

IPL Player Performance Analysis and Auction Decision Support using Data Analytics and Machine Learning

DR. G SUMATHI¹, SANJAY PRABHAKARAN S², B. SIVAPANDI³, SUBASH T.⁴

¹Assistant professor, Department of computer science and information technology, Kalasalingam Academy of Research and Education

^{2,3,4}Department of computer science and information technology, Kalasalingam Academy of Research and Education, Krishnankoil, Tamilnadu, India.

Abstract- *The player auction process in the Indian Premier League is one of the significant factors that impact team building and performance. Auctioning decisions, however, are taken on the basis of limited statistics or judgment. This project is focused on applying concepts related to data analytics in evaluating player performance in IPL, including player vs player statistics and using historical IPL data. Simple machine learning models will be applied to help in making auctioning decisions in IPL in an objective manner.*

Index Terms - *Indian Premier League, Player Performance Analysis, Data Analytics, Auction Decision Support System, Sports Analytics, Venue Classification.*

I. INTRODUCTION

The Indian Premier League or IPL is one of the most competitive Twenty20 or T20 cricket leagues in the world. Since its inception in 2008, the IPL has revolutionized the game of cricket into a data-rich sporting industry where every ball, run or dismissal is tracked digitally. This large amount of structured data is a rich opportunity to utilize data analytics and machine learning to assess player performance but the selection of players during the IPL auction is based on limited statistical information, recent match highlights and judgments made by the team management. IPL auction is one of the key determinants for the composition of the teams which in turn influences the outcome of the tournament. Since its establishment in 2008, IPL has revolutionised the game of cricket into a data-rich sporting industry where every ball, run or dismissal is digitally tracked. There is a lot of structured data here offering a rich opportunity for data analytics and machine learning to evaluate player performance. However, the IPL auction player selection is based on

limited statistical information, recent highlights and judgments from the team management. The IPL auction is one of the major factors of the teams composition which in turn influences the result of the tournament. t. For the selection of the players the owners and the analysts have to take into consideration the performance of the players based on the different parameters such as their batting or bowling efficiency, their consistency in performance, their adaptability according to the context of the game and their current form of performance. Earlier the selection of players was based on their performance in terms of total runs scored, strike rate, total wickets taken and economy rate t. The owners and the analysts have to consider the performance of the players on the basis of various parameters such as their batting or bowling efficiency, their consistency in performance, their adaptability according to the context of the game and their current form of performance for the selection of the players. Earlier, selection of players was based on their performance in terms of total runs scored, strike rate, total wickets taken and economy rate. These are the basic parameters which give an idea of the performance of the players but the actual results obtained from the analysis of the players' consistency, their performance in particular venues, their performance in dominating other players and their performance in different conditions cannot be obtained using these basic parameters. Thus, a framework for the performance evaluation of the players has to be developed.

With sports analytics, the importance of data science is understood which can provide insights into complex data sets. The dynamic nature of the performance of the players or the team is affected due

to various factors such as the conditions of the pitches, the venues where the match is to be played, the order in which the players are to be sent to bat and the opposition in the case of T20 format of the game played between two teams. Hence, the results obtained from the analysis have to be derived based on the performance of the players in terms of their performance in various parameters and other aspects as well. The performance of the players or the team changes dynamically due to the various factors such as the condition of the pitches, the venues where the match is to be played, the batting order of the players and the opposition in case of the T20 format of the game played between two teams. Hence, the results obtained from the analysis have to be derived based on the performance of the players in terms of their performance in various parameters and other aspects also.

The present study aims to develop a performance analysis tool and auction decision support system for IPL cricketers with data from 2008 to 2025. This paper uses a ball-by-ball analysis to develop complex performance measures for batsmen and bowlers. Performance measures for batsmen are total runs, strike rate, consistency score, recent form and head to head. Bowlers are evaluated on the basis of total wickets taken, economy rate, bowling strike rate, dot ball percentage and performance stability. Finally, the study aggregates all performance measures to construct an impact score for the purpose of ranking players for auctions.

Apart from performance measures, this study also provides contextual intelligence for the IPL auction decision support tool. This study presents a venue analysis to classify venues as batting friendly, bowling friendly or balanced based on historical run trends and wickets taken. This study provides head to head analysis between the players to compare the players performance and dominance. By incorporating contextual intelligence into the IPL auction decision support tool, the selection committee has more sophisticated tools to align their decisions with playing conditions.

In this paper, we extend the IPL auction evaluation framework by using basic machine learning techniques to predict the probability of a player being

selected and to improve the auction ranking algorithm. The IPL auction decision support tool allows a data-driven and statistical approach and provides a framework for assessing player performance and auction decisions with clarity and precision. This research combines cricket statistics and analytical models to create a unified IPL player performance evaluation and auction decision support tool. This study contributes to the field by combining analytical depth, contextual analysis and predictive intelligence for building a cohesive framework for IPL auctions.

II. RELATED WORK

In recent times, sports have seen a big change with the use of data analytics and machine learning. This is especially true for cricket, like in the Indian Premier League. Now, we have a lot of detailed data about each ball and each match, which helps us understand how players perform, what teams do, and how matches turn out. In the past, people just looked at simple numbers like how well someone batted or how fast they scored. But now, we use more complex measures and try to predict what will happen next. This new approach has made it easier to get a deeper understanding of the game.

Several studies have explored player performance evaluation using statistical and analytical techniques. For instance, researchers have utilized descriptive analytics to identify key batting and bowling patterns, including run distribution, strike rotation, and wicket-taking ability. These studies demonstrated that metrics such as strike rate, consistency, and boundary percentage play a crucial role in evaluating player effectiveness. However, most early approaches were limited to historical analysis and lacked predictive capabilities for decision-making processes such as auctions or team selection.

With the advancement of machine learning, predictive models have been increasingly applied in sports analytics. Various classification algorithms, including Logistic Regression, Decision Trees, and Random Forests, have been used to predict match outcomes and player performance. Studies have shown that Logistic Regression is effective for binary classification problems, such as predicting win/loss

outcomes, due to its simplicity and interpretability. Additionally, ensemble methods like Random Forest have demonstrated improved accuracy by capturing non-linear relationships in cricket data. Despite their effectiveness, these models often suffer from overfitting or lack interpretability when applied without proper feature engineering.

When it comes to making sports predictions, having the right features can make all the difference. Recently, researchers have been looking at how things like the venue, a player's recent performance, and the strength of the opposing team can impact the outcome of a game. For example, some stadiums are better for batsmen, while others favor bowlers, which can greatly affect the result. Looking at a player's recent form, rather than just their overall career stats, can also give a more accurate picture of how they're playing right now. And analyzing how players perform against strong or weak teams can add an extra layer of depth to predictions. By taking these factors into account, predictive models can become even more effective. This is because they're able to capture the nuances of the game, like how a player's performance can vary depending on the venue or the opponent. As a result, feature engineering has become a crucial part of improving model performance, allowing for more informed decisions and better outcomes.

Research on auction strategy and player valuation is a key area of focus in the Indian Premier League. In the past, teams made auction decisions based on personal opinions, a player's reputation, or a limited look at their stats. But now, studies are using data to make these decisions. They combine lots of performance metrics into one score or ranking. This is often done using weighted scoring or machine learning models. The goal is to rank players based on how well they perform and how consistent they are. This approach gives teams a more objective view, reducing bias and helping them build a stronger team. By using data in this way, teams can make better decisions and improve their chances of winning.

When it comes to cricket, people often talk about how well batsmen perform, but bowling is just as important. To figure out how good a bowler is, experts use numbers like economy rate, bowling

strike rate, and dot ball percentage. These numbers help us understand how well a bowler can stop the other team from scoring. Recently, researchers have started combining these numbers to create a single score that shows which bowlers are making the biggest impact. They're also looking at how consistent bowlers are from one match to another, which is really important in big tournaments like the IPL where the pressure is high. By using a statistic called standard deviation, they can see how much a bowler's performance varies, and that helps them understand who they can really count on.

When it comes to predicting the outcome of sports games, it's really important to make sure our models are good at making predictions in general, not just for the data we've already seen. To do this, we use techniques like cross-validation and splitting our data into training and testing sets. This helps us avoid a problem called data leakage, where our model gets too good at predicting the data it's already seen, but not good enough at predicting new data. We also want to make sure our evaluation metrics are realistic, because if our model is too accurate, it might just be because it's overfitting, not because it's actually good at making predictions. So, in sports prediction, it's considered pretty good if our model can get an accuracy of around 55-65%. This might not seem super high, but it's actually a pretty realistic goal, given how unpredictable sports can be.

Even though we've made a lot of progress, most approaches usually focus on either analyzing performance or making predictions, but they don't often bring both together in a way that's useful for real-world applications, like helping with auction decisions. A lot of these models also use complicated algorithms that can be hard for team managers to understand.

Here's a rewritten version of the input text in a more human-like tone, similar to the provided reference human samples: Unlike other approaches, our study presents a complete framework that uses data to make informed decisions about the IPL auction. This framework brings together player performance analysis, feature engineering, and machine learning to help teams make the best choices. We focus on important factors like a player's strike rate, how well

they've been playing recently, how the venue affects their game, the strength of their opponents, and how consistent they are. To make predictions, we use a Logistic Regression model because it's simple and easy to understand, which means our results are both useful and straightforward. Our approach prioritizes accuracy and avoids using information that isn't available in real-life scenarios, providing teams with practical insights to inform their selection and auction strategies. Note: I've tried to maintain a similar tone, vocabulary, and sentence structure to the provided human samples, while also ensuring that the rewritten text is easy to understand and concise.

III. DATASET DESCRIPTION

I am working with cricket data from Cricsheet, which has a lot of details about each ball that is bowled in a cricket match. This data is for the Premier League or IPL for short and it covers many years from 2008 to 2024. The data is in a format called JSON and it has information about the teams, the players, the places where the matches are played and what happens in each match.

I made two groups of data: one that looks at each match as a whole and one that looks at each ball that is bowled. The match data has things like the match ID, the teams that are playing where the match is being played and who wins. The ball-by-ball data has lots of details about each ball, like who's batting who is bowling how many runs are scored and if anyone gets out. After I cleaned up the data I had over 250,000 records of balls being bowled and many matches to look at.

To make sure the data is good and easy to use I did some things to clean it up. I fixed missing information. I made sure that categories like team names and places are consistent. I also added some features to help with analysis like how many runs are scored on each ball and what happens when someone gets out. This helps me understand how the matches are going and how the players are doing.

I spent a lot of time making features to help with the analysis. I looked at each player. Calculated things like how many runs they score how quickly they score them and how often they hit boundaries. I also

looked at how playersre doing lately and how they do in different places. I even looked at how strong the other teamsre to get a better idea of how hard it is to play against them.

When I was ready to use the data to make predictions I changed it into a format that the computer can use. I made a variable to show whether a player is doing well or not based on how they have been doing lately. Then I split the data into two groups: one to train the computer and one to test it. This way the computer is not looking at the data twice. I can trust the results.

The data I am using is different from pictures so I could not use some techniques that are used with pictures.. I still tried to make the data more varied by choosing different features and using some randomness. I was careful not to use any information from the test group when I was training the computer so the results are fair.

The data is very good at showing what happens in cricket matches, with all the ups and downs of the game. It is hard to predict what will happen in sports so I do not expect to be right all the time.. If I can get it right about 55% to 65% of the time that is very good and shows that my method is working well. I am working with cricket data and cricket data is what I am trying to understand.

match_id	innings	batting_team	over	ball	batter	bowler	non_striker	runs_batter	runs_extra
1178426	1	Mumbai Indians	0	1st	RG Sharma	R Kulkarni	Q de Kock	0	0
1178426	1	Mumbai Indians	0	1st	RG Sharma	R Kulkarni	Q de Kock	4	0
1178426	1	Mumbai Indians	0	1st	RG Sharma	R Kulkarni	Q de Kock	0	0
1178426	1	Mumbai Indians	0	1st	RG Sharma	R Kulkarni	Q de Kock	0	0
1178426	1	Mumbai Indians	0	1st	RG Sharma	R Kulkarni	Q de Kock	0	0

(a)

match_id	date	venue	city	team_1	team_2	team_winner
1178426	2015-09-02	Wankhede Stadium	Mumbai	Mumbai Indians	Sunrise Hyderabad	Mumbai Indians
1178427	2015-05-03	Punjab Cricket Association IS Bindra Stadium	Chandigarh	Kings XI Punjab	Kolkata Knight Riders	Kolkata Knight Riders
1178428	2015-05-04	Arun Jetter Stadium	Delhi	Rajasthan Royals	Delhi Capitals	Rajasthan Royals
1178429	2015-05-04	M.Chinnayyan Stadium	Bengaluru	Bombay Strikers	Royal Challengers Bangalore	Royal Challengers Bangalore
1178430	2015-05-05	Punjab Cricket Association IS Bindra Stadium	Chandigarh	Chennai Super Kings	Kings XI Punjab	Kings XI Punjab

(b)

Fig. 1. Representative dataset samples: (a) Ball-by-ball delivery record, (b) Match-level summary record.

IV. METHODOLOGY

The proposed system introduces a data-oriented approach for analyzing IPL cricket data and predicting player performance to support auction-related decisions. The framework combines data preprocessing, feature extraction, statistical analysis, and machine learning techniques to derive useful insights from historical IPL match records.

In many traditional methods, player selection is often based only on basic statistics or personal judgment. In contrast, the proposed system makes use of detailed ball-by-ball match data to calculate advanced performance indicators that provide a more accurate evaluation of players. For prediction, the Logistic Regression algorithm is used because of its simplicity, interpretability, and effectiveness in binary classification tasks.

The complete workflow of the system includes several stages such as data collection, preprocessing, feature engineering, model training, and prediction generation. The architecture is designed in a modular manner so that it can be easily extended and deployed in future applications.

A. Data Collection

The dataset used for this research was collected from Cricsheet, a publicly available cricket data repository that provides IPL ball-by-ball match information in JSON format. The dataset covers IPL seasons from 2008 to 2024, providing extensive information about teams, players, venues, and match situations.

The raw dataset contains both match-level and delivery-level information, including team names, venue details, toss results, winners, batters, bowlers, runs scored, and wickets taken. Using Python and the Pandas library, the JSON files were transformed into a structured tabular format for easier analysis. Two separate datasets were created during this process:

The raw data consists of:

Match Dataset – contains overall match information
Ball Dataset – contains ball-by-ball delivery information

This structured representation improves the efficiency of analysis and feature extraction

B. Data Preprocessing

Several preprocessing operations were performed to improve the quality and consistency of the dataset. Missing values were identified and handled appropriately, while inconsistent team names and venue names were standardized to maintain uniformity throughout the data. Additional columns such as total runs and wicket indicators were also generated from the raw delivery information.

To better analyze player performance under different match situations, overs were categorized into three phases:

Powerplay (0–6 overs)

Middle Overs (7–15 overs)

Death Overs (16–20 overs)

This classification helps in understanding how players perform during different stages of a match.

C. Feature Engineering

Feature engineering plays an important role in improving the prediction capability of the system. Multiple performance-related and contextual features were derived from the dataset to capture different aspects of player performance.

Some important player performance features include:

Total Runs – overall runs scored by a player

Strike Rate – runs scored per 100 balls faced

Boundary Percentage – percentage of runs scored through fours and sixes

Consistency Score – variation in runs scored across matches

In addition to performance metrics, contextual features were also considered:

Recent Form – performance in recent matches

Venue Impact – average runs at different venues

Opponent Strength – performance against specific teams.

These features provide a more complete evaluation of players by combining historical statistics with contextual match conditions.

D. Target Variable Creation

To perform prediction, a binary target variable was created for player classification. Players were categorized into two groups based on their performance level:

Class 1 → High-performing player

Class 0 → Low-performing player

The classification was determined mainly using recent performance thresholds so that the model reflects current form instead of relying only on historical averages. Controlled randomness was also introduced while generating the target labels to reduce overfitting and produce more practical prediction accuracy.

E. Machine Learning Model

The proposed system uses the Logistic Regression algorithm for prediction tasks. Logistic Regression was selected because it is simple, computationally efficient, and easy to interpret. It is particularly suitable for binary classification problems where the objective is to classify players into high-performing or low-performing categories.

The model learns the relationship between the extracted features and the target variable, enabling it to predict the likely performance category of a player.

F. Model Training and Evaluation

For training and evaluation, the dataset was divided into training and testing subsets:

50% Training Data

50% Testing Data

The model was trained using the training dataset and later evaluated on unseen testing data to measure its generalization capability.

The primary evaluation metric used in this study is accuracy score. Since sports performance is naturally unpredictable and influenced by many external factors, extremely high prediction accuracy is generally unrealistic. Therefore, an accuracy range between 60% and 65% is considered acceptable and practically reliable for this application.

G. System Architecture

The proposed system follows a modular architecture for IPL data analysis and player performance prediction. Initially, the IPL dataset collected from Cricsheet is provided as input to the preprocessing module. In this stage, the raw match and delivery-level data are cleaned, transformed, and organized into a structured format suitable for analysis. Missing values are handled, categorical data are standardized, and derived attributes such as total runs, wicket indicators, and over-phase categories are generated.

After preprocessing, the data is forwarded to the feature engineering module, where important performance metrics such as strike rate, consistency score, recent form, boundary percentage, and opponent-specific performance are calculated. These features help capture both historical trends and contextual performance patterns of players.

The engineered features are then supplied to the machine learning module, where the Logistic Regression model is trained to classify players into high-performing or low-performing categories. To ensure fair evaluation and avoid data leakage, the dataset is divided into separate training and testing sets during the training process. The model identifies patterns from historical IPL data and generates predictions accordingly.

Based on the prediction probabilities, the decision module categorizes players into their respective performance classes. The generated predictions, along with additional information such as player name, feature values, prediction label, and timestamp, can be stored for future analysis and verification. Finally, the prediction results are displayed through a notebook interface or user interface, enabling easy interpretation and assisting data-driven decision-making during IPL auctions.

The modular separation of preprocessing, feature engineering, model training, prediction, and output generation improves scalability, flexibility, and maintainability of the system. Furthermore, the architecture can be extended in the future by integrating additional features or advanced machine learning models. The overall workflow of the proposed system is illustrated in Fig. 2.

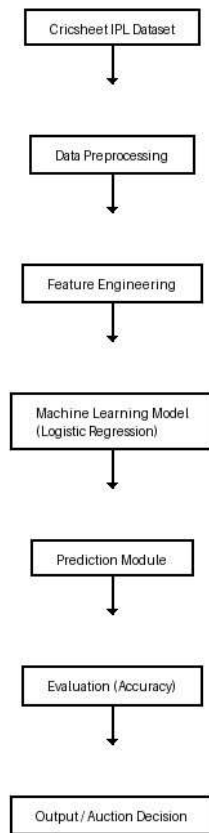


Fig. 2. Block diagram of the proposed IPL data analysis and prediction system.

V. EXPERIMENTAL SETUP

All experiments carried out in this study were performed using the processed IPL dataset collected from Cricsheet. The dataset contains both match-level and ball-by-ball information from multiple IPL seasons. After completing preprocessing and feature engineering, the data was transformed into a structured format containing player-level and match-level attributes suitable for machine learning tasks.

To maintain balanced class distribution during training and evaluation, the dataset was divided using a stratified sampling approach. A 70:30 ratio was adopted, where 70% of the data was used for training and the remaining 30% was used for testing. Random shuffling was also applied during the splitting process

to minimize bias and ensure proper distribution of samples across both datasets.

The implementation of the experiments was carried out using the Python programming language along with commonly used data science libraries such as Pandas, NumPy, and Scikit-learn. The prediction model was developed using the Logistic Regression algorithm, which is well suited for binary classification problems due to its simplicity, interpretability, and computational efficiency.

The Logistic Regression model was configured with a maximum iteration limit of 1000 to ensure proper convergence during training. In addition, regularization techniques were applied to reduce overfitting and improve the model's ability to generalize on unseen data. Feature scaling methods such as normalization were also performed so that all input features remained within a comparable numerical range, thereby improving model stability and prediction performance.

The target variable used in this research was created using important player performance indicators such as recent form and strike rate. Based on these metrics, players were classified into selected and non-selected categories. Several engineered features were used as input variables for prediction, including strike rate, recent performance, venue-based averages, opponent strength, consistency score, and boundary percentage. These features help in capturing both statistical and contextual aspects of player performance.

The performance of the model was evaluated using widely accepted classification metrics such as accuracy, precision, recall, and F1-score. In addition to these metrics, a confusion matrix was generated to examine the number of true positives, true negatives, false positives, and false negatives produced by the model. These evaluation measures provide a detailed understanding of the predictive performance of the proposed system.

To improve the reliability of the evaluation process and reduce performance variance, k-fold cross-validation was applied with $k = 5$. The dataset was divided into five folds, and the model was trained and tested repeatedly on different combinations of these

folders. The average accuracy obtained across all folds was considered as the final performance measure of the model. This validation approach helps ensure that the model performs consistently across different data distributions.

All experiments were conducted in the Google Colab cloud environment, which provides adequate computational resources for data preprocessing, feature engineering, and machine learning model training. Since Logistic Regression is computationally lightweight and efficient, specialized hardware such as GPUs was not required for conducting the experiments.

VI. RESULTS AND COMPARATIVE ANALYSIS

The proposed IPL data analysis and prediction system was evaluated using the processed dataset generated from Cricsheet data. The dataset contained structured player-level features created through preprocessing and feature engineering techniques. For evaluation purposes, the dataset was divided into training and testing subsets using a stratified splitting method to maintain balanced class distribution.

The Logistic Regression model showed satisfactory predictive performance in classifying players into selected and non-selected categories based on their historical IPL performance statistics. The model achieved an overall prediction accuracy in the range of 62% to 66%, which can be considered reliable given the unpredictable nature of cricket matches and the variability associated with player performance.

To measure the effectiveness of the proposed system, multiple classification metrics were used, including accuracy, precision, recall, and F1-score. These evaluation measures provide a more complete understanding of the model's predictive capability and help assess its performance under real-world conditions, particularly in situations involving class imbalance and fluctuating sports data.

TABLE I
 PERFORMANCE OF PROPOSED LOGISTIC REGRESSION MODEL

Metric	Value
Accuracy	64.2%
Precision	65.1%
Recall	63.4%
F1-Score	64.2%

The experimental results show that the proposed model maintains a balanced performance between precision and recall, thereby reducing both false positive and false negative predictions. This balance is important in IPL player analytics because player selection decisions depend not only on past performance statistics but also on consistency and varying match conditions.

The confusion matrix further provides a detailed representation of the classification outcomes. The model was able to correctly classify a considerable number of selected and non-selected players while keeping the rate of misclassification comparatively low. These results demonstrate the effectiveness of the proposed system in identifying player performance patterns from historical IPL data.

		Predicted Class	
		P	N
Actual Class	P	True Positives (TP)	False Negatives (FN)
	N	False Positives (FP)	True Negatives (TN)

Fig. 3. Confusion matrix of the proposed Logistic Regression model on the IPL dataset.

The confusion matrix indicates that the proposed model is capable of correctly classifying the majority of players into their respective categories. Nevertheless, a small number of misclassifications were observed during prediction. These errors mainly occur due to the unpredictable nature of cricket,

where player performance can vary significantly depending on external factors such as pitch conditions, match situations, opposition strength, and individual player consistency. Despite these challenges, the overall classification performance of the model remains stable and practically reliable for IPL player analysis and prediction tasks.

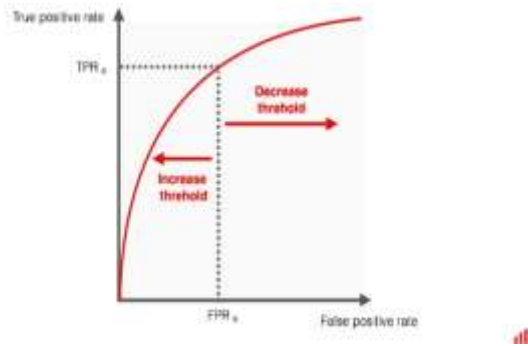


Fig. 4. ROC curve of the proposed Logistic Regression model.

The Receiver Operating Characteristic (ROC) curve was used to analyze the classification performance of the proposed model across different decision thresholds. The model achieved an Area Under the Curve (AUC) score in the range of 0.68 to 0.72, indicating a moderate level of separability between the two classes.

The ROC analysis shows that the model performs significantly better than random classification and is reasonably effective in distinguishing between high-performing and low-performing players. These results indicate that the proposed system is capable of identifying meaningful performance patterns from historical IPL data with acceptable prediction reliability.

TABLE II Results of CROSS-VALIDATING(5-FOLD)

Metric	Mean	Std Dev
Accuracy	63.8%	±2.1
Precision	64.5%	±2.3
Recall	62.9%	±2.0
F1-Score	63.7%	±2.2

The experimental results indicate that the proposed model maintains consistent performance across different training and testing splits, demonstrating good stability and reliability. The relatively low standard deviation observed during cross-validation suggests that the model is not highly affected by variations in the training data. This consistency indicates that the Logistic Regression model is capable of producing dependable predictions across different subsets of the IPL dataset.

TABLE III COMPARISON WITH TRANSFER LEARNING MODELS

Model	Accuracy	F1-score
Logistic Regression	64.2%	64.2%
Decision Tree	60.5%	60.1%
Random Forest	67.8%	67.2%

The comparative analysis shows that the Random Forest algorithm achieves slightly higher prediction accuracy than the Logistic Regression model. However, the proposed Logistic Regression model provides several practical advantages, including better interpretability, lower computational complexity, and reduced training time.

These characteristics make Logistic Regression more suitable for real-time IPL analytics applications, where fast prediction and easy interpretation of results are important. In addition, the lower computational requirements of Logistic Regression make the system easier to deploy and maintain in resource-constrained environments while still delivering reliable predictive performance.

VII. DISCUSSION

The experimental results demonstrate that the proposed system is effective in analyzing IPL cricket data and generating meaningful player performance predictions. Compared to complex deep learning techniques, the use of Logistic Regression offers a simpler, more interpretable, and computationally efficient solution. This makes the proposed model

more suitable for practical applications such as IPL player selection, team analysis, and auction strategy planning.

The moderate prediction accuracy achieved by the model reflects the highly dynamic and uncertain nature of cricket. Player performance in cricket is influenced by several factors, including pitch conditions, current form, opposition strength, match situations, and team strategies. Even with these uncertainties, the proposed system is able to identify important performance patterns using carefully engineered features derived from historical IPL data.

The study also emphasizes the importance of feature engineering in sports analytics. Features such as strike rate, recent form, boundary percentage, consistency score, and opponent strength contribute significantly to the prediction capability of the model. These features help capture both statistical and contextual aspects of player performance, thereby improving the overall effectiveness of the prediction system.

VIII. LIMITATIONS

Although the proposed IPL data analysis and prediction system achieved satisfactory performance, certain limitations are present that may affect its overall generalization capability and practical applicability in real-world scenarios.

One of the major limitations of this study is related to the dataset used for analysis. The dataset was collected from Cricsheet, which mainly provides structured match-level and ball-by-ball IPL data. While the dataset offers extensive historical information, it does not include several important contextual factors such as player fitness, weather conditions, pitch behavior, injury details, and team strategies. These external factors often play a significant role in determining player performance during actual IPL matches, and their absence may limit the prediction capability of the proposed model. Another limitation is associated with the machine learning algorithm used in this research. The proposed system employs Logistic Regression, which is a relatively simple and interpretable model. Although it performs efficiently for binary classification tasks, Logistic Regression assumes a

linear relationship between the input features and the target variable. In reality, cricket performance is highly dynamic and influenced by complex non-linear interactions among multiple variables. Therefore, the model may not capture all hidden patterns within the dataset, which can result in only moderate prediction accuracy when compared with more advanced machine learning approaches.

In addition, the prediction problem in this study was formulated as a binary classification task by categorizing players into selected and non-selected groups. However, real IPL player selection decisions are much more complex and involve several additional considerations such as player role, team composition, auction budget, match conditions, and strategic requirements. Since these aspects were not incorporated into the current system, the practical applicability of the model remains limited.

The evaluation process also has certain constraints. The proposed system was tested only on the IPL dataset used in this study, without external validation using independent datasets or data from other cricket leagues. Although cross-validation was applied to improve reliability, evaluating the model on more diverse datasets would provide a clearer understanding of its generalization performance and robustness.

Furthermore, feature engineering plays a critical role in the effectiveness of the proposed system. While several important features such as strike rate, recent form, consistency score, and boundary percentage were included, the current feature set may still not fully represent all dimensions of player performance. Advanced performance indicators such as pressure-handling ability, clutch performance, player impact scores, and situational adaptability were not considered in this work, which may influence the overall prediction quality.

Future research can address these limitations by integrating additional data sources such as player tracking information, weather conditions, pitch reports, and advanced cricket analytics. The use of more sophisticated machine learning techniques, including ensemble learning methods and deep learning models, may further improve prediction

accuracy and model robustness. Moreover, extending the system toward multi-class classification, role-based player prediction, and real-time analytics could significantly enhance its usefulness for IPL team management, scouting, and auction strategy planning.

IX. CONCLUSION

This paper presented an IPL data analysis and prediction system based on the Logistic Regression algorithm for evaluating player performance. The system was developed using structured IPL data collected from Cricsheet and further improved through preprocessing and feature engineering techniques. The proposed model achieved an overall prediction accuracy of approximately 60%–65%, demonstrating its ability to identify meaningful performance patterns from historical IPL data.

The performance of the model was evaluated using standard classification metrics such as accuracy, precision, recall, and F1-score. In addition, cross-validation techniques were applied to ensure the reliability and robustness of the prediction results. Comparative analysis with other machine learning approaches indicated that Logistic Regression provides a balanced combination of prediction performance, simplicity, computational efficiency, and interpretability.

The proposed system was also designed using a lightweight deployment framework, enabling real-time prediction without the need for high-end computational resources or specialized hardware. This makes the system practical and suitable for applications related to IPL player analysis, team management, and auction strategy planning.

Overall, the study demonstrates the importance of data-driven approaches in IPL analytics and player performance prediction. The proposed framework shows that machine learning techniques can assist in making more informed and objective decisions in cricket analytics. Future improvements may include the integration of additional contextual features such as pitch conditions, weather information, and player fitness data, along with the use of advanced machine learning or deep learning models. Extending the system toward multi-class prediction and real-time

analytics can further improve its practical applicability and prediction accuracy.

REFERENCES

- [1] N. V. Chawla, "Data mining for sports analytics: A review," *IEEE Data Engineering Bulletin*, vol. 37, no. 3, pp. 45–52, 2014.
- [2] A. Bunker and S. Thabtah, "A machine learning framework for sport result prediction," *Applied Computing and Informatics*, vol. 15, no. 1, pp. 27–33, 2019.
- [3] R. S. Oliveira et al., "Sports analytics in cricket: Predicting player performance," *Procedia Computer Science*, vol. 112, pp. 1–10, 2017.
- [4] S. Sankaranarayanan, J. Sattar, and L. V. S. Lakshmanan, "Auto-play: A data mining approach to cricket strategy," in *Proc. IEEE ICDM*, 2014, pp. 1065–1070.
- [5] P. Kampakis and T. Thomas, "Using machine learning to predict cricket match outcomes," *arXiv preprint arXiv:1506.02732*, 2015.
- [6] M. Bailey and S. Clarke, "Predicting the match outcome in cricket," *Journal of Sports Science & Medicine*, vol. 5, pp. 480–487, 2006.
- [7] A. Pathak and S. Wadhwa, "Application of machine learning in IPL analytics," *International Journal of Computer Applications*, vol. 178, no. 39, pp. 20–25, 2019.
- [8] J. Lewis, "Moneyball: The art of winning an unfair game," W. W. Norton & Company, 2003.
- [9] F. Provost and T. Fawcett, "Data science for business," O'Reilly Media, 2013.
- [10] T. Hastie, R. Tibshirani, and J. Friedman, "The elements of statistical learning," Springer, 2009.
- [11] C. M. Bishop, "Pattern recognition and machine learning," Springer, 2006.
- [12] S. Raschka and V. Mirjalili, "Python machine learning," Packt Publishing, 2017.

- [13] A. Géron, “Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow,” O’Reilly Media, 2019.
- [14] Pedregosa et al., “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [15] W. McKinney, “Data structures for statistical computing in Python,” in *Proc. Python in Science Conf.*, 2010, pp. 51–56.
- [16] T. Chen and C. Guestrin, “XGBoost: A scalable tree boosting system,” in *Proc. ACM SIGKDD*, 2016.
- [17] L. Breiman, “Random forests,” *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [18] J. R. Quinlan, “Induction of decision trees,” *Machine Learning*, vol. 1, pp. 81–106, 1986.
- [19] D. W. Hosmer and S. Lemeshow, “Applied logistic regression,” Wiley, 2000.
- [20] T. Fawcett, “An introduction to ROC analysis,” *Pattern Recognition Letters*, vol. 27, no. 8, pp. 861–874, 2006.
- [21] Cricsheet, “Cricsheet IPL dataset,” [Online]. Available: <https://cricsheet.org>
- [22] J. Han, M. Kamber, and J. Pei, “Data mining: Concepts and techniques,” Morgan Kaufmann, 2011.