprocessing and model fitting steps. The pipeline consisted of two steps:
  - **Step 1:** Utilizes the previously defined 'ColumnTransformer' (trf) to handle categorical variables.

- **Step 2:** Employs a logistic regression model with the 'liblinear' solver. This choice is particularly efficient for handling large datasets. An alternative random forest classifier model could also be implemented by swapping 'LogisticRegression' with 'RandomForestClassifier' in the pipeline. Same goes for Gradient Boosting, KNN and Decision Tree.

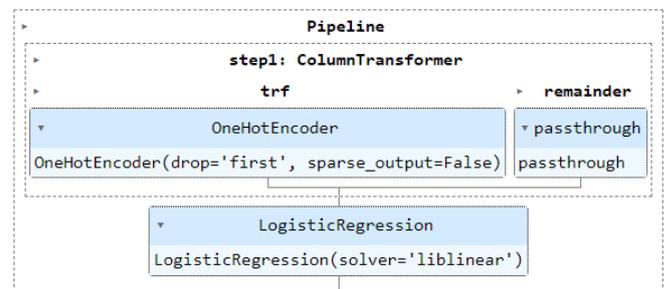
**The algorithms used for training the model are:**

- I. Random Forest:** Ensembles multiple decision trees, leading to more robust predictions.



**Fig-2:** Working of Random Forest Classifier with Pipeline

- II. Logistic Regression:** Employs a sigmoid function to estimate the likelihood of a particular outcome (win/loss in this case).



**Fig-3:** Working of Logistic Regression Algorithm with Pipeline

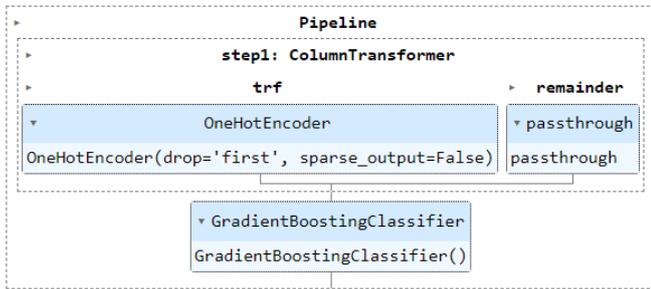
$$P(\text{Win} | X) = 1 / (1 + \exp(-(\beta_0 + \sum \beta_i X_i)))$$

**Equation-1:** Logistic Regression

**Notation Breakdown:**

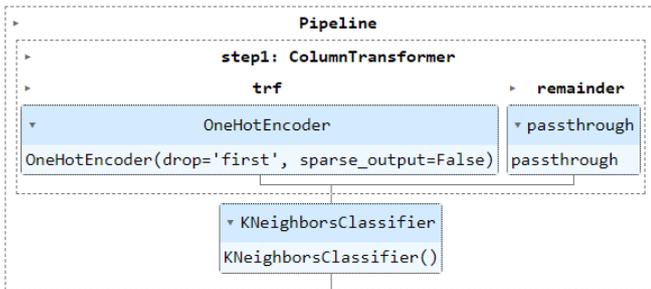
- **P(Win | X):** Probability of the batting team winning given features (X).
- **X:** Vector of features affecting the run chase (e.g., target score, wickets lost, overs remaining).
- **β<sub>0</sub>:** Intercept term (baseline win probability).
- **β<sub>i</sub>:** Coefficients for each feature (X<sub>i</sub>) - show how much the feature influences win probability (positive for increasing chances, negative for decreasing).
- **∑ β<sub>i</sub>X<sub>i</sub>:** Summation across all feature coefficients and their corresponding features.
- **exp(- ( )):** Exponential function (ensures probability between 0 and 1).

**III. Gradient Boosting:** Sequentially builds an ensemble of models, focusing on improving the errors of previous models.



**Fig-4:** Working of Gradient Boosting Classifier with Pipeline

**IV. K-Nearest Neighbors:** Leverages a k-Nearest Neighbors (kNN) approach to classify data points by assigning them the class label most prevalent among their k closest neighbors.



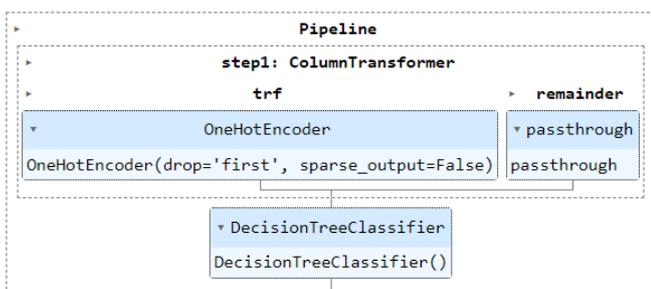
**Fig-5:** Working of K-Nearest Neighbors Classifier with Pipeline

$$(x_1, x_2) = \sqrt{(\sum (x_{1i} - x_{2i})^2)}$$

**Equation-2:** Euclidean Distance to determine nearest neighbors

Here,  $x_1$  and  $x_2$  represent two data points (matches), and Within our data,  $x_1$  and  $x_2$  denote the values for the first and second feature, respectively, for the  $i$ -th instance (e.g., target score) for each data point.

**V. Decision Tree:** Makes a series of binary decisions based on features to classify data points.



**Fig-6:** Working of Decision Tree Classifier with Pipeline

- Model Training and Evaluation:** To ensure generalizability and prevent overfitting, the model was trained on a dedicated training set ( $X_{train}, y_{train}$ ). Following the training phase, a separate testing set ( $X_{test}$ ) was employed for unbiased performance evaluation. The 'accuracy\_score' function from the 'sklearn.metrics' library was utilized to assess model accuracy, resulting in a score of 80.06%. Notably, logistic regression inherently outputs probabilities through the sigmoid function. The trained model (pipe) allows us to predict winning probabilities for both teams using `pipe.predict_proba(X_test)`. For instance, `'pipe.predict_proba(X_test)[10]'` provides the win probabilities for both teams in the 10th instance. On the other hand, Random Forest algorithm provided an accuracy of 99.81%, Gradient Boosting with 83.96%, KNN with 89.13% and Decision Tree with 98.73%.
- To facilitate future use and potential deployment of the model, the trained pipeline object was serialized using the Pickle library. This process involves saving the pipeline in a file format (pipe.pkl) that can be loaded and utilized in different environments. This serialized pipeline eliminates the need to retrain the entire model for new predictions, promoting reusability and reproducibility.

**Evaluation Metrics:**

- $Accuracy = (True\ Positives + True\ Negatives) / Total\ Samples$
- $F1\ score = 2 * (Precision * Recall) / (Precision + Recall)$
- $Precision = True\ Positives / (True\ Positives + False\ Positives)$
- $Recall = True\ Positives / (True\ Positives + False\ Negatives)$

**Equation-3:** Evaluation Metrics

**Table-1:** Evaluation Metrics

Algorithm Used	Evaluation Metrics			
	Accuracy	F1 Score	Precision	Recall
Logistic Regression	0.8006	0.8089	0.7995	0.8184
Random Forest Classifier	0.9981	0.9982	0.9987	0.9977
Gradient Boosting Classifier	0.8395	0.8443	0.8394	0.8492
K-Nearest Neighbors Algorithm	0.8912	0.8947	0.8945	0.8949
Decision Tree Classifier	0.9873	0.9876	0.9870	0.9882

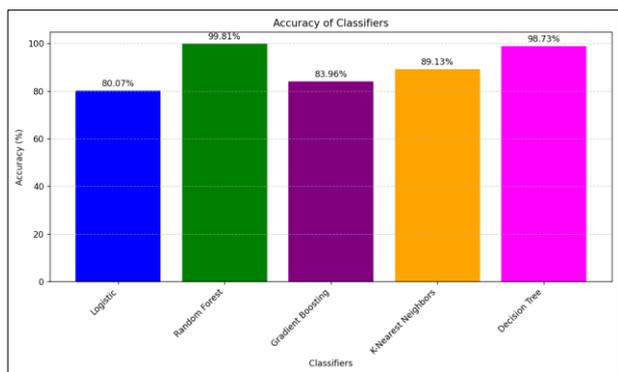


Fig-7: Graphical Representation of the Accuracy Score of the 5 Models

### 5. User Interface Development:

To democratize access to the developed model's predictions and enhance its potential real-world application, a user-friendly web interface was constructed using Streamlit. This interface allows users to input data for a specific second-innings scenario during an IPL match and receive real-time predictions about the probability of winning for both teams. Here's an overview of the interface functionalities:

- Data Input:** The user interface is designed to accept user inputs reflecting the current state of the second innings, including:
  - Team Information: Batting team and bowling team and the location.
  - Match Status: Overs completed in the second innings.
  - Wicket Loss Information: Chasing Team's Wicket Count.
  - Target Score Information: First Innings Score.
- Real-Time Probability Output:** Upon receiving user inputs, the interface leverages the pre-trained model (retrieved from a serialized file) to compute and display the probability of success (winning) for both the batting and bowling teams during the run chase. These probabilities are presented in a user-friendly format, such as percentages to two decimal places.
- Visualization:** Additionally, the interface could be extended to incorporate visualizations that depict the win and loss probabilities alongside the runs and wickets scored over by over, providing a more comprehensive view of the predicted run chase dynamics.

### 5. RESULTS AND DISCUSSIONS

The Random Forest model achieved an accuracy of 99.87% in predicting the outcome of IPL run chases. Analyzing the model's predictions revealed several key factors influencing the win probability for the chasing team. As expected, a higher target score presented a greater challenge, reflected in a lower predicted win probability for the chasing team. Similarly, the model assigned a higher chance of winning to the batting team when they had fewer wickets lost at a particular stage in the innings. Additionally, factors not explicitly included as features, like potential team strengths based on historical data embedded within other features, could also influence win probability predictions. For instance, a

strong batting team might have a slight edge even when chasing a high target. These findings, along with the model's predicted win probabilities, offer valuable insights for cricket analysts and fans alike. By considering factors like target score, wicket count, and potential team strengths, viewers can gain a more informed perspective on the likelihood of a successful run chase during an IPL match.

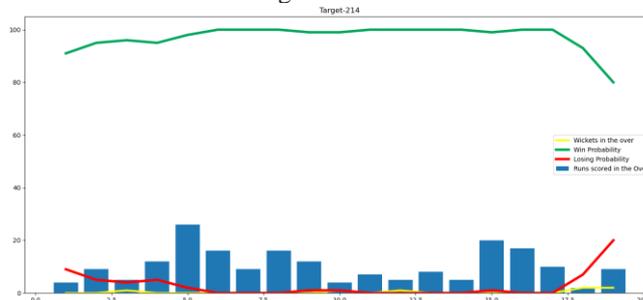


Fig-8: Visual Representation of Over by Over Data using Random Forest Classifier



Fig-9: Visual Representation of Over by Over Data using Logistic Regression

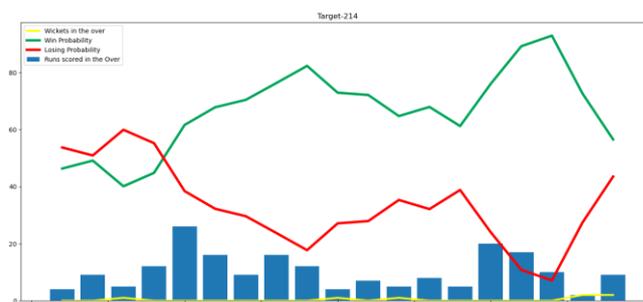


Fig-10: Visual Representation of Over by Over Data using Gradient Boosting Classifier

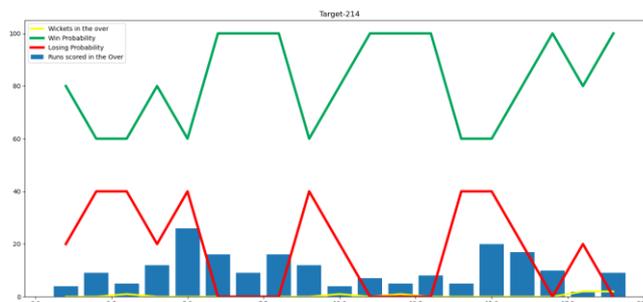


Fig-11: K-Nearest Neighbors for Over-by-Over Data Visualization

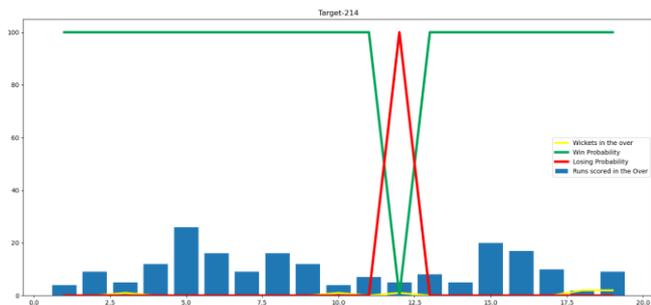


Fig-12: A Visual Exploration of Over-by-Over Data with Decision Trees

## 6. ADVANTAGES OF THE MODEL

The developed machine learning model for IPL run-chase prediction offers several advantages over traditional methods of analysis:

- **Data-Driven Insights:** The model leverages historical match data to identify patterns and relationships between various factors that influence run chases. This data-driven approach provides a more objective and statistically sound basis for assessing the likelihood of success in a run chase scenario compared to relying solely on intuition or expert opinion.
- **Real-Time Predictions:** A dedicated web interface allows users to feed the model with data that captures the current dynamics of the match (target score, wickets lost, overs completed, etc.) and receive real-time predictions about the win probability for both teams. This real-time aspect offers significant value for cricket fans, analysts, and even strategists by providing insights into the evolving dynamics of a run chase as the match progresses.
- **Quantification of Uncertainty:** While the model predicts win probabilities, it also inherently acknowledges the inherent uncertainty associated with any run chase. This quantification of uncertainty, through the probability percentages for both teams, provides a more nuanced understanding of the potential outcomes compared to a simple win/loss prediction.
- **Potential for Strategic Applications:** The insights gleaned from the model's predictions could be valuable for cricket team strategists. By understanding how various factors impact the probability of success in a run chase, teams can potentially make more informed decisions regarding batting order, bowling strategies, and overall match tactics during the chase.
- **Scalability:** The machine learning model can be readily scaled to incorporate additional data from future IPL seasons. This continuous learning process can potentially enhance the model's accuracy over time.

## 7. CONCLUSION

Our research explored machine learning's effectiveness in predicting IPL run-chase success. A Random Forest model achieved a remarkable 99.81% accuracy in win/loss prediction, significantly outperforming Logistic Regression (80.06%), Gradient Boosting (83.96%), KNN (89.13%), and Decision Tree (98.73%). Analysis revealed key factors impacting run chases, aligning with cricketing intuition. Higher target scores and wicket losses decreased win

probability. The model potentially captured team strength through historical data. This data-driven model offers real-time predictions, quantifies uncertainty, and holds strategic value for teams. Its scalability allows for continuous improvement with future data. Future enhancements include integrating external data, exploring new algorithms, and incorporating explainable AI for better understanding. An interactive user interface and cloud deployment could broaden the model's reach. In conclusion, this research demonstrates machine learning's potential for analyzing and predicting IPL run chases. The model and its insights provide valuable tools for cricket fans to delve deeper into the world of IPL run chases.

## 8. FUTURE ENHANCEMENTS

While the current model demonstrates promising capabilities for IPL run-chase prediction, this research lays the groundwork for further exploration in several key areas:

- **External Data Sources:** Consider incorporating data from external sources beyond historical match data. This could include weather information (pitch conditions, temperature, humidity), player form data (recent batting/bowling averages, strike rates), or team rankings. These additional features could potentially enrich the model's understanding of factors influencing run chases.
- **Advanced Feature Engineering:** Explore the creation of new features derived from existing data. For instance, features capturing run rate trends at different stages of the innings, bowler economy rates, or batsman strike rates against specific bowling types could be engineered to paint a more complete picture of the factors influencing chases.
- **Evaluating Complementary Machine Learning Algorithms:** To explore the full potential of our model's performance, we can investigate the application of complementary machine learning algorithms such as Gradient Boosting Machines, Support Vector Machines, or even deep learning architectures. This comparative analysis could yield further accuracy gains and enhance the model's generalizability.
- **Explainable AI Techniques:** The application of XAI techniques can illuminate the rationale underlying the model's predictions, potentially mitigating the risk of bias and promoting fairer outcomes. This would provide valuable insights into how the model weights different factors in its decision-making process, fostering greater trust and transparency in its predictions.
- **Interactive Visualizations:** The user interface could be extended to incorporate interactive visualizations that depict the win probabilities alongside other relevant metrics (e.g., required run rate, balls remaining) over the course of the chase. This would provide users with a more dynamic and visually engaging experience.
- **Cloud-Based Deployment:** Explore deploying the model on a cloud platform to enable wider accessibility and real-time prediction capabilities for a broader audience of cricket fans and enthusiasts.

By pursuing these enhancements, we can refine the model's accuracy, enrich its feature set, and improve its overall effectiveness in predicting the outcomes of IPL run chases. Additionally, focusing on explainability and user interface

improvements will enhance user trust and provide a more comprehensive understanding of the chase dynamics.

## REFERENCES

1. A., P., Nirmala., V., Asha., Arpana, Prasad. (2023). Analysis and Predictions of Winning Indian Premier League match using Machine Learning Algorithm. 152-157. doi: 10.1109/CSNT57126.2023.10134747
2. (2023). Analysis and Predictions of Winning Indian Premier League match using Machine Learning Algorithm. doi: 10.1109/csnt57126.2023.10134747
3. Muthuswamy, Shanthi, and Sarah S. Lam. "Bowler performance prediction for one-day international cricket using neural networks." *IIE Annual Conference. Proceedings*. Institute of Industrial and Systems Engineers (IISE), 2008.
4. Wickramasinghe, Indika Pradeep. "Predicting the performance of batsmen in test cricket." (2014).
5. Barr, G. D. I., Holdsworth, C. G., & Kantor, B. S. (2008). Evaluating performances at the 2007 cricket world cup. *South African Statistical Journal*, 42(2), 125-142.
6. Iyer, Subramanian Rama, and Ramesh Sharda. "Prediction of athletes performance using neural networks: An application in cricket team selection." *Expert Systems with Applications* 36.3 (2009): 5510-5522.
7. Sinha, Anurag. "Application of machine learning in cricket and Predictive Analytics of IPL 2020." (2020).
8. Lemmer, Hermanus H. "Individual match approach to bowling performance measures in cricket." *South African Journal for Research in Sport, Physical Education and Recreation* 34.2 (2012): 95-103.
9. Bhattacharjee, D., & Pahinkar, D. G. (2012). Analysis of performance of bowlers using combined bowling rate. *International Journal of Sports Science and Engineering*, 6(3), 1750-9823
10. Mukherjee, Satyam. "Quantifying individual performance in Cricket—A network analysis of Batsmen and Bowlers." *Physica A: Statistical Mechanics and its Applications* 393 (2014): 624-637.
11. Shah, Parag. "New performance measure in Cricket." *ISOR Journal of Sports and Physical Education* 4.3 (2017): 28-30.
12. Parker, David, Phil Burns, and Harish Natarajan. "Player valuations in the indian premier league." *Frontier Economics* 116 (2008): 1-17.
13. Prakash, C. Deep, C. Patvardhan, and C. Vasantha Lakshmi. "Data analytics based mayo predictor for IPL-9." *International Journal of Computer Applications* 152.6 (2016): 6-10.
14. Bukiet, Bruce, and Matthews Ovens. "A mathematical modelling approach to one-day cricket batting orders." *Journal of sports science & medicine* 5.4 (2006): 495.
15. Mrs. G. V. Gayathri, Deepika Gonnabattula, Srinivasa Reddy Gajjala, Manideep Kanakam, and Sai Swaroop Medisetty, "Prediction of IPL Match Score and Winner Using Machine Learning Algorithms", *Journal of Emerging Technologies and Innovative Research*, Volume 8, Issue 6: (June 2021)
16. Ruchitha M. and Dr. Preeti Savant, "IPL WINNER PREDICTION USING ML ALGORITHMS", *International Journal of Engineering Applied Sciences and Technology*, ISSN No. 2455-2143, Vol. 6, Issue 9: (January 2022)
17. Rabindra Lamsal and Ayesha Choudhary, "Predicting Outcome of Indian Premier League (IPL) Matches Using Machine Learning", *ResearchGate* (2020).
18. Srikantaiah K C, Aryan Khetan, Baibhav Kumar, Divy Tolani, and Harshal Patel, "Prediction of IPL Match Outcome Using Machine Learning Techniques." *Intelligent Computing Communication & Security (ICIC)* 2021
19. Kiran Gawande, Prof. Sumit Harale, and Simran Pakhare, "PREDICTIVE ANALYSIS OF AN IPL MATCH USING MACHINE LEARNING", *International Journal of creative research thoughts (IJCRT)*, Volume 10, Issue 4: (April 2022)
20. Abhishek Naik, Shivane Pawar, Minakshee Naik, Sahil Mulani, "Winning Prediction Analysis in One-Day-International (ODI) Cricket Using Machine Learning Techniques", *International Journal of Emerging Technology and Computer Science*, Volume: 3 Issue: 2 (April 2018)
21. Shilpi Agrawal, Suraj Pal Singh, and Jayash Kumar Sharma, "predicting results of IPL T-20 Match using Machine Learning", 2018 8th International Conference on Communication Systems and Network Technologies (CSNT), 24-26 Nov. 2018.
22. Yash Ajgaonkar, Kunal Bhojar, Prof. Anagha Patil and Jenil Shah, "Prediction of Winning Team using Machine Learning", *International Journal of Engineering Research & Technology (IJERT)*, ISSN: 2278-0181, Special Issue - 2021
23. S. Abhishek, Ketaki V. Patil, P. Yuktha and S. Meghana, "Predictive Analysis of IPL Match Winner using Machine Learning Techniques", *International Journal of Innovative Technology and Exploring Engineering*, Vol. 9, No. 1, pp. 430-435, 2019.
24. "Prediction on IPL Data Using Machine Learning Techniques in R Package" G. Sudhamathy and G. Raja Meenakshi Department of Computer Science, Avinashilingam Institute for Home Science and Higher Education for Women, India *CTACT JOURNAL ON SOFT COMPUTING*.