

Data-Driven Cricket: A Machine Learning Approach to IPL Score Prognostication

Anurag Das, Prasanth Ayitapu, Rohit S Nair, Satvik Raghav, S. Lalitha

Department of Electronics and Communication Engineering

Amrita School of Engineering, Bengaluru

Amrita Vishwa Vidyapeetham, India

{anuragdas070, payitapu, rohitsnair2004, satvikraghav007}@gmail.com

*Corresponding Author: s_lalitha@blr.amrita.edu

Abstract—The Indian Premier League (IPL) is a very popular and globally recognized cricket tournament, attracting many spectators worldwide. The proposed work aims to predict cumulative scores of teams in IPL using attributes such as batting and bowling team data, total runs, wickets, overs, runs scored in the last 5 overs, batting taken in the last 5 overs, and Total Score. This work employs Machine Learning (ML) techniques such as Linear Regression (LR), Random Forest (RF), and Support Vector Machine (SVM) models to predict IPL scores. The database considered in this work was carefully obtained from well-known websites such as ESPN and Cricbuzz, ensuring the reliability and credibility of the research. The evaluation of the models reveals good performance, with a mean absolute error (MAE) of 12.12, a root mean square error (RMSE) of 15.84, and an R-squared value (R^2) of 0.75.

In addition to the prediction models, a user-friendly graphical user interface (GUI) is implemented for easy interaction with the prediction system. This GUI allows the user to enter the necessary match data and displays the predicted score. Uniquely, this system offers a range of points that can be extended up to 5 runs from the estimated 5-run below. This range recognizes the unpredictable nature of cricket matches and provides more realistic and practical predictions. The combination of powerful ML models and an intuitive GUI makes this system a powerful tool for IPL score prediction.

Index Terms—Linear Regression, Random Forest, Support Vector Machine, Pandas, One-hot Encoding

I. INTRODUCTION

Cricket, as a sport, has achieved unprecedented global popularity, with leagues such as the Indian Premier League (IPL) emerging as a cultural phenomenon that transcends geographical boundaries. The IPL, known for its exciting matches and star-studded lineups, has become a cornerstone of the world of cricket, captivating audiences across the globe [1].

The study is rooted in the dynamic field of the IPL, focusing on predictive modeling of cumulative team scores – an important aspect of the league that contributes to its intense and unpredictable nature. Leveraging unique feature vectors like batting and bowling team data, total runs, wickets, overs, runs scored in the last 5 overs, wickets taken in the last 5 overs, and total score, a robust system for IPL score prediction will be developed.

The methodology seamlessly integrates RF and SVM models, underscoring a commitment to meticulous analysis and

predictive precision. Drawing upon inspiration from diverse applications like Emotion Variation Detection [2], Mixed-lingual Affective State Recognition System [3], and Analysis of Heart Disease Prediction [4], these models have established themselves as robust frameworks in the field. Moreover, their utilization in IPL score prediction further emphasizes their efficacy and reliability. The dataset employed in this study is carefully curated from esteemed platforms like ESPN and Cricbuzz, ensuring the trustworthiness and credibility of the data utilized in the predictive model.

The primary motivation for conducting this research stems from the limited existing literature in the area of IPL score prediction. Recognizing the dearth of research in this domain, this work aims to fill this gap by providing valuable insights and contributing to the evolving landscape of game analytics. This study not only increases the understanding of IPL matches but also shares its findings and methodology with the wider academic and cricketing communities.

II. LITERATURE SURVEY

In the proposed work of IPL score prediction, this literature review evaluates seminal studies using different ML techniques in cricket analytics. Through this brief survey, the goal is to explore concepts, identify effective methodologies, and highlight areas that require further research while contributing to the growing discourse in sports analytics.

Ishan Jain et al. [5] suggest ML algorithms like RF Classifier and Multivariate Polynomial Regression (MPR) to predict the final score and winner of the match. The accuracy of the prediction model was reported to be 67.3% for MPR and 55% for RF classification. A distinctive feature of this work is the integration of a data mining module with a corresponding visualization that improves the predictive and analytical capabilities of the system.

Indika Wickramasinghe et al. [6] pointed out that various ML techniques have been used in cricket. These include SVM, K nearest neighbors (KNN), Regression, RF, and XGBoost. The accuracy of this method was determined using the Fuzzy matrix-based method, F-score, RMSE, ROC, Cohen's kappa statistics, MAE, MCC, and R2. Feature selection techniques such as correlation-based feature selection, iterative elimination, and Chi-Square are commonly used. These findings

distinguish this work from others by providing a comprehensive review of ML methods, accuracy, and feature selection methods used in cricket analysis.

Manoj Ishi et al. [7] recommended the use of ML techniques such as LR, NB Bayes, KNN, SVM, decision tree, RF, GBM, XGBoost, and CatBoost. The accuracy achieved varies, with SVM reaching 93.54% and the highest when combined with the LR feature optimization reaching 94.28%. A distinctive feature of this work is the use of nature-inspired algorithms for feature selection leading to improved accuracy compared to other works.

Yogesh Kumar et al. [8] stated that ML techniques such as Multinomial NB Bayes, LR, Ridge Classifier, RF, SGD Classifier, and Decision Tree were used to predict sentiment in IPL using Twitter data. The accuracy achieved by this method, Ridge Classifier reached the highest accuracy of 90.27%. A distinctive feature of this work is the use of a large dataset of 7 million tweets and the use of an optimization methodology to improve classification accuracy.

Parmeet Kaur et al. [9] proposed ML methods, including the KNN algorithm, to predict the outcome of IPL cricket matches. The accuracy of the KNN algorithm with $k=4$ was observed to be around 71%. The distinguishing feature of this work lies in its dynamic approach to the prediction of match results and the use of the non-relational HBase database for scalability.

Saranya G et al. [10] utilized ML techniques to predict the performance of players and select the top 15 players for IPL. The work achieved an average accuracy rate of 94% for analyzing IPL match results. The distinctive features of this work include the use of data cleaning techniques, the study of batting and bowling performances, and special data extraction for analysis.

Muhammad Sajjad et al. [11] presented a model that employs CNN-LSTM architecture for cricket activity recognition. It achieves an F1 score of 91%, a recall of 91%, and a precision of 92%, outperforming several recurrent neural network variations. Multiple convolutional layers with 64 feature maps, rectified linear unit (ReLU) activation functions, and 2D max-pooling with a 4x4 filter size and 2x2 step size are some of the distinguishing characteristics. The model's accuracy beats that of current approaches, indicating that it is a potential tool for recognizing athletic activities.

Amitabha Chakrabarty et al. [12] centered the use of SVM with a linear kernel. For the linear kernel SVM, the model's accuracy is reported to be 92%. Important characteristics are chosen using feature selection algorithms like similarity selection and duplicate feature reduction.

E. L. Lekamge et al. [13] proposed a prediction model of cricket player performance using LR, SVM with linear and polynomial kernel. The accuracy of this model was 91.5% for Tamim, the batter, and 75.3% for Mahmudullah, the bowler. Relevant attributes are extracted using feature selection methods like similarity selection and duplicate feature deletion.

Afsana Khan et al. [14] utilized ML techniques such as R-CNN and YOLO v3. The proposed model achieved an

accuracy of 84.71% in the classification of cricket shots. A distinctive feature of this work is the use of instant segmentation and the creation of special databases. No external sources were used in this literature review.

From the literature review performed, it is observed that in cricket analytics, a spectrum of ML algorithms has been used to predict scores, player performances, and match outcomes. Some models include LR, Naive Bayes, KNN, SVM with linear and polynomial kernels, Decision Tree, RF, and Gradient Boosting Machine, along with other advanced techniques like Ridge Classifier, SGD Classifier, Multinomial NB Bayes, MPR, GNLM, K-Means Clustering, CNN, and MLP Neural Network. These algorithms are suitable for a wide range of analytical problems, using optimization methodologies, extensive databases, and feature selection techniques to improve accuracy in cricket analytics. The models of KNN, CNN, SVM- that are being applied for IPL prediction are also found to play a major role in the detection of Parkinson Speech [15], Multilingual and Mixed-lingual Digit Speech Recognition System [16], and Prediction of Coronary Artery Disease [17].

The current literature on cricket analysis reveals several research gaps. First, there are inconsistencies in the reported accuracy of ML models used for tasks such as score prediction and match outcome prediction, suggesting the need for more reliable assessment methods. The prediction models lack real-time data integration, ignoring the potential benefits of timely information during live matches. There is not enough emphasis is placed on the interpretability and explanatory power of the model, which hinders the understanding and acceptance of the prediction model by cricket stakeholders. The current research lacks the integration of a user-friendly GUI website to improve the accessibility and usability of prediction models. Closing these gap can lead to deeper and more reliable insights into cricket analysis, thereby increasing the effectiveness and application of predictive models in real-world scenarios.

III. METHODOLOGY

Fig. 1 shows the detailed workflow of this proposed architecture which is explained below.

A. Dataset

The dataset used in this study has been carefully collected from reputable sources including ESPN [18] and Cricbuzz [19], known for their comprehensive cricket statistics. The dataset consists of a significant size of 76,015 samples and consists of 15 attributes which include Mid, Date, Venue, Bat_Team, Bowl_Team, Batsman, Bowler, Runs, Wickets, Overs, Run_Last_5, Wickets_Last_5, Striker, Non-Striker, and Total. The dataset captures various match-related parameters and player performance metrics. Spanning the years from 2008 to 2017, this comprehensive dataset serves as a rich repository of information for the analysis. Specifically, data up to the year 2016 is designated as the training set, while data after 2016 constitutes the testing set. This strategic division allows for a robust investigation of the temporal evolution of the dataset

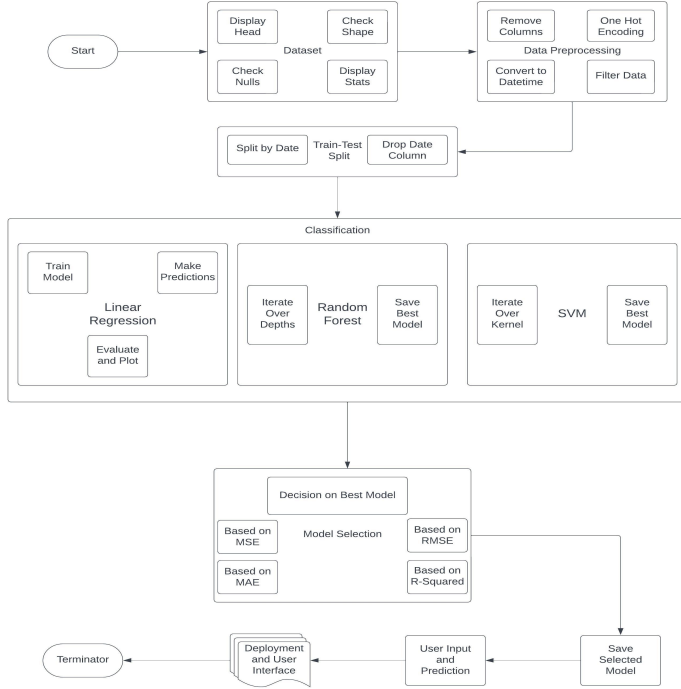


Fig. 1. Architecture of the Proposed Model

and enables the development and evaluation of predictive models over different time frames.

B. Data Pre-processing

In the data processing stage, the database is expanded to filter and organize the data for optimal analysis. Unnecessary features have been removed, and focus on important features like Bat_Team, Bowl_Team, Runs, Wickets, Overs, Runs_Last_5, Wickets_Last_5, and Total as it is represented in TABLE I. This selective approach ensures that the database captures the key elements for the predictive model in this research work. To improve the categorical variables, the process of hot coding [20] is used, changing the categorical characteristics into a suitable format for quantitative analysis. One-hot Encoding is used to convert categorical features ('bat team' and 'bowl team') into binary columns, improving the model's understanding of team data. This process results in a 'cat df' dataframe containing one-hot-encoded columns alongside essential features like 'runs', 'wickets', and 'overs', preparing the data for ML model training. Additionally, to maintain consistency and relevance in the analysis, the database was filtered based on consistent teams across the league. The successive teams selected are Kolkata Knight Riders, Chennai Super Kings, Rajasthan Royals, Mumbai Indians, Kings XI Punjab, Royal Challenger Bangalore, Delhi Daredevils, and Sunrisers Hyderabad. This refined database, now stripped of redundancies and enriched with relevant features, forms the basis for further analyzes and predictive models as shown in Table I.

TABLE I
SELECTION OF FEATURES BEFORE AND AFTER
PRE-PROCESSING

Features In The Dataset	Features After Selection
Mid	Bat_Team
Date	Bowl_Team
Venue	Runs
Bat_Team	Wickets
Bowl_Team	Overs
Batsman	Runs_Last_5
Bowler	Wickets_Last_5
Runs	Total
Wickets	
Overs	
Run_Last_5	
Wickets_Last_5	
Striker	
Non-Striker	
Total	

C. Train-Test Split

This dataset is strategically divided into training and test sets based on temporal criteria. Cases up to and including 2016 were selected for the training set, and from 2017 onwards for the test set. This chronological separation allows efficient evaluation of forecast models in different time periods. As a result, after the splitting is completed, the attribute "date" column is redundant for the modeling process and is deleted as a result. This approach guarantees that the temporal evolution of the database during training and testing is considered correctly and contributes to the reliability of the analysis.

D. Classification

In the classification phase of the analysis, three different ML models are used - LR, RF, and SVM.

For LR, this process includes building a model in the training set, making predictions using the built model, and then evaluating the results. The results are carefully evaluated and illustrated using several plots that provide insight into the model's performance and predictive ability.

The RF model undergoes a more complex iterative process, exploring a range of tree depths to determine the optimal configuration. Performance metrics such as RMSE, MSE, MAE, and R^2 are used to determine the effectiveness of the model at each depth. The best model, as determined by these criteria, is selected and retained for further analysis and comparison.

SVM models systematically explore different types of kernels, including linear, radial basis function (RBF), polynomial, and sigmoid kernels. By continuously evaluating the performance of the model with each type of kernel, it is identified, and the most effective kernel is identified based on performance metrics such as RMSE, MSE, MAE, and R^2 . This sensible approach ensures that the chosen kernel type improves the prediction accuracy and forms the basis for a detailed evaluation of the SVM's classification capabilities.

E. Model Selection And Saving Best Model

In the critical stage of model selection, it is systematically evaluated and compared the performance of three ML models—LR, RF, and SVM. The decision-making process involves a thorough evaluation based on key performance metrics such as MSE, MAE, RMSE, and R^2 . By examining the results of the model against these criteria, the most effective and reliable model is determined. The selected model, which is deemed superior in prediction accuracy and overall performance, is then retained. This systematic approach ensures that the chosen model, whether it is LR, RF, or SVM, represents the optimal solution to the IPL score prediction problem.

F. Real-Time Prediction

Designing a model to receive input data for the last 5 overs in real time in a live IPL match is a complex task but important. The model should contain some dynamic factors like batting team, bowling team, current score, the highest score, the score in the last 5 overs, and the number of wickets in the last 5 overs. These real-time inputs are essential for accurate predictions as the game progresses. It is necessary to implement a continuous model update mechanism with the latest data to reflect the importance of model predictions and game-changing dynamics. This mechanism must continuously integrate incoming real-time data, allowing the model to adapt and refine its predictions based on the latest game conditions. By constantly updating the model with the latest data, it can maintain accuracy and provide valuable insight into the expected score based on the team's latest performance.

IV. EVALUATION METRICS

The evaluation prominently used for the cricket score prediction are MAE, MSE, RMSE, R^2 which are explained as follows [21]:

A. Mean Absolute Error (MAE)

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

B. Mean Squared Error (MSE)

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

C. Root Mean Squared Error (RMSE)

$$\text{RMSE} = \sqrt{\text{MSE}} \quad (3)$$

D. R-squared (R^2)

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (4)$$

V. IMPLEMENTATION AND ANALYSIS

The evaluation in this work is simulated using Python programming. Initially, the feature vector consists of 15 attributes, which are reduced to 8 after pre-processing. The pre-processed data includes important elements like Bat_Team, Bowl_Team, Runs, Wickets, Overs, Runs_Last_5, Wickets_Last_5, and Total. The work is carried out in two phases:

Phase 1 deals with model implementation and analysis. The three models—LR, RF, and SVM—are applied to predict IPL scores based on selected attributes after pre-processing. The scatter plots and line plots are also being inclusive to have better analysis of the models. The models are compared with respect to year-wise score average and their performance metrics.

Phase 2 involves Graphical User Interface (GUI) Integration to improve user interaction. The GUI provides a user-friendly interface for visualizing and interpreting predictions with LR, RF, and SVM models for better exploration and understanding of the analysis results through the GUI, facilitating a seamless and interactive experience. GUI integration increases the practical utility of the model's predictions, making the results more accessible and convenient for stakeholders.

PHASE 1:

A. Proposed IPL score prediction framework using LR

For the feature vectors of the IPL, is given as input to the LR model and the predicted values. The MAE, a measure of the average absolute difference between the predicted value and the actual value, was calculated as 12.12. This means, on average, the model predictions deviate by approximately 12.12 units from the actual values. The MSE, which measures the average squared difference between the predicted value and the actual value, gave a value of 251.01. The RMSE, calculated as the square root of the MSE, was found to be 15.84, indicating the average size of the forecast error. The R^2 value, a statistical measure of model fit, was determined to be 0.75. This R^2 value shows that approximately 75.23% of the variation in the dependent variable (IPL score, in this case) can be explained by a LR model. To have a further in-depth analysis, the line plot and scatter plot are shown in Fig. 2 and Fig. 3.

TABLE II
PREDICTED EVALUATION METRIC VALUES USING LR MODEL

Evaluation Metrics	Values
MAE	12.118617546193295
MSE	251.00792310417438
RMSE	15.843229566732106
R^2	0.7522633566350527

From Fig. 2, it is observed that the range of actual values in the x-axis tested for the model is between 60 and 220, and the range of predicted values for the y-axis is between 75 and 250. This plot shows that the range of predicted values that are generally quite small, which indicates consistent predictions. In particular, the range of predicted values is smaller within the

range of runs from 135 to 200, a common scoring range in IPL matches. Within this range, where most groups usually enter, the model shows heightened accuracy and minimal variation in predictions. In particular, outliers, which are still captured by the model, can show a rather larger range in the predicted score, indicating that the model adapts to changes in the scoring pattern.

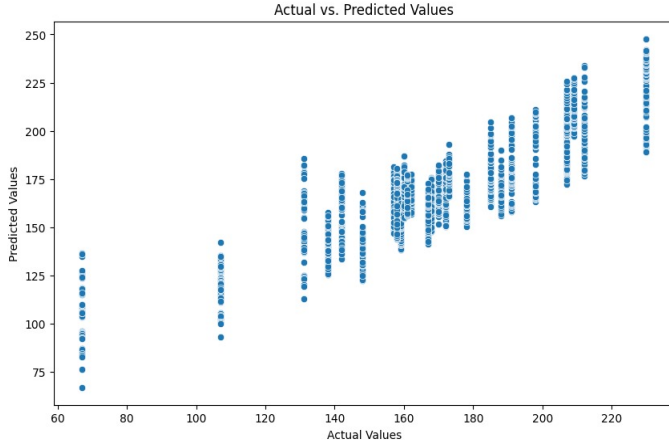


Fig. 2. Graph showing analysis of line plot using LR model

From Fig. 3, it is observed that the line plot shows the range of predicted scores between 60 and 240 and shows the ability of the model to provide diverse predictions across the spectrum of possible scores. In particular, the predicted line consistently matches the actual score, even for outliers, indicating the robust performance of the model in capturing different scenarios. Closer inspection reveals a cluster of data points and prediction lines between the range 140 and 175, indicating a higher frequency of accurate predictions in this range of scores. This clustering confirms the effectiveness of the model in consistently predicting scores, especially for IPL matches in the range of frequent observations.

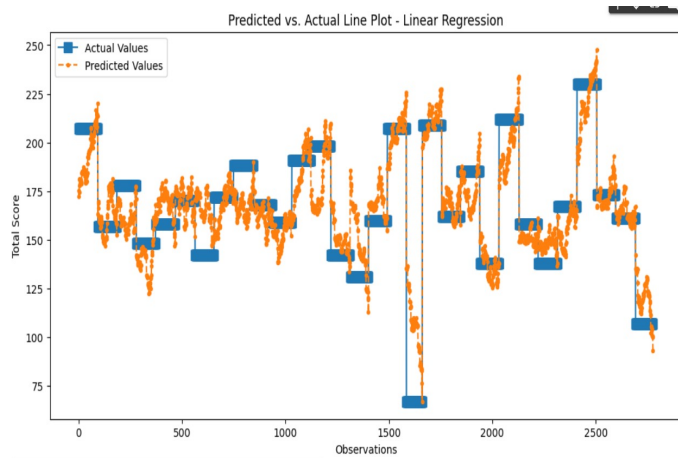


Fig. 3. Graph showing analysis of line plot using LR model

B. Proposed IPL score prediction framework using RF

Evaluation metrics for the RF model are calculated at different max depth values to determine the optimal depth that gives the best prediction performance as shown in Table II. For max depth = 9, this model shows an MAE of 13.43, an MSE of 297.56, RMSE of 17.25, and an R^2 value of 0.71. As max depth increased, the performance of the model improved with max depth = 11 showing the best results: MAE 13.06, MSE 291.40, RMSE 17.07, and R^2 value 0.71. This means that about 71.24% of the variation in the IPL score is explained by the RF model with max depth 11. In particular, the R^2 value for the best-performing model closely matches the LR model and shows that it is comparable prediction efficiency. Using the optimized RF model with max depth = 11, the final predicted score for user input is 193.74. This comprehensive measure provides a detailed assessment of the accuracy and robustness of the RF model in various deep configurations.

TABLE III
PREDICTED EVALUATION METRIC VALUES USING RF MODEL

Evaluation Metrics	Max Depth	Values
MAE MSE RMSE R^2	9	13.431570993070723 297.56426740564103 17.250051229072945 0.706313761411314
MAE MSE RMSE R^2	10	13.179801689575434 293.9452995976107 17.14483302915519 0.7098855647477176
MAE MSE RMSE R^2	11	13.061816283946026 291.5606972565823 17.075148528097266 0.7122390895784115
MAE MSE RMSE R^2	12	13.052005121403399 295.2223824819888 17.182036622065173 0.7086251255426039
MAE MSE RMSE R^2	13	13.192187916308637 305.2494186875832 17.47138857353883 0.6987287674446095

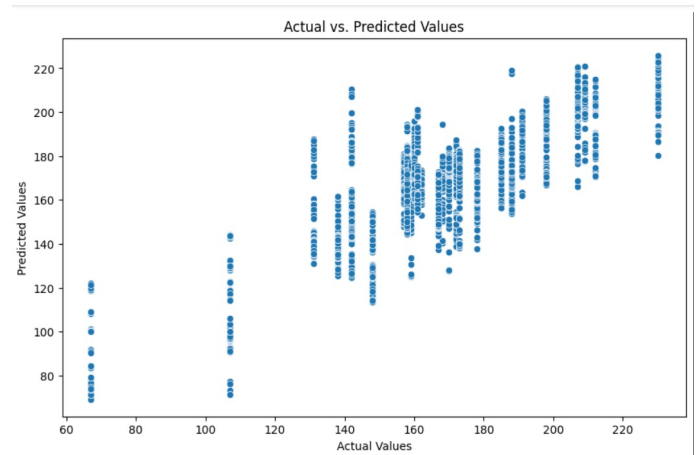


Fig. 4. Graph showing analysis Of scatter plot using RF model

From Fig. 4, it is observed that the range of actual values in the x-axis tested for the model is between 60 and 220, and the range of predicted values for the y-axis is between 80 and 220. It reveals a greater range in predicted values and more variability in model predictions. The distribution of data points on the x-axis shows that most team scores cluster between 130 and 210. However, the scatter plot shows that the model is not consistent in its predictions because the distribution of the points varies. The predicted values represent the range determined by the input attributes, indicating the adaptability of the model to different scenarios, but also reflecting the inherent unpredictability in some cases.

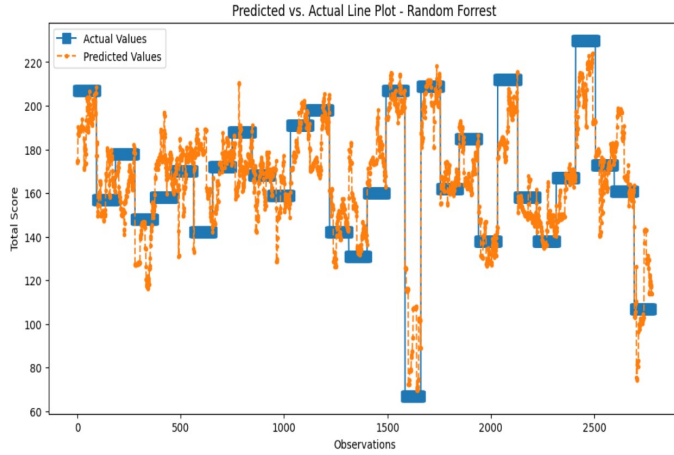


Fig. 5. Graph showing analysis of line plot using RF model

From Fig. 5, it is observed that the plot shows a broader range of predicted values, between 70 and 220, indicating that the model makes diverse predictions in various score scenarios. Unlike the LR plot, the predicted line for RF appears to be less consistent, showing higher variability and inaccuracies in predicting the actual score. The wider the spread of predicted scores for a given actual score, the more inconsistent the model's predictions. This variation indicates that the RF model may struggle to consistently capture the underlying patterns in the data, resulting in a less accurate representation of the actual score.

C. Proposed IPL score prediction framework using SVM

Various kernel functions, namely linear, RBF, sigmoid, and 3rd-degree polynomial are used to evaluate the performance of the SVM model. Table III illustrates the corresponding RMSE values for each kernel. The linear kernel exhibited competitive performance with an RMSE of 16.44. In contrast, the 3rd-degree polynomial kernel yields an RMSE of 18.16, while the RBF kernel shows an RMSE of 16.96. However, the sigmoid kernel showed a high RMSE of 64.54, indicating a less favorable fit for the given dataset. Further analysis identified the linear kernel as the optimal choice, predicting a score of 187.62 for user input, as detailed in Table IV. Evaluation metrics for the loaded SVM model using a linear kernel, MAE,

MSE, RMSE, and R^2 , highlight the effectiveness of the linear kernel in achieving accurate IPL score predictions.

TABLE IV
PREDICTED EVALUATION METRIC VALUES USING SVM MODEL

Evaluation Metrics	Values
SVM Metrics for kernel-linear	
Root Mean Squared Error	16.4385896162757
SVM Metrics for kernel-poly	
Root Mean Squared Error	18.16494478983923
SVM Metrics for kernel-RBF	
Root Mean Squared Error	16.962782579599097
SVM Metrics for kernel-sigmoid	
Root Mean Squared Error	64.5390812299391

TABLE V
BEST PREDICTION MODEL BASED ON SVM MODEL

Best SVM Model (kernel-linear)	
SVR (kernel='linear')	
SVM Predicted Score for User Input	187.62165038

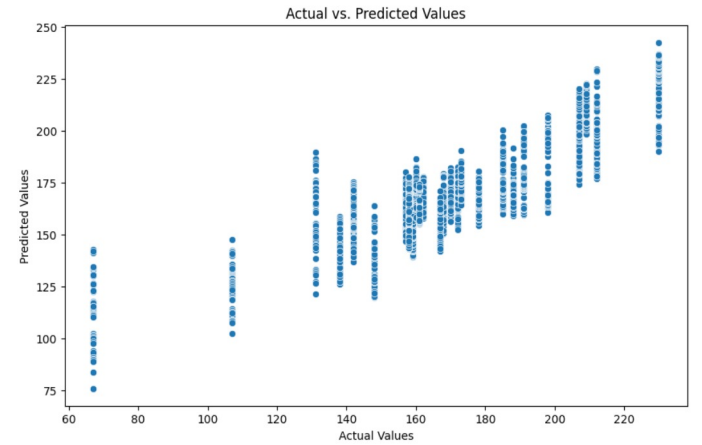


Fig. 6. Graph showing analysis of scatter plot using SVM model

From Fig. 6, it is observed that the range of actual values in the x-axis tested for the model is between 60 and 220, and the range of predicted values for the y-axis is between 60 and 250. This plot shows similar characteristics to an LR scatter plot, except that there are small ranges in the predicted values that represent consistent predictions. According to IPL scoring trends, the X-axis has a concentration of data points between 130 and 210, which represents the common scoring range of most teams. This shows that the SVM obtains a pattern that matches the typical IPL scoring scenarios. Although the scatter plot shows the ability of the model to provide more stable predictions, some variations may still be observed, especially in accommodating outliers or nuanced situations.

From Fig. 7, it is observed that this plot repeats the characteristics seen in the LR model scenario with predicted values from 60 to 240. The predicted line consistently matches the actual score, showing the reliability of the model in

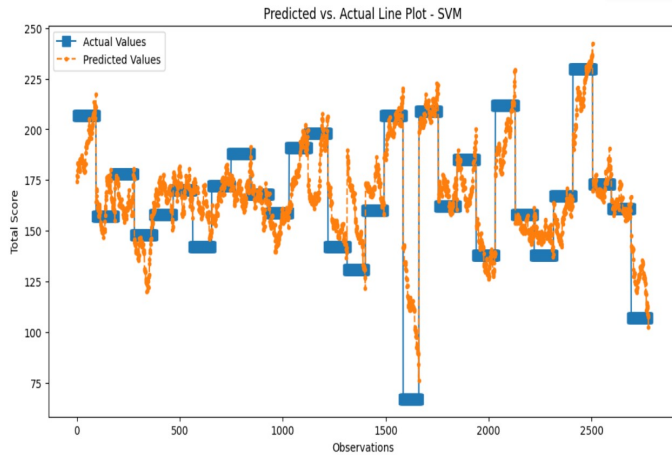


Fig. 7. Graph showing analysis of line plot using SVM model

capturing different scenarios. Similar to LR, even outliers are consistently predicted, confirming the robust performance of the model. In particular, there is a cluster of data points and prediction lines between 140 and 175, indicating a higher frequency of accurate predictions in this scoring interval. Sharing similarities with LR, the unique properties of SVM contribute to its effectiveness of consistently predicting scores across a range of common IPL scores, especially in unique situations.

D. Comparative Analysis of the Proposed Framework Across All the Models

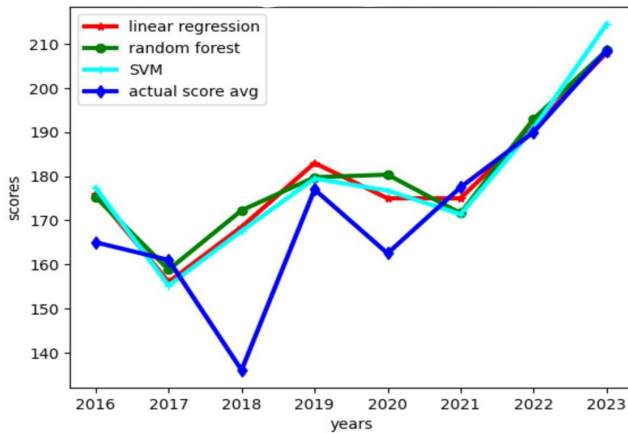


Fig. 8. Comparison of various models with respect to year-wise real score average

In the examination of IPL matches in 2023, featuring Mumbai Indians vs Royal Challenger Bangalore and Rajasthan Royals vs Sunrisers Hyderabad, have been reviewed in detail as viewed in Fig. 8. Mumbai Indians were impressive from 99 to 174 for 2 wickets between 10 and 15 overs, scoring 75 runs in the last 5 overs to reach a total of 200 for 4 wickets. On the other hand, Sunrisers Hyderabad scored 217 for 6, with a substantial 64 runs in the last 5 overs. The average

score in the two games was 208.5. Predictions were made for each game using LR, RF, SVM models, resulting in scores of 208, 208.58, and 214.5, respectively. This prediction is then compared with 217 points, providing a comprehensive assessment of the model's performance from 2016 to 2023. In the same way, random matches were chosen from each edition of the IPL from 2016 to 2023, and a similar analysis was done to get the points, which were plotted on the graph.

E. Analysis of Proposed Framework Using Evaluation Metrics With Respect to All the Models

Evaluating three machine learning models - LR, RF, and SVM - the main performance metrics are MAE, MSE, RMSE, and R^2 . LR emerges as a top performer. LR shows the lowest MSE at 251.01, surpassing both SVM (MSE: 270.23) and RF (MSE: 291.56). LR has a higher MAE (12.12) and RMSE (15.84) compared to RF and SVM, which shows slightly higher error. When evaluating the goodness of fit, LR has the highest R^2 value of 0.75, better than RF (0.71) and SVM (0.73). The results show LR as a robust model, providing accurate and robust predictions for IPL scores compared to other models.

TABLE VI
COMPARISON OF PERFORMANCE METRICS FOR DIFFERENT ML TECHNIQUES

Metrics	LR	RF	SVM
MAE	12.1186	13.0618	12.1651
MSE	251.0079	291.5607	270.2272
RMSE	15.8432	17.0751	16.4386
R^2	0.7523	0.7122	0.7333

Phase 2: Deployment and User Interface

In creating an online cricket match final score prediction app, the user interface will allow users to enter various details such as the batting team, bowling team, current runs, wickets, score in the last 5 overs, and wickets fallen in the last 5 overs. This user interface will be designed using HTML, CSS, and possibly JavaScript to provide users with a visually appealing and interactive experience.

The back-end of the application will be powered by Flask, a web framework for Python. In Flask, the stack file model that was previously trained and saved as a .pkl file will be loaded using the pickle module. This model will be responsible for processing the user inputs and predicting the final score based on the details provided.

Once the user submits a form with the required details, Flask processes the form submission, processes the input data, and passes it to the loaded stack file model for prediction. The predicted final score is then displayed back to the user, either on a new page or on the same page, giving them the valuable insight they are looking for.

This approach not only provides an intuitive and user-friendly experience but also uses the power of machine learning to provide accurate predictions, increasing overall user engagement and satisfaction with the app.

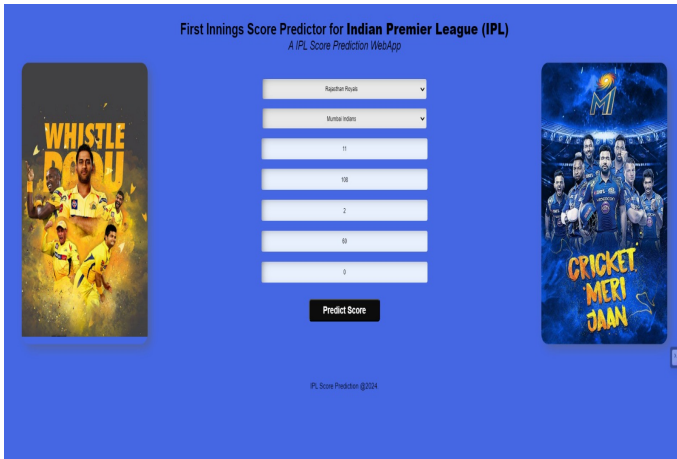


Fig. 9. Website User Interface of Proposed Model

An innovative approach has been taken to improve the clarity of predicted IPL scores in the website layout. The final score is displayed as a 10-digit range rather than as a single point estimate. This range is achieved by offering a wider range of potential outcomes, including +5 and -5 to the predicted score. By presenting that section, the placement strategy recognizes the uncertainty involved in predicting cricket scores and gives users a more complete understanding of possible score changes. This thoughtful presentation not only provides the model's forecast but also provides transparency and informed decision-making for users interested in the website's score forecast, as well as the level of uncertainty.

of the final score by the model. This instance showcases the model's ability to provide valuable insights into match outcomes based on given match parameters.

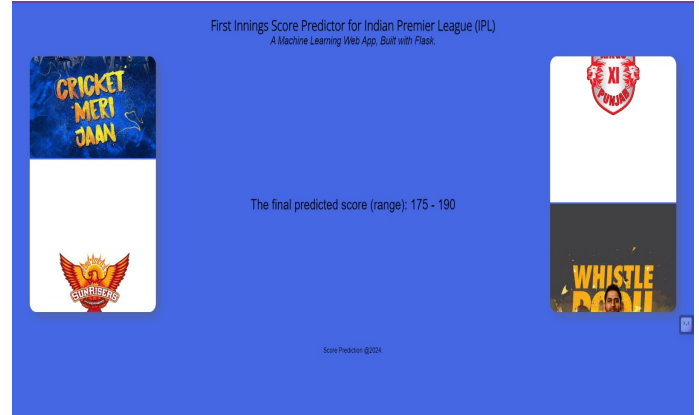


Fig. 11. Output of the Predicted Proposed Model

VI. CONCLUSION

In conclusion, this work serves as a demonstration of ML application for predictive insights in sports data. While LR has made a significant contribution to sports analytics, the consideration of more complex models such as RF and SVM suggests the preference of LR in this context. This discovery not only enriches the field of sports analytics but also shows the integral role of machine learning in shaping the dynamic landscape of sports analytics, creating a solid foundation for future research.

This work demonstrates the significant potential of ML, especially LR, in sports analytics, with a special focus on IPL score prediction. While the accuracy of the predictions shows promise, there is continual room for improvement. The incorporation of additional features such as player shape, team composition, and game conditions is an important way to increase the efficiency of the model. In addition, the exploration of more sophisticated ML models, including decision trees, RF, and neural networks, promises to capture complex relationships in the data, thus increasing the accuracy of predictions. Overall, this research contributes to the development of sports analytics, laying the foundation for future progress and methodology.

VII. FUTURE SCOPE

The future scope of this work involves advancements in cricket predictive analytics. Continuous improvement, realistic predictions during live IPL matches, and improving model performance by integrating additional data sources offer opportunities for development. A user-friendly interface improves accessibility and extends the model to provide insight into player choices, adding valuable strategic information. By exploring this area, this work aims to make a significant contribution to cricket predictive analytics, team management, sports analysts, and cricket enthusiasts.

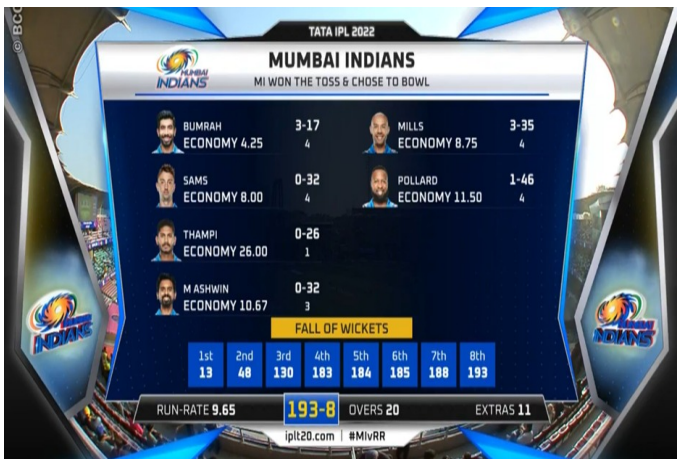


Fig. 10. Output of the Actual Proposed Model

In this work, specific input values were provided to the trained ML model, including details such as the batting team being Rajasthan Royals, bowling team as Mumbai Indians, 11 overs played, 108 runs scored, 2 wickets taken, 60 runs in the last 5 overs, and no wickets in the last 5 overs. The model generated a predicted score ranging between 175 and 190. Upon comparison with the actual score of 193, it reveals a reasonably accurate prediction, with a slight underestimation

REFERENCES

- [1] Indian Premier League (IPL), <https://www.iplt20.com/>
- [2] Lalitha S., Gupta D., Zakariah, M., Alotaibi Y.A., "Mental Illness Disorder Diagnosis Using Emotion Variation Detection from Continuous English Speech", *ICTACT Journals*, 2021.
- [3] Lalitha S., Gupta D., "Investigation of Automatic Mixed-lingual Affective State Recognition System for Diverse Indian Languages", 2021.
- [4] Joshi K., Reddy G.A., Kumar S., Anandaram H., Gupta A., Gupta H., "Analysis of Heart Disease Prediction using Various Machine Learning Techniques: A Review Study", 2023.
- [5] Eeshan Mundhe, Ishan Jain, Sanskar Shah, "Live Cricket Score Prediction Web Application using Machine Learning", *IEEE*, 2021.
- [6] Indika Wickramasinghe, "Applications of Machine Learning in cricket: A systematic review", *Elsevier*, 2022.
- [7] Manoj Ishi, Jayantrao Patil, Vaishali Patil, "An efficient team prediction for one day international matches using a hybrid approach of CS-PSO and machine learning algorithms", *Elsevier*, 2022.
- [8] Yogesh Kumar, Harendra Sharma, Ritu Pal, "Popularity Measuring and Prediction Mining of IPL Team Using Machine Learning", *IEEE*, 2021.
- [9] Shubhra Singh, Parmeet Kaur, "IPL Visualization and Prediction Using HBase", *Elsevier*, 2017.
- [10] Saranya G, Aravind Swaminathan, Joel Benjamin J, "IPL Data Analysis and Visualization for Team Selection and Profit Strategy", *IEEE*, 2023.
- [11] Waqas Ahmad, Muhammad Munsif, Habib Ullah, Mohib Ullah, Alhanouf Abdulrahman Alsuwailam, Abdul Khader Jilani Saudagar, Khan Muhammad, Muhammad Sajjad, "Optimized deep learning-based cricket activity focused network and medium scale benchmark", *Elsevier*, 2023.
- [12] Aminul Islam Anik, Sakif Yeaser, A.G.M Imam Hossian, Amitabha Chakrabarty, "Player's Performance Prediction in ODI Cricket Using Machine Learning Algorithms", *IEEE*, 2018.
- [13] E. L. Lekamge, K. R. Wickramasinghe, S. E. Gamage, T. M. K. L. Thennakoon, Mr. Prasanna S. Haddela, Ms. Sandamini Senaratne, "Cric-squad: A System To Recommend Ideal Players To A Particular Match And Predict The Outcome Of The Match", *IEEE*, 2023.
- [14] Afsana Khan, Fariha Haque Nabila, Masud Mohiuddin, Mahadi Mollah, Md Ashraful Alam, Md Tanzim Reza, "An Approach to Classify the Shot Selection by Batsmen in Cricket Matches Using Deep Neural Network on Image Data", *IEEE*, 2022.
- [15] Lalitha S., Jayashree R.J., Ganesh S., Karanth S.C., Automatic Detection of Parkinson Speech Under Noisy Environment, 2021
- [16] Lalitha S., Rachana N., Vinay Bhargav J., Multilingual and Mixed-lingual Digit Speech Recognition System for Indian Context, 2023
- [17] V. Parthasarathy, L. Pallavi, H. Anandaram, M. Praveen, S. Arun, and R. Krishnamoorthy, Prediction of Coronary Artery Disease by Adapting Hybrid Approach of Machine Learning Methods, 2022
- [18] ESPN, <https://www.espnricinfo.com/>.
- [19] Cricbuzz, <https://www.cricbuzz.com/>.
- [20] Scikit-Learn One Hot Encoding Documentation, <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.OneHotEncoder.html>.
- [21] C. M. Bishop, "Pattern Recognition and Machine Learning", First Edition, Springer, 2006.