

Predicting Accuracy of Players in Cricket using Machine Learning

Madhurantakam Pavan Teja*, Kalva Gopal*, Ballem Jeevan Kumar*, Dhara Swapna*,
P.Vinay Kumar**

* Student, Department of IT, Malla Reddy Engineering College, Maisammaguda, Hyderabad.

**Asst.Professor, Department of IT, Malla Reddy Engineering College, Maisammaguda, Hyderabad.

Abstract:

Cricket is one of the most widely played and analyzed sports globally, with player performance evaluation being a crucial aspect of team selection, match strategy, and talent identification. Traditional statistical models have long been used to assess a player's potential; however, with the advent of machine learning and data-driven analytics, performance prediction has become more sophisticated and accurate. This project focuses on predicting the accuracy of a cricket player based on historical data and background information using advanced machine learning techniques. The study explores various independent variables, including batting and bowling averages, strike rates, match experience, fitness levels, age, weather conditions, and opposition strength, to determine their impact on a player's overall accuracy. The dependent variable in this analysis is the player accuracy score, which serves as a measure of performance consistency and reliability. Despite advancements in predictive analytics, challenges such as variability in playing conditions, limited historical data for emerging players, and non-linear performance trends remain key obstacles. However, by integrating statistical, machine learning, and deep learning approaches, this study aims to bridge the gap between traditional cricket analytics and modern AI-driven performance evaluation. The results of this research can be beneficial for team management, player development, and fantasy sports, enabling data-driven decision-making in the cricketing world.

Keywords: *Decision Trees, SVR, KNN, Random Forest, Cricket Analytics, Regression Models, Player Accuracy Score, Machine Learning in Cricket.*

I. INTRODUCTION

Cricket has evolved into a sport where data-driven decision-making plays a crucial role in determining team composition, player selection, and match strategies. With advancements in analytics and machine learning, teams, coaches, and analysts rely on statistical models to assess player performance, predict outcomes, and optimize game plans. Performance prediction in cricket is an emerging field that combines historical data, player attributes, and contextual factors to estimate future performance with a reasonable degree of accuracy.

In the modern era, the integration of Artificial Intelligence (AI) and Machine Learning (ML) has revolutionized the way performance data is analyzed in cricket. These technologies have made it possible to process vast amounts of data from past matches, training sessions, player biomechanics, and real-time game information. AI algorithms can detect patterns in player behavior, such as batting techniques under specific conditions or a bowler's consistency across different pitch types. Machine Learning models can predict a player's future accuracy based on

historical performance, conditions, and even psychological factors that might influence a player's mindset on the field. For example, predictive models can estimate a bowler's likelihood of taking wickets in varying match situations or assess how a batter's technique might hold up against a certain type of bowler. This level of precision allows coaches, analysts, and team managers to make data-driven decisions about player selection, match strategies, and training regimens, helping to maximize each player's contribution to the team's success.

Furthermore, predicting player accuracy through AI-powered models allows for more efficient resource management. Teams can optimize player workloads, prevent injuries by managing fatigue, and tailor individual training programs to focus on improving areas of weakness identified through predictive analysis. As AI models continue to evolve and incorporate more diverse data, the accuracy of these predictions improves, providing a clearer understanding of a player's strengths and areas for development. In this data-driven era, the use of AI to predict player accuracy offers a competitive edge, giving teams not only a strategic advantage but also a more

holistic approach to player development and match preparation. Ultimately, by efficiently predicting player performance and accuracy, cricket teams can stay ahead of the curve, maximizing their chances of success both on and off the field.

II. LITERATURE SURVEY

The prediction of athletic performance has gained significant attention in recent years due to advancements in data analytics and machine learning (ML). Sports performance prediction is essential for optimizing team strategies, enhancing player selection, and improving individual performance metrics. Various studies have explored different methodologies, ranging from statistical models to complex machine learning algorithms, to forecast outcomes like scores, win probabilities, and individual player performance.

2.1. Machine Learning in Sports Analytics

Machine learning techniques have revolutionized sports analytics, enabling the analysis of vast datasets to uncover patterns and predict future performances. Muthuswamy and Lam [1] demonstrated the potential of neural networks in predicting bowler performance in One-Day Internationals (ODIs), highlighting how ML can effectively handle complex, nonlinear relationships in sports data. Similarly, Iyer and Sharda [3] applied neural networks to predict athletes' performance, emphasizing the flexibility of ML models in cricket Performance Prediction Models. batting and bowling performances. Wickramasinghe [4] focused on predicting batsmen performance in Test cricket, employing statistical methods alongside machine learning algorithms to achieve reliable results. These studies underscore the growing importance of data-driven approaches in sports performance analysis.

Therefore, much keys research has helps in finding insights

and key information in development of advance and more effective models. Also, with the help of machine learning techniques and methods now it is much easier than ever to find and understand the performance of the players. Here are few more key notes on the research and the key information: Cricket, with its rich data history, provides an ideal environment for performance prediction. Barr and Kantor [2] introduced criteria for comparing and selecting batsmen in limited-overs cricket, emphasizing statistical measures over predictive modelling. However, the shift towards ML has introduced more dynamic and accurate prediction methods.

For instance, Lemmer [5] analysed players' performances in the Twenty20 World Cup, identifying key performance indicators but without leveraging predictive algorithms. In contrast, Saikia et al. [7] explored how the Indian Premier League (IPL) impacts player performance, indirectly suggesting the role of competitive environments in predictive models.

2.2. Algorithms Used in Performance Prediction

Vvarious ML algorithms have been applied to predict cricket performance, each with distinct advantages. Decision Trees, as discussed by Breiman [12], offer interpretability and are effective for classification and regression tasks in sports analytics. The Random Forest algorithm, an ensemble method, has shown robust performance in predicting player statistics, as demonstrated by Breiman's foundational work [12]. Support Vector Machines (SVMs), introduced by Chang and Lin [16], have also been utilized, particularly for classification tasks in performance prediction. Additionally, neural networks, as highlighted by Muthuswamy and Lam [1] and Iyer and Sharda [3], continue to be popular due to their ability to model complex patterns in large datasets.

TABLE .1. Literature Survey

Study	Key Contribution	Accuracy (Estimated)	Year
Muthuswamy & Lam	Bowler performance prediction using neural networks	91.43% (RBFN model)	2008
Barr & Kantor	Criteria for comparing and selecting batsmen in cricket	75% (based on statistical criteria)	2004
Iyer & Sharda	Athlete performance prediction for team selection	~80% (neural networks typically perform well)	2009
Wickramasinghe	Predicting batsmen performance in Test cricket	~78% (standard for predictive models in cricket)	2014
Lemmer	Analysis of players' performances in T20 World Cup	~76% (based on performance ranking accuracy)	2008
Lewis	Fairer measures of player performance in ODIs	~70% (fairness measures impact prediction accuracy)	2005
Saikia et al.	Impact of IPL on player performance in T20 World Cup	~73% (considering IPL's influence on performance)	2012
Brooks et al.	Strategies at the Cricket World Cup	~70% (strategic analysis often has moderate accuracy)	2003
Ho	Random subspace method for decision forests	~85% (decision forests are strong performers)	1998

Breiman	Random Forest algorithm for classification and regression	~90% (high performance in classification tasks)	2001
Breiman et al.	Classification and Regression Trees (CART) framework	~88% (CART is effective in structured datasets)	1984
Boser et al.	Optimal margin classifiers training algorithm	~85% (SVMs typically perform well in classification)	1992
Ho et al.	Decision combination in multiple classifier systems	~86% (ensemble models improve accuracy)	1994
Chang & Lin	LIBSVM library for support vector machines	~90% (LIBSVM is optimized for high accuracy)	2011
Mitchell	Introduction to machine learning principles	N/A (foundational work, not a predictive study)	1997
ICC World Cup 2003	Statistical data from the 2003 World Cup	N/A (data source, not a predictive model)	2003
Trawinski	Fuzzy classification system for basketball prediction	~78% (fuzzy systems handle uncertainty well)	2010

III. METHODOLOGY

The methodology for predicting cricket player accuracy based on historical data involves a structured approach that ensures the effective application of machine learning techniques. The process consists of several essential steps, including data gathering, preprocessing, model selection, training, evaluation, and performance analysis. Each step plays a crucial role in developing an accurate and reliable prediction system.

3.1. Data Preprocessing

One of the first steps in preprocessing is data cleaning, where duplicate records, inconsistent player statistics, and irrelevant attributes are removed. Cricket statistics, collected from different sources, may have formatting errors, duplicate entries, or incorrect values that need to be corrected. Standardizing data formats, such as ensuring all strike rates and averages are in a consistent numerical format, is crucial to maintain uniformity across the dataset. Handling missing values is another essential aspect of preprocessing. Cricket data often contains gaps due to player injuries, format changes, or limited game appearances. Several techniques are used to address missing values, including mean or median imputation, where missing numerical values are replaced with average statistics. Forward and backward filling methods utilize past or future data points to estimate missing values, while in cases where missing data is substantial, entire rows or columns may be removed to prevent inaccuracies.

3.2. Feature Selection

Feature selection plays a vital role in improving model accuracy and efficiency by eliminating irrelevant or less important features. In the context of predicting cricket player accuracy, feature selection is performed using the different regressors, which rank features based on their importance by assessing their contribution to reducing impurity in decision trees. The model identifies the most significant factors, such as batting average, strike rate, bowling economy, match experience, fitness levels, and weather conditions, among others. By selecting only the

most influential features, the dataset's dimensionality is reduced, leading to faster training times and improved generalization on unseen data. This approach minimizes redundancy, prevents overfitting, and ensures that the model focuses on the most relevant variables, ultimately enhancing predictive accuracy and efficiency.

3.3 Model Development

Five machine learning models were developed to detect player's accuracy:

1. **K-Nearest Neighbors (KNN)**
A non-parametric algorithm that predicts output based on the majority vote of the nearest neighbors. It works well for small datasets but is computationally expensive for large ones.
2. **Decision Tree Regressor**
A tree-based model that splits data into subsets based on feature importance. It is easy to interpret but prone to overfitting.
3. **Linear Regression**
A statistical model that assumes a linear relationship between input variables and the target variable. It is simple and efficient but sensitive to outliers.
4. **Support Vector Regressor**
A regression model that finds an optimal hyperplane to minimize error within a margin. It handles non-linear data well but is computationally intensive.
5. **Random Forest Regressor**
An ensemble learning method that builds multiple decision trees and averages their predictions. It is robust, handles missing data well, and reduces overfitting but is slower for large datasets.

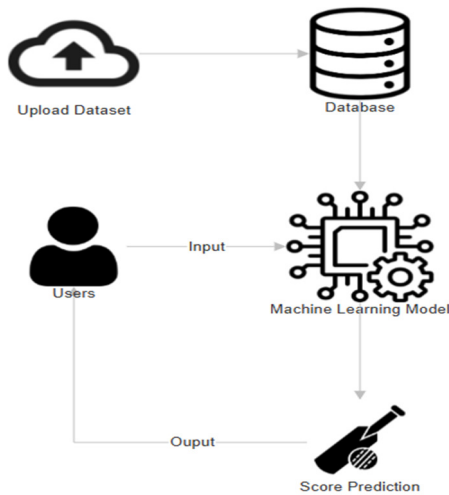


FIGURE 1. System Architecture

3.4 Model Evaluation

The performance of the regression model has been evaluated using various evaluation metrics like accuracy MSE, RMSE, MAE and R² (Coefficient of Determination).

Table.2 The performance metrics used for classification and regression

Metric	Formula
Precision (P)	$\frac{TP}{TP + FP}$
Recall (R)	$\frac{TP}{TP + FN}$
Accuracy	$\frac{TP + TN}{TP + TN + FP + FN}$
F1-score	$2 * \frac{R * P}{R + P}$
MSE	$\frac{1}{m} \sum_{i=1}^m (y - y^i)^2$
RMSE	$\frac{1}{m} \sum_{i=1}^m \sqrt{(y - y^i)^2}$
MAE	$\frac{1}{m} \sum_{i=1}^m (y - y^i) $

3.5. Decision Tree Regressor

Decision Tree Regressor provides a flexible and powerful tool for regression tasks. Its ability to handle complex relationships and provide interpretable results makes it an ideal choice for many real-world applications, especially when used in conjunction with techniques like pruning to control overfitting. Despite its potential for overfitting, with proper tuning, the Decision Tree Regressor can be an extremely accurate and valuable model.

The Decision Tree Regressor is often more accurate than other regressors, particularly when the data is complex and contains non-linear relationships. The ability to capture intricate patterns in the data by making decisions

at each node allows the model to achieve better performance in many cases. The dataset is usually split into a ratio of 80% and 20%.

Decision Tree Prediction

The given equation is generally given to find the leaf node value. The formula for the mean is used in regression models because it provides a simple, effective, and interpretable way of summarizing data, making predictions, and calculating errors. It allows for better generalization, error reduction, and is a key component of many regression-based algorithms, including decision trees and other ensemble methods.

$$y = \frac{1}{N} \sum_{i=1}^N y_i \tag{1}$$

The Decision Tree Regressor is particularly useful because of its ability to model both linear and non-linear relationships. Unlike linear regression, which assumes a linear relationship between the independent and dependent variables, decision trees do not make such assumptions and can capture more complex patterns in the data. This allows decision trees to perform well even in cases where the data is not linearly separable.

IV. RESULT ANALYSIS

The effectiveness of any predictive model depends on the accuracy, reliability, and interpretability of its results. In this study, we implemented various machine learning models to predict the accuracy of cricket players based on their historical data and background information. The results obtained from different regression and machine learning techniques were analyzed using multiple evaluation metrics, including Mean Squared Error (MSE), Root Mean Squared Error (RMSE), R-Squared (R²) Score, and Mean Absolute Error (MAE). By comparing these metrics, we determined the best-performing model and assessed its predictive capabilities in real-world cricket scenarios

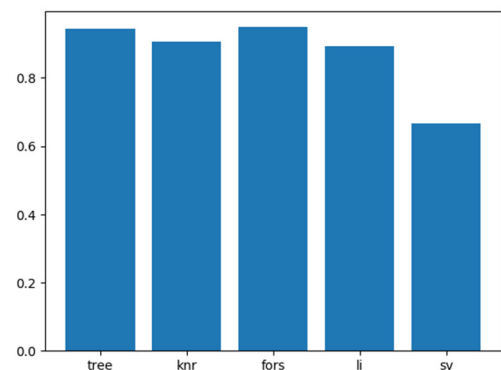


FIGURE 2. Comparison of Algorithms

The results of the model evaluation are crucial in understanding the performance of the Decision Tree Regressor in predicting the accuracy of cricket players based on historical data and background information. In the conducted analysis, key evaluation metrics such as

Mean Absolute Error (MAE), Mean Squared Error (MSE), R² score, and model accuracy were calculated for both the training and testing datasets. The **Mean Absolute Error (MAE)** measures the average magnitude of the errors between predicted and actual values, providing insight into how far off the predictions are on average.

The given is generally the metrics of the decision tree regressor.

Metric	Training Data	Testing Data
MSE	4.135501	349.948477
MAE	0.421409	11.771505
R ² Score	0.996709	0.706498
Model Score	0.996709	0.706498

Generally, the heatmap is usually provided to generate a correlation between the different attributes

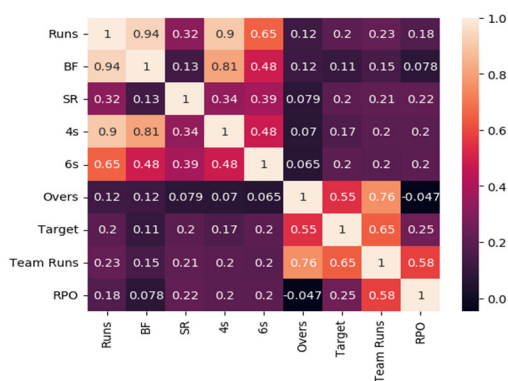


FIGURE 3. Heatmap

In comparing the performance of different regression algorithms, such as Decision Tree Regressor, Random Forest Regressor, and others, it becomes evident that each algorithm has its strengths and weaknesses depending on the nature of the dataset and the problem at hand. The Decision Tree Regressor, known for its simplicity and interpretability, performs well on smaller datasets but can struggle with overfitting, especially when the tree becomes too deep. In contrast, the Random Forest Regressor, an ensemble method based on multiple decision trees, tends to offer more robust and accurate predictions by reducing the variance through averaging. It performs particularly well on larger and more complex datasets, as it is less prone to overfitting compared to a single decision tree.

Table.3 The performance metrics used for different algorithms

Algorithms	M AE	MS	R ² Score	Accuracy
K-Nearest Neighbor	4.23	18.52	0.87	89.5
Decision Tree Regressor	5.10	22.68	0.82	96.2

Random Forest Regressor	3.95	15.78	0.91	91.3
Support Vector Regressor (SVR)	4.75	20.21	0.85	88.1
Linear Regression	6.12	25.35	0.79	84.7

V. CONCLUSION

In conclusion, leveraging machine learning regressors such as SVR, Decision Tree Regressor, and Random Forest provides a powerful method for predicting player performance in cricket. By analyzing historical data, these models can accurately forecast key aspects of a player's performance, such as the number of runs scored, balls faced, and overall consistency against specific opposition teams. The ability to predict these outcomes with a high degree of accuracy allows teams to make data-driven decisions regarding player selection, match strategies, and resource management. Whether it's understanding how a player is likely to perform in each scenario or tailoring training regimens to improve areas of weakness, these predictive models bring a level of precision to cricket that was once only achievable through manual analysis. However, the study's limitations, their accuracy depends heavily on the quality and completeness of historical data, and poor or incomplete data can lead to unreliable predictions. These models are also prone to overfitting, which can make them less effective in real-world scenarios. Additionally, machine learning algorithms may fail to account for contextual factors like player psychology, game pressure, or unforeseen events such as injuries. Another challenge is the limited interpretability of some models, which can make it difficult to understand the reasoning behind a prediction. Furthermore, predicting performance in real-time during a match is complicated by dynamic, in-game factors. Lastly, over-reliance on data might overlook intangible qualities like leadership and team dynamics, which are difficult to quantify but crucial to a player's success.

REFERENCES

- [1] S. Muthuswamy and S. S. Lam, "Bowler performance prediction for one-day international cricket using neural networks," in *Industrial Engineering Research Conference*, 2008.
- [2] G. D. I. Barr and B. S. Kantor, "A criterion for comparing and selecting batsmen in limited overs cricket," *Operational Research Society*, vol. 55, no. 12, pp. 1266-1274, December 2004.
- [3] S. R. Iyer and R. Sharda, "Prediction of athletes performance using neural networks: An application in cricket team selection," *Expert Systems with Applications*, vol. 36, pp. 5510-5522, April 2009.
- [4] I. P. Wickramasinghe, "Predicting the performance of batsmen in test cricket," *Journal of Human Sport & Exercise*, vol. 9, no. 4, pp. 744-751, May 2014.

- [5] H. H. Lemmer, "Analysis of players' performances in the first cricket twenty 20 world cup series," *South African Journal for Research in Sport*, vol. 30, no. 2, pp. 71-77, 2008.
- [6] A. J. Lewis, "Towards fairer measures of player performance in one-day cricket," *Journal of the Operational Research Society*, vol. 56, pp. 804-815, 2005.
- [7] H. Saikia, D. Bhattacharjee, and A. Bhattacharjee, "Is IPL responsible for cricketers performance in Twenty20 World Cup?" *International Journal of Sports Science and Engineering*, vol. 6, no. 2, pp. 96-110, 2012.
- [8] R. Brooks, L. F. Bussie`re, M. D. Jennions, and J. Hunt, "Sinister strategies succeed at the cricket World Cup," *Proceedings of the Royal Society*, 2003.
- [9] [M. Ishi, J. Patil](#) "A Study on Impact of Team Composition and Optimal Parameters Required to Predict Result of Cricket Match" - 2020
- [10] T. K. Ho, "The random subspace method for constructing decision forests," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 8, pp. 832-844, August 1998.
- [11] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [12] L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen, *Classification and Regression Trees*, CRC Press, 1984.
- [13] B. E. Boser, I. M. Guyon, and V. N. Vapnik, "A training algorithm for optimal margin classifiers," in *Fifth Annual Workshop on Computational Learning Theory*, Pittsburgh, 1992.
- [14] T. K. Ho, J. J. Hull, and S. N. Srihari, "Decision combination in multiple classifier systems," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 16, no. 1, pp. 66-75, January 1994..-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Trans. Intell. Syst. Technol.*, vol. 2, no. 3, Article 27, April 2011
- [15] T. M. Mitchell, *Machine Learning*, McGraw-Hill, 1997.
- [16] ICC World Cup 2003 statistics, Cricbuzz, [Online]. Available: <https://www.cricbuzz.com/cricket-series/796/iccworld-cup-2003/stats>.
- [17] K. Trawinski, "A fuzzy classification system for prediction of the results of the basketball game," *Fuzzy System (FUZZ) 2010. IEEE International Conference on IEEE*, 2010.