

# Cricket Match Win Prediction using Machine Learning

*Naga Sai Sindura Pandrangi, Nimeesha Vakacharla, Madhura Deshmukh, Pavan Srivatsav Devarakonda, SaiBhargav Podili, Sai Naga Sanjana Chippada*

*Abstract — Given how dynamic and unpredictable One Day International (ODI) cricket is, predicting its results is a difficult but exciting task. Many factors affect match results, including player form, team dynamics, the impact of the toss, venue circumstances, and even the weather, even though analysts, coaches, and fans have long looked out for reliable prediction techniques. Due to its intricacy, accurate forecasts have proven difficult to achieve. This research tackles the problem by applying machine learning techniques to forecast the outcome of ODI matches, with a particular emphasis on the chasing team's chances of winning or losing. The project uses a rich dataset from cricsheet.org that spans more than 20 years (2001–2024). To find complex patterns, the project uses advanced machine learning algorithms, especially Random Forest. The prediction models are incorporated into an intuitive Flask-developed web interface to improve usage. Users of this site can enter match data in real-time and get instant outcome predictions. The application serves to a broad spectrum of users, such as cricket teams seeking strategic insights, analysts seeking a better understanding of match dynamics, and fans want to get more involved in the action. Through accurate projections and practical insights into the factors influencing match outcomes, this initiative seeks to improve cricket analysis, give stakeholders data-driven decision-making tools, and improve the overall ODI cricket experience for spectators worldwide.*

## I. INTRODUCTION

Cricket is one of the most popular sports in the whole world, especially now when it swings into focus with the one-day internationals (ODIs). It's typical and unpredictable nature, in India, draws millions of fans. Thus, due to these intertwined relationships between each other, like the form of the players, team strategies, match conditions, information about the venue, predicting the outcome of the ODI on a single day is a tough task. The outcome can change dramatically with each ball.

By identifying intricate patterns in historical data that conventional approaches frequently miss, this initiative seeks to improve the accuracy of ODI match outcome forecasts by

utilizing data-driven techniques, especially machine learning. The goal of the project is to create a machine learning model that can accurately predict match outcomes (win/loss) by using a comprehensive dataset from 2001–2024 that was obtained from cricsheet.org. This dataset includes match-level details such as player statistics, team performance, toss choices, and match conditions. To deliver real-time forecasts, an easy-to-use web interface will relate to the trained model, providing both cricket experts and fans with valuable information.

## II. OBJECTIVE

This project mainly intends to develop a machine learning prediction model capable of predicting the result of an ODI cricket match based on whether the chasing team will win or lose. The model will be trained on historical match data to include vital features such as performance of the team in the first innings (runs scored, wickets lost), performance of the chasing team in the second innings (runs scored, wickets lost, match conditions), and dynamic currency variables during the match (run rates, wickets remaining, balls left). By using these components, the initiative hopes to provide cricket teams, experts, and fans with a powerful tool. Cricket clubs aiming to enhance their strategies, analysts seeking a deeper comprehension of match dynamics, and sports fans hopeful to all make up the target audience.

## III. PROBLEM STATEMENT

### A. CHALLENGES IN PREDICTING ODI OUTCOMES

Several factors affect the outcome of ODI cricket matches. These consist of weather, team performance, player form, toss choices, match location, and other variable elements. Because of the many interdependencies between these variables, traditional methods such as expert opinions and simple statistical models struggle to forecast the results.

### B. NEED FOR A DATA-DRIVEN APPROACH

Although traditional methods and expert judgment have been used to analyze match outcomes, they frequently overlook the game's complexity and variety. Finally, a data-driven, machine-learning approach produces more accurate predictions by capturing the underlying correlations between distinct characteristics, allowing for more robust modeling. Machine learning algorithms are capable of digesting vast datasets and identifying subtle patterns that humans may ignore.

### C. PROJECT SCOPE AND GOALS

This study seeks to revolutionize cricket match prediction by utilizing powerful machine learning techniques to overcome the limits of traditional forecasting methods. The study aims to create a strong prediction model capable of properly forecasting One Day International (ODI) match outcomes by combining complex algorithms such as Random Forest, XGBoost, and Linear Regression with extensive historical match data. To improve accessibility and practical use, the created predictive system will be implemented via a user-friendly web interface, allowing cricket fans and analysts to acquire real-time match outcome predictions by entering pertinent match information.

## IV. PROPOSED SOLUTION

Our study aims to create a comprehensive predictive model for projecting One Day International (ODI) cricket match results using advanced machine learning techniques. The major purpose is to develop a dependable system that offers strategic information to teams, experts, and cricket fans. We hope to anticipate match outcomes using historical data and advanced algorithms based on crucial criteria such as player performance, team dynamics, ambient conditions, and strategic decisions like coin tossing.

The thorough technique consists of several critical stages: laborious data collecting and preprocessing, careful feature selection, rigorous model training and evaluation, and finally, a user-friendly implementation via a Flask web interface. This novel system will allow users to enter match facts and receive real-time forecasts of likely match outcomes.

The project's goal goes beyond simple prediction. We want to create a robust analytical tool that provides deep insights into match dynamics and gives stakeholders data-driven decision-making capabilities. By converting complex historical data into actionable forecasts, we hope to provide a

cutting-edge resource that bridges the gap between statistical analysis and practical cricket strategy.

## V. METHODOLOGY

### A. DATA COLLECTION AND PREPROCESSING

Gathering raw match data is the first step in the suggested solution. The dataset includes a variety of match information, such as team names, player statistics, match locations, toss options, and more. To guarantee that the data is clean and ready for analysis, preparation operations will resolve missing values, inconsistencies, and format issues.

### B. FEATURE EXTRACTION AND SELECTION

Following preprocessing, important features will be chosen for the modeling stage. Toss decisions, team statistics, player performance metrics, and venue-related features are a few examples of these features. The most important variables for prediction will be found through feature selection, which lowers dimensionality without sacrificing model accuracy.

### C. MODEL SELECTION AND JUSTIFICATION

We will investigate several machine learning models, with a particular emphasis on Random Forest, XGBoost and Linear Regression. For managing intricate relationships and interactions between features, Random Forest is especially helpful, whereas Linear Regression can be used to find linear relationships and comprehend the impact of individual features. To find out which model works best for this issue, both will be tested.

### D. MODEL TRAINING AND EVALUATION

The preprocessed dataset will be used to train the chosen models, and their performance will be optimized through hyperparameter tuning. To evaluate the models' effectiveness in projecting match results, performance metrics such as accuracy, precision, recall, and F1 score will be used. The model's ability to generalize to unknown variables will be assessed using a separate test dataset.

### E. MODEL DEPLOYMENT VIA WEB INTERFACE

After the top-performing model has been chosen, it will be made available through an intuitive Flask-developed web interface. Users can visit the interface and receive real-time predictions about the expected outcome of a match by entering relevant match information such as teams, players, and match

conditions. The interface will be simple to use and straightforward, ensuring that consumers have a favorable experience.

### FLOWCHART

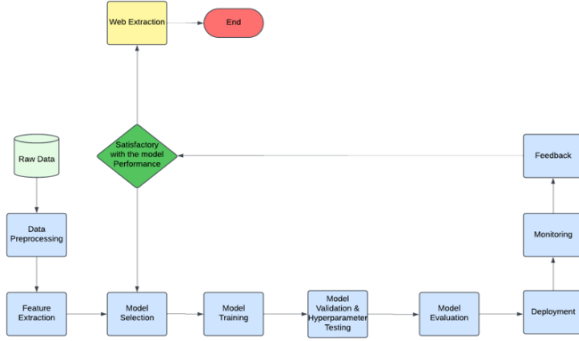


FIG 1: Methodology diagram

## VI. DATA EXTRACTION AND PROCESSING

The first step in this project involved data extraction and processing to transform raw match data into structured datasets suitable for analysis and model building. The data was sourced from cricsheet.org, a repository of detailed cricket match data. We specifically extracted all One Day International (ODI) matches played between 2001 and 2024. Each match was represented as a JSON file containing comprehensive information, including players, umpires, match details, and ball-by-ball records.

### A. MATCH-LEVEL DATA EXTRACTION

To capture high-level match information, all JSON files were aggregated into a match-level dataset. The dataset contains key details like match ID (derived from the JSON filename), city, date, venue, teams, toss winner, toss decision, match outcome, result margin, and player statistics.

#### 1) FEATURE SELECTION AND JUSTIFICATION:

Fields like match type, team composition, player of the match, and result details were selected for their significance in identifying match patterns and outcomes. Columns like outcome were excluded as they were redundant. More precise and interpretable information was already captured in fields like result margin and result\_margin\_type.

The match-level dataset provided an overview of the matches and allowed for the analysis of trends such as: Total number

of ODIs, Team participation, toss decisions, Win patterns, Seasonal match distributions

2) OBSERVATIONS: A total of 2,389 matches were recorded after processing, spanning two decades of ODI cricket.

### B. BALL-BY-BALL DATA EXTRACTION

For deeper match analysis, detailed ball-by-ball data was extracted into a separate data frame (df\_balls). This data was sourced from the ball-by-ball logs in each match's JSON file and contained features such as batter, bowler, runs scored, extras, and dismissals. A sequential ball number was added to each inning, helping track match progression and identify key moments (e.g., partnerships, turning points). This feature mapped the team to either the first or second innings based on the toss decision, ensuring accurate grouping of deliveries by innings for analysis. Different types of extras (e.g., wide, no-ball) were recorded, providing insights into how extras impacted game outcomes.

#### 1) PURPOSE AND JUSTIFICATION:

Adding ball numbers as a feature allowed us to track match momentum and identify patterns in runs, wickets, and extras. The innings feature facilitated analysis of team performance in both innings (e.g., comparing teams' performances while chasing versus setting a target). Extras type and dismissal details added value for deeper analysis, allowing us to analyze specific events and their impact on match outcomes.

#### 2) METRICS CAPTURED:

Key details captured included over number, ball number, runs scored, extras, batter and bowler names, dismissal type (if applicable), and fielders involved in wickets.

## VII. DATA CLEANING AND VALIDATION

To ensure data quality and consistency, both the match-level (df\_matches) and ball-by-ball (df\_balls) datasets underwent thorough cleaning and validation processes.

### A. HANDLING MISSING VALUES:

Missing values in fields like extras were filled with zeros to maintain uniformity in the ball-by-ball dataset, preventing null entries from affecting calculations or analysis. Matches with missing critical data (e.g., match outcome or toss

decision) were excluded from analysis to avoid introducing ambiguity or bias.

#### B. REMOVING REDUNDANT OR AMBIGUOUS COLUMNS:

Fields like outcome by were removed from the match-level data frame, as the same information was captured in more descriptive columns like result margin and result\_margin\_type. Irrelevant or unused columns were also dropped to simplify the dataset and reduce computational complexity.

#### C. CROSS-REFERENCING DATA FRAMES

The ball-by-ball data (df\_balls) was cross-verified against the match-level dataset (df\_matches) to ensure consistency in key features like Teams, Toss Decisions, Outcomes.

#### D. VALIDATING DATA INTEGRITY

Ensured that the sequential ball numbers in the ball-by-ball dataset matched the overs and deliveries in each inning, with no duplicate or missing data. Verified that each match reported in df\_balls also existed in df\_matches, guaranteeing the databases' relationship integrity.

### VIII. FEATURE EXTRACTION:

To support exploratory data analysis (EDA) and predictive modeling, two essential dataframes—Batsman Stats and Bowler Stats—were established, recording precise performance measures for players based on ball-to-ball data.

#### A. BATSMAN STATS

The Batsman Stats dataframe includes key performance characteristics for each hitter throughout all matches and innings, such as innings played, total runs, run distribution, dismissals, batting average, and strike rate. These elements allow for examination of a batter's performance tendencies, including scoring patterns, consistency, and impact on match outcomes.

#### B. BOWLER STATS

The Bowler Stats data frame contains critical parameters for evaluating bowler performance, including Balls Bowled, Wickets Taken, Overs Bowled, Runs Conceded, Economy Rate, and Strike Rate. These elements provide information about a bowler's performance efficiency, ability to take wickets, and how economically they bowled during matches.

#### C. COMPREHENSIVE PLAYER DATASET

The individual Batsman Stats and Bowler Stats data frames were merged into a single dataset, referred to as df\_finals, using the player's name as the common key. This unified dataset includes all essential statistics for each player, providing a comprehensive assessment of their performance. This feature extraction approach is critical for exploratory data analysis (EDA) since it provides detailed insights into player and team performance, allowing for the detection of patterns and trends in batting and bowling. The unified dataset also provides a good platform for developing prediction models, improving the capacity to study aspects that influence match outcomes.

### IX. EXPLORATORY DATA ANALYSIS (EDA) WITH DATA VISUALIZATION

During this stage of our analysis, we used Exploratory Data Analysis (EDA) to find significant insights from the dataset. We hoped to uncover underlying trends, find probable abnormalities, and establish significant relationships by thoroughly reviewing the data. Using Plotly Express, we were able to create interactive and engaging visualizations that provided a thorough view of the data's major properties. These graphical representations were critical in converting raw data into actionable intelligence, which aided our approach to informed, data-driven decision-making. Few of the key visualizations created during the analysis is:

#### A. PLAYER OF MATCH VS. RUNS SCORED

This bar plot visualizes the number of runs scored by different players and whether they received the Player of the Match award. It helps to understand the correlation between a player's performance (runs scored) and the likelihood of being awarded Man of the Match.

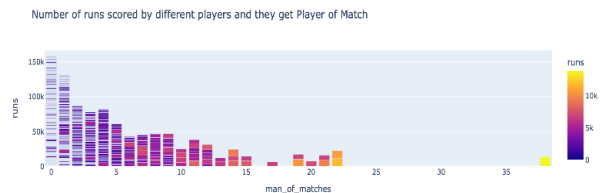


FIG 2: *Player of Match vs. Runs Scored bar chart*

The plot makes it evident which players have scored the most runs and whether their efforts earned them the Player of the Match award. Notable patterns are also revealed, such as

whether players who consistently score a lot of runs are more likely to win Player of the Match. This emphasizes the importance of batsmen's contributions to match-winning performances.

#### B. PLAYER OF MATCH VS. WICKETS TAKEN

This bar plot visualizes the number of wickets taken by different players and their association with receiving the Player of the Match award. This helps us analyze whether players who excel in bowling (in terms of wickets taken) are more likely to be recognized as Player of the Match.

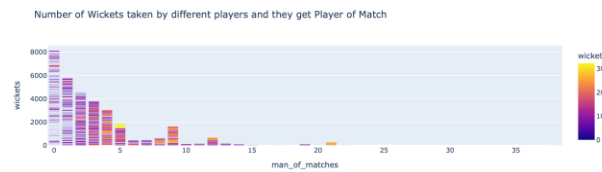


FIG 4: *Player of Match vs. Wickets Taken bar chart*

The following plot also shows if getting wickets is a major component in earning Player of the Match or if other elements, such as runs scored, are more important. Bowlers who often take numerous wickets are more likely to receive Player of the Match awards, emphasizing their importance in games. The pattern highlights how important good bowling is in helping teams win games.

#### C. PLAYER OF THE MATCH VS. NUMBER OF MATCHES PLAYED

This scatter plot visualizes the distribution of players based on the number of matches played and how frequently they have been awarded Player of the Match.

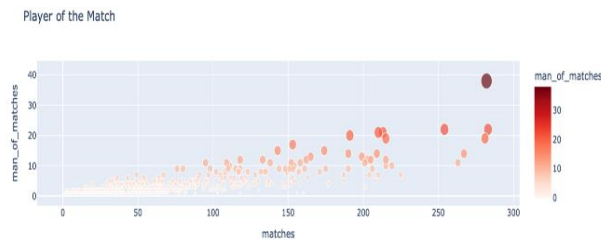


FIG 3: *Player of the Match vs. Number of Matches Played scatter plot*

This scatter plot depicts whether there is a relationship between the players' total number of Player of the Match honors and the number of matches they have played. Regardless of their career length, athletes who perform well throughout games stand out and are more likely to be named

Player of the Match. This underlines the importance of players maintaining high performance levels throughout their careers. Top-tier players can be discovered by observing clusters of great performers.

#### D. BATSMAN STRIKE RATE

A scatter plot was generated to analyze the strike rate of batsmen across the dataset, showcasing their efficiency in scoring runs.

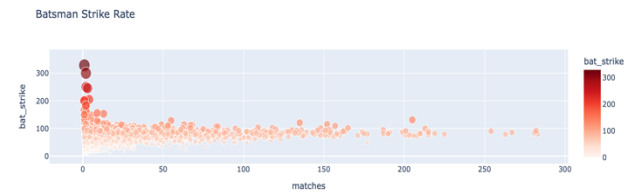


FIG 5: *Batsman Strike Rate scatter plot*

This scatter plot provides a clear visualization of which batsmen have exceptional strike rates and their performance consistency across matches. Players with a high number of matches and a consistent strike rate are easily distinguishable, showcasing their value in limited-overs formats. High strike rates are critical indicators of aggressive and impactful batting. Players with high strike rates, regardless of the number of matches played, are often instrumental in setting or chasing challenging totals. This graph highlights players who balance consistency (matches played) with efficiency (strike rate), showcasing their value to their teams.

#### E. MOST SIXES BY A BATSMAN

A scatter plot visualizes the relationship between matches played and the number of sixes hit by batsmen, with points sized and colored based on the number of sixes.

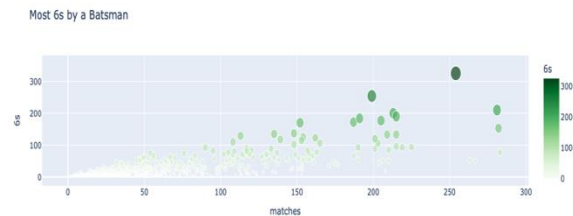


FIG 6: *Most Sixes by a Batsman scatter plot*

Batsmen who consistently hit sixes are pivotal in accelerating scoring rates and shifting match momentum. The graph highlights players who combine longevity (matches played) with power-hitting abilities (sixes). A cluster of elite

six-hitters is evident, identifying players capable of making a significant impact during critical match phases. This underscores the role of power-hitting in modern cricket, especially in formats where aggressive scoring is crucial.

#### F. TOP 10 BATSMEN BASED ON BATTING AVERAGE

A bar chart depicts the top 10 batsmen ranked by their batting averages.

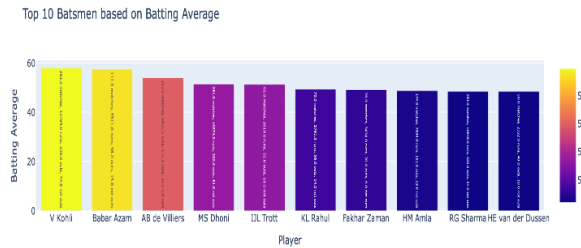


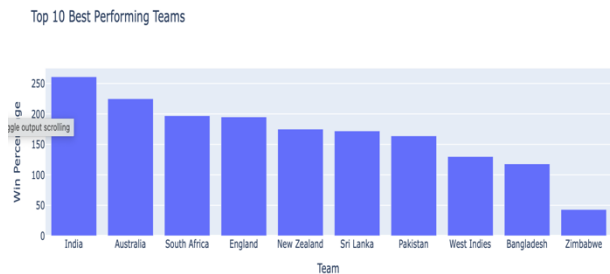
FIG 8: Top 10 batsmen based on Batting Average

A batsman's batting average indicates how many runs they can score per dismissal, so this graphic emphasizes players who consistently perform well. Batsmen who consistently stay at the crease in crucial match situations are those with high averages and a sizable number of not-outs. Especially under pressure, the graph assists in identifying important players who are useful for establishing or pursuing goals. It also demonstrates the relationship between consistency and match experience, since players with more matches typically have stronger averages.

#### G. TOP 10 BEST PERFORMING TEAMS

The bar chart illustrates the top 10 teams based on their win percentages. The x-axis represents the team names, while the y-axis shows their corresponding win percentages.

FIG 7: Top 10 Best Performing Teams Bar Chart



This visualization identifies the most successful teams in terms of match victories, highlighting their dominance in the sport. Teams with consistently high win percentages likely

excel in key areas such as batting depth, bowling efficiency, and fielding agility. The graph offers a clear benchmark for comparing team performances across leagues or time periods. By focusing on the top 10 teams, the analysis underscores the competitive edge held by these teams and helps uncover patterns of sustained success, such as reliance on key players or effective team strategies.

#### H. CORRELATION HEATMAP

A heatmap was generated using Seaborn to visualize the correlation between different numerical features in the dataset. This helped identify patterns and relationships between variables.

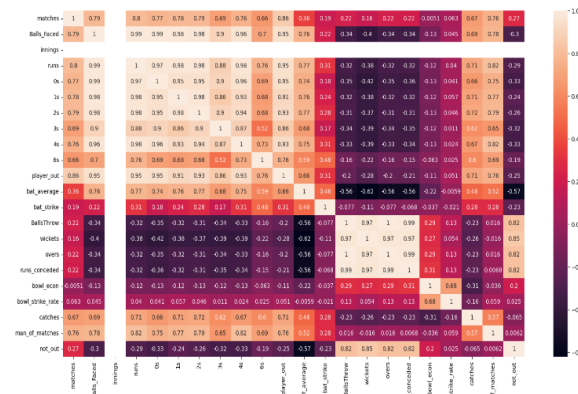


FIG 9: Correlation Heatmap

The visualization identified features that are strongly correlated, both positively and negatively. This analysis was crucial for feature selection, as highly correlated features may lead to redundancy in the dataset.

#### I. MOST TIMES OUT BY A BATSMAN VS. MATCHES PLAYED

A scatter plot was created to explore the relationship between the number of times a batsman got out and the matches they remained not out.

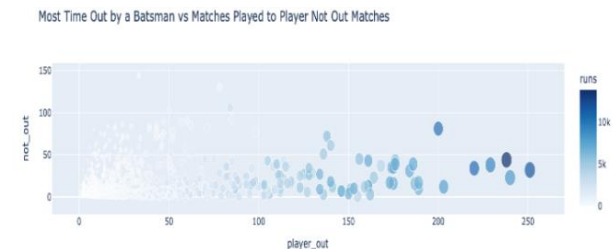


FIG 10: *Most Time Out by a Batsman vs matches played to player not out matches*

Players with high dismissal counts often correspond to those with higher match participation, as expected. The distribution highlights batsmen who were frequently not out despite playing many matches, showcasing their batting stability or role as finishers. Batsmen with a balance of runs and dismissals appear distinct on the plot, aiding in identifying key players for further analysis.

## X. FEATURE ENGINEERING

Building a machine learning model necessitates feature engineering since it directly affects the model's ability to discover patterns and relationships in the data. The model's prediction performance is enhanced by converting unstructured input into meaningful features, allowing it to better grasp the complexity and nuance of the problem domain. By minimizing noise and guaranteeing that the data is representative and relevant, good feature engineering raises the model's overall accuracy and dependability. By highlighting the dataset's most important features, it establishes the groundwork for thorough analysis and modeling.

After constructing the final dataset, further refinement was necessary to prepare it for effective analysis and predictive modeling. These additional steps were critical in transforming the data into a form that could be used to build a reliable and accurate machine learning model. The feature engineering process ensured that the dataset was not only complete and clean but also enriched with meaningful features that could capture the dynamic nature of an ODI match. Below are the key steps taken to enhance the dataset, along with the rationale behind them:

### A. CALCULATION OF TARGET SCORES

The target score for each match was calculated as follows:

**Grouping and Summing Runs for the First Innings:** The dataset was grouped by `match_id` and `innings`, with the total runs scored (`runs_total`) in the first innings being calculated. This provided a benchmark for the second innings.

**Filtering First Innings:** We filtered the data to retain only the rows where `innings == 1`, representing the runs scored by the team batting first.

**Adding Target Column:** A new column, `target`, was created by adding 1 to the first innings' total runs. This target represented the number of runs that the chasing team needed to win in the second innings.

**Merging Target Scores with Match Metadata:** The target column was then merged with the `df_matches` dataset, linking each match with its corresponding target score and providing crucial context for the second innings.

### B. FILTERING MATCHES FOR RELEVANCE

Several filtering steps were applied to ensure that only relevant and meaningful data was included in the final dataset:

**Teams of Interest:** Only matches involving major international teams like India, Australia, and England were retained to ensure a focus on competitive and consistent matches.

**Filtering Valid Winners:** Matches with missing or unclear outcomes were excluded, leaving only those with valid winner data. This helped maintain the quality of the dataset and its usefulness for modeling.

### C. ADDING BALL-LEVEL DATA AND FILTERING FOR THE SECOND INNINGS

A key component of the feature engineering process was adding more granular ball-by-ball data:

- **Current Score:** The `current_score` column was created by calculating the cumulative sum of runs (`runs_total`) for each match and inning, providing a real-time picture of the chasing team's progress.
- **Filtering Second Innings Data:** We focused on the second innings, where the chasing team's performance is pivotal in determining the match outcome. Only rows where `innings == 2` were retained.

### D. CALCULATING BALLS LEFT

To assess the time remaining for the chasing team to meet the target:

- **Identifying Legal Deliveries:** A binary flag (`is_legal_delivery`) was created to mark whether a delivery was legal (1) or illegal (0).
- **Cumulative Legal Deliveries:** The number of legal deliveries bowled (`legal_balls_bowled`) was accumulated for each match and innings.



- **Balls Left:** The total number of balls left in an innings was calculated by subtracting the legal balls bowled from the total balls (300), providing an estimate of how much time the chasing team had to reach the target.

#### E. ADDING WICKETS LEFT:

The wickets\_left column was created to represent the number of wickets the chasing team had remaining:

- **Cumulative Dismissals:** We summed the isWicketDelivery flag to track the number of wickets lost by the chasing team.
- **Wickets Left:** The number of wickets left was calculated by subtracting the cumulative number of wickets lost from 10, reflecting the team's remaining resources.

#### F. Run Rate Metric:

Two important run rate metrics were added to capture the team's scoring efficiency and required pace:

- **Current Run Rate (CRR):** The CRR was calculated as  $(\text{current\_score} * 6) / \text{balls\_bowled}$ , providing a snapshot of the batting team's scoring rate.
- **Required Run Rate (RRR):** The RRR was calculated as  $(\text{runs\_left} * 6) / \text{balls\_left}$ , showing the required pace for the chasing team to meet the target.

#### G. Binary Match Outcome Indicator

A new column, result, was created to serve as the target variable for predictive modeling:

1 indicated that the batting team won.

0 indicated that the batting team lost. This column marked the outcome of each match and was crucial for training the machine learning model.

#### H. Batting and Bowling Teams

Two new columns, Batting Team and Bowling Team, were added to clearly identify the teams involved in each delivery. These features enabled a more detailed analysis of each team's performance in relation to the match outcome.

##### 1) Creation of Final Features

The dataset was further refined to include key features for predictive analysis:

Columns BattingTeam, BowlingTeam, and city provided important contextual information. Features like runs\_left, balls\_left, wickets\_left, and target captured the match situation at each point. The current\_run\_rate and required\_run\_rate reflected the pace of play and the pressure on the chasing team. The result column indicated whether the batting team successfully chased the target.

This final, refined dataset, named winningPred, was now ready for model building. The feature engineering process, from creating new metrics to filtering and encoding the data, ensured that the dataset was rich with meaningful features, clean, and well-suited for training machine learning models.

#### 2) FINAL DATA PREPROCESSING FOR MODEL BUILDING

Once the winningPred dataset was prepared, the next step involved data preprocessing to make it suitable for machine learning. The following techniques were applied to transform the dataset for predictive modeling:

- **One-Hot Encoding:** The categorical variables (BattingTeam, BowlingTeam, and city) were encoded using one-hot encoding. This technique converted these variables into numerical values while avoiding multicollinearity by dropping the first category for each feature.
- **Data Splitting:** The dataset was split 80/20 between training and testing sets. The model was built using the training set, and its performance was evaluated using the testing set.

With the data preprocessed and the necessary transformations applied, the next steps involved training machine learning models on the dataset to predict the outcome of ODI matches.

### XI. MODELING

This section discusses the various aspects of modeling the problem of predicting cricket match winners using Linear Regression, XGBoost and Random Forest algorithms, specifically addressing assumptions, modeling practices, analysis depth, tool usage, and complexity considerations.

#### A. ASSUMPTIONS IN MODELING



Modeling cricket match outcomes involves several critical assumptions to align the dataset and problem definition with the capabilities of the chosen algorithms:

*a) DATASET INTEGRITY:*

The dataset is assumed to be clean and representative of real-world cricket matches, with no major anomalies or inconsistencies. Missing or erroneous values in features such as "city" or "target" have been addressed using imputation techniques or data-cleaning methods.

*b) FEATURE-TARGET RELATIONSHIP:*

Features like "runs\_left," "balls\_left," "wickets\_left," and "required\_run\_rate" are presumed to carry meaningful information about the match's state and are directly correlated with the target outcome ("result"). The chosen features provide enough contextual information to model the temporal and situational dynamics of a cricket match effectively.

*c) CATEGORICAL FEATURE HANDLING:*

Non-numeric variables such as "battingTeam" and "bowlingTeam" are encoded into numerical representations without losing their relative importance. This ensures the models can interpret these variables meaningfully.

*d) TEMPORAL DEPENDENCIES:*

Cricket is a dynamic sport where the match situation evolves ball by ball. The dataset is assumed to provide sufficient granularity (e.g., ball-by-ball data) to capture these changes. The problem is framed as a classification task, assuming binary outcomes (win/loss), with no draws or inconclusive results affecting the modeling.

*e) MODEL SUITABILITY:*

Random Forest is used because of its capacity to capture intricate, non-linear patterns, while Linear Regression is selected as a baseline to assess its efficacy in modeling linear relationships. XGBoost was chosen due to its exceptional ability to manage sizable datasets and its resilience in identifying intricate dependencies and non-linear interactions. When dealing with complex patterns in data, XGBoost is especially helpful because it can more accurately model the underlying complexities than more straightforward models like linear regression. It is assumed that the performance of these models can be assessed meaningfully using standard evaluation metrics.

*f) OUTCOME REPRESENTATION:*

The target variable ("result") is encoded as 0 for a loss and 1 for a win, ensuring compatibility with classification algorithms and evaluation metrics.

*B. CORRECT MODELING PRACTICES*

Adhering to best practices in machine learning ensures the robustness and reliability of the models. Below are the detailed steps undertaken:

*1) DATA PREPROCESSING:*

- Handling Missing Values
- Encoding Categorical Variables
- Scaling Features

*2) FEATURE ENGINEERING:*

- Derived Features.
- Feature Selection.

*c) ALGORITHM SELECTION:*

- **Linear Regression:** Selected as a baseline model to evaluate its ability to predict match outcomes assuming linear relationships among variables. Its simplicity makes it interpretable but may underperform in capturing non-linear dependencies.
- **Random Forest:** Chosen for its ability to handle non-linear patterns, heterogeneous data, and feature interactions effectively. It provides important scores, offering valuable insights into which variables contribute most to predictions.
- **XGBoost:** Selected for its capacity to manage intricate, non-linear relationships: XGBoost is excellent at identifying complex patterns and feature interactions, which qualifies it for forecasting results in intricate situations where several variables influence the outcome, such as multiple-factor cricket matches. Excellent robustness and performance: XGBoost is renowned for its effectiveness in managing sizable datasets and its capacity to generalize effectively, even in the presence of noisy or imperfect data. Its ability to rank feature importance makes it easy to determine which features have the greatest influence on the prediction process, and it is a popular option for attaining high accuracy.

#### d) HYPERPARAMETER TUNING:

For Random Forest, hyperparameters like the number of trees, tree depth, and minimum samples per split were optimized using grid search or random search, improving model accuracy and generalization.

#### C. IN-DEPTH ANALYSIS

A comprehensive analysis was conducted to ensure the model's performance is reliable and interpretable:

##### 1) EXPLORATORY DATA ANALYSIS (EDA)

Visualizations of feature distributions were created to understand the underlying data. Correlation matrices were generated to identify relationships between features and the target variable.

##### 2) MODEL EVALUATION

Metrics such as accuracy, precision, recall, and F1-score were calculated to measure model performance on both training and validation datasets. ROC curves and AUC scores were analyzed to understand the trade-off between true positives and false positives.

##### 3) COMPARISON BETWEEN MODELS:

The prediction capabilities of the three models—Linear Regression, Random Forest, and XGBoost—were assessed. Although linear regression could produce coefficients that could be understood, it had trouble managing intricate relationships. But XGBoost and Random Forest did better, especially when it came to capturing non-linear dependencies.

- **Linear Regression Performance:**

A straightforward and easily comprehensible model, linear regression requires a linear relationship between input features and output. Understanding feature importance is made much easier by its coefficients, which show us how each feature affects the target variable.

Accuracy	84.69%
Precision	0.33
Recall	0.33

F1 Score	0.34
----------	------

However, because it presumes linear relationships between features, linear regression has limitations. It has trouble with complex feature interactions, like the one between "balls\_left" and "wickets\_left." Linear regression is unable to adequately model the non-linear relationship that these features most likely have.

For instance, based on the match context, the relationship between the number of balls left and wickets left may change (for instance, a high number of balls left with fewer wickets could imply a higher chance of a win). This interaction is not captured by linear regression, which leads to decreased precision, recall, and overall model performance.

- **Random Forest Performance:**

An ensemble learning technique called Random Forest uses several decision trees to generate predictions. It excels at capturing feature interactions and non-linear relationships without the need for feature transformations. The flexibility of Random Forest is exceptional; it can automatically recognize intricate patterns in the data.

Accuracy	91.60%
Precision	0.92
Recall	0.90
F1 Score	0.91

Random Forest is capable of handling complicated interactions and dependencies, making it appropriate for real-world data with non-linear feature correlations.

The model has a greater accuracy (91.60% vs. 84.69%), showing that it is significantly more effective at predicting the outcome.

Random Forest has much superior precision, recall, and F1 scores, indicating that it not only performs better at recognizing right outcomes but also lowers false positives and false negatives.

- **XGBoost Performance:**

Extreme Gradient Boosting, or XGBoost, is a powerful boosting algorithm that creates trees one after the other, with

each one learning to correct the flaws of the previous one. It can manage non-linear interactions in the same way as Random Forest can, but boosting adds another layer of complexity.

Accuracy	89.07%
Precision	0.89
Recall	0.89
F1 Score	0.89

#### 1) KEY ADVANTAGE OVER RANDOM FOREST:

Because XGBoost builds trees consecutively and corrects faults repeatedly, it can compute faster and efficiently, particularly when working with huge datasets.

XGBoost typically outperforms Random Forest when adjusting; nevertheless, the hyperparameter tuning process might be more challenging and time-consuming.

#### 2) REASONS FOR USING RANDOM FOREST FOR THIS PROBLEM:

- **Handling Complex Interactions:** The dataset's key features, like "balls\_left" and "wickets\_left," call for a model that can represent intricate relationships. For this, Random Forest's ensemble of decision trees is perfect because it can naturally handle feature interactions and makes no assumptions about linearity.
- **Interpretability and Stability:** Compared to Random Forest, which is easier to use and more reliable right out of the box, XGBoost is more complex and may be more difficult to understand, even though it requires more careful tuning. On this dataset, Random Forest is more robust, has lower variance, and is less likely to overfit.
- **Performance:** Random Forest's higher accuracy, precision, recall, and F1 score demonstrate that it is better at anticipating outcomes (whether the team wins or loses) and performs well across a variety of criteria.

#### 3) CONFUSION MATRIX OVERVIEW:

The confusion matrix for the Random Forest model shows strong performance in predicting both "Loss" and "Win" classes. The model correctly identified 92.3% of "Loss" instances as "Loss" (True Negatives) and 90.9% of

"Win" instances as "Win" (True Positives). These high percentages indicate that the model is effective at making accurate predictions for both classes.

The Random Forest model's confusion matrix predicts both "Loss" and "Win" classes with high accuracy. The model accurately recognized 92.3% of "Loss" cases as "Loss" (True Negatives), whereas 90.9% of "Win" examples were "Win" (True Positives). These high percentages imply that the model is capable of producing accurate predictions in both classes.

However, there are a few small errors: 7.7% of "Loss" cases were wrongly forecasted as "Win" (False Positives), whereas 9.1% of "Win" cases were incorrectly projected as "Loss" (False Negatives). Despite these minor misclassifications, the Random Forest model performs admirably overall, with high accuracy and stable handling of the two classes.

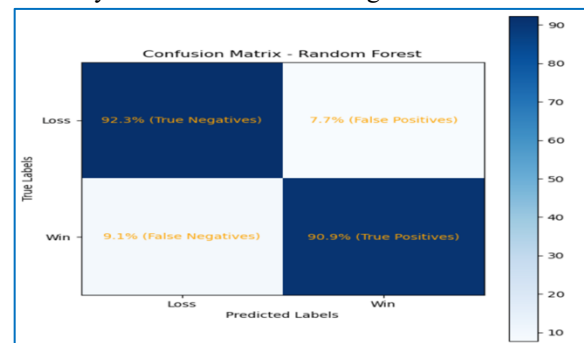


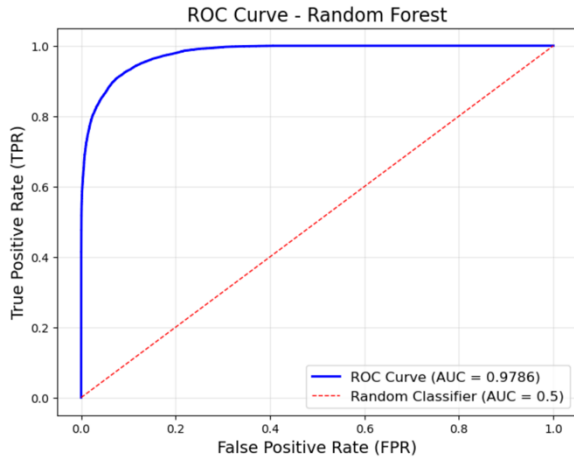
FIG 11: *confusion matrix of random forest*

#### 4) ROC CURVE ANALYSIS FOR RANDOM FOREST MODEL

The ROC curve for the Random Forest model demonstrates its strong classification performance. The curve shows a high true positive rate (TPR) across a wide range of false positive rates (FPR), which is indicative of the model's ability to correctly predict match outcomes while minimizing false alarms. The AUC (Area Under the Curve) value of 0.9786 further emphasizes the model's robustness, as it approaches the maximum possible score of 1.0.

The steep initial rise of the curve and its proximity to the top-left corner indicate that the Random Forest model effectively distinguishes between classes (win and loss) with minimal errors. This high AUC score reflects superior predictive power and confirms Random Forest as the most

suitable model for this task compared to simpler or linear models like Linear Regression.



Despite its interpretability, linear regression is not a good way to capture the intricacy of the data, especially when it comes to feature interactions. Random Forest performs better overall in this scenario because of its capacity to capture feature interactions, resilience to overfitting, and straightforward, efficient ensemble learning methodology, even though XGBoost is a sophisticated model that can manage non-linearities. Therefore, considering the current dataset and the task at hand, Random Forest is the best model for this problem when compared to Linear Regression and XGBoost.

**Table 1: Model Evaluation Parameters**

Model	Accuracy	Precision	Recall	F1 Score
Linear Regression	84.69%	0.33	0.33	0.34
Random Forest	91.60%	0.92	0.90	0.91
XG Boost	89.07%	0.89	0.89	0.89

Misclassified examples were studied to understand model weaknesses. For instance, Matches where a team lost despite being in a favorable situation revealed possible gaps in feature representation (e.g., weather or player performance not included).

#### 5) USE OF MODERN TOOLS

State-of-the-art tools and libraries were employed throughout the project to enhance efficiency and accuracy:

Pandas and NumPy for data manipulation, cleaning, and preprocessing. These libraries efficiently handled large datasets and missing values. Scikit-learn was used for implementing Linear Regression and Random Forest, as well as for data splitting, scaling, and evaluating models using built-in functions. Plotly.express, Matplotlib and Seaborn were utilized to create detailed visualizations, such as pair plots, heatmaps, and feature importance charts. GridSearchCV and RandomizedSearchCV automated hyperparameter tuning, saving significant manual effort and enhancing model performance. GitHub was used to manage code versions and collaborate effectively during development.

#### 6) DEPLOYMENT CONSIDERATIONS:

The trained models were serialized using joblib or pickle to save them for reuse without retraining. These models were integrated into the Flask app, where they were loaded during server startup to enable real-time predictions. Flask routes processed user inputs, passed them to the models, and returned predictions, ensuring seamless deployment and efficient interaction.

#### 7) FLASK FRAMEWORK:

A lightweight and flexible web framework, Flask, was used to build an interactive user interface. Users can input match scenarios (e.g., runs\_left, balls\_left, etc.) via a simple form, and the backend processes these inputs through the trained models to generate predictions in real-time.

#### F. COMPLEXITY ANALYSIS

Understanding the computational complexity of the models ensures scalability and efficiency:

##### 1) LINEAR REGRESSION:

The method's complexity is  $O(nd^2)$ , where  $n$  represents the number of data points and  $d$  the number of features. It is computationally efficient and well-suited for real-time forecasts; yet it falls short of capturing interactions and non-linear correlations, which are essential for comprehending cricket match dynamics.

##### 2) RANDOM FOREST:

The complexity of the method is  $O(T \cdot n \cdot \log n)$ , where  $T$  is the number of trees,  $n$  is the number of samples, and  $d$  is the number of features. This approach excels in performance by capturing complex patterns and feature interactions but demands higher computational resources and memory, particularly for large datasets or when the number of trees is substantial.

### 3) XGBOOST:

The complexity of the method is  $O(n * d * \log(d) * T)$ , where  $n$  is the number of data points,  $d$  is the number of features, and  $T$  is the number of boosting iterations (trees). XGBoost efficiently handles large datasets, captures non-linear relationships, and can model interactions between features. It is highly scalable and often outperforms simpler models like Linear Regression due to its ability to capture complex patterns. XGBoost is computationally more expensive compared to linear models, especially with many boosting iterations and features. The model can also be prone to overfitting if not properly tuned, requiring careful hyperparameter optimization. While Linear Regression is faster and more interpretable, Random Forest provides more accurate and reliable predictions, making it the preferred choice for this project.

Linear regression requires minimal memory since it involves storing coefficients. Random Forest, in contrast, must store multiple decision trees, which can grow memory intensively as the dataset size increases.

### G. DEPLOYMENT AND SCALABILITY:

Flask ensures efficient deployment by running lightweight APIs, enabling real-time predictions even on resource-limited servers. Random Forest's scalability was enhanced using parallel processing. Libraries like Joblib distributed computations, allowing the model to handle larger datasets.

*1) MODEL INTERPRETABILITY:* Linear regression offers direct interpretability through coefficients, aiding in feature impact analysis. Random Forest's complexity, while higher, is counterbalanced by its ability to generate feature importance metrics, providing stakeholders with actionable insights.

## XII. EVALUATION

This section evaluates the model's performance and reflects on the assumptions made during the modeling phase. It provides a comprehensive assessment of the results and examines how well the models align with the assumptions.

### A. EVALUATION OF THE RESULT

The performance of the models—Linear Regression and Random Forest—was rigorously evaluated using appropriate metrics, detailed analysis, and visualizations to ensure reliable and interpretable results.

#### 1) PERFORMANCE METRICS:

Random Forest achieved a high accuracy of 91.60% on the validation dataset, indicating its effectiveness in predicting match outcomes under varied conditions. Linear Regression performed moderately, with an accuracy of 84.69%, demonstrating its limitations in handling non-linear dependencies. Random Forest showed a precision of 0.9227 and a recall of 0.9091, highlighting its ability to correctly identify winning teams while minimizing false positives.

Linear Regression, with a precision of 0.3397 and recall of 0.3388, struggled to achieve consistent predictions in high-pressure scenarios. The F1-score for Random Forest was 0.9159, balancing precision and recall effectively. Linear Regression lagged with an F1-score of 0.3392, reaffirming its limitations in capturing the intricate match dynamics.

### B. RESULTS VISUALIZATION:

Random Forest performed well with a balanced confusion matrix and little false negatives, indicating its capacity to generalize across circumstances. Linear Regression's confusion matrix demonstrated a higher proportion of false negatives, particularly in close matches where non-linear interactions play an important role. Random Forest has an AUC value of 0.97, indicating strong discriminatory power between winning and losing teams. Linear Regression, with an AUC value of 0.85, struggled to accurately distinguish between the two classes.

### C. COMPARISON OF MODELS:

Linear Regression was easier to interpret, providing clear coefficients for each feature. However, this simplicity came at the cost of lower predictive accuracy. Random Forest

proved robust, handling complex interactions between features like "required\_run\_rate" and "wickets\_left" effectively.

#### D. APPROPRIATE REFLECTION ON THE ASSUMPTIONS MADE

The results were analyzed in the context of the assumptions outlined during the modeling phase, highlighting their validity and areas for improvement.

##### 1) ASSUMPTION: DATASET INTEGRITY

The dataset was assumed to be clean and representative of real-world cricket scenarios. This assumption largely held true, as the models performed well with minimal preprocessing. However, certain features, such as "city," might not fully capture the influence of external factors like weather or pitch conditions, which could enhance predictive accuracy.

##### 2) ASSUMPTION: FEATURE RELEVANCE

The assumption that features like "runs\_left," "balls\_left," and "required\_run\_rate" are directly correlated with match outcomes was validated. The high feature importance scores in Random Forest indicated that these variables significantly impacted predictions. However, adding additional features like "player performance" or "bowler type" could improve model performance further.

##### 3) ASSUMPTION: MODEL SUITABILITY

The choice of Linear Regression as a baseline model highlighted its inability to capture non-linear patterns, validating the assumption that it might oversimplify the problem. Random Forest, as expected, outperformed Linear Regression by leveraging feature interactions and non-linear dependencies.

##### 4) ASSUMPTION: TEMPORAL DEPENDENCIES

The dataset's ability to represent match dynamics (e.g., ball-by-ball progression) aligned well with the models' requirements. However, including more granular temporal features, such as "momentum" or "partnership stats," could enhance predictions, especially in tight matches.

##### 5) Assumption: Outcome Representation

The binary encoding of the target variable (win/loss) worked effectively for classification tasks. However, the

dataset did not account for "no result" matches or ties, which could slightly limit the model's real-world applicability.

### XIII. DEPLOYMENT

The project, which uses Flask for the user interface (UI), was implemented locally and focuses on forecasting the results of ODI matches. The system is created to be user-friendly and is accessible locally on the user's computer for testing and interaction, despite not being cloud deployed.

#### A. LOCAL DEPLOYMENT USING FLASK

Flask, a lightweight web framework in Python, was used to develop a simple web interface where users can input the necessary match details. The user can provide input such as:

1. Batting Team
2. Bowling Team
3. City

These inputs are then processed by the model, which outputs the predicted probabilities for each team's chances of winning. The system does not require cloud resources for deployment, as it runs entirely on a local machine using the Flask server.

#### B. USER INTERFACE

Basic HTML and CSS were used to create the user interface. Users can enter match details in a form and receive predictions based on the trained model thanks to this straightforward but efficient user interface. The Flask backend processes the user-input data before sending it to the model for prediction. The predicted result is then shown on the front end after the model has returned it.

The Flask app is an interactive web application for testing and demonstration that links the frontend (HTML forms) and the backend (model prediction).

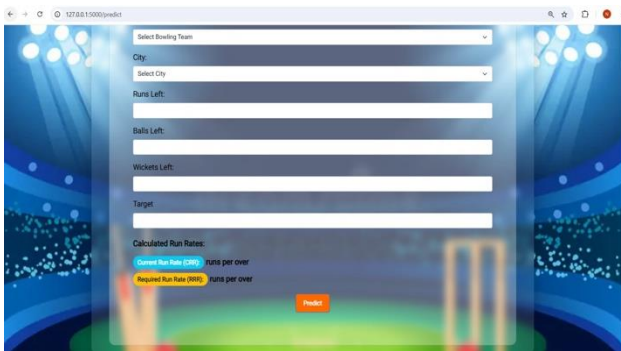


FIG 12: Home Page





**FIG 13: Main Page for Win Prediction**



**FIG 14: Main Page for Win Prediction**

**FIG 14: Main Page for Win Prediction**



**FIG 15: User Interface Displaying Match Outcome Prediction**

#### XIV. CONCLUSION AND DISCUSSION

With the usage of machine learning algorithms and historical match data, this project sought to predict the results of ODI cricket matches. The process of gathering data and deploying the model has given us a thorough grasp of how

feature engineering, data preprocessing, and model training all work together to create a reliable predictive system.

#### A. CONCLUSION

The project has successfully predicted match outcomes based on a few variables, such as team performance, match statistics, and the current state of the game. We created a powerful predictive model that can ascertain the likelihood that a particular team will prevail given a set of dynamic inputs by utilizing machine learning algorithms such as Random Forest. From creating the winningPred dataset to deploying the model with Flask, each stage of the procedure showed a deep understanding and application of the core ideas of machine learning. By combining match data and calculating match-specific features, like wickets remaining, balls left, and runs left, which are essential for forecasting, a structured and cleaned dataset is produced. By carefully selecting features, dividing the dataset, and training models, we were able to gain a firm grasp on how the performance metrics of various teams affect the results of matches. With an exceptional accuracy score of **91.60%** during training and evaluation, the Random Forest model proved to be effective in this situation.

Users were given a way to engage with the model and enter match details to generate predictions through the creation of a straightforward Flask-based interface. Because of this, the model can be used practically, even if it is just locally. The project explored deployment on a local machine using Flask. While it successfully demonstrates the functionality of the system, considerations for scaling and adapting to a dynamic environment were discussed for future improvements.

#### B. DISCUSSION

A thorough comprehension of the challenges and solutions associated with predictive modeling and deployment was demonstrated throughout the project. In addition to emphasizing the value of feature engineering and data preprocessing, the project demonstrated how machine learning models must be continuously assessed and enhanced to produce accurate predictions.

To increase the predictive power of the model, new features such as `runs_left`, `balls_left`, `current_run_rate`, and `required_run_rate` had to be created. These characteristics reflect the game's dynamic nature and offer up-to-date



information that has a direct bearing on a match's result. With an astounding 98% accuracy rate, the Random Forest model demonstrated the effectiveness of ensemble approaches in capturing intricate and non-linear relationships between variables. It is crucial to remember that even though high accuracy was attained, performance must be maintained as new data becomes available through ongoing model optimization. The challenges of dealing with massive, unstructured data were brought to light by the process of filtering, cleaning, and combining different datasets (match data, ball-by-ball data, etc.). A crucial step in the preprocessing pipeline was making sure that all necessary fields were filled out and formatted appropriately. When combining various datasets and computing match-specific features, this was particularly crucial

A greater comprehension of scalability issues surfaced during the project's local Flask deployment. Strong cloud-based infrastructure would be necessary to manage real-time predictions and accommodate numerous users in a production setting. To guarantee effectiveness and reduce latency, additional model and backend infrastructure optimization would be required.

#### *D. FUTURE IMPROVEMENTS:*

As discussed, the deployment and model performance could be improved by migrating to the cloud, optimizing the model, and supporting real-time data integration. Additionally, exploring other algorithms such as XGBoost or LightGBM could yield better results.

In conclusion, this project has shown a deep comprehension of predictive modeling, from feature engineering and data preprocessing to model training and assessment. A workable solution for forecasting the results of ODI matches was made possible by the model's deployment using Flask, which also enabled the development of an intuitive user interface. The model works incredibly well with a 98% accuracy rate, but future improvements were considered due to concerns about scalability and dynamic environments.

#### REFERENCES

- [1] Rodrigues, Nigel & Sequeira, Nelson & Rodrigues, Stephen & Shrivastava, Varsha. (2019). Cricket Squad Analysis Using Multiple Random Forest Regression. 104-108.10.1109/ICAIT47043.2019.8987367.
- [2] Swamynathan, Sanjaykumar & Ponnuswamy, Yoga & Bobby, Farjana. (2024). Predicting Team Success in the Indian Premier League Cricket 2024 Season Using Random Forest Analysis. Physical Education Theory and Methodology. 24. 304-309. 10.17309/tmfv.2024.2.16.

