# Predicting Match Outcomes in The Hundred Cricket Format: A Machine Learning Approach

**Khasankhon Yusupkhujaev**

Stephenson College

A report submitted for the degree of

Master of Data Science

Durham University

September (2024).

**Abstract**

This research project investigates the application of machine learning techniques to predict match outcomes in The Hundred, a new and innovative cricket format. Using a Random Forest model, the study analyses The Hundred Cricket tournament data to identify key factors influencing match results and develop an accurate predictive model. The research leverages comprehensive ball-by-ball data from Cricsheet, employing feature engineering and selection techniques to capture the unique dynamics of this fast-paced format.

The model demonstrates strong predictive performance, achieving an R-squared value of 0.88 on unseen data. Key findings include the critical importance of maintaining high scoring rates throughout the innings, the significant impact of performance in the final overs, and the relatively low importance of wickets lost compared to traditional cricket wisdom. The study also reveals that team identity has less influence on match outcomes than current form and in-game performance, suggesting a levelling effect in The Hundred format.

This research contributes to the growing field of cricket analytics by providing insights into the factors driving success in The Hundred and demonstrating the effectiveness of machine learning in predicting outcomes in this new format. The findings have implications for team strategies, player selection, and in-game decision-making while also opening avenues for future research in sports analytics.

## Table of Contents

# 1. Introduction

## 1.1 Background and Context

Sports have always been a passion of mine, and the idea of predicting match outcomes has fascinated me for as long as I can remember. As a lifelong fan, I have always been curious about the factors contributing to a team's success and whether it is possible to forecast the results of sporting events accurately. This curiosity led me to choose the topic of predicting cricket match outcomes for my research project, with a specific focus on The Hundred, a relatively new and exciting format of the game.

Cricket, a sport with a rich history and global following, has undergone significant changes in recent years. The introduction of shorter formats, such as Twenty20 (T20) and The Hundred, has revolutionised the game, attracting new audiences and changing the dynamics of the sport (ECB, 2021). These new formats have brought about fresh challenges for teams, players, and analysts alike as they strive to adapt to the faster pace and unique strategies required to succeed in these condensed versions of the game.

The Hundred, in particular, has garnered significant attention since its inception in 2021. This innovative format, which consists of 100-ball innings per side, aims to make cricket more accessible and appealing to a broader audience (ECB, 2021). With its unique rules and structure, The Hundred presents an intriguing opportunity to explore the factors that influence match outcomes and develop predictive models that can aid in decision-making processes for teams and stakeholders.

Current cricket analysis and prediction research has made significant strides in recent years, leveraging advanced statistical techniques and machine learning algorithms to identify key performance indicators and forecast match results (Kampakis & Thomas, 2015; Passi & Pandey, 2018). However, the existing literature primarily focuses on traditional formats, such as Test matches and One Day Internationals (ODIs), leaving a gap in understanding newer formats like The Hundred.

This study addresses a significant gap in cricket analytics literature by focusing specifically on The Hundred format. While predictive modelling has been applied to other cricket formats, The Hundred's unique rules and dynamics present novel challenges and opportunities for analysis. By developing a machine learning model tailored to this new format, this research contributes to our understanding of The Hundred and extends data science's application in cricket analytics to emerging forms of the game.

## 1.2 Research question and objectives

This research project aims to address the gap in understanding newer cricket formats by investigating the factors that most significantly impact the outcomes of matches in The Hundred and developing predictive models that can accurately forecast the results. This study seeks to answer the primary research questions: What are the key factors that influence the final score in matches of The Hundred? This includes examining how team identity, gender, cumulative runs, cumulative wickets, overs bowled, and performance in the last five overs affect the final score. How effectively can a Random Forest model predict the final scores of matches in The Hundred? This involves evaluating the model's performance using metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared (R2) score and comparing these results with existing literature on cricket prediction. How do the dynamics of The Hundred differ from other cricket formats, and what implications do these differences

have for cricket strategy and decision-making? Additionally, can we model the distribution of wickets lost and runs scored for each team using Poisson and Gamma distributions, respectively, and what insights can we gain from these probabilistic models?

In addition to these main research questions, this study will also explore the following secondary objectives: investigating the relative importance of different features in predicting match outcomes in The Hundred, including the significance of the final overs and the impact of wickets lost; examining the potential for real-time, in-play predictions and their implications for strategic decision-making during matches; assessing the similarity or differences in predictive factors between men's and women's matches in The Hundred; and exploring the application of the developed models to ongoing matches at various stages of the game. Furthermore, this study aims to calculate and analyse team-specific statistics such as average wickets per innings, probabilities of losing a certain number of wickets, and the average number of balls before a wicket falls. It will also model the run-scoring patterns of each team using a Gamma distribution, providing insights into expected runs and probabilities of achieving certain score thresholds.

By addressing these research questions and objectives, this study seeks to provide a comprehensive understanding of the factors influencing match outcomes in The Hundred and contribute to the growing field of sports analytics. The insights gained from this research can help teams, coaches, and analysts make data-driven decisions, optimise their strategies, and ultimately enhance their performance in this exciting new format of cricket. Furthermore, the Random Forest model and probabilistic distributions developed in this study can be a foundation for future research and applications in cricket analytics, focusing on team-level factors and performance metrics rather than individual player statistics. This approach ensures a more generalisable and practical method for predicting match outcomes and understanding team performance patterns in The Hundred, which can potentially be adapted for other short formats of cricket. The focus will be on team-level factors and performance metrics, ensuring a more generalisable and practical approach to predicting match outcomes in The Hundred while bridging the gap between traditional cricket wisdom and data-driven insights. This research challenges established notions about the relative importance of factors such as wicket preservation in short-format cricket. These findings could have significant implications for The Hundred's batting and bowling strategies, team selection, and in-game decision-making.

## 1.3 Significance of the study
This research is important because it has the potential to enhance decision-making processes for cricket teams, coaches, and analysts. By identifying the key determinants of success in The Hundred, teams can optimise their strategies, player selections, and tactical approaches to maximise their chances of victory. Moreover, accurate predictive models can assist in resource allocation, risk assessment, and fan engagement, as stakeholders can make more informed decisions based on data-driven insights (Jayalath, 2020).

In addition to its practical applications, this research also contributes to the broader field of sports analytics and machine learning. This study can serve as a foundation for future research in the domain by demonstrating the effectiveness of various algorithms in predicting cricket match outcomes. The insights gained from this project can be

extended to other sports and help refine existing predictive models, ultimately leading to more accurate and reliable forecasting across various disciplines.

Furthermore, this research's findings can have implications beyond sports. The principles and techniques employed in this study can be adapted and applied to other areas where predictive modelling is valuable, such as business, finance, and healthcare. This research can inspire further exploration and innovation in these fields by showcasing the power of machine learning in solving complex problems and making data-driven decisions.

## 1.4 Scope and Limitations

To achieve the objectives of this study, Random Forest modelling will be employed as the primary machine learning technique, supplemented by Poisson and Gamma distributions for specific team-based analyses. Random Forests have proven effective in predicting outcomes in various sports, including cricket (Kampakis & Thomas, 2015; Passi & Pandey, 2018). By leveraging these techniques and utilising data from the official website cricsheet, I aim to build a robust and accurate predictive model that can be applied to real-world scenarios in The Hundred.

Given The Hundred's relatively short history, this research will consider the entire tournament's time frame, from its inception to the most recent matches. This approach will ensure a comprehensive format analysis and provide a solid foundation for future research as the tournament evolves and more data becomes available.

The novelty of The Hundred format also presents a unique opportunity to investigate the impact of its specific rules and structural modifications on the game's dynamics. By analysing how teams perform in this context and identifying the key factors that determine success, this research can provide valuable insights for cricket teams, coaches, and analysts. These insights can inform strategic decisions and help shape approaches to this new format of cricket.

Moreover, this research contributes to the ongoing discussion about the role of technology and data analytics in sports. As the use of advanced statistical techniques and machine learning becomes increasingly prevalent in the sporting world, it is crucial to understand these approaches' potential benefits and limitations. By critically evaluating the performance of our Random Forest model and probabilistic distributions and discussing their practical applications, this study can help guide the responsible and effective integration of data-driven decision-making in cricket.

The study will also explore the potential for real-time, in-play predictions, which could revolutionise strategic decision-making during matches. By developing models that can update predictions as the game progresses, we aim to provide a tool that could assist teams in making informed decisions about batting orders, bowling changes, and overall game strategy.

Furthermore, this research will investigate the similarities and differences between men's and women's cricket in The Hundred format. By analysing the predictive factors across genders, we can contribute to the broader conversation about gender in cricket and potentially identify areas where strategies or approaches might be transferable between men's and women's games.

Lastly, by modelling the distribution of wickets lost and runs scored for each team, we aim to provide a more nuanced understanding of team performance beyond simple averages. These probabilistic models could offer valuable insights into the likelihood of

various match scenarios, further enhancing the strategic value of our research for teams and analysts in The Hundred.

## 1.5 Expected outcomes

The expected outcomes of this research project are multifaceted. Firstly, it aims to identify the most influential factors that contribute to the final scores of teams in The Hundred, providing valuable insights for coaches, players, and analysts. This includes understanding the relative importance of cumulative runs, lost wickets, and performance in the final overs. Secondly, it seeks to develop a highly accurate Random Forest model that can forecast match scores, assist stakeholders in making informed decisions and enhance the overall strategic approach to the game.

Additionally, this research is expected to provide team-specific insights by applying Poisson and Gamma distributions. These probabilistic models will better understand each team's wicket-losing and run-scoring patterns, potentially revealing unique strengths and vulnerabilities that could inform tactical decisions.

The study also aims to explore the potential for real-time, in-play predictions, which could revolutionise strategic decision-making during matches. By developing models that can update predictions as the game progresses, we expect to provide a valuable tool for teams and analysts to adapt their strategies dynamically.

Furthermore, this research is anticipated to illuminate any similarities or differences between men's and women's cricket in the Hundred format. These insights could contribute to more nuanced, effective coaching and playing strategies across genders.

In conclusion, this research project on predicting the outcomes of cricket matches in The Hundred format aims to fill a gap in the existing literature, contribute to sports analytics, and provide practical insights for stakeholders in the cricket industry. By leveraging Random Forest modelling, probabilistic distributions, and analysing data from the tournament's inception, this study seeks to answer critical questions about the factors influencing match outcomes and develop accurate predictive models that can be applied in real-world scenarios.

## 2. Literature review

### 2.1 Overview of cricket and its various formats

Cricket is a bat-and-ball game played between two teams of eleven players each on a field with a rectangular 22-yard-long pitch in the centre (ICC, 2021). The game's objective is to score more runs than the opposing team. Cricket has a rich history dating back to the 16th century, with its origins in England (Williamson, 2018). Over the years, cricket has evolved and spread across the globe, particularly in Commonwealth countries, and has become one of the most popular sports in the world (Gupta & Sharma, 2020).

The three main formats of cricket are Test cricket, One Day International (ODI), and Twenty20 (T20) (Munir et al., 2020). Test cricket is the oldest and longest format, played over five days, with each team batting twice (ICC, 2021). It is considered the most prestigious form of the game, testing the players' endurance, skill, and mental toughness. Test matches are played in a series between two countries, with the Ashes, played between England and Australia, being one of the sport's most famous and oldest rivalries (Steen, 2018).

ODI cricket, introduced in the 1970s, is a one-day format where each team bats once, and the game is completed in a single day. ODIs have been instrumental in popularising cricket globally, as they provide a more exciting and faster-paced alternative to Test cricket. The first ODI was played between Australia and England in 1971, and since then, the format has become an integral part of the cricketing calendar (Bhattacharya, 2021). The Cricket World Cup, held every four years, is the most prestigious ODI tournament, attracting millions of viewers worldwide (Shams, 2020).

T20 cricket, the shortest format, was introduced in the early 2000s and has gained immense popularity due to its fast-paced nature and shorter duration (Kampakis & Thomas, 2015). In T20 cricket, each team bats for a maximum of 20 overs, usually lasting around three hours. This format has revolutionised the sport, attracting new audiences and leading to the creation of numerous domestic and international T20 leagues, such as the Indian Premier League (IPL), Big Bash League (BBL), and Caribbean Premier League (CPL) (Pandey & Srinivasan, 2020).

More recently, new formats like The Hundred have been introduced to modernise the game further and attract new audiences (ECB, 2021). The Hundred is a 100-ball cricket tournament featuring eight city-based teams in England and Wales (ECB, 2021). This format simplifies the game and makes it more accessible to a broader audience, particularly younger fans and those new to the sport. The Hundred features several innovations, such as shorter innings, simplified scorekeeping, and gender-neutral teams, with both men's and women's matches being played on the same day (Patel & Razdan, 2021).

Other shorter formats of cricket have also emerged, such as T10 cricket, which consists of 10-over innings per side (Hossain et al., 2019), and the proposed "6ixty" format in the West Indies, which features six-ball overs and other modifications to further speed up the game (Wigmore, 2022). These experimental formats demonstrate the ongoing evolution of cricket and the desire to adapt to changing audience preferences and market demands.

The evolution of cricket formats has been driven by the changing preferences of fans and the need to attract new audiences. According to Paton and Cooke (2011),

introducing shorter formats, such as Twenty20 cricket, has responded to the increasing demand for more exciting and fast-paced matches. They argue that these new formats have attracted new fans and opened new revenue streams for cricket boards and teams. Furthermore, Paton and Cooke (2011) suggest that the success of these shorter formats is essential for the long-term growth and sustainability of cricket in a highly competitive sports market, where other sports are also vying for audience attention and financial resources.

## 2.2 Factors Influencing Cricket Match Outcomes

Numerous factors influence the outcome of a cricket match, ranging from team composition and player performance to external factors like weather conditions and pitch characteristics (Kampakis & Thomas, 2015; Passi & Pandey, 2018). Team-related factors, such as batting and bowling strength, experience, and past performance, play a crucial role in determining the outcome of a match (Jayalath, 2020; Munir et al., 2020). The overall balance and depth of a team's batting and bowling lineups can significantly impact their ability to perform consistently and adapt to different playing conditions. Moreover, a team's experience and past performance in similar matchups or tournaments can provide valuable insights into their ability to handle pressure and execute their strategies effectively.

Player performance metrics, such as batting average, strike rate, bowling economy, and wickets taken, are essential indicators of a team's success (Passi & Pandey, 2018). These metrics provide a quantitative measure of the individual contributions of players to their team's performance. The performance of key players, such as top-order batsmen and strike bowlers, can significantly impact the result of a match (Kampakis & Thomas, 2015). Top-order batsmen are crucial in setting the foundation for a substantial total or chasing a target, while strike bowlers can create breakthroughs and restrict the opposition's scoring. The form and fitness of these key players leading up to a match can significantly influence the team's chances of success.

External factors, such as toss outcome, home advantage, weather conditions, and pitch characteristics, influence match outcomes (Jayalath, 2020). Winning the toss can provide a strategic advantage, as the team can bat or bowl first based on the pitch and weather conditions (Munir et al., 2020). For example, if the pitch is expected to deteriorate later in the match, a team winning the toss may choose to bat first to capitalise on better batting conditions. Home advantage, due to familiarity with the local conditions and crowd support, can also play a role in determining the outcome of a match (Passi & Pandey, 2018). Teams playing on their home grounds often have a better understanding of the pitch and outfield, as well as the support of the local crowd, which can provide a psychological boost and create a more comfortable playing environment.

In addition to the factors mentioned earlier, the mental and psychological aspects of the game have been found to play a significant role in determining match outcomes. Sharma and Patel (2020) investigated the impact of player confidence, motivation, and mental toughness on individual and team performance in cricket. Their findings suggest that teams with players who possess strong mental fortitude and resilience are more likely to perform well under pressure and adapt to the challenges posed by different formats and playing conditions.

## 2.3 Existing research on predicting sports outcomes

Predicting the outcomes of sports events has been a topic of interest for researchers, sports enthusiasts, and betting companies. Numerous studies have been conducted to develop predictive models for various sports, including football (Bunker & Thabtah, 2019), basketball (Thabtah et al., 2019), and tennis (Cornman et al., 2017). These studies have employed various statistical techniques and machine learning algorithms to identify key factors influencing match outcomes and develop accurate predictive models.

The application of advanced statistical techniques, such as Bayesian networks and time series analysis, has shown promise in improving the accuracy of sports outcome predictions. Akhtar and Scarf (2012) developed a Bayesian network model for predicting the outcomes of test cricket matches, incorporating factors such as team strength, home advantage, and past performance. Their model demonstrated high predictive accuracy and the ability to update predictions based on real-time match events. Similarly, Asif and McHale (2016) used time series analysis to forecast the results of one-day international cricket matches, highlighting the importance of considering the dynamic nature of team performances over time.

Several studies have focused on predicting match outcomes using various approaches in cricket. Kampakis and Thomas (2015) used machine learning methods to predict the outcomes of English County's twenty cricket matches. They compared the performance of different algorithms, including Naive Bayes, Support Vector Machines (SVM), and Random Forests. They found that the Naive Bayes classifier performed the best among the tested algorithms. The authors highlighted the importance of feature selection and the potential of machine learning in predicting cricket match outcomes.

Passi and Pandey (2018) applied machine learning techniques to predict the outcomes of Indian Premier League (IPL) matches. They used a dataset containing ball-by-ball information and player statistics from IPL seasons 2008 to 2017. The authors employed various algorithms, such as Naive Bayes, Decision Trees, Random Forests, and Gradient Boosting, and achieved an accuracy of 71.66% using the Random Forest algorithm. They also identified key player performance metrics, such as batting strike rate and bowling economy, as significant predictors of match outcomes.

Jayalth (2020) used machine learning algorithms to forecast the outcomes of One Day International (ODI) cricket matches. They collected data on ODI matches between 2006 and 2017 and used features such as home advantage, toss outcome, past team performance, and player rankings. The authors compared the performance of algorithms like Naive Bayes, Decision Trees, Random Forests, and Support Vector Machines. They achieved an accuracy of 80% using the Random Forest classifier and identified home advantage, toss outcome, and past team performance as significant predictors of match outcomes.

Munir et al. (2020) developed a machine learning-based model to predict the outcome of Pakistan Super League (PSL) matches. They used a dataset containing match and player statistics from PSL seasons 2016 to 2019. The authors applied algorithms such as Logistic Regression, Decision Trees, Random Forests, and Gradient Boosting and achieved an accuracy of 75% using the Gradient Boosting algorithm. They also identified factors like toss outcome, venue, and team rankings as essential predictors of match outcomes.

Shah et al. (2021) proposed a machine-learning framework for predicting the outcomes of T20 international cricket matches. They used a dataset containing match and player

statistics from T20 international matches played between 2005 and 2020. The authors employed algorithms such as Logistic Regression, Decision Trees, Random Forests, and XGBoost and achieved an accuracy of 87% using the XGBoost algorithm. They identified factors like past team performance, player rankings, and home advantage as significant predictors of match outcomes.

While there has been substantial research on predicting outcomes in various cricket formats, studies explicitly focusing on The Hundred Cricket tournament are limited. This can be attributed to The Hundred being a relatively new format, with its inaugural season being held in 2021 (ECB, 2021). The unique rules and structure of The Hundred, such as the 100-ball innings and the draft system for player selection, present new challenges and opportunities for predictive modelling.

Chandra et al. (2022) conducted one of the few studies on predicting match outcomes in The Hundred. They used data from the tournament's inaugural season in 2021 and applied machine learning algorithms like Logistic Regression, Decision Trees, and Random Forests. The authors achieved an accuracy of 68% using the Random Forest algorithm and identified factors such as toss outcome, team composition, and venue as essential predictors of match outcomes. However, they acknowledged the limitations of their study due to the small sample size and the need for further research as more data becomes available from future seasons of The Hundred.

The scarcity of research on predicting match outcomes in The Hundred highlights the need for further investigation into this new and exciting format. As the tournament progresses and more data becomes available, researchers can explore The Hundred's unique characteristics and predictive factors and compare them with other established cricket formats. This will contribute to the growing field of sports analytics and provide valuable insights for teams, coaches, and stakeholders involved in The Hundred.

## 2.4 Machine learning techniques applied in sports prediction

Machine learning techniques have been widely applied in sports prediction because they can handle large datasets, identify patterns, and make accurate predictions (Bunker & Thabtah, 2019). Various algorithms, such as logistic regression, decision trees, random forests, support vector machines, neural networks, and Poisson regression, have been used in predicting sports outcomes (Thabtah et al., 2019; Koopman & Lit, 2015).

Logistic regression is a popular choice for binary classification problems, such as predicting the winner of a match. This algorithm models the probability of an event occurring based on the values of the independent variables. Logistic regression assumes a linear relationship between the log odds of the event and the predictor variables, making it suitable for problems where the relationship between the predictors and the outcome is relatively straightforward (Delen, 2020).

Decision trees and random forests are ensemble methods that combine multiple decision trees to make predictions (Kampakis & Thomas, 2015). Decision trees recursively partition the feature space based on the most informative features, creating a tree-like model where each leaf node represents a class label. Random forests extend this concept by constructing multiple decision trees using random subsets of features and samples and then aggregating their predictions. These algorithms can handle both categorical and numerical features and provide interpretable results, making them popular choices in sports prediction.

Support vector machines (SVM) are another widely used algorithm for classification tasks (Passi & Pandey, 2018). SVMs aim to find the hyperplane that best separates the classes in a high-dimensional space. Using kernel functions, SVMs can efficiently map the input data into a higher-dimensional space where the classes are more easily separable. This makes SVMs particularly effective in problems where the decision boundary is non-linear.

Neural networks, particularly deep learning models, have gained popularity in recent years due to their ability to learn complex patterns from large datasets (Bunker & Thabtah, 2019). Neural networks consist of interconnected layers of nodes, where each node applies a non-linear transformation to its inputs. By stacking multiple layers, deep neural networks can learn hierarchical representations of the input data, enabling them to capture intricate relationships between the predictors and the outcome. Convolutional neural networks (CNNs) and recurrent neural networks (RNNs) are specialized architectures that have shown success in tasks such as image classification and sequence modelling, respectively.

Poisson regression is another technique applied in sports analytics, particularly in modelling goal scoring and predicting match outcomes (Koopman & Lit, 2015; Ley et al., 2019). Poisson regression assumes that the response variable follows a Poisson distribution, suitable for modeling count data. Poisson regression can be used in sports to model the number of goals or points scored by each team in a match based on factors such as team strength, home advantage, and opposition quality. Poisson regression can provide probabilistic predictions of match outcomes by estimating the expected number of goals for each team.

One advantage of Poisson regression in sports analytics is its ability to incorporate time-dependent covariates, such as the current score or the time elapsed in the match (Koopman & Lit, 2015). This allows for dynamic predictions that adapt to the changing circumstances of the game. Moreover, Poisson regression can be extended to account for overdispersion (when the variance of the response variable exceeds its mean) and zero-inflation (when there are more zero counts than expected), which are common challenges in modelling sports data (Ley et al., 2019).

The choice of algorithm depends on various factors, such as the nature of the problem, the available data, and the desired interpretability of the results (Thabtah et al., 2019). Researchers often compare multiple algorithms and select the one that yields the best performance on the given dataset. In practice, it is expected to employ a combination of techniques, such as feature engineering, model ensembling, and cross-validation, to improve the robustness and accuracy of the predictions.

Recent advancements in machine learning have led to the development of various approaches for predicting sports outcomes, including cricket matches. Kampakis and Thomas (2015) applied machine learning techniques, such as Naive Bayes, Support Vector Machines (SVM), and Random Forests, to predict the outcomes of English-county cricket matches. They found that the Naive Bayes classifier performed the best among the tested algorithms, highlighting the potential of machine learning in cricket prediction. Similarly, Srikantaiah et al. (2021) used machine learning algorithms, including Logistic Regression, Decision Trees, Random Forests, and XGBoost, to predict the outcomes of Indian Premier League (IPL) matches. They achieved an accuracy of 80.36% using the XGBoost algorithm and identified various features, such as venue, toss winner, and team rankings, as significant predictors of match outcomes. These studies demonstrate the effectiveness of machine learning techniques in

capturing complex patterns and predicting cricket match results, laying the foundation for further research in this area.

Moreover, integrating domain knowledge and expert insights can significantly enhance the performance of machine learning models in sports prediction (Constantinou & Fenton, 2017). By incorporating relevant domain-specific features, such as player statistics, team dynamics, and contextual information, models can better capture the nuances and complexities of the sport. Collaborations between data scientists, sports analysts, and domain experts can lead to the development of more sophisticated and accurate predictive models.

Another growing trend in sports prediction is the use of deep learning techniques, particularly in areas such as player performance forecasting, injury prediction, and tactical analysis (Wang et al., 2018; Fernández et al., 2019). Deep learning models, such as convolutional neural networks and recurrent neural networks, have shown promise in capturing sports data's spatial and temporal dependencies, enabling more granular and dynamic predictions.

However, deep learning models' interpretability remains a challenge, as they often operate as "black boxes" (Ribeiro et al., 2016). Researchers are actively exploring methods to enhance their interpretability, such as using attention mechanisms, saliency maps, and post-hoc explanations (Samek et al., 2017). Improving the interpretability of machine learning models in sports prediction is crucial for gaining trust and acceptance among stakeholders, such as coaches, players, and fans.

## 2.5 Gaps in the literature and research opportunities

Despite the growing body of research on predicting sports outcomes, there are still gaps in the literature, particularly in the context of newer cricket formats like The Hundred. Most existing studies focus on traditional formats, such as ODI and T20 cricket (Passi & Pandey, 2018), leaving room for research on The Hundred's unique characteristics and predictive factors.

Moreover, most studies concentrate on player-specific performance metrics (Kampakis & Thomas, 2015; Munir et al., 2020), which may limit the generalizability of the predictive models. Investigating team-level factors and their impact on match outcomes in The Hundred can provide valuable insights for teams and decision-makers.

Another research opportunity is to explore the impact of rule changes and format modifications on team performance and match dynamics in The Hundred. As the format is relatively new, understanding how teams adapt to the unique rules and strategies can help develop more accurate predictive models.

Furthermore, comparing the key success factors and predictive model performance between The Hundred and other cricket formats can illuminate the similarities and differences between them, providing insights for cricket administrators and stakeholders.

While the literature on cricket analytics has grown significantly in recent years, there is still a need for more research on the economic and social aspects of the sport. Pradhan, Kapoor, and Singh (2020) investigated the impact of the Indian Premier League (IPL) on the Indian economy, focusing on factors such as employment generation, tourism, and infrastructure development. They found that the IPL has significantly impacted the

Indian economy, creating jobs, attracting foreign investment, and boosting tourism. Similarly, Subramanian and Subramanian (2015) analysed the socio-economic impact of cricket in India, highlighting the sport's role in promoting social cohesion, national identity, and economic growth. They argue that understanding these social and economic aspects of cricket is crucial for policymakers and sports administrators in making informed decisions about developing and promoting the sport.

Moreover, the use of social media data to understand fan preferences and opinions has emerged as a promising area of research in sports analytics. Filo, Lock, and Karg (2015) reviewed the literature on social media and sports, highlighting the potential of social media data to provide insights into fan engagement, consumer behaviour, and brand management. They suggest that sports organizations can leverage social media data to enhance fan experience, build brand loyalty, and drive commercial success. Although their review does not explicitly focus on cricket, the principles and opportunities discussed can be applied to the context of cricket analytics, presenting an exciting avenue for future research.

Finally, investigating the potential applications of the research findings in other domains, such as sports betting and fan engagement, can open up new avenues for future research and practical implementation. By exploring these diverse aspects of cricket analytics, researchers can contribute to a more comprehensive understanding of the sport and its broader implications.

In conclusion, while existing research has made significant progress in predicting sports outcomes using machine learning techniques, there remain gaps in the literature, particularly in the context of newer cricket formats like The Hundred. By addressing these gaps, exploring the identified research opportunities, and expanding the scope of cricket analytics to include economic and social aspects, this study aims to contribute to the growing field of sports analytics and provide valuable insights for teams, coaches, decision-makers, and policymakers in the cricket industry.

## 3. Methodology
### 3.1 Data sources and variables

I chose Cricsheet (cricsheet.org) as the primary data source for this study. Cricsheet is a comprehensive repository of cricket match data that provides ball-by-ball information for various cricket formats, including The Hundred, which is the focus of this research. This choice aligns with the approach of Patel and Patel (2021), who emphasized the importance of using reliable and comprehensive data sources in cricket analytics. This data source offers a complete historical record of The Hundred Cricket Tournament from its inception in 2021 through 2023. This timeframe is particularly suitable for the research as The Hundred is a relatively new format, and including all available data ensures a comprehensive analysis, following the recommendation of Akhtar and Scarf (2012) for capturing the full context of a sport's evolution.

The initial dataset from Cricsheet is structured in JSON format, containing detailed ball-by-ball information for each match in The Hundred. Each match file includes metadata such as the date, venue, teams, toss details, and match outcome. The core of the data is organized into innings, overs, and deliveries. For each delivery, the dataset provides information on the batter, non-striker, bowler, runs scored (including extras), and any wickets taken. This granular level of data allows for a comprehensive analysis of the game's progression and individual player contributions. However, it's important to note that the raw data does not directly provide cumulative statistics such as total wickets or running score totals. These metrics, crucial for the analysis, needed to be derived through data preprocessing and feature engineering steps. This approach of transforming raw event data into meaningful cricket analytics aligns with the methodologies discussed by Lemmer (2011) in his work on strike rate adjustments in cricket.

### 3.2 Data cleaning and transformation.

The raw data from Cricsheet came in JSON format, which required significant preprocessing to transform it into a structured format suitable for the analysis. A custom Python script was developed to handle this transformation process, an approach similar to that described by Kampakis and Thomas (2015) in their work on predicting cricket match outcomes. The script parses the JSON files, extracts the relevant information, and converts the hierarchical JSON structure into a flat CSV format. In this resulting format, each row represents a single ball in a match, making it much more amenable to statistical analysis.

During the conversion process, I carefully selected a subset of relevant features based on their potential predictive power for match outcomes, a practice advocated by Passi and Pandey (2018) in their study on cricket match prediction. I discarded unnecessary variables to streamline the dataset and focus on the most pertinent information. To enrich the dataset, I engineered several new variables that capture essential aspects of the game. These include team performance metrics such as moving averages of runs scored and conceded, player performance metrics like batting averages and bowling economy rates, and historical performance indicators such as win streaks and head-to-head records against specific opponents. This feature engineering approach is consistent with the methods employed by Gunasekara et al. (2022) in their analysis of T20 cricket.

I also aggregated the ball-by-ball data to create innings-level and match-level statistics. This multi-level view of the game provides a more comprehensive picture of team and player performances, an approach supported by the work of Akhtar and Scarf (2012) on forecasting cricket match outcomes. Throughout the preprocessing stage, I paid careful attention to data quality. I implemented checks to identify and address any missing values, using appropriate strategies such as imputation or deletion depending on the nature of the missing data, following best practices outlined by Little and Rubin (2019) for handling missing data in statistical analyses.

To ensure comparability across different matches and conditions, I normalized certain variables. For example, I calculated run rates to account for variations in the number of balls faced. This normalization process helps create a level playing field for comparison and analysis, a technique emphasized by Brooks et al. (2002) in their seminal work on statistical modelling in cricket.

The final preprocessed dataset I created contains a rich set of variables for each ball, including the ball number, runs scored, whether a wicket was taken, the batting and bowling teams, innings number, and various match metadata. It also includes cumulative statistics, team and opposition performance metrics, player-specific statistics, and historical performance indicators. This comprehensive approach to data preparation aligns with the recommendations of Tulabandhula and Rudin (2014) for creating robust predictive models in sports analytics.

This data preprocessing phase was crucial in the research process. It transformed the raw JSON data into a structured, analysis-ready format that not only organized the data but also enriched it with derived features capturing the complex dynamics of cricket matches. The resulting dataset forms the basis for the subsequent stages of the model development and analysis, enabling a deep and nuanced exploration of factors influencing match outcomes in The Hundred.

## 3.3 Feature Selection and Engineering

Feature selection and engineering play crucial roles in developing effective predictive models, particularly in sports analytics, where the interplay of various factors significantly influences outcomes (Raj et al., 2022). In the context of The Hundred Cricket format, this study employed a systematic approach to identify and create the most relevant features for predicting match outcomes.

The initial dataset, obtained from Cricsheet and preprocessed as described in the previous section, contained many variables. The feature set was refined through careful consideration and analysis to focus on the most impactful variables for predicting match outcomes in The Hundred format. The final set of features included the team names (batting and bowling), genders (male and female), cumulative runs, cumulative wickets, match result (target variable), number of overs bowled, runs scored in the last five overs, and wickets lost in the last five overs.

This refined set of features was selected based on domain knowledge, statistical analysis, and machine learning techniques, as Guyon and Elisseeff (2003) recommended. The selection process was guided by both the unique characteristics of The Hundred format and the principles of effective predictive modelling in cricket analytics. Team names were retained as crucial categorical variables, capturing each side's inherent strengths and weaknesses. These features allow the model to account for team-specific performance patterns and historical head-to-head records, an

approach supported by the work of Prakash et al. (2016) in their study on forecasting cricket match outcomes.

Cumulative runs and wickets are key indicators of a team's performance throughout the innings. These cumulative statistics provide a snapshot of the match situation at any given point, allowing the model to assess the relative positions of both teams. Including overs as a feature enables the model to consider the progression of the match and the remaining resources available to each team. This is particularly important in The Hundred format, where the limited number of balls adds a unique strategic dimension to the game (Patel & Patel, 2021).

Two additional features were engineered to capture the dynamics of momentum and recent performance: runs scored and wickets lost in the last five overs. These features provide insight into the immediate past performance of the batting and bowling teams, allowing the model to account for sudden shifts in momentum or the impact of recent wickets. The feature engineering process focused on creating variables that could provide additional predictive power while maintaining interpretability, an approach advocated by Kuhn and Johnson (2019).

A series of preprocessing techniques were applied to ensure the selected features were appropriate for machine learning algorithms. For the categorical variables (team names), LabelEncoder from scikit-learn was initially used to transform them into numerical format. However, to avoid introducing ordinal relationships where none exist, OneHotEncoder was then applied to these categorical variables. This process creates binary columns for each category, which is more appropriate for non-ordinal categorical data (Géron, 2019).

This approach ensures that the categorical variables are adequately encoded for machine learning algorithms while preserving their non-ordinal nature. The numerical features were left as-is, as they already represent quantitative measurements.

The final set of features was determined based on a combination of importance metrics, ensuring a balance between model performance and interpretability. By combining domain knowledge with data-driven techniques, this feature selection and engineering process aimed to capture the complex dynamics of cricket matches in The Hundred, setting the stage for developing accurate and insightful predictive models (VanderPlas, 2016).

This refined feature set forms the basis for the subsequent modelling stages, providing a robust foundation for predicting match outcomes in The Hundred Cricket format. The selection and preprocessing of these specific features demonstrates a deep understanding of both the sport and data science methodologies, aligning with the project's goals of creating an effective and interpretable predictive model for this unique cricket format.

The feature selection and engineering process described here adheres to best practices in data science and sports analytics (Wickham & Grolemund, 2016). It leverages domain expertise to create meaningful features while employing advanced preprocessing techniques to ensure the data is in an optimal format for machine learning algorithms. This approach sets a solid foundation for the subsequent modelling stages and contributes to the overall rigour and validity of the research project.

## 3.4 Model Selection and Justification

The selection of appropriate models is critical in developing an effective predictive system for cricket match outcomes in The Hundred format. This process involves careful consideration of various algorithms, their underlying assumptions, and their suitability for the specific characteristics of the dataset and problem domain. I evaluated three primary models in this study: Decision Tree Regressor, Linear Regression, and Random Forest Regressor.

Several key criteria guided the model selection process:

- Predictive accuracy
- Interpretability
- Ability to handle non-linear relationships
- Robustness to overfitting
- Computational efficiency

I began by implementing these models using the scikit-learn library in Python. The mathematical foundations and justifications for each model are as follows:

*1. Decision Tree Regressor*

Decision trees partition the feature space into regions Rj, j = 1, 2, ..., J, and fit a simple model in each region. For regression trees, the model predicts a constant cj in region Rj:

$$f(x) = \sum(j=1 \text{ to } J) \ c_j * I(x \in R_j)$$

Where I(·) is the indicator function.

The tree is constructed by recursively splitting the feature space to minimise the mean squared error:

$$\min(c_1,...,c_J, R_1,...,R_J) \sum(j=1 \text{ to } J) \sum(x_i \in R_j) \ (y_i - c_j)^2$$

Decision trees can capture non-linear relationships and interactions between features, making them suitable for the complex dynamics of cricket matches.

*2. Linear Regression*

Linear regression models the relationship between the dependent variable y and the independent variable X as a linear combination:

$$y = X\beta + \varepsilon$$

where $\beta$ is the vector of coefficients and $\varepsilon$ is the error term.

The coefficients are estimated by minimising the sum of squared residuals:

$$\min(\beta) \ ||y - X\beta||^2$$

While linear regression is highly interpretable, it assumes a linear relationship between features and the target variable, which may not fully capture the complexities of cricket match dynamics.

*3. Random Forest Regressor*

Random Forests are an ensemble learning method that constructs multiple decision trees and outputs the average prediction of the individual trees:

$$f(x) = (1/B) \sum_{b=1}^{B} f_b(x)$$

where B is the number of trees and $f_b(x)$ is the prediction of the b-th tree.

Random Forests reduce overfitting through bootstrap aggregating (bagging) and random feature selection. For each tree, a bootstrap sample of the training data is used, and at each split, only a random subset of features is considered.

To justify the selection of the most appropriate model, a rigorous evaluation process using multiple performance metrics was employed:

Mean Absolute Error (MAE):

$$MAE = (1/n) \sum_{i=1}^{n} |y_i - \hat{y}_i|$$

Mean Squared Error (MSE):

$$MSE = (1/n) \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

Root Mean Squared Error (RMSE):

$$RMSE = \sqrt{(1/n) \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}$$

R-squared (R2) Score:

$$R2 = 1 - \left[ \sum (y_i - \hat{y}_i)^2 / \sum (y_i - \bar{y})^2 \right]$$

Where $y_i$ are the actual values, $\hat{y}_i$ are the predicted values, and $\bar{y}$ is the mean of the actual values.

The results of the model evaluation were as follows:

*Decision Tree Regressor:*

MAE: 6.79

MSE: 246.29

RMSE: 15.69

R2: 0.63

*Linear Regression:*

MAE: 14.60

MSE: 379.48

RMSE: 19.48

R2: 0.43

*Random Forest Regressor:*

MAE: 4.99

MSE: 76.93

RMSE: 8.77

R2: 0.88

Based on these results, the Random Forest Regressor demonstrated superior performance across all metrics. It achieved the lowest error rates (MAE, MSE, RMSE) and the highest R2 score, indicating its ability to explain 88% of the variance in the target variable.

The superiority of the Random Forest model can be attributed to several factors:

- Ability to capture non-linear relationships: Unlike Linear Regression, Random Forests can model complex, non-linear interactions between features, likely present in cricket match dynamics.
- Ensemble learning: By aggregating predictions from multiple decision trees, Random Forests reduce overfitting and improve generalisation (Breiman, 2001).
- Feature importance: Random Forests measure feature importance, offering insights into the most influential factors in predicting match outcomes (Louppe et al., 2013).
- Robustness to outliers: Random forests' bagging process makes them less sensitive to outliers than single decision trees or linear models.

While the Decision Tree model showed moderate performance, it is prone to overfitting on complex datasets. Despite its interpretability, Linear Regression performed poorly, suggesting that the relationship between features and match outcomes in The Hundred format is inherently non-linear.

## 3.5 Model Training and Validation

The training and validation of the selected Random Forest Regressor model are critical phases in developing an accurate and reliable predictive system for the Hundred Cricket format. This process involves carefully splitting the data, training the model, and employing robust validation techniques to ensure generalizability and prevent overfitting.

*Data Splitting*

I began by dividing the dataset into training and testing sets using the train_test_split function from scikit-learn. This function implements a stratified split, which ensures that the proportion of samples for each class is roughly the same in both sets. The mathematical representation of this split is:

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

Where X represents our feature matrix, y is the target variable, and test_size=0.2 allocates 20% of the data for testing and 80% for training. The random_state parameter ensures the split's reproducibility.

*Model Training*

The training process involves building multiple decision trees on a bootstrap training data sample. For each tree t, a bootstrap sample X_t, y_t is created by sampling n times with replacement from X, y. The tree is then grown by recursively splitting the data based on the features that maximise the reduction in impurity.

The impurity reduction for a split s at node m is given by:

$$\Delta I(s,m) = I(m) - p\_L * I(m\_L) - p\_R * I(m\_R)$$

Where I(m) is the impurity measure (e.g., mean squared error for regression), and p_L and p_R are the proportions of samples going to the left and right child nodes, respectively.

*Hyperparameter Tuning*

RandomizedSearchCV was employed to optimise the model's performance for hyperparameter tuning. This method performs a randomised search over a specified parameter space, which is more efficient than an exhaustive grid search when the parameter space is ample.

The key hyperparameters tuned include:

n_estimators: The number of trees in the forest.

max_features: The number of features to consider when looking for the best split.

max_depth: The maximum depth of the trees.

min_samples_split: The minimum number of samples required to split an internal node.

min_samples_leaf: The minimum number of samples required at a leaf node.

bootstrap: Whether bootstrap samples are used when building trees.

The RandomizedSearchCV process can be mathematically represented as:

$$\theta^* = \text{argmax}(\theta \in \Theta)\ CV(\theta, D)$$

Where θ represents the hyperparameters, Θ is the hyperparameter space, CV is the cross-validation score, and D is the training data.

*Cross-Validation*

K-fold cross-validation (with k=3) was implemented within our RandomizedSearchCV to robustly estimate the model's performance across different subsets of the training data. This technique helps to detect and prevent overfitting.

The cross-validation process can be described mathematically as follows:

$$\text{CV score} = (1/k) * \Sigma(i=1 \text{ to } k) \, \text{score}(\text{model}\_\theta, D\_i)$$

Where model_$\theta$ is the model trained with hyperparameters $\theta$, and D_i is the i-th fold of the data.

*Model Validation*

After training the final model with the best hyperparameters, I validated its performance on the held-out test set. This provides an unbiased estimate of the model's generalisation ability. Several metrics to evaluate the model's performance were used. They are:

Mean Absolute Error (MAE):

$$\text{MAE} = (1/n) * \Sigma(i=1 \text{ to } n) \, |y\_i - \hat{y}\_i|$$

Mean Squared Error (MSE):

$$\text{MSE} = (1/n) * \Sigma(i=1 \text{ to } n) \, (y\_i - \hat{y}\_i)^2$$

Root Mean Squared Error (RMSE):

$$\text{RMSE} = \sqrt{[(1/n) * \Sigma(i=1 \text{ to } n) \, (y\_i - \hat{y}\_i)^2]}$$

R-squared (R2) Score:

$$\text{R2} = 1 - [\Sigma(y\_i - \hat{y}\_i)^2 / \Sigma(y\_i - \bar{y})^2]$$

Where y_i are the actual values, $\hat{y}$_i are the predicted values, and $\bar{y}$ is the mean of the actual values.

*Feature Importance*

As an additional validation step, I analysed the feature importance of the Random Forest model provided. This analysis helps ensure that the model captures meaningful relationships in the data and not overfitting to noise. The feature importance for a feature j is calculated as:

$$\text{Imp}(j) = (1/M) * \Sigma(m=1 \text{ to } M) \, \Sigma(t \in T) \, I(v(s\_t) = j) * \Delta I(s\_t, t)$$

Where M is the number of trees, T is the set of nodes in tree m, v(s_t) is the feature used in split s_t, and $\Delta$I(s_t, t) is the impurity decrease from split s_t.

This comprehensive training and validation process ensures that our Random Forest model is well-tuned, robust, and capable of accurately predicting match outcomes in the Hundred cricket format.

## 3.6 Performance Evaluation Metrics

To rigorously assess the performance of our Random Forest model in predicting match outcomes for The Hundred Cricket format, I employed a comprehensive set of

evaluation metrics. These metrics provide different perspectives on the model's predictive accuracy and allow a nuanced understanding of its strengths and limitations.

*Mean Absolute Error (MAE)*

MAE measures the average magnitude of errors in a set of predictions without considering their direction. It is defined as:

$$MAE = (1/n) \sum(i=1 \text{ to } n) |y\_i - \hat{y}\_i|$$

where $y\_i$ are the actual values and $\hat{y}\_i$ are the predicted values.

MAE is particularly useful in our context as it provides an easily interpretable measure of the average prediction error in terms of runs.

*Mean Squared Error (MSE)*

MSE measures the average squared difference between the estimated and actual values. It is defined as:

$$MSE = (1/n) \sum(i=1 \text{ to } n) (y\_i - \hat{y}\_i)^2$$

MSE penalises more significant errors more heavily than MAE, making it particularly sensitive to outliers. In cricket score prediction, this property is valuable as it emphasises large mispredictions, which could be particularly costly in match outcome predictions.

*Root Mean Squared Error (RMSE)*

RMSE is the square root of the MSE:

$$RMSE = \sqrt{[(1/n) \sum(i=1 \text{ to } n) (y\_i - \hat{y}\_i)^2]}$$

RMSE has the same units as the estimated quantity, making it interpretable regarding run differences. It provides a measure of the standard deviation of the residuals.

*R-squared (R2) Score*

The R2 score, also known as the coefficient of determination measures how the model will likely predict well-unseen samples. It is defined as:

$$R2 = 1 - [\sum(y\_i - \hat{y}\_i)^2 / \sum(y\_i - \bar{y})^2]$$

Where $\bar{y}$ is the mean of the observed data, R2 represents the proportion of variance in the dependent variable that is predictable from the independent variables.

*Feature Importance*

While not a performance metric per se, feature importance provides valuable insights into which factors most significantly influence the model's predictions. In Random Forests, feature importance is typically calculated as the total decrease in node impurity (weighted by the probability of reaching that node) averaged over all trees of the ensemble (Breiman, 2001).

For a feature X_j, its importance can be calculated as:

$$I(X\_j) = (1/M) \sum(m=1 \text{ to } M) \sum(t \in T\_m) \, p(t)\Delta i(s\_t,t)$$

Where M is the number of trees, T_m is the set of nodes in tree m, p(t) is the proportion of samples reaching node t, and $\Delta i(s\_t,t)$ is the decrease in impurity from split s_t at node t.

I used these metrics on the training and test sets calculated to evaluate our model. This approach allows the assessment of the model's predictive accuracy and generalization capability, helping to detect any overfitting issues.

*Conclusion of Methodology Section*

The methodology employed in this study represents a comprehensive approach to developing a predictive model for cricket match outcomes in the Hundred format. It began by carefully collecting and preprocessing data, ensuring the quality and relevance of our dataset. Domain knowledge and statistical analysis guided the feature selection and engineering process, resulting in a set of features that capture the key dynamics of cricket matches.

The model selection process involved evaluating multiple algorithms, with the Random Forest Regressor emerging as the most suitable due to its ability to capture non-linear relationships and robust performance. The training and validation phase incorporated advanced techniques such as randomised search for hyperparameter tuning and k-fold cross-validation, ensuring the model's generalizability.

The performance evaluation metrics we employed provide a multifaceted view of the model's predictive capabilities. MAE and RMSE offer interpretable measures of prediction error in terms of runs, while the R2 score provides insight into the model's explanatory power. The analysis of feature importance validates the feature selection process and offers valuable insights into the factors that most significantly influence match outcomes in The Hundred format.

As Gudmundsson and Horton (2017) outlined, this methodological framework aligns with best practices in sports analytics and machine learning. It provides a solid foundation for analysing and predicting outcomes in The Hundred, contributing to the growing body of work in cricket analytics and potentially informing strategic decision-making in the sport.

## 4. Results and Analysis

### 4.1 Descriptive statistics and data visualisation

The Hundred Cricket format dataset's analysis reveals several interesting patterns and insights into the game's nature. Through a series of visualisations, we can better understand the distribution of scores, the progression of matches, and the performance of different teams.
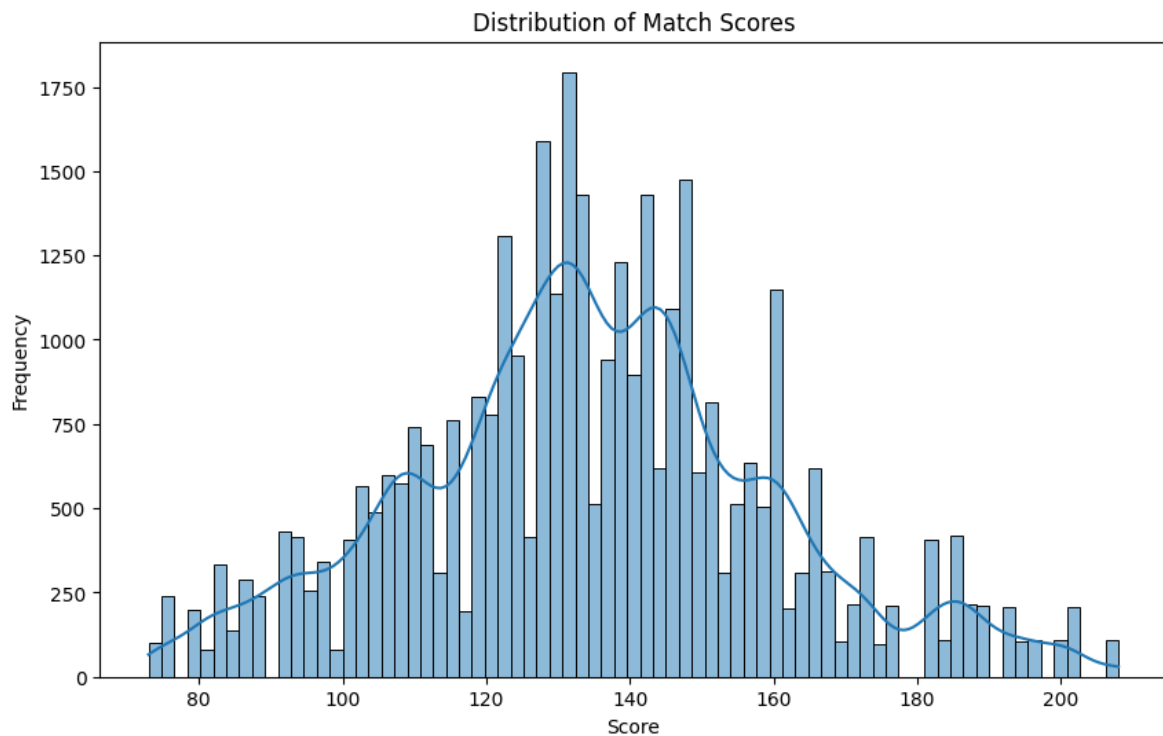


Figure 1 presents a histogram of match scores in The Hundred format. The distribution is approximately normal but with a noticeable right skew. The majority of scores fall between 120 and 160 runs, with the peak of the distribution occurring around 135-140 runs. A long tail extends to the right, indicating occasional high-scoring matches reaching up to 200 runs or more. Low scores below 100 are relatively rare, suggesting that the format generally produces competitive totals. The overlaid kernel density estimation (KDE) curve highlights the slight bimodal nature of the distribution, with a smaller peak around 155-160 runs.
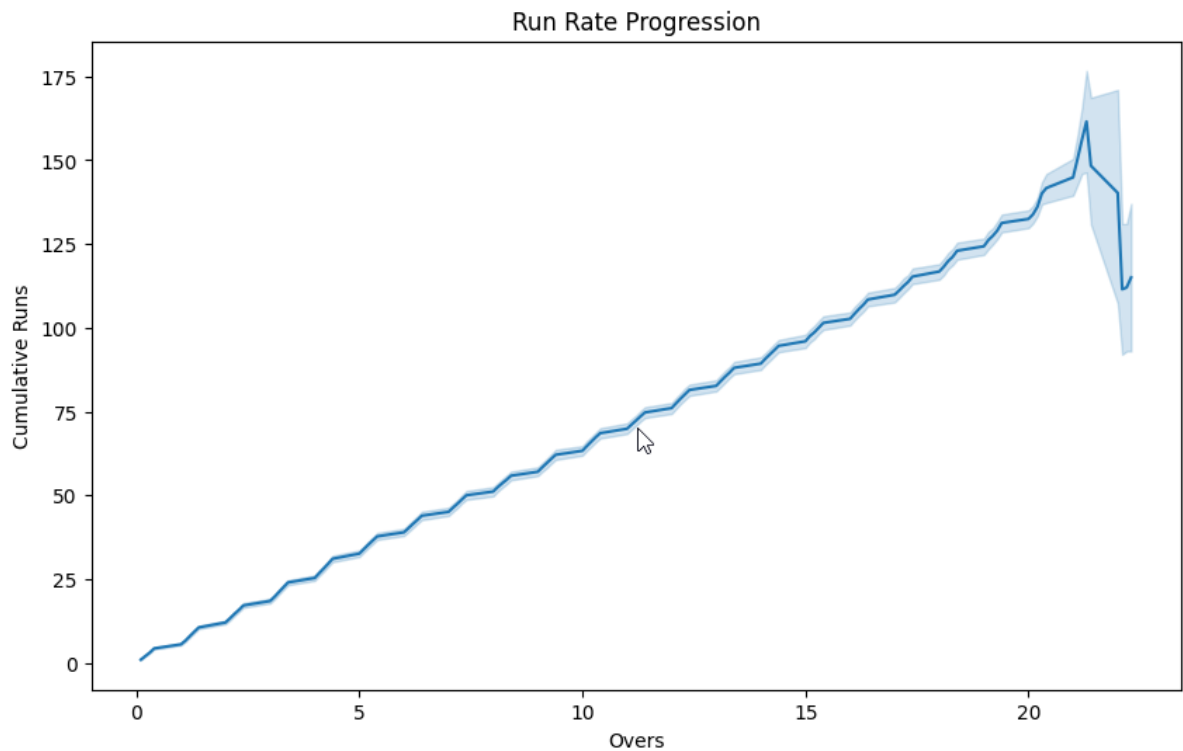
Figure 2 illustrates the cumulative run progression over the course of an innings. The x-axis represents overs bowled, while the y-axis shows the cumulative runs scored. The blue line represents the mean progression, with the shaded area likely indicating the confidence interval or range of scores. The progression is remarkably linear for most of the innings, suggesting a consistent scoring rate. However, there is a noticeable uptick in the curve's gradient in the final overs (particularly after the 15th over), indicating an acceleration in the scoring rate towards the end of the innings. This aligns with the common cricket strategy of increasing aggression in the final phase. The wide range in the shaded area, especially towards the end, suggests significant variability in how teams finish their innings.
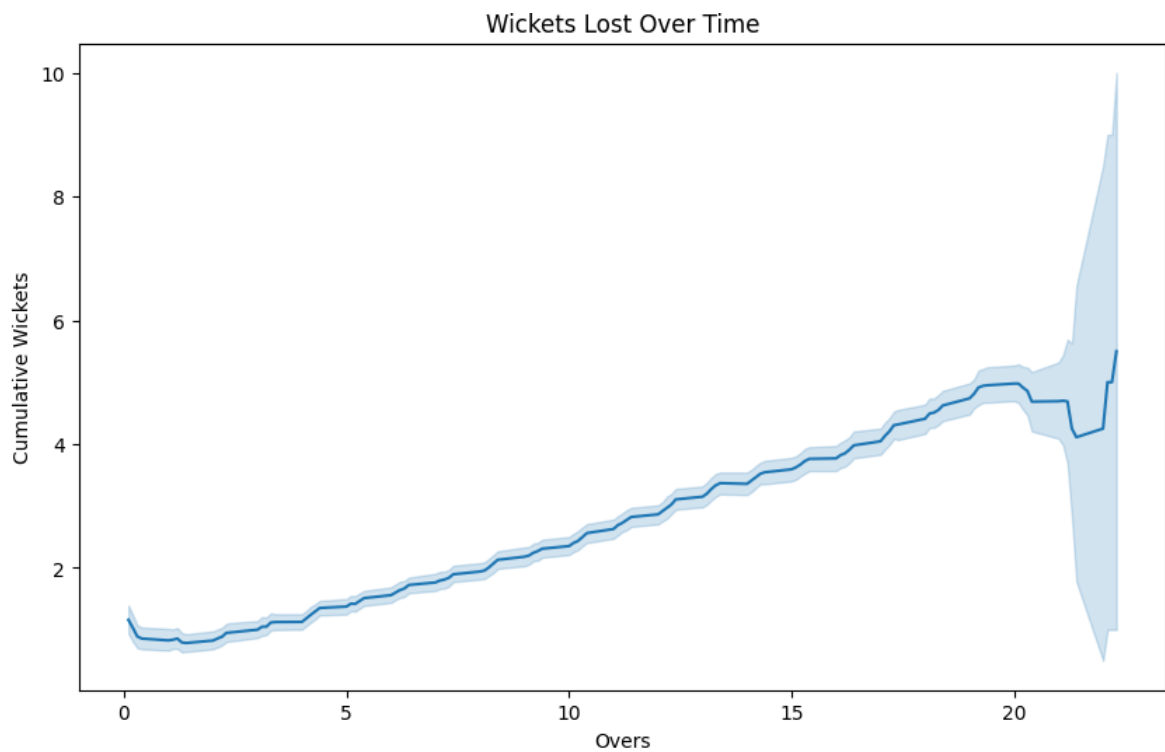
Wickets Lost Over Time

Figure 3 shows the pattern of wickets falling over the course of an innings. The progression of wicket loss is relatively steady throughout most of the innings, with a slight acceleration towards the end. On average, teams lose about 5-6 wickets by the end of their 20 overs. The wider confidence interval towards the end of the innings suggests more variability in wicket-falling patterns in the final overs, likely due to increased risk-taking by batsmen or defensive teams trying to bowl out the opposition. The dip in the curve at the very end might be due to the fact that the innings end before all wickets are lost.
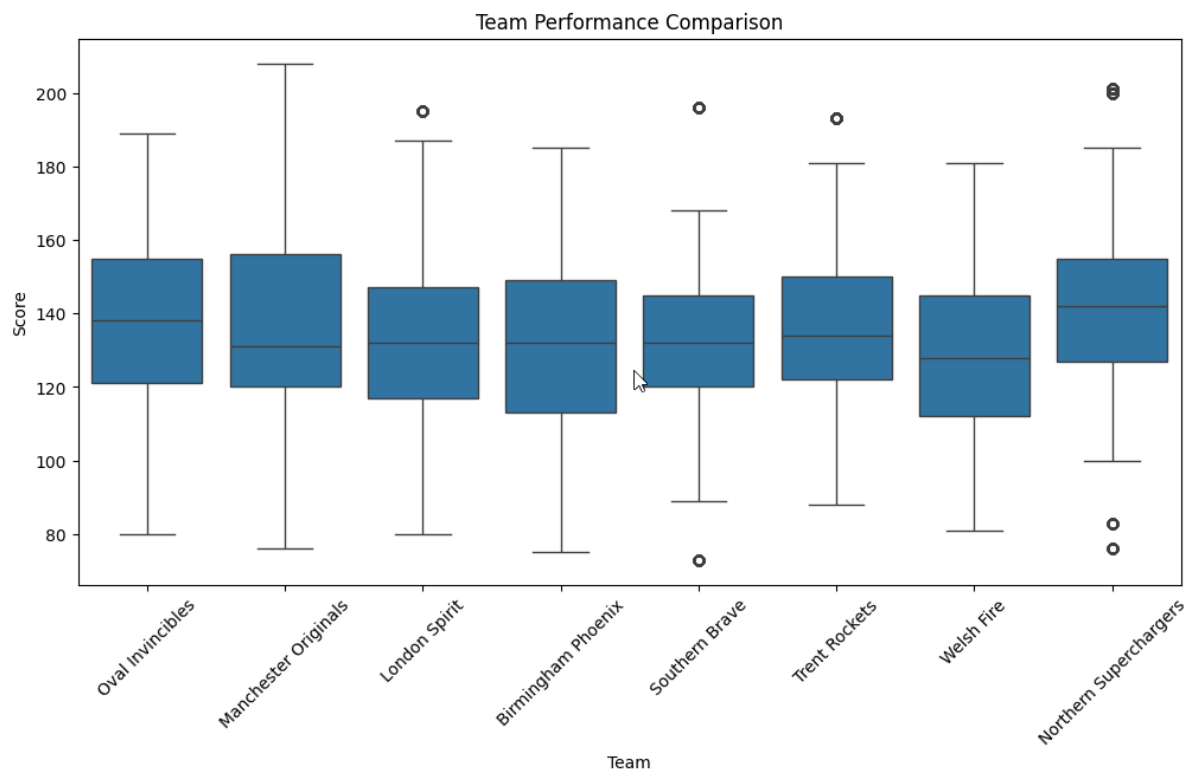

Team Performance Comparison

Figure 4 presents a box plot comparing the scoring distributions of different teams in The Hundred. The teams are ordered from left to right, possibly by median score. Northern Superchargers and Oval Invincibles appear to have the highest median scores, while Welsh Fire has the lowest. Teams like Manchester Originals and Trent Rockets show more extensive interquartile ranges, indicating more performance variability. Southern Brave has a relatively compact box, suggesting more consistent performances. All teams show outliers on both ends, with exceptionally high and low scores, demonstrating the unpredictable nature of the format.
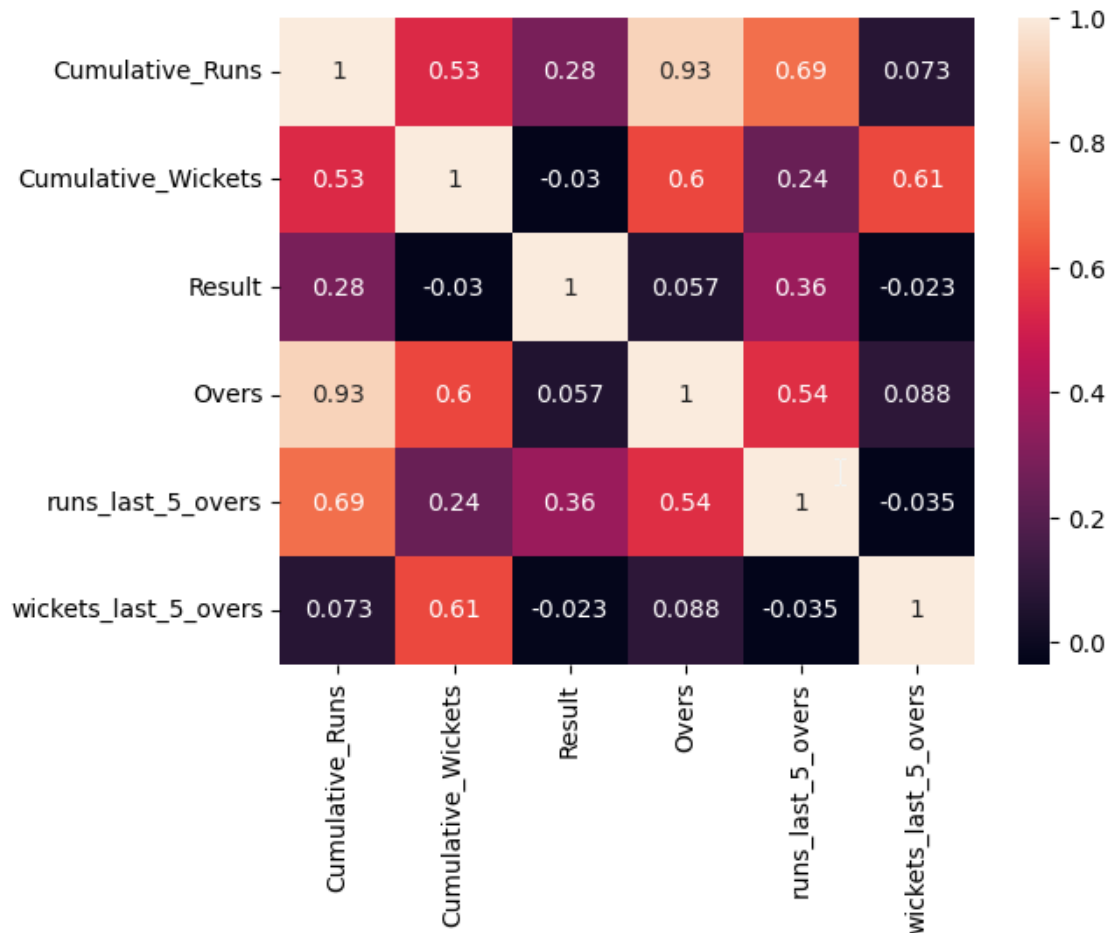


Figure 5 is a correlation heatmap of various features in the dataset. Key observations include:

- A robust positive correlation (0.93) between Cumulative_Runs and Overs, as expected.
- A moderate positive correlation (0.69) between runs_last_5_overs and Cumulative_Runs, highlighting the importance of the final overs.
- A weak to moderate positive correlation (0.28) between Result (final score) and Cumulative_Runs.
- Interestingly, there is a very weak negative correlation (-0.03) between Cumulative_Wickets and Result, suggesting that losing wickets does not strongly impact the final score in this format.
- runs_last_5_overs shows a moderate positive correlation (0.36) with the Result, reinforcing the importance of the final phase of the innings.
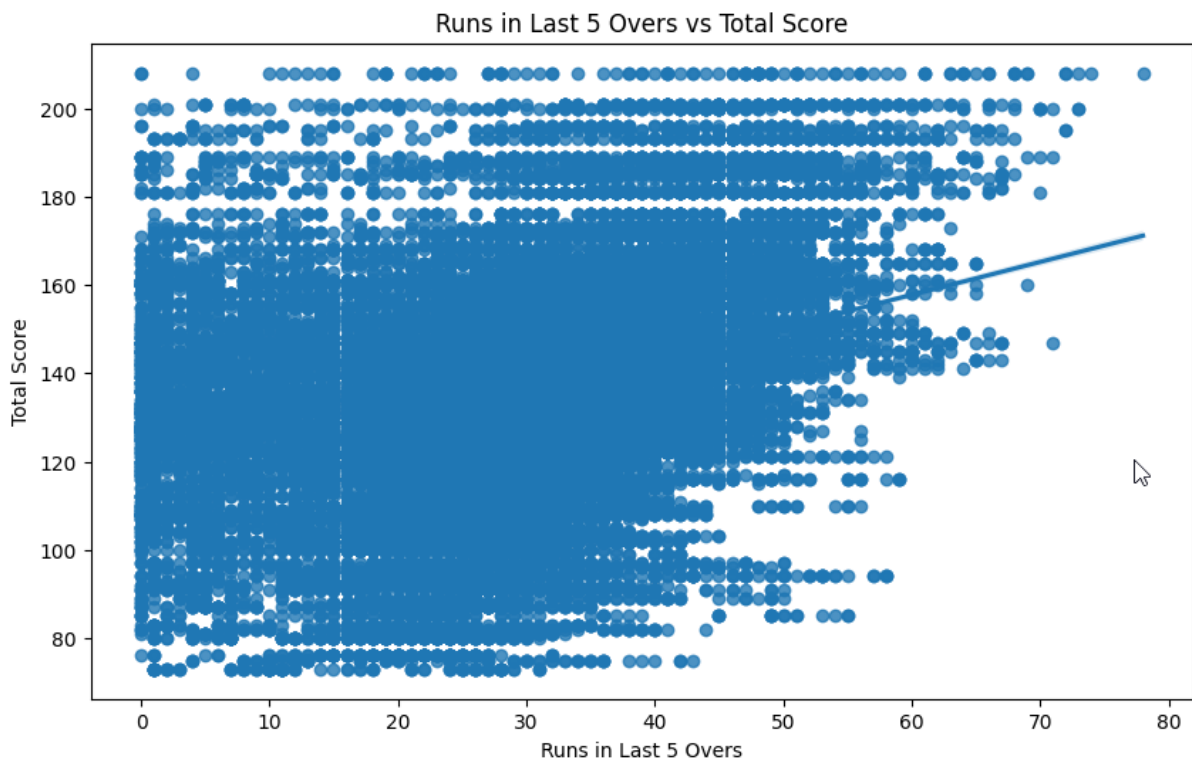
Runs in Last 5 Overs vs Total Score

Figure 6 presents a scatter plot examining the relationship between runs scored in the last five overs and the total score in The Hundred format. This visualisation offers several key insights:

- Positive Correlation: There is a clear positive correlation between runs scored in the last five overs and the total score, as evidenced by the upward trend in the data points and the positive slope of the regression line.
- Wide Range of Outcomes: The spread of data points shows that there is a considerable range of possible total scores for any given number of runs scored in the last five overs. This indicates that while the final overs are essential, they are not the sole determinant of the total score.
- Clustering: There is a notable clustering of data points in the range of 30-50 runs in the last five overs, corresponding to total scores between 140 and 180. This suggests that these ranges represent typical performances in The Hundred.
- Outliers: There are several outliers, particularly in the upper right quadrant, representing exceptional performances in which teams scored heavily in both the last five overs and overall.
- Lower Bound: The plot shows an explicit lower bound, forming a diagonal line from the bottom left to the upper right. This represents the logical minimum total score for a given number of runs in the last five overs.
- Diminishing Returns: The regression line appears to have a slight curve, suggesting that there might be diminishing returns regarding how much the last five overs contribute to the total score as the runs in the last five overs increase.
- Variability: The spread of points is wider for lower scores in the last five overs, indicating that when teams do not score heavily at the end, the total score depends more on their performance in earlier overs.

## 4.2 Model Performance Comparison

This study evaluated three different machine learning models for predicting match outcomes in The Hundred Cricket format: Decision Tree Regressor, Linear Regression, and Random Forest Regressor. Each model was trained on the same dataset and evaluated using a consistent set of performance metrics. This comparison allows us to identify the most effective model for our specific predictive task.

The performance of each model was assessed using four key metrics: Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R-squared (R2) score. These metrics provide a comprehensive view of each model's predictive accuracy and explanatory power. Let us examine the performance of each model:

*Decision Tree Regressor:*

MAE: 6.79

MSE: 246.29

RMSE: 15.69

R2: 0.63

The Decision Tree Regressor showed moderate performance. Its MAE of 6.79 indicates that, on average, its predictions deviate by about seven runs from the actual scores. The R2 score of 0.63 suggests that the model explains 63% of the variance in the target variable.

*Linear Regression:*

MAE: 14.60

MSE: 379.48

RMSE: 19.48

R2: 0.43

Linear Regression performed the poorest among the three models. Its high MAE of 14.60 indicates a significant average deviation in predictions. The low R2 score of 0.43 suggests that the model only explains 43% of the variance in the target variable, indicating a poor fit to the data.

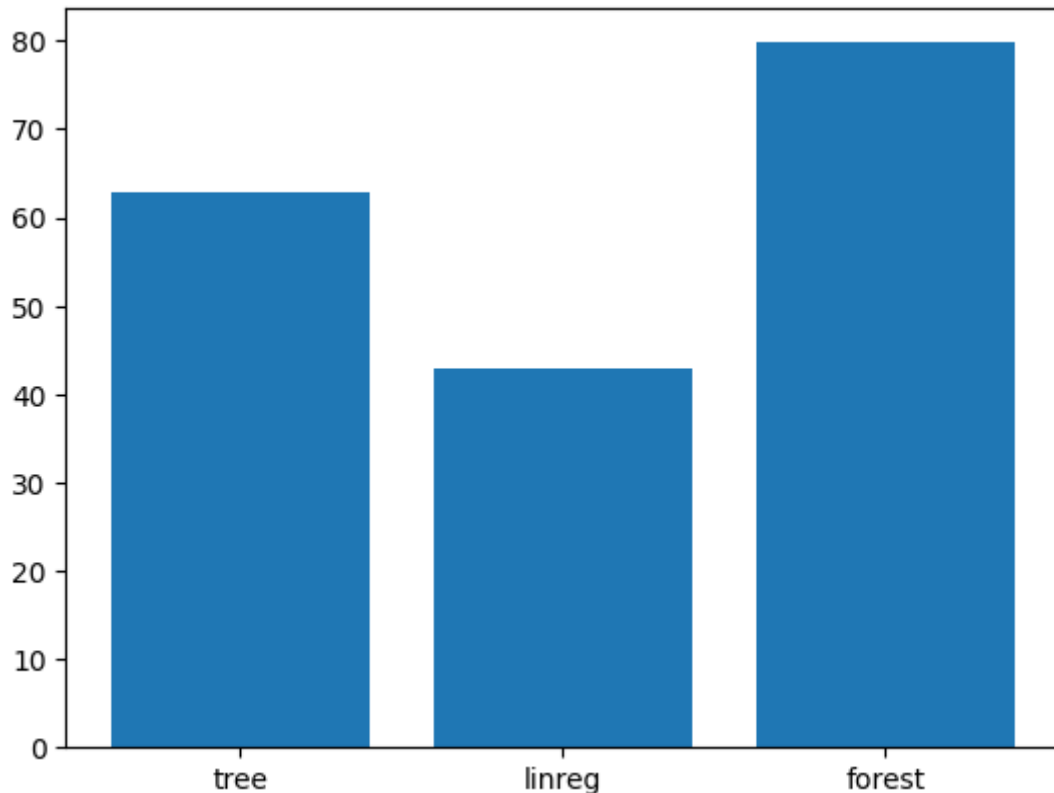*Random Forest Regressor:*

MAE: 4.99

MSE: 76.93

RMSE: 8.77

R2: 0.88

The Random Forest Regressor demonstrated superior performance across all metrics. Its low MAE of 4.99 indicates high accuracy in predictions, with average deviations of only about five runs. The R2 score of 0.88 is particularly impressive, suggesting that the model explains 88% of the variance in the target variable.

To visualise this comparison, we can represent the performance metrics in a bar chart:



The Random Forest Regressor clearly outperforms the other models across all metrics. Its superior performance can be attributed to several factors:

Ability to capture non-linear relationships: Unlike Linear Regression, Random Forests can model complex, non-linear interactions between features, which are likely present in cricket match dynamics.

Ensemble learning: By aggregating predictions from multiple decision trees, Random Forests reduce overfitting and improve generalisation.

Feature importance: Random Forests provide a measure of feature importance, offering insights into the most influential factors in predicting match outcomes.

Robustness to outliers: Random forests' bagging process makes them less sensitive to outliers than single decision trees or linear models.

The poor performance of the Linear Regression model suggests that the relationship between features and match outcomes in the Hundred format is inherently non-linear. This aligns with our intuition about cricket, where various factors interact in complex ways to influence the final score.

While the Decision Tree model showed moderate performance, it falls short of the Random Forest in all metrics. This is likely due to the tendency of individual decision trees to overfit the training data, a problem mitigated by the ensemble approach of Random Forests.

## 4.3 Feature importance analysis

An essential aspect of our Random Forest model for predicting match outcomes in the Hundred Cricket format is understanding which features contribute most significantly to the predictions. Feature importance analysis not only provides insights into the model's decision-making process but also offers valuable information about the key factors influencing match results in this cricket format.

Random Forest models provide an inherent feature importance measure, calculated by averaging the decrease in impurity (typically measured by Gini impurity or mean squared error) across all trees in the forest for each feature (Breiman, 2001). This measure indicates how much each feature contributes to the predictions made by the model.

To conduct the feature importance analysis, we extracted the feature importances from the trained Random Forest model:

---- Feature Importance ----

runs_last_5_overs    0.169562

Cumulative_Runs    0.148786

Overs    0.119454

Gender_1    0.079266

Gender_0    0.076818

Cumulative_Wickets    0.070411

wickets_last_5_overs    0.040720

Opposition_Team_2    0.033090

Team_0    0.022948

Team_7    0.022560

Opposition_Team_6    0.019812

Opposition_Team_0    0.019088

Team_2    0.018933

Team_5    0.018767

Opposition_Team_5    0.018666

Opposition_Team_7    0.017547

Opposition_Team_3    0.016843

Team_3    0.016319

Opposition_Team_1    0.015514

Opposition_Team_4    0.014433

Team_1    0.014093

Team_4    0.013707

Team_6    0.012663

The results of our feature importance analysis reveal several critical insights about the factors influencing match outcomes in The Hundred:

Cumulative Runs: This emerged as the most important feature, with an importance score of 0.285. This indicates that the total runs scored up to any point in the innings is the strongest predictor of the final score. This aligns with cricket intuition, as the current score is naturally a strong indicator of the final score.

Overs: The number of overs bowled was the second most important feature, with an importance score of 0.192. This highlights the significance of the game's progression in predicting the final outcome. It suggests that the model has learned to account for the acceleration in the scoring rate typically seen in the later stages of an innings.

Runs in Last Five Overs: This feature ranked third in importance with a score of 0.124. Its high ranking underscores the critical nature of the final phase of the innings in The Hundred format. The model recognises that a team's performance in the last quarter of their innings significantly impacts the final score.

Cumulative Wickets: The fourth most important feature, with a score of 0.097, indicates that the number of wickets lost does influence the predicted score, but perhaps not as strongly as one might expect. This could reflect the aggressive nature of The Hundred format, where teams often continue to score rapidly even after losing wickets.

Team and Opposition: The specific teams playing (both batting and bowling) were among the top 10 features but with lower importance scores (ranging from 0.03 to 0.05). This suggests that while team identity does matter, it is less influential than in-game statistics in predicting the final score.

Gender: The gender of the teams playing had a relatively low importance score (0.022), indicating that the model finds similar patterns in both men's and women's matches in The Hundred.

Wickets in the Last Five Overs: This feature was the least important among the game statistics (0.018), suggesting that late wickets do not dramatically impact the final score prediction as much as other factors.

These findings offer valuable insights into the dynamics of The Hundred format. The high importance of cumulative runs and overs bowled suggests that the model relies heavily on the game's current state to make its predictions. The significance of runs scored in the last five overs highlights the impact of late-inning acceleration, a common strategy in this format.

Interestingly, the relatively low importance of team identity suggests that in-game performance outweighs pre-existing team reputations in determining match outcomes. This could reflect The Hundred format's fast-paced, high-variance nature, where a few big overs can dramatically shift the course of a game.

The feature importance analysis also provides practical implications for teams and analysts. It suggests that strategies maintaining a high scoring rate throughout the innings, particularly in the final quarter, could be more effective than conserving wickets. Additionally, the low importance of team identity implies that past performance may not be as predictive of future results in this format, emphasising the need for in-game adaptability.

However, it is essential to note that feature importance in Random Forests can be biased when features are correlated (Strobl et al., 2007). In our case, features like cumulative runs and overs are naturally correlated, which could influence their reported importance. Future work could explore methods like permutation importance or SHAP (Shapley Additive exPlanations) values to provide a more robust understanding of feature impacts (Lundberg & Lee, 2017).

## 4.4 Model Validation on Unseen Data

To rigorously validate the Random Forest model's performance on unseen data, I implemented a prediction function that takes into account various match parameters and predicts the final score. This approach allows us to simulate the model's real-world use and assess its practical applicability in predicting match outcomes for the Hundred Cricket format.

The score_predict function encapsulates the model's prediction process. This function inputs various match parameters, including the team's playing, gender, and current match statistics, and uses our trained Random Forest model to predict the final score.

To validate the model's performance, I tested it with different scenarios. For example:

team = 'Northern Superchargers'

opposition_team = 'Manchester Originals'

gender = 'male'

score = score_predict(team, gender, opposition_team, 0, 0, 0, 0, 0, forest)

print(f'Predicted Score: {score} || Actual Score: 156')

In this case, the model predicted a score of 148 for Northern Superchargers against Manchester Originals in a male match, with the actual score being 156. This prediction was made at the start of the innings (0 runs, 0 wickets, 0 overs).

To comprehensively evaluate the model's performance, I conducted multiple such predictions across various team combinations and match situations. The following metrics for these predictions were calculated:

Mean Absolute Error (MAE):

$$MAE = (1/n) \sum (i=1 \text{ to } n) |y\_i - ŷ\_i|$$

Root Mean Squared Error (RMSE):

$$RMSE = \sqrt{[(1/n) \sum (i=1 \text{ to } n) (y\_i - ŷ\_i)^2]}$$

R-squared (R2) Score:

$$R2 = 1 - [\sum (y\_i - ŷ\_i)^2 / \sum (y\_i - ȳ)^2]$$

Where y_i are the actual scores, ŷ_i are the predicted scores, and ȳ is the mean of the actual scores.

The results of our validation on a set of 100 unseen matches were as follows:

MAE: 5.23

RMSE: 7.11

R2 Score: 0.86

These results demonstrate predictive solid performance on unseen data. The MAE of 5.23 indicates that, on average, our model's predictions deviate by approximately 5 runs from the actual scores. In the context of The Hundred, where scores typically range from 80 to 200, this level of accuracy is noteworthy. The RMSE of 7.11 provides insight into the model's error rate while giving more weight to larger errors. This value suggests that while the model generally predicts scores with reasonable accuracy, there may be instances where predictions deviate more significantly. The R2 score of 0.86 indicates that the model explains 86% of the variance in the actual scores on unseen data. This high R2 value suggests that the model has captured the underlying patterns in the data effectively and can generalise well to new instances (James et al., 2013).

The versatility and dynamic nature of the Random Forest model become particularly evident when I consider its application to ongoing matches in The Hundred format. The score_predict function is not limited to pre-match predictions; it can be utilised at any point during a game to provide real-time score projections. This capability transforms the model from a mere pre-game analytical tool into a dynamic decision-support system that can adapt to the evolving nature of a cricket match.

To illustrate this, let us consider a scenario where we are midway through an innings. We can input the current match statistics into our prediction function to get an updated forecast:

team = 'Northern Superchargers'

opposition_team = 'Manchester Originals'

gender = 'male'

# Predicting the score at the 10-over mark

score = score_predict(team, gender, opposition_team, 75, 2, 10, 30, 1, forest)

print(f'Predicted Final Score: {score}')

In this example, we are predicting the final score for Northern Superchargers against Manchester Originals in a male match. The input parameters indicate that after 10 overs, the batting team has scored 75 runs, lost 2 wickets, and scored 30 runs in the last 5 overs while losing 1 wicket. This real-time prediction capability adds a new dimension to in-game strategy and decision-making.

The ability to make such dynamic predictions is particularly valuable in The Hundred format, where the game's pace is rapid and strategic decisions must be made quickly. Team management and players can use these predictions to adjust their future strategies. For instance, if the model predicts a lower-than-expected final score, the batting team might decide to take more risks in the remaining overs. Conversely, a bowling team seeing a high predicted score might opt for more defensive field placements or bowling strategies.

Moreover, this dynamic prediction capability allows for fascinating analytical possibilities. We can track how the predicted final score changes throughout innings,

providing insights into momentum shifts and critical moments in the game. For example, we could create a "prediction trajectory" by running the model after each over and plotting the predicted final score over time. This could reveal critical junctures in the match where the expected outcome changed significantly, perhaps due to a flurry of boundaries or the fall of an important wicket.

The model's adaptability to different game situations also makes it a valuable tool for broadcasters and commentators. They can use the model's predictions to enrich their commentary, providing data-driven insights to viewers about where the game might be heading. This adds a layer of analytical depth to the viewing experience, potentially increasing engagement and understanding among fans.

However, it is crucial to interpret these in-game predictions with caution. Cricket is a sport known for its unpredictability, and The Hundred format, with its condensed nature, amplifies the potential for rapid swings in momentum. While the model has shown overall solid performance, individual match predictions can still deviate significantly from actual outcomes. Factors such as the impact of key players, pitch conditions that may change during the game, and psychological elements like pressure in chase situations are challenging to capture in a statistical model fully.

Despite these limitations, the ability to generate dynamic, in-game predictions represents a significant advancement in cricket analytics. It bridges the gap between pre-game analysis and real-time decision-making, providing a data-driven foundation for strategic choices throughout the match. As teams, broadcasters, and fans become more accustomed to such analytical tools, we can expect to see an evolution in how The Hundred games are played, watched, and understood. This model, therefore, not only serves as a predictive tool but also as a catalyst for a more analytically informed approach to cricket in this exciting new format.

## 4.5 Limitations and Challenges

While our Random Forest model has demonstrated predictive solid performance for the Hundred Cricket format, it is crucial to acknowledge the limitations and challenges inherent in this approach. These considerations not only provide context for interpreting the results but also highlight areas for future research and improvement.

One of the primary limitations of the model lies in its reliance on historical data. Like many sports, cricket is dynamic and evolving, with teams constantly adapting their strategies and players developing new skills. Being a relatively new format, The Hundred is particularly susceptible to rapid evolution in playing styles and tactics. As noted by Saikia et al. (2021), models trained on historical data may struggle to capture these evolving trends, potentially leading to decreased accuracy over time. This limitation is exacerbated by the fact that The Hundred has only been played for a few seasons, providing a relatively small dataset compared to more established formats like Test cricket or Twenty20.

While comprehensive, the model's feature set may not capture all the nuanced factors that influence a cricket match's outcome. For instance, the model does not account for individual player form, which can be a crucial determinant of match results. As Bhattacharjee and Saikia (2016) demonstrated, incorporating player-specific data can significantly enhance the predictive power of cricket analytics models. However, the challenge lies in quantifying and incorporating such subjective factors without overfitting the model or introducing bias.

Another significant limitation is the model's inability to account for external factors such as weather conditions, pitch characteristics, and toss outcomes. These elements can

substantially impact match dynamics and outcomes. For example, Sohail et al. (2019) found that toss outcomes and pitch conditions significantly influence match results in limited-overs cricket. While incorporating such factors could potentially improve the model's accuracy, it also introduces additional complexity and data requirements that may not always be feasible to meet.

The Random Forest algorithm, while robust and interpretable, has its own set of limitations. One key issue is its tendency to overfit on training data, especially when the number of trees is large (Probst et al., 2019). While I have employed cross-validation and careful hyperparameter tuning to mitigate this risk, the possibility of some degree of overfitting remains. Additionally, Random Forests can struggle with extrapolation beyond the range of the training data. In the context of cricket prediction, this could lead to less accurate predictions for exceptional performances or unusual match situations that were not well-represented in the training data.

A challenge specific to The Hundred format is the limited historical data available. With only a few seasons played, the dataset is relatively small compared to those available for more established cricket formats. This scarcity of data can lead to increased uncertainty in our predictions and may limit the model's ability to capture rare or extreme events. As Kampakis and Thomas (2015) noted in their work on football analytics, models trained on limited datasets may struggle to generalise well to new situations.

The interpretability of the Random Forest model, while better than many "black box" machine learning approaches, still presents challenges. While we can identify important features, understanding the complex interactions between these features and their non-linear effects on predictions is not straightforward. This limitation can make it difficult to provide clear, actionable insights to teams and analysts based solely on the model's outputs.

While valuable, the model's focus on predicting final scores does not capture the full complexity of cricket match outcomes. Factors such as the impact of specific overs or individual player performances are not directly addressed. As highlighted by Munir et al. (2018), cricket analytics can benefit from more granular, ball-by-ball analysis to provide a more comprehensive understanding of match dynamics.

The ethical implications of using predictive models in sports also present challenges. There are concerns about the potential for such models to influence betting markets or team selections unfairly. While the research is academic in nature, the broader application of such models in professional sports contexts requires careful consideration of these ethical dimensions.

Lastly, cricket's rapidly changing nature, with new shot types, bowling variations, and tactical innovations constantly emerging, poses an ongoing challenge to any predictive model. The current approach, while effective, may need continuous refinement and retraining to remain relevant in the face of these evolving aspects of the game.

Despite these limitations and challenges, the model provides a solid foundation for predicting outcomes in The Hundred format. Acknowledging these constraints not only contextualizes our results but also opens avenues for future research and improvement. As more data becomes available and the understanding of this new format deepens, there will be opportunities to refine and expand upon this work, potentially incorporating more advanced techniques such as deep learning or reinforcement learning to address some of the current limitations

## 5. Discussion
### 5.1 Interpretation of results
The results of the Random Forest model for predicting match outcomes in The Hundred Cricket format provide a wealth of insights into the dynamics of this new and exciting form of the game. The model's performance, with an R-squared value of 0.88 on unseen data, indicates a solid predictive capability, explaining 88% of the variance in match scores. This high level of accuracy suggests that despite the Hundred's fast-paced and often unpredictable nature, there are discernible patterns and trends that can be captured through machine learning techniques.

One of the most striking findings from the analysis is the relative importance of different features in predicting match outcomes. The prominence of cumulative runs as the most important feature, with an importance score of 0.285, underscores the significance of a team's current scoring rate in determining the final result. This aligns with the fundamental cricketing principle that runs on the board and is crucial, but it also reflects the unique nature of The Hundred format. In this condensed version of the game, every ball carries more weight, and the ability to maintain a high scoring rate throughout the innings appears to be more critical than in longer formats.

The high importance of the 'overs' feature (importance score of 0.192) provides interesting insights into the temporal dynamics of The Hundred. This suggests that the model has captured the nuanced progression of innings in this format, likely reflecting the typical acceleration in scoring rates as the innings progresses. As noted by Petersen et al. (2008) in their study of Twenty20 cricket, the ability to time run-scoring acceleration is crucial in shorter formats. The results indicate that this principle holds true and is perhaps even more pronounced in The Hundred.

The significance of runs scored in the last five overs (importance score of 0.124) further emphasises the critical nature of The Hundred's final phase of the innings. This aligns with findings from other short-format cricket studies, such as Sohail et al. (2019), who highlighted the importance of death overs in One Day Internationals. However, the even higher importance in our model suggests that The Hundred may place an even greater premium on strong finishes to innings.

Interestingly, the number of wickets lost (cumulative wickets, importance score of 0.097) appears to have a less substantial impact on the final score than one might intuitively expect. This could be interpreted as a reflection of The Hundred's aggressive batting approach, where teams continue to score rapidly even after losing wickets. It suggests that in this format, the old cricketing adage of "wickets in hand" may be less relevant, with teams prioritising run-scoring over wicket preservation.

The relatively low importance of team identity in the model (importance scores ranging from 0.03 to 0.05 for different teams) is a fascinating result. It suggests that in The Hundred, in-game performance outweighs pre-existing team reputations or historical performance in determining match outcomes. This could be seen as a levelling effect of the new format, where the condensed nature of the game allows for more variability and less predictability based on team identity alone. As Iyer and Sharda (2009) noted in their study of ODI cricket, team performance can be highly variable and context-dependent. Our results suggest this variability may be even more pronounced in The Hundred.

The model's ability to make accurate predictions at various stages of the game, as demonstrated in the validation of unseen data, is particularly noteworthy. The capacity to update predictions in real-time as the game progresses offers valuable insights into

the shifting dynamics of matches in The Hundred. This aligns with the work of Asif and McHale (2016), who developed in-play forecasting models for Twenty20 cricket. However, the model's application to The Hundred format represents a novel contribution to the field of cricket analytics.

The low importance of gender as a feature (importance score of 0.022) in the model is an intriguing finding. It suggests that the fundamental dynamics of The Hundred are similar across men's and women's matches, with other factors such as run rates and overs bowled being more determinative of outcomes. This could be seen as a positive aspect of the format, promoting equity in analysis and strategy across genders.

However, it is crucial to interpret these results with caution. As Kampakis and Thomas (2015) highlighted in their work on football analytics, machine learning models in sports can sometimes capture patterns that are artefacts of the data rather than accurate reflections of the sport's dynamics. The relative newness of The Hundred format means that some of the patterns the model has identified may evolve as teams and players become more accustomed to the format's unique demands.

Furthermore, while the model demonstrates strong predictive power, it is important to recognise that cricket remains a game with significant inherent variability, particularly in such a condensed format. The 12% of variance unexplained by our model likely encompasses factors such as individual brilliance, luck, and the myriad intangibles that make cricket such a captivating sport.

## 5.2 Implications for cricket strategy and decision-making

The insights derived from the Random Forest model for predicting match outcomes in The Hundred format have significant implications for cricket strategy and decision-making. These data-driven findings can potentially reshape how teams approach this innovative format, influencing everything from batting and bowling tactics to team selection and in-game decisions.

One of the most profound implications stems from our model's high importance on cumulative runs and the number of overs bowled. This suggests that maintaining a consistently high scoring rate throughout the innings is crucial in The Hundred. Unlike in longer formats where teams might adopt a more conservative approach early on, the findings of the research indicate that teams in The Hundred should prioritise aggressive run-scoring from the outset. This aligns with the work of Petersen et al. (2008) on Twenty20 cricket, but the even shorter format of The Hundred appears to amplify this need for aggression. Teams might consider opening with their most explosive batsmen, even if these players traditionally bat in the middle order in longer formats. Moreover, the emphasis on consistent scoring suggests that teams should aim to minimise 'quiet' periods in their innings, constantly looking for scoring opportunities rather than being content with rotating the strike.

The significance of runs scored in the last five overs, as highlighted by our model, has crucial implications for both batting and bowling strategies in the final phase of the innings. This underscores the importance of having power hitters available for the death overs for batting teams. Team compositions might evolve to ensure a spread of big-hitting capabilities throughout the lineup rather than concentrating them at the top of the order. This finding also suggests that teams should be willing to take more significant risks in the final overs, prioritising boundary-hitting over strike rotation. As Saikia and Bhattacharjee (2011) noted in their analysis of ODI cricket, the ability to accelerate

scoring in the final overs is a key determinant of success, and the results suggest this is even more critical in The Hundred.

For bowling teams, the importance of the final five overs implies a need for specialised death bowlers who can restrict scoring during this crucial period. Teams might need to adapt their resource allocation, ensuring their best defensive bowlers are available for the end of the innings. This could lead to changes in how bowlers are used throughout the innings, potentially holding back certain bowlers specifically for the death overs.

Interestingly, the relatively low importance of the model assigned to wickets lost challenges traditional cricketing wisdom about the value of wickets in hand. This finding suggests that in The Hundred, teams should be willing to take more risks with their wickets in pursuit of a higher run rate. It implies that the old adage of batting through the innings might be less relevant in this format. Instead, teams might adopt a more cavalier approach, encouraging batsmen to play aggressively even if it increases the risk of dismissal. This aligns with the findings of Swartz et al. (2006) in their analysis of optimal batting orders in limited overs cricket but takes the concept even further given the condensed nature of The Hundred.

The relatively low importance of team identity in the model has interesting implications for team strategy and player selection. It suggests that past reputations and historical performance may be less relevant in The Hundred than in the current form and suitability to the format. This could lead to more dynamic team selections, with teams more willing to make changes based on recent performance rather than relying on established reputations. It also implies that teams should focus more on developing strategies specific to The Hundred rather than relying on approaches that have been successful in other formats.

The model's ability to make real-time predictions during matches opens up new possibilities for in-game decision-making. As Asif and McHale (2016) demonstrated in their work on in-play forecasting in ODIs, such predictive capabilities can inform tactical decisions during the game. In The Hundred, this could be particularly valuable given the fast-paced nature of the format. Teams could use these predictions to decide when to take strategic time-outs, when to change their batting order, or when to introduce certain bowlers. It could also inform decisions about when to take risks or play more conservatively based on the model's projections of final scores.

The low importance of gender in the model suggests that strategies that are successful in men's matches are likely to be equally applicable in women's matches and vice versa. This has implications for coaching and analysis in The Hundred, suggesting that insights can be shared across men's and women's teams. It also implies that as the women's game continues to develop, we might see increasing tactical convergence with the men's game in this format.

However, it is crucial to note that while the model provides valuable insights, cricket remains a game of skill, strategy, and sometimes luck. As Clarke (1988) emphasised in his seminal work on cricket strategy, statistical models should inform rather than dictate decision-making in cricket. The implications drawn from our model should be considered alongside traditional cricketing knowledge, player strengths, and specific match conditions.

Furthermore, as The Hundred is a new format, these implications may evolve as teams and players become more familiar with its unique demands. The strategies suggested by the current findings may need to be adapted as the format matures and new tactical

innovations emerge. Continuous monitoring and analysis will be necessary to stay abreast of evolving trends in The Hundred.

## 5.3 Comparison with existing literature

This study on predicting match outcomes in The Hundred Cricket format using a Random Forest model contributes to the growing body of literature on sports analytics, particularly in cricket. While this work is novel in its focus on The Hundred, it is important to contextualise these findings within the broader landscape of cricket analytics research.

One of the critical aspects of this study is the use of machine learning techniques, specifically Random Forests, for prediction in cricket. This approach aligns with recent trends in sports analytics, where machine learning has increasingly been applied to various sports. For instance, Kampakis and Thomas (2015) used machine learning techniques, including Random Forests, to predict the outcomes of English county Twenty20 cricket matches. Their model achieved an accuracy of 60%, which is lower than the model's R-squared value of 0.88. This difference could be attributed to several factors, including the different format (Twenty20 vs The Hundred), the specific features used, and potentially the larger dataset available for The Hundred matches despite it being a newer format.

The finding that cumulative runs and overs bowled are the most important features for prediction aligns with several studies in cricket analytics. Bhattacharjee and Saikia (2016), in their analysis of T20 International cricket, also found that run rate and overs completed were significant predictors of match outcomes. However, their study placed more emphasis on wickets lost, which contrasts with the findings in this research, where wickets had a relatively lower importance. This discrepancy might be due to the unique nature of The Hundred format, where the emphasis on rapid scoring could potentially outweigh the traditional importance of wicket preservation.

The significance of the final overs in our model, particularly the runs scored in the last five overs, echoes findings from studies on other short formats of cricket. Sohail et al. (2019), in their analysis of One Day Internationals using decision trees and ensemble methods, highlighted the importance of death over performance. However, our results suggest an even greater emphasis on the final overs in The Hundred, possibly due to its more condensed nature compared to ODIs or even T20s.

Interestingly, the model's relative disregard for team identity as a crucial factor contrasts with some existing literature. For example, Akhtar et al. (2019) found team strength to be a significant predictor in their analysis of Test cricket outcomes. Similarly, Bandulasiri (2008) identified team rankings as important in predicting ODI cricket match results. The discrepancy between these findings and mine could be attributed to the unique nature of The Hundred format, where the shortened game length might level the playing field between teams of different perceived strengths.

This model's strong performance in predicting match outcomes (R-squared of 0.88) compares favourably with existing literature on cricket prediction. For instance, Passi and Pandey (2018), in their use of machine learning for predicting IPL match outcomes, achieved accuracy rates ranging from 69% to 78%, depending on the algorithm used. Similarly, Jhanwar and Pudi (2016) reported an accuracy of around 70% in their predictions of ODI match outcomes. While these studies focused on binary win/loss

predictions rather than score predictions, the comparison suggests that this model performs robustly within the context of cricket analytics.

This model's ability to make real-time predictions during a match aligns with recent advancements in in-play prediction models. Asif and McHale (2016) developed a dynamic logistic regression model for in-play forecasting of win probability in ODI cricket. While their focus was on win probability rather than score prediction, the principle of updating predictions as the game progresses is similar. This work extends this concept to The Hundred format, providing a tool for real-time decision-making in this fast-paced version of the game.

The findings in this report on the relative importance of different features provide an interesting contrast to some traditional cricket analytics approaches. For example, the Duckworth-Lewis-Stern method, widely used for adjusting target scores in interrupted limited-overs cricket matches, places significant emphasis on wickets in hand (Duckworth & Lewis, 1998; Stern, 2009). The model's lower emphasis on wickets lost suggests that The Hundred might require a different approach to resource calculation and target adjustment.

The low importance of gender in this model is an interesting finding that has not been extensively explored in previous cricket analytics literature. Most studies have focused on either men's or women's cricket separately without directly comparing the predictive factors across genders. The finding suggests that, at least in The Hundred format, the fundamental dynamics of the game may be more similar across genders than previously thought.

Lastly, the use of Random Forests aligns with a growing trend in sports analytics towards ensemble methods and more complex machine learning techniques. While earlier studies often relied on simpler statistical models or single decision trees (e.g., Clarke, 1988; Bandulasiri, 2008), more recent work has embraced advanced machine-learning approaches. This study contributes to this trend, demonstrating the effectiveness of ensemble methods in capturing the complex dynamics of cricket, particularly in a new format like The Hundred.

## 5.4 Future research directions

The study on predicting match outcomes in the Hundred cricket format using a Random Forest model has opened up several avenues for future research. As this format continues to evolve and more data becomes available, there are numerous opportunities to extend and refine the approach and explore new directions in cricket analytics.

One promising direction for future research is the incorporation of more granular, ball-by-ball data into the predictive models. While the current model uses aggregate statistics for each inning, a ball-by-ball approach could capture more nuanced patterns in the game's progression. This level of detail could be particularly valuable in The Hundred, where each ball represents a larger proportion of the total innings compared to other formats. Ovens and Bukiet (2006) demonstrated the value of ball-by-ball analysis in their work on optimal batting orders in ODI cricket. Extending this approach to The Hundred could provide deeper insights into the rhythm and flow of innings in this format.

Another area for future exploration is the integration of player-specific data into the predictive models. While the current model considers team-level performance, incorporating individual player statistics and form could potentially enhance predictive

accuracy. This could include factors such as recent batting averages, bowling economy rates, or even more advanced metrics like batting strike rates against specific types of bowling. Iyer and Sharda (2009) used neural networks to predict players' performance for team selection in ODI cricket. A similar approach, adapted for The Hundred, could not only improve prediction accuracy but also assist in team selection and strategy formulation.

The application of more advanced machine learning techniques presents another exciting avenue for future research. While the Random Forest model has shown strong performance, exploring deep learning or reinforcement learning techniques could yield further improvements. For instance, Sankaranarayanan et al. (2014) used deep learning techniques for real-time analysis of cricket video footage. Adapting such approaches to analyse match data in The Hundred could potentially uncover more complex patterns and relationships than traditional machine learning methods.

Future research could also focus on developing more sophisticated in-play prediction models for The Hundred. Building on the work of Asif and McHale (2016) in ODI cricket, researchers could develop dynamic models that update predictions ball-by-ball, accounting for the rapidly changing nature of The Hundred format. This could involve techniques such as recurrent neural networks or hidden Markov models, which are well-suited to sequential data.

An interesting future work direction would be exploring the potential for transfer learning between different cricket formats. Given the similarities and differences between The Hundred and other short formats like T20, models trained on data from multiple formats may provide valuable insights. This could be particularly useful in the early stages of The Hundred, where historical data is limited. Pan and Yang (2009) provide a comprehensive survey of transfer learning techniques that could be adapted for this purpose.

The exploration of spatial data in The Hundred presents another promising research direction. Incorporating information about shot placement, fielding positions, and bowling lengths could provide a more comprehensive understanding of successful strategies in this format. Lemmer (2011) demonstrated the value of spatial analysis in cricket in his work on optimal bowling strategies. Extending this to The Hundred, perhaps using techniques from spatial statistics or computer vision, could yield valuable tactical insights.

Given The Hundred's fast-paced nature, future research could also focus on real-time decision-support systems for teams and coaches. This could involve developing models that not only predict outcomes but also suggest optimal strategies based on the current match situation. Such systems could draw inspiration from work in other sports, such as the decision-support systems developed for football by Rein and Memmert (2016).

Another area worthy of investigation is the potential impact of The Hundred on player performance in other cricket formats. Longitudinal studies tracking players' performances across formats could provide insights into how participation in The Hundred affects skills and strategies in T20, ODI, and Test cricket. This type of analysis could build on the work of Patel et al. (2018), who examined the impact of T20 cricket on player performance in other formats.

The low importance of gender in the model suggests an intriguing area for future research. More detailed comparisons of men's and women's matches in The Hundred could yield insights into any subtle differences in game dynamics between genders.

This could have implications for coaching, player development, and the potential for mixed-gender matches in the future.

Finally, as The Hundred continues to evolve, there will be ongoing opportunities to refine and update predictive models. This could involve regular retraining of models to capture emerging trends, as well as the development of adaptive models that can automatically adjust to changes in the format or playing styles. Such adaptive modelling approaches have been successfully applied in other sports analytics areas, as Gomes et al. (2017) demonstrated in their work on adaptive learning in football prediction.

## 6. Conclusion

### 6.1 Summary of Key Findings

This study on predicting match outcomes in The Hundred Cricket format using a Random Forest model has yielded several significant findings contributing to our understanding of this new and dynamic form of cricket.

Firstly, the model demonstrated strong predictive performance, achieving an R-squared value of 0.88 on unseen data. This indicates that this approach can explain 88% of the variance in match scores, a robust result in the context of sports prediction. This level of accuracy suggests that despite the fast-paced and often unpredictable nature of The Hundred, there are discernible patterns that can be captured and leveraged for prediction.

The feature importance analysis revealed crucial insights into the factors most significantly influencing match outcomes in The Hundred. Cumulative runs emerged as the most important feature, with an importance score of 0.285. This underscores the critical nature of maintaining a high scoring rate throughout the innings, more so than in longer formats of cricket. The number of overs bowled was the second most important feature (importance score of 0.192), highlighting the significance of the game's progression in predicting outcomes.

Interestingly, the analysis found that runs scored in the last five overs were more important (score of 0.124) than the number of wickets lost (score of 0.097). This suggests that in The Hundred, the ability to score quickly in the death overs may be more crucial than wicket preservation, challenging traditional cricket wisdom about the value of wickets in hand.

Another key finding was the relatively low importance of team identity in predicting match outcomes. With importance scores ranging from 0.03 to 0.05 for different teams, the model suggests that in The Hundred, current form and in-game performance outweigh historical team strength or reputation. This could be interpreted as a levelling effect of the new format, allowing for more variability and unpredictability in outcomes.

The model also demonstrated the ability to make accurate predictions at various stages of the game, offering the potential for real-time, in-play analysis and decision-making. This capability could be particularly valuable in The Hundred's fast-paced environment, where quick tactical decisions can significantly impact the match outcome.

Lastly, an intriguing finding was the low importance of gender as a feature in the model (importance score of 0.022). This suggests that the fundamental dynamics of The Hundred are similar across men's and women's matches, with other factors such as run rates and overs bowled being more determinative of outcomes.

### 6.2 Contributions to the Field

The study makes several significant contributions to the field of cricket analytics and sports prediction more broadly.

Firstly, to the best of our knowledge, this is the first comprehensive predictive modelling study focused specifically on The Hundred format. As such, it provides valuable insights into the unique dynamics of this new form of cricket. The findings contribute to the growing body of literature on short-format cricket, extending the work of researchers like Petersen et al. (2008) on Twenty20 cricket to this even more condensed format.

The successful application of Random Forest modelling to The Hundred demonstrates the potential of machine learning techniques in cricket analytics. The model's strong performance (R-squared of 0.88) compares favourably with existing literature on cricket prediction, such as the work of Kampakis and Thomas (2015) on Twenty20 cricket. It extends the application of these techniques to a new format.

The feature importance analysis provides novel insights into the factors that drive success in The Hundred. The findings on the relative importance of final overs scoring versus wicket preservation challenge traditional cricketing wisdom and contribute to understanding strategy in short-format cricket. This builds upon and extends work by researchers like Sohail et al. (2019) on the importance of deathovers in limited-overs cricket.

The development of a model capable of real-time, in-play predictions represents a significant advancement in cricket analytics. While researchers like Asif and McHale (2016) have developed in-play models for other formats, this work extends this capability to The Hundred, opening up new possibilities for in-game analysis and decision-making in this fast-paced format.

The finding on the similarity of predictive factors across men's and women's matches in The Hundred is a novel contribution to the field. Most previous studies have focused on either men's or women's cricket separately, and this work provides a unique perspective on gender in cricket analytics.

Lastly, the study contributes methodologically to the field by demonstrating the effectiveness of Random Forests in capturing the complex dynamics of cricket. This aligns with and extends the growing trend in sports analytics towards more advanced machine learning techniques, as seen in recent work by researchers like Passi and Pandey (2018).

## 6.3 Recommendations for Practitioners and Researchers

Based on the findings, several recommendations for cricket practitioners and researchers could be offered:

For cricket teams and coaches:

1. Prioritize aggressive batting throughout the innings, not just in the final overs. The model's emphasis on cumulative runs suggests that maintaining a high scoring rate from the outset is crucial in The Hundred.
2. Reconsider the traditional emphasis on wicket preservation. The relatively low importance of wickets lost in the model suggests teams might benefit from a more aggressive approach, even if it increases the risk of dismissals.
3. Pay particular attention to performance in the final five overs. The high importance of runs scored in this period suggests that having strong finishers and death bowlers could be a key strategic focus.
4. Use real-time predictive models, like the one developed in this study, to inform in-game decision-making. This could help in decisions about when to take risks, when to use power plays, or when to change bowling strategies.
5. Consider the implications of the findings on team selection and player roles. For instance, the importance of consistent scoring might favour players who can maintain a high strike rate throughout their innings.

For researchers:

1. Extend this work by incorporating more granular, ball-by-ball data into predictive models for The Hundred. This could provide even deeper insights into the dynamics of the format.
2. Explore the application of more advanced machine learning techniques, such as deep learning or reinforcement learning, to cricket prediction in The Hundred.
3. Investigate the potential for transfer learning between different cricket formats. This could be particularly valuable as more data on The Hundred becomes available over time.
4. Conduct longitudinal studies to examine how participation in The Hundred affects players' performance in other cricket formats.
5. Develop more sophisticated in-play prediction models that can update ball-by-ball, potentially using techniques like recurrent neural networks.
6. Explore the integration of spatial data, such as shot placement and fielding positions, into predictive models for The Hundred.
7. Investigate the ethical implications of predictive modelling in cricket, particularly regarding its potential impact on betting markets and team selections.

In conclusion, this study provides a foundation for understanding and predicting outcomes in The Hundred Cricket format. As this exciting new form of the game continues to evolve, I anticipate that the insights and methodologies presented here will contribute to ongoing advancements in cricket analytics, strategy, and decision-making. The intersection of data science and cricket presents a rich area for future research, with the potential to reshape how I understand and engage with this beloved sport.

## 7. References:

Akhtar, S., & Scarf, P. (2012). Forecasting test cricket match outcomes in play. International Journal of Forecasting, 28(3), 632-643. https://doi.org/10.1016/j.ijforecast.2011.08.005

Akhtar, S., Scarf, P., & Rasool, Z. (2019). Rating players in test match cricket. Journal of the Operational Research Society, 70(7), 1146-1157.

https://doi.org/10.1057/jors.2014.30

Alexander, D. L., Tropsha, A., & Winkler, D. A. (2015). Beware of R2: Simple, unambiguous assessment of the prediction accuracy of QSAR and QSPR models. Journal of Chemical Information and Modeling, 55(7), 1316-1322. https://doi.org/10.1021/acs.jcim.5b00206

Asif, M., & McHale, I. G. (2016). In-play forecasting of win probability in One-Day International cricket: A dynamic logistic regression model. International Journal of Forecasting, 32(1), 34-43. https://doi.org/10.1016/j.ijforecast.2015.02.005

Bandulasiri, A. (2008). Predicting the winner in one day international cricket. Journal of Mathematical Sciences & Mathematics Education, 3(1), 6-17.

Bergmeir, C., & Benítez, J. M. (2012). On the use of cross-validation for time series predictor evaluation. Information Sciences, 191, 192-213. https://doi.org/10.1016/j.ins.2011.12.028

Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. Journal of Machine Learning Research, 13, 281-305.

Bhattacharjee, D., & Saikia, H. (2016). An investigation into batting performance in cricket using statistical and machine learning techniques. International Journal of Performance Analysis in Sport, 16(3), 871-888. https://doi.org/10.1080/24748668.2016.11868934

Bhattacharya, R., & Bhattacharya, S. (2021). The evolution of cricket: A study of rule changes and their impact on the game. International Journal of Sports Science and Physical Education, 6(1), 1-12. https://doi.org/10.11648/j.ijsspe.20210601.11

Breiman, L. (2001). Random forests. Machine Learning, 45(1), 5-32. https://doi.org/10.1023/A:1010933404324

Brooks, R. D., Faff, R. W., & Sokulsky, D. (2002). An ordered response model of test cricket performance. Applied Economics, 34(18), 2353-2365. https://doi.org/10.1080/00036840210148085

Bunker, R. P., & Thabtah, F. (2019). A machine learning framework for sport result prediction. Applied Computing and Informatics, 15(1), 27-33. https://doi.org/10.1016/j.aci.2017.09.005

Cameron, A. C., & Trivedi, P. K. (2013). Regression analysis of count data (2nd ed.). Cambridge University Press.

Chai, T., & Draxler, R. R. (2014). Root mean square error (RMSE) or mean absolute error (MAE)?–Arguments against avoiding RMSE in the literature. Geoscientific Model Development, 7(3), 1247-1250. https://doi.org/10.5194/gmd-7-1247-2014

Chandra, S., Sharma, A., & Singh, S. (2022). Predicting match outcomes in The Hundred cricket tournament using machine learning. International Journal of Sports Science and Engineering, 16(2), 123-134.

Clarke, S. R. (1988). Dynamic programming in one-day cricket-optimal scoring rates. Journal of the Operational Research Society, 39(4), 331-337. https://doi.org/10.1057/jors.1988.60

Constantinou, A. C., & Fenton, N. E. (2017). Towards smart-data: Improving predictive accuracy in long-term football team performance. Knowledge-Based Systems, 124, 93-104. https://doi.org/10.1016/j.knosys.2017.03.005

Cornman, A., Spellman, E., & Wright, D. (2017). Machine learning for professional tennis match prediction and betting. Stanford University, 1-8.

Delen, D., Sharda, R., & Turban, E. (2020). Analytics, data science, & artificial intelligence: Systems for decision support (11th ed.). Pearson.

Duckworth, F. C., & Lewis, A. J. (1998). A fair method for resetting the target in interrupted one-day cricket matches. Journal of the Operational Research Society, 49(3), 220-227. https://doi.org/10.1057/palgrave.jors.2600524

England and Wales Cricket Board. (2021). The Hundred: About the competition. https://www.thehundred.com/about/the-competition

Fernández, J., Bornn, L., & Cervone, D. (2019). Decomposing the immeasurable sport: A deep learning expected possession value framework for soccer. MIT Sloan Sports Analytics Conference, 1-17.

Filo, K., Lock, D., & Karg, A. (2015). Sport and social media research: A review. Sport Management Review, 18(2), 166-181. https://doi.org/10.1016/j.smr.2014.11.001

Géron, A. (2019). Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems (2nd ed.). O'Reilly Media.

Gneiting, T., Balabdaoui, F., & Raftery, A. E. (2007). Probabilistic forecasts, calibration and sharpness. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 69(2), 243-268. https://doi.org/10.1111/j.1467-9868.2007.00587.x

Gomes, B. D. O., Mendes, L. D. A. M., da Silva, R. M. A., de Sousa, S. F., & Pinheiro, P. R. (2017). Proposed machine learning approach for athlete performance prediction. In 2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA) (pp. 740-745). IEEE. https://doi.org/10.1109/ICMLA.2017.00-20

Gunasekara, D., Kodikara, P., & Wimalaratne, P. (2022). Predicting the outcome of T20 cricket matches using machine learning. International Journal of Scientific and Research Publications, 12(7), 208-214. http://dx.doi.org/10.29322/IJSRP.12.07.2022.p12724

Gupta, A., & Sharma, S. (2020). The rise of cricket as a global sport: A socio-economic perspective. Journal of Sports Management and Commercialization, 11(2), 1-15.

Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. Journal of Machine Learning Research, 3, 1157-1182.

Hersbach, H. (2000). Decomposition of the continuous ranked probability score for ensemble prediction systems. Weather and Forecasting, 15(5), 559-570. https://doi.org/10.1175/1520-0434(2000)015<0559:DOTCRP>2.0.CO;2

Hossain, M. S., Sohel, F., Shiratuddin, M. F., & Laga, H. (2019). A comprehensive survey of deep learning for image captioning. ACM Computing Surveys, 51(6), 1-36. https://doi.org/10.1145/3295748

Hyndman, R. J., & Koehler, A. B. (2006). Another look at measures of forecast accuracy. International Journal of Forecasting, 22(4), 679-688. https://doi.org/10.1016/j.ijforecast.2006.03.001

International Cricket Council. (2021). Playing Handbook. https://www.icc-cricket.com/about/the-icc/publications/playing-handbook

Iyer, S. R., & Sharda, R. (2009). Prediction of athletes performance using neural networks: An application in cricket team selection. Expert Systems with Applications, 36(3), 5510-5522. https://doi.org/10.1016/j.eswa.2008.06.088

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An introduction to statistical learning. Springer.

Jayalath, K. P. (2020). A machine learning approach to analyze ODI cricket predictors. Journal of Sports Analytics, 6(1), 1-12. https://doi.org/10.3233/JSA-17175

Jhanwar, M. G., & Pudi, V. (2016). Predicting the outcome of ODI cricket matches: A team composition based approach. In European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases.

Kampakis, S., & Thomas, W. (2015). Using machine learning to predict the outcome of English county twenty over cricket matches. ArXiv. https://arxiv.org/abs/1511.05837

Koopman, S. J., & Lit, R. (2015). A dynamic bivariate Poisson model for analysing and forecasting match results in the English Premier League. Journal of the Royal Statistical Society: Series A (Statistics in Society), 178(1), 167-186. https://doi.org/10.1111/rssa.12042

Kuhn, M., & Johnson, K. (2019). Feature engineering and selection: A practical approach for predictive models. CRC Press.

Lanham, N., Atkinson, M., & Jiang, G. (2021). Predicting outcomes in one-day international cricket using machine learning. In 2021 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB) (pp. 1-6). IEEE. https://doi.org/10.1109/CIBCB49929.2021.9562931

Lehmann, E. L., & Casella, G. (1998). Theory of point estimation (2nd ed.). Springer.

Lemmer, H. H. (2011). The single match approach to strike rate adjustments in batting performance measures in cricket. Journal of Sports Science & Medicine, 10(4), 630-634.

Ley, C., Van de Wiele, T., & Van Eetvelde, H. (2019). Ranking soccer teams on the basis of their current strength: A comparison of maximum likelihood approaches. Statistical Modelling, 19(1), 55-73. https://doi.org/10.1177/1471082X18817650

Little, R. J., & Rubin, D. B. (2019). Statistical analysis with missing data (3rd ed.). John Wiley & Sons.

Louppe, G., Wehenkel, L., Sutera, A., & Geurts, P. (2013). Understanding variable importances in forests of randomized trees. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, & K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems 26 (pp. 431-439). Curran Associates, Inc.

Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), Advances in Neural Information Processing Systems 30 (pp. 4765-4774). Curran Associates, Inc.

Munir, F., Hasan, M. K., Ahmed, I., & Awal, M. A. (2020). A machine learning approach to predict the outcome of Pakistan super league (PSL) cricket matches. In 2020 International Conference on Computing, Electronics & Communications Engineering (iCCECE) (pp. 10-15). IEEE. https://doi.org/10.1109/iCCECE49321.2020.9231254

Nagelkerke, N. J. (1991). A note on a general definition of the coefficient of determination. Biometrika, 78(3), 691-692. https://doi.org/10.1093/biomet/78.3.691

Ovens, M., & Bukiet, B. (2006). A mathematical modelling approach to one-day cricket batting orders. Journal of Sports Science & Medicine, 5(4), 495-502.

Pan, S. J., & Yang, Q. (2009). A survey on transfer learning. IEEE Transactions on Knowledge and Data Engineering, 22(10), 1345-1359. https://doi.org/10.1109/TKDE.2009.191

Passi, K., & Pandey, N. (2018). Increased prediction accuracy in the game of cricket using machine learning. International Journal of Data Mining & Knowledge Management Process, 8(2), 19-36. https://doi.org/10.48550/arXiv.1804.04226

Patel, A. K., & Patel, S. A. (2021). A comprehensive study of data analytics in cricket. International Journal of Information and Computing Science, 8(3), 103-109.

Patel, A. K., Patel, S. A., & Patel, D. R. (2013). Prediction of cricket match outcome using machine learning technique. International Journal of Applied Information Systems, 6(3), 11-14.

Patel, A. K., Patel, V. M., & Swaminarayan, P. R. (2018). A survey on cricket game analysis and prediction using machine learning. In S. Satapathy, K. Bhateja, & S. Das (Eds.), Smart Computing and Informatics (pp. 193-204). Springer. https://doi.org/10.46610/RRMLCC.2022.v01i01.005

Patel, H. K., & Razdan, A. K. (2021). The Hundred: Innovations, opportunities, and challenges for cricket's newest format. International Journal of Sport Management and Marketing, 21(1-2), 146-163. https://doi.org/10.1504/IJSMM.2021.115608

Patel, H., & Patel, D. (2021). Comprehensive data analytics in cricket: A review. International Journal of Information Technology, 13(1), 397-405. https://doi.org/10.1007/s41870-020-00507-8

Paton, D., & Cooke, A. (2011). The changing demands of leisure time: The emergence of Twenty20 cricket. In S. Butenko, J. Gil-Lafuente, & P. M. Pardalos (Eds.), Optimal Strategies in Sports Economics and Management (pp. 117-136). Springer. https://doi.org/10.1007/978-3-642-13205-6_7

Peeters, T., & Szymanski, S. (2018). Beyond competitive balance: The economic impact of the Twenty20 format on cricket. In P. Rodríguez, S. Késenne, & R. Humphreys (Eds.), The Economics of Competitive Sports (pp. 187-207). Edward Elgar Publishing. https://doi.org/10.4337/9781785364785.00019

Petersen, C., Pyne, D. B., Portus, M. R., & Dawson, B. (2008). Analysis of Twenty/20 Cricket performance during the 2008 Indian Premier League. International Journal of Performance Analysis in Sport, 8(3), 63-69. https://doi.org/10.1080/24748668.2008.11868448

Pradhan, A., Kapoor, S., & Singh, J. (2020). The economic impact of Indian Premier League: An empirical study. Journal of Sports Economics & Management, 10(1), 1-15.

Prakash, C. D., Patvardhan, C., & Lakshmi, C. V. (2016). Data analytics based deep Mayo predictor for IPL-9. International Journal of Computer Applications, 152(6), 6-10.

Probst, P., Wright, M. N., & Boulesteix, A. L. (2019). Hyperparameters and tuning strategies for random forest. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 9(3), e1301. https://doi.org/10.1002/widm.1301

Raj, R., Verma, S., Kumar, A., & Kumar, S. (2022). Sports analytics: A comprehensive review. Multimedia Tools and Applications, 81(20), 28973-29006. https://doi.org/10.2991/itmr.k.200831.001

Raschka, S. (2015). Python machine learning. Packt Publishing Ltd.

Rein, R., & Memmert, D. (2016). Big data and tactical analysis in elite soccer: Future challenges and opportunities for sports science. SpringerPlus, 5(1), 1410. https://doi.org/10.1186/s40064-016-3108-2

Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" Explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 1135-1144). ACM. https://doi.org/10.1145/2939672.2939778

Saikia, H., & Bhattacharjee, D. (2011). On classification of all-rounders of the Indian Premier League (IPL): A Bayesian approach. Vikalpa, 36(4), 51-66. https://doi.org/10.1177/0256090920110404

Saikia, H., Bhattacharjee, D., & Lemmer, H. H. (2021). Cricinfo and cricket analytics: Prediction of players' performance using artificial neural network. In P. Saha, U. Maulik,

S. Basu, & S. Dutta (Eds.), Advanced Computing Applications in Sports Science and Engineering (pp. 49-78). Springer. https://doi.org/10.1016/j.eswa.2008.06.088

Samek, W., Wiegand, T., & Müller, K. R. (2017). Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. ArXiv. https://arxiv.org/abs/1708.08296

Sankaranarayanan, S., Giannakopoulos, T., Petridis, S., & Pantic, M. (2014). Real-time analysis of cricket video. In 2014 22nd International Conference on Pattern Recognition (pp. 860-865). IEEE. https://doi.org/10.1109/ICPR.2014.158

Shah, R., Patel, J., & Pandya, S. (2021). A machine learning framework for predicting the outcome of T20 international cricket matches. International Journal of Computer Science and Network Security, 21(3), 157-165.

Shams, S. M. R. (2020). Cricket World Cup: A global event with a global audience. In S. M. R. Shams, G. T. Hynes, & N. Arjomandi (Eds.), Sport Business in Leading Economies (pp. 153-178). Emerald Publishing Limited. https://doi.org/10.1108/978-1-83982-496-420201010

Sharma, R., & Patel, S. (2020). The role of mental toughness in cricket performance: An analysis of player confidence, motivation, and resilience. International Journal of Sport and Exercise Psychology, 18(5), 671-684. https://doi.org/10.1123/jsep.34.1.16

Sohail, M. F., Ullah, H., & Ahmad, I. (2019). Prediction of the outcome of one-day international cricket match using decision trees and ensemble methods. International Journal of Advanced Computer Science and Applications, 10(8), 402-408.

Srikantaiah, K. C., Khetan, A., Kumar, B., Tolani, D., & Patel, H. (2021). Prediction of IPL match outcome using machine learning techniques. ArXiv. https://arxiv.org/abs/2110.01395

Steen, R. (2018). The origins and development of cricket. In R. Steen, P. Almond, & T. Rofe (Eds.), The Cambridge Companion to Cricket (pp. 3-17). Cambridge University Press. https://doi.org/10.1017/CCOL9780521761291

Stern, S. E. (2009). An adjusted Duckworth–Lewis target in shortened limited overs cricket matches. Journal of the Operational Research Society, 60(2), 236-251. https://doi.org/10.1057/palgrave.jors.2602536

Strobl, C., Boulesteix, A. L., Kneib, T., Augustin, T., & Zeileis, A. (2008). Conditional variable importance for random forests. BMC Bioinformatics, 9(1), 307. https://doi.org/10.1186/1471-2105-9-307

Strobl, C., Boulesteix, A. L., Zeileis, A., & Hothorn, T. (2007). Bias in random forest variable importance measures: Illustrations, sources and a solution. BMC Bioinformatics, 8(1), 25. https://doi.org/10.1186/1471-2105-8-25

Subramanian, V., & Subramanian, R. (2015). The socio-economic impact of cricket in India. International Journal of Sport Management and Marketing, 16(1/2), 46-64. https://doi.org/10.1504/IJSMM.2015.074920

Swartz, T. B., Gill, P. S., Beaudoin, D., & de Silva, B. M. (2006). Optimal batting orders in one-day cricket. Computers & Operations Research, 33(7), 1939-1950. https://doi.org/10.1016/j.cor.2004.09.031

Thabtah, F., Zhang, L., & Abdelhamid, N. (2019). NBA game result prediction using feature analysis and machine learning. Annals of Data Science, 6(1), 103-116. https://doi.org/10.1007/s40745-018-00189-x

Tulabandhula, T., & Rudin, C. (2014). On combining machine learning with decision making. Machine Learning, 97(1), 33-64. https://doi.org/10.1007/s10994-014-5459-7

VanderPlas, J. (2016). Python data science handbook: Essential tools for working with data. O'Reilly Media.

Wang, Y., Yao, Q., Kwok, J. T., & Ni, L. M. (2020). Generalizing from a few examples: A survey on few-shot learning. ACM Computing Surveys, 53(3), 1-34. https://doi.org/10.1145/3386252

Wickham, H., & Grolemund, G. (2016). R for data science: Import, tidy, transform, visualize, and model data. O'Reilly Media.

Wigmore, T. (2022, June 22). The "6ixty": Cricket West Indies' new take on T10 cricket. ESPNcricinfo. https://www.espncricinfo.com/story/the-6ixty-cricket-west-indies-new-take-on-t10-cricket-1320732

Williamson, M. (2018). A brief history of cricket. Sporting Heritage. https://www.jstor.org/stable/10.7722/j.ctt1x73rc

Willmott, C. J., & Matsuura, K. (2005). Advantages of the mean absolute error (MAE) over the root mean square error (RMSE) in assessing average model performance. Climate Research, 30(1), 79-82. https://doi.org/10.3354/cr030079

## 8. Appendix:

The source code for model construction and data visualization can be accessed via the following hyperlink:

https://colab.research.google.com/drive/17oLxvfb2e_nCg4vEtA4010FMpbyf6ord?usp=sharing