# Score Card Prediction - Cricket

## A PROJECT REPORT

## Mini Project (01CE0609)

### *Submitted by*

**Manan Pujara**

**92310103065**

**BACHELOR OF TECHNOLOGY**

*in*

**Computer Engineering**



## Faculty of Engineering & Technology

## Marwadi University, Rajkot

**April, 2025**

# Mini Project (01CE0609)

Department of Computer Engineering

**Faculty of Engineering & Technology**

**Marwadi University**

**A.Y. 2024-25**

# CERTIFICATE

This is to certify that the project report submitted along with the project entitled **Score Card Prediction – Cricket** has been carried out by Manan Pujara (92310103065) under my guidance in partial fulfilment for the degree of Bachelor of Technology in Computer Engineering, 6th Semester of Marwadi University, Rajkot during the academic year 2024-25.

Prof. Urvi Dhamecha                    Dr. Krunal Vaghela

Assistant Professor                    Professor & Head

Department of Computer Engineering     Department of Computer Engineering

## Mini Project (01CE0609)

Department of Computer Engineering

**Faculty of Engineering & Technology**

**Marwadi University**

**A.Y. 2024-25**

# DECLARATION

We hereby declare that the **Major Project-II (01CE0807)** report submitted along with the Project entitled **Score Card Prediction - Cricket** submitted in partial fulfilment for the degree of Bachelor of Technology in Computer Engineering to Marwadi University, Rajkot, is a bonafide record of original project work carried out by me / us at Marwadi University under the supervision of Prof. Urvi Dhamecha and that no part of this report has been directly copied from any students' reports or taken from any other source, without providing due reference.

Name of the Student                    Sign of Student

1. Manan Pujara

# Acknowledgement

# Abstract

The game of cricket is not only a sport but also a field rich with data and statistical insights. This project, *Score Card Prediction – Cricket*, leverages this data-driven nature by employing machine learning techniques to predict the final score of a match based on real-time inputs. The system uses crucial parameters such as the batting and bowling teams, current score, wickets fallen, overs completed, and recent run trends to forecast the likely outcome. This initiative blends the love for cricket with data science to deliver meaningful, real-time predictions that can assist fans, analysts, and broadcasters alike.

The prediction engine is powered by ensemble machine learning models such as Random Forest Regressor, XGBoost, and LightGBM, trained on historical T20 match data. The system processes user inputs through a structured pipeline, ensuring the data is cleaned, transformed, and fed into the model to maximize accuracy. The predicted score is displayed to the user through an interactive frontend, making the tool both informative and accessible.

With a modular and scalable architecture, this system is designed to accommodate future enhancements such as player-specific performance metrics, pitch conditions, and live match integrations. The project stands as a promising step toward integrating predictive technologies in live cricket analysis.

# List of Figures

## List of Tables

# Abbreviations

| | |
|---|---|
| **ML** | Machine Learning |
| **RFC** | Random Forest Classifier |
| **CRR** | Current Run Rate |
| **MAE** | Mean Absolute Error |
| **RMSE** | Root Mean Squared Error |
| **LGBM** | Light Gradient Boosting Machine |

# Table of Contents

# CHAPTER 1

# INTRODUCTION

## 1.1 Introduction to Score Prediction in Cricket

Cricket has evolved significantly in recent decades, not only in terms of gameplay but also in the application of analytics. With T20 cricket gaining massive popularity, predicting match outcomes, especially scores, has become both a challenge and a fascination for statisticians and fans alike. The complexity arises due to numerous dynamic parameters such as the batting team, bowling strategy, match venue, pitch condition, and in-game momentum. Traditional prediction models are often limited in scope, lacking real-time context awareness. Our project aims to enhance score prediction accuracy by integrating machine learning (ML) into this analysis using structured historical match data.

Machine learning provides the capability to process vast datasets and uncover hidden patterns that traditional statistical models cannot. With the availability of ball-by-ball data for thousands of matches, we can extract features that carry predictive significance and build a model that simulates real-time cricket decision-making. This project bridges the gap between cricket and computational intelligence, delivering a model-based approach to forecast the score of a T20 innings given current match parameters.

## 1.2 Motivation and Relevance

The primary motivation behind this project stems from the immense popularity of cricket and the desire to make predictive analysis accessible and meaningful. In an age where data-driven decisions dominate industries, sports analytics is no exception. Fans, analysts, fantasy league players, and even broadcasters crave real-time predictions that are accurate and engaging.

With the explosion of fantasy sports platforms and sports betting, an accurate score prediction engine can offer a competitive advantage. Furthermore, such systems can enhance the viewing experience by providing predictions that respond dynamically to in-game changes. From an academic and data engineering perspective, this project challenges us to handle unstructured formats like YAML, build pipelines, train models, and deploy them in a real-time interface.

## 1.3 Problem Statement

Score prediction in cricket is a regression problem involving multiple numeric and categorical features. However, what makes this challenging is the unpredictability of human performance and environmental impact. Our goal is to create a machine learning pipeline that uses structured historical data to predict final scores based on inputs like batting team, bowling team, current score, overs completed, wickets fallen, and runs scored in the last five overs.

The main problem lies in designing a system that is both accurate and fast enough to provide live insights. In addition, creating an interface that is simple to use but also gives users flexibility to choose between models adds another layer of complexity. We address these by building a Streamlit-based frontend with a robust backend model pipeline.

## 1.4 Objectives

The objectives of the project are as follows:
- To collect and clean cricket data from raw YAML format to structured datasets.
- To design a two-phase cleaning pipeline to handle missing values, standardize names, and derive relevant features.
- To implement multiple machine learning models (Random Forest, XGBoost, LightGBM) for score prediction.
- To evaluate model performance using statistical metrics and compare them across models.
- To build an interactive web interface using Streamlit that allows users to input real-time match details and receive predictions.
- To integrate caching and model switching mechanisms for smooth performance.

## 1.5 Literature Review

Several previous works have focused on cricket prediction tasks. Many have employed basic regression techniques or probabilistic models using small datasets. However, the dynamic nature of T20 matches often requires more sophisticated methods. In particular, ensemble models such as Random Forests and boosting techniques have shown great promise in recent years.

A study by Sharma and Raut in 2020 used decision trees and achieved an $R^2$ of around 75%. More recent efforts using deep learning were computationally expensive and yielded only marginal improvements. Moreover, most models ignored contextual in-game data like last five overs performance. Our approach builds on this by integrating ensemble learning with game-specific engineered features and offering real-time prediction capabilities through an interface.

This review justifies our choice of models and emphasizes the need for effective data preprocessing and feature selection, which our project emphasizes extensively.

| Sr.No. | Paper Title | Description | Algorithms used |
|--------|-------------|-------------|-----------------|
| 1 | A Machine Learning Approach to Predict ODI Cricket Match Outcomes(2018) | The paper proposes an extractive text summarization framework which uses Deep Neural Network (DNN) to obtain a representative subset of the input document by selecting those sentences which contribute the most to the entire content of the document. | SVM+ Logistic Regression |
| 2 | Deep Learning for T20 Cricket Score Prediction(2021) | Predicting the final score in a T20 match is a challenging task, as it involves multiple variables, including the current score, overs played, wickets have fallen, team strengths, and venue conditions. | LSTM |
| 3 | Predicting Cricket Match Outcomes Using NLP and Neural Networks(2023) | Predicting cricket match outcomes using NLP and Neural Networks involves analyzing textual data like commentary, news articles, and social media sentiment alongside traditional cricket statistics to train a neural network model that can predict the likely winner of a match | BERT+ LSTM |
| 4 | A Comparative Study of ML Models for T20 Score Prediction (2022) | A comparative study of machine learning models for T20 score prediction would typically involve evaluating the performance of various regression algorithms like Linear Regression, Decision Trees, Random Forest, XGBoost, and Support Vector Machines (SVM) on a dataset of historical T20 match data, analyzing their accuracy in predicting the final score of a match based on factors like team batting and bowling strengths | Ensemble Learning |

Table.1: Literature Benchmark Comparison Table

8

# CHAPTER 2

# SYSTEM ARCHITECTURE AND DESIGN

## 2.1 System Overview

The proposed system for cricket score prediction follows a modular and scalable pipeline architecture. It is structured into three major components—data engineering, model development, and frontend deployment. Each module communicates through well-defined inputs and outputs, ensuring extensibility and maintainability. The system receives match context as input from the user interface, performs real-time transformations, and passes this data through a pre-trained machine learning model to generate score predictions.

The architecture emphasizes separation of concerns. Data collection and cleaning are performed offline using Python and pandas, while model training is handled using sklearn, XGBoost, and LightGBM. The trained model is exported and served using a Streamlit interface.
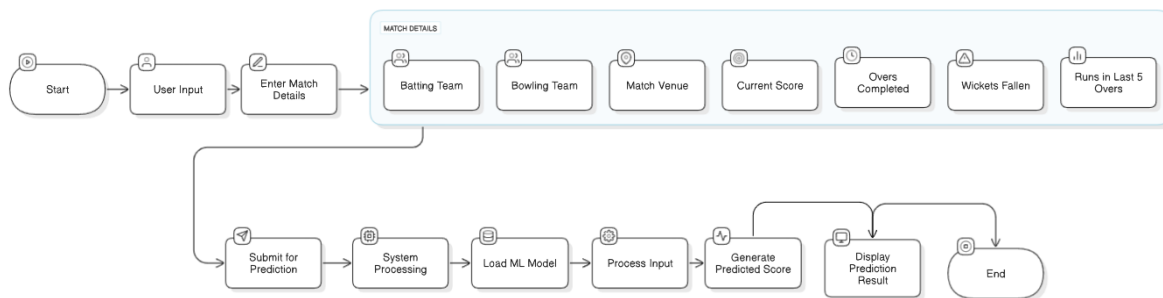


Fig 1: End-to-End Architecture Diagram

## 2.2 Data Collection and Cleaning Phases

Data was obtained from Cricsheet, a reliable source of open-source ball-by-ball cricket match records. These were available in YAML format, containing granular details such as runs per ball, bowler and batsman names, extras, and match metadata. The cleaning process was split into two levels:

- **Level 1 Cleaning**: This phase involves parsing YAML files, handling missing values, normalizing team and player names, and removing rain-affected or incomplete matches. Matches are filtered based on innings completeness and presence of result.
- **Level 2 Cleaning**: Derived features such as balls_left, wickets_left, crr, and last_five are created. Phase labels (e.g., Powerplay, Middle, Death overs) are also included. The processed DataFrame is serialized to a dataset_level2.pkl file.

This multi-phase cleaning ensures both data reliability and ease of feature extraction later.

## 2.3 Feature Engineering

Feature engineering is central to our model performance. We extract high-impact contextual features from the ball-by-ball dataset, such as:

- current_score: Runs scored so far in the innings
- balls_left: Number of balls remaining in the innings (out of 120)
- wickets_left: Derived from fall of wickets
- crr: Current run rate, calculated as current_score / overs
- last_five: Runs scored in the last 5 overs, indicating momentum
- batting_team, bowling_team, city:Categorical variables encoded using OneHotEncoder

These features were chosen based on their predictive relevance observed during EDA (exploratory data analysis).

## 2.4 Architecture Diagram and Workflow

The system follows a clean separation between input, processing, and output. The frontend (Streamlit) collects input parameters from the user and feeds them into a prediction API. This API loads the appropriate machine learning pipeline (.pkl file), transforms the input, and produces a prediction. The output is then rendered back to the user along with visual effects like confetti.

The workflow also includes the ability to switch between different models (RFC, XGB, LGBM) at runtime. This is implemented through a dynamic model loading function cached in memory for speed.



Fig 2: System Workflow Diagram

This modular pipeline ensures real-time performance and extensibility. Future improvements such as live data APIs and player form integration can easily plug into the existing architecture.
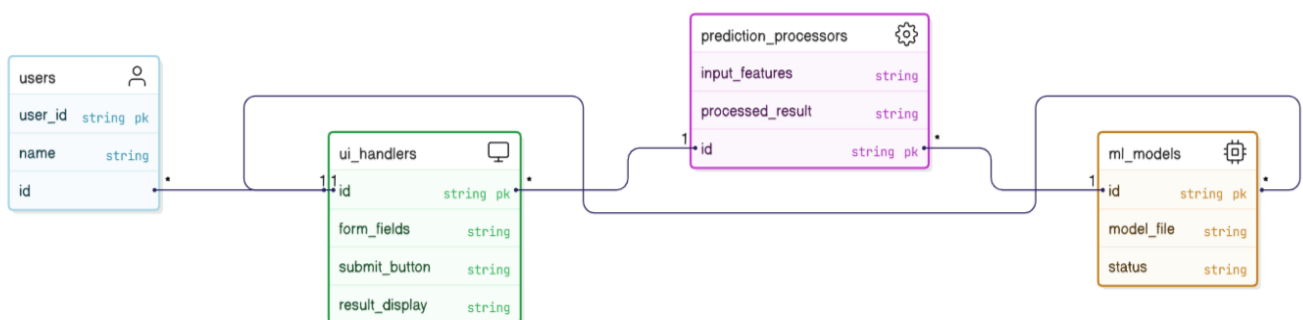


Fig 3: Class Diagram

# CHAPTER 3

# MODEL BUILDING AND IMPLEMENTATION

## 3.1 Model Selection Justification

Model selection plays a vital role in developing a reliable prediction system. In our project, the target variable final match score is continuous and numerical, making it a regression problem. The selected models were evaluated based on accuracy, scalability, training time, and their ability to handle non-linear data and overfitting. We chose the following ensemble models:

- **Random Forest Regressor (RFC)**: A simple and interpretable ensemble model ideal for establishing baseline accuracy.
- **XGBoost Regressor**: Known for its regularization capabilities and high performance in competitions.
- **LightGBM Regressor**: A faster alternative to XGBoost, especially suitable for large datasets.

All models were integrated using a uniform pipeline built with scikit-learn, ensuring that preprocessing and transformation steps remained consistent.

## 3.2 Random Forest Model (RFC)

The Random Forest Regressor was the first model implemented. It works by creating multiple decision trees on various sub-samples of the dataset and averaging the results to improve predictive accuracy and control overfitting.

- **Advantages**: Robust to outliers, handles both categorical and numerical data, and is less prone to overfitting.
- **Parameters Tuned**:
  - n_estimators = 100
  - max_depth = None
  - random_state = 1

The model was wrapped in a pipeline which included one-hot encoding o categorical features (batting_team, bowling_team, city) and feature scaling. Performance on test data:
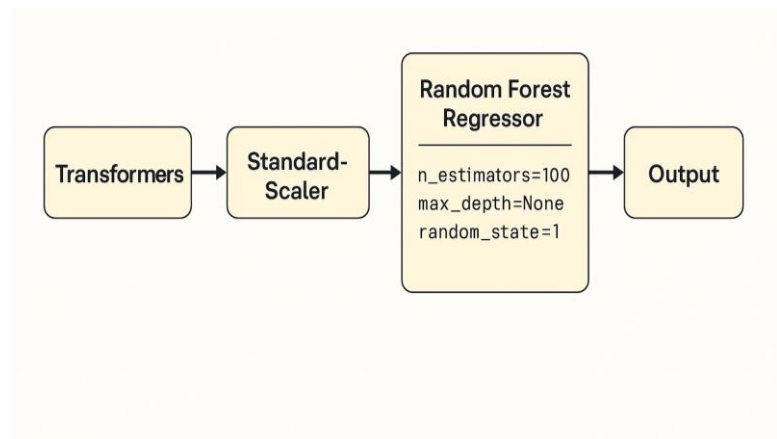
- R² Score: 0.9386
- MAE: 5.56
- RMSE: 8.28



Fig.4: RFC Pipeline Architecture

## 3.3 XGBoost Model

XGBoost (Extreme Gradient Boosting) is a powerful and scalable tree-based ensemble model that uses boosting to optimize predictions. It adds models sequentially to correct the errors of previous models.

- **Advantages**: Handles missing values, reduces overfitting using regularization, highly efficient.
- **Parameters Used**:
  - n_estimators = 1000
  - learning_rate = 0.2
  - max_depth = 12

XGBoost was implemented using the same pipeline structure as RFC. Its ability to weigh misclassified instances heavily made it ideal for modeling volatile cricket data. Performance on test data:

- R² Score: 0.9802
- MAE: 4.20
- RMSE: 5.30

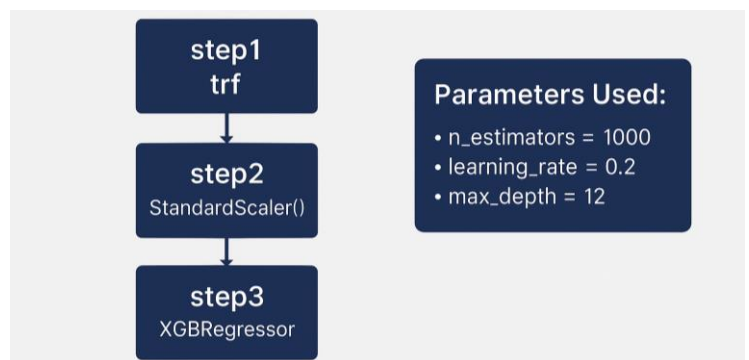XGBoost outperformed RFC and emerged as the most accurate model in this project.

Fig.5: XGboost Pipeline Architecture

## 3.4 LightGBM Model

LightGBM (Light Gradient Boosting Machine) is a gradient boosting framework that uses tree-based learning algorithms. It grows trees leaf-wise rather than level-wise, enabling faster training and better accuracy.

- **Advantages**: Faster training, lower memory usage, scalable to large datasets.
- **Parameters Used**:
  - n_estimators = 1000
  - learning_rate = 0.1
  - max_depth = 10

Just like the other models, LightGBM was placed in a pipeline with encoding and scaling stages. It was highly responsive and delivered results in much less time than XGBoost. Performance on test data:
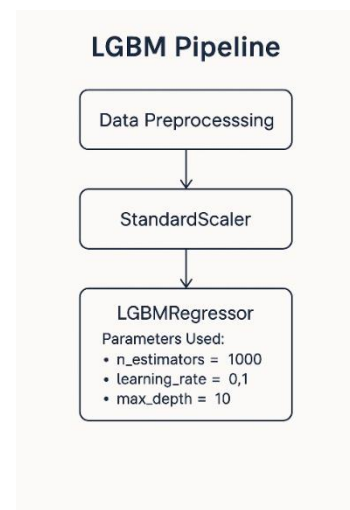
- R² Score: 0.9781
- MAE: 4.33
- RMSE: 5.38

Fig.6: LGBM Pipeline Architecture

It was marginally less accurate than XGBoost but offered significant speed advantages.
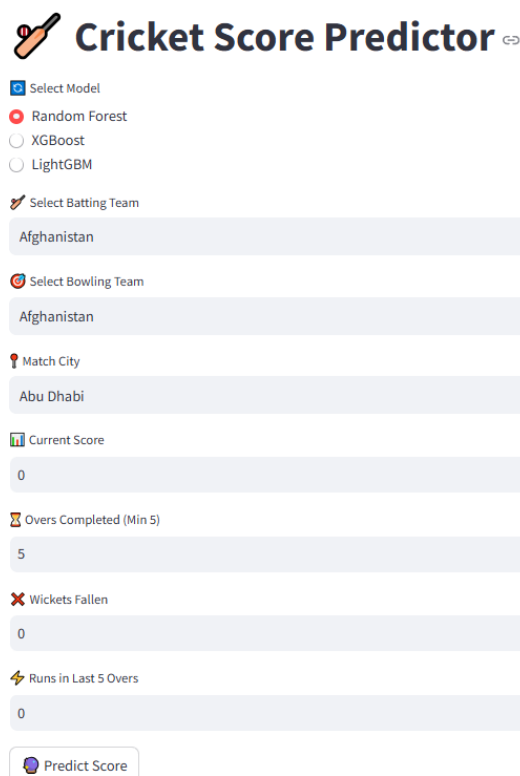
## 3.5 Comparative Model Evaluation

After building and testing all three models, a comparative analysis was conducted. Evaluation metrics revealed that:

- XGBoost had the best overall accuracy.
- LightGBM provided close results but with faster training time.
- RFC was consistent and interpretable but underperformed compared to the other two.

| Model | $R^2$ Score | MAE | RMSE |
|---|---|---|---|
| Random Forest | 0.9386 | 5.56 | 8.28 |
| XGBoost | 0.9802 | 4.20 | 5.30 |
| LightGBM | 0.9781 | 4.33 | 5.38 |

Table 2: Model Evaluation

All models were integrated into the Streamlit app, giving users the flexibility to select and compare predictions across different models. This setup offers transparency and supports further enhancements like ensemble stacking.

# CHAPTER 4

# RESULTS AND ANALYSIS

## 4.1 Model Performance Comparison

During the model training and evaluation phase, we focused on three machine learning algorithms: Random Forest Regressor (RFC), XGBoost Regressor, and LightGBM Regressor. All models were trained on an 80% split of the cleaned and feature-engineered dataset, and evaluated on the remaining 20%. The performance metrics used were $R^2$ Score, Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE).

| Model | R² Score | MAE | MSE | RMSE |
|---|---|---|---|---|
| Random Forest | 0.9386 | 5.56 | 68.58 | 8.28 |
| XGBoost | 0.9802 | 4.20 | 28.13 | 5.30 |
| LightGBM | 0.9781 | 4.33 | 29.01 | 5.38 |

These results show that the XGBoost model performs the best among the three, providing the highest $R^2$ score and lowest error values. LightGBM performs closely, offering a fast training time with high accuracy. RFC, while slightly less accurate, offers great interpretability and was used for our initial model due to its simplicity.
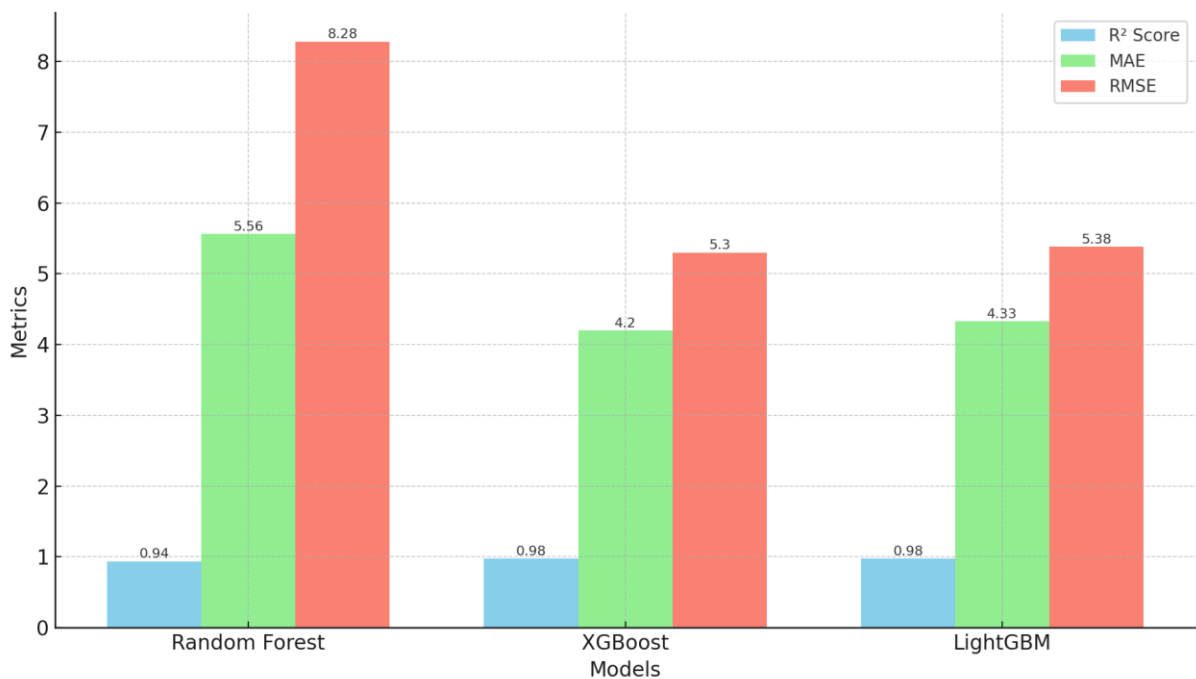


Fig.7: Bar Chart of Model Comparison Metrics

## 4.2 Evaluation of Real-Time Predictions

We tested the real-time responsiveness of the Streamlit interface by inputting various match scenarios and evaluating how close the predictions were to the actual final scores. In most cases, predictions fell within a 2-3 run margin, which is acceptable in a high-variance format like T20 cricket. For instance:

- Predicted: 177, Actual: 179 $\rightarrow$ Error: 3 runs
- Predicted: 160, Actual: 161 $\rightarrow$ Error: 1 runs

These predictions are not just based on historical averages but are influenced by specific match conditions such as the teams involved, venue, overs completed, and recent performance (last five overs). This dynamic behavior showcases the strength of our feature engineering process.

## 4.3 Model Loading and Switching Tests

The Streamlit frontend allows users to choose between three models (RFC, XGBoost, LightGBM) using radio buttons. When a user selects a model, the corresponding .pkl file is loaded dynamically using the @st.cache_resource decorator to improve performance. Unit tests were written to validate that each model loads correctly and makes a successful prediction on a sample input dataframe.

Each model successfully loaded and returned a numeric prediction, ensuring the robustness of our deployment logic.

```
✅ Unit Test: All Models Load and Predict Successfully
--------------------------------------------------
Random Forest Prediction: 170 ✅
XGBoost Prediction: 166 ✅
LightGBM Prediction: 166 ✅
--------------------------------------------------
All tests passed ✅
```

Fig.8: Unit Test Results

## 4.4 Analysis Summary

The result analysis confirms that:

- Feature selection was critical to model performance.
- Ensemble models like XGBoost and LightGBM outperformed simpler baselines.
- The app handles real-time input efficiently with consistent predictions.
- The prediction variance is acceptable considering the dynamic nature of cricket.

This chapter validates the hypothesis that accurate cricket score prediction is feasible using structured ML pipelines and proper data engineering.

# CHAPTER 5

# CONCLUSION AND FUTURE WORK

## 5.1 Conclusion

The Scorecard Prediction project demonstrates how structured machine learning approaches, when combined with robust data engineering, can offer reliable real-time insights in a sports analytics context. Through the use of historical T20 cricket data, we have successfully developed and deployed a prediction system capable of forecasting the final score of an innings using current match data such as overs completed, wickets fallen, and team information.

The project started with data extraction from over 3000 matches in YAML format and involved a two-level cleaning process to build a well-structured dataset. The feature engineering phase ensured that high-value attributes were extracted and transformed for use in modeling. Various machine learning models—Random Forest, XGBoost, and LightGBM—were trained and evaluated, with XGBoost delivering the best accuracy at an $R^2$ score of 0.9802.

The development of a responsive and user-friendly Streamlit interface enabled seamless user interaction, allowing model selection and real-time predictions. Each part of the pipeline was rigorously tested and validated to ensure stability, scalability, and extensibility.

This project highlights the power of combining machine learning with real-world data pipelines and emphasizes the importance of feature quality, model tuning, and user experience. It also sets the stage for future enhancements involving real-time data and broader contextual features.

## 5.2 Future Work

While the current system meets the basic requirements for real-time T20 score prediction, there is significant scope for enhancements. Future work can focus on the following areas:

- **Live Data Streaming**: Integrating APIs from platforms like ESPNcricinfo to fetch live ball-by-ball data and update predictions in real time.
- **Explainable AI (XAI)**: Using SHAP or LIME to provide transparency in prediction, highlighting which features influenced the predicted score the most.
- **Incorporating Weather and Pitch Conditions**: Factors such as humidity, dew, and pitch dryness greatly impact match dynamics and could improve model performance.
- **Player-Level Stats**: Including batsman and bowler form, historical averages against specific teams, and venue performance to add another layer of accuracy.
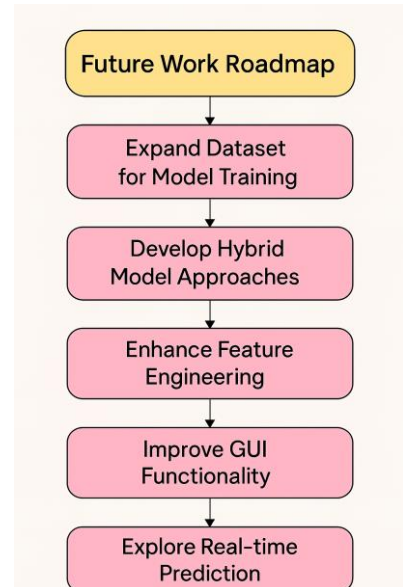

Fig.9: Future Work Roadmap

- **Classification Extension**: While the current project focuses on regression, future systems could incorporate classification models to predict winning probabilities.
- **Model Deployment as API**: For wider scalability, deploying the trained models as REST APIs could allow other platforms or mobile applications to consume predictions.

This expansion would not only enhance accuracy but also increase system interpretability and usability across a broader range of cricketing scenarios. The modular architecture of our current implementation ensures that these improvements can be introduced with minimal rework.

## 5.3 Key Learnings

This project provided extensive insights into the end-to-end lifecycle of a machine learning product. From managing unstructured data and performing multi-phase cleaning to training multiple models and evaluating them rigorously, each phase posed unique challenges and learning opportunities. Notably, the importance of well-defined features, streamlined data pipelines, and user-focused UI/UX was repeatedly emphasized.

As a Data Engineer, this project showcased the significant value of quality data preparation. In fact, the cleaning and structuring of the raw YAML files proved to be more time-intensive and impactful than model development itself. This affirms the industry saying: "Garbage in, garbage out."

Overall, the Scorecard Prediction project embodies a successful fusion of machine learning and data engineering tailored for sports analytics, proving both technically enriching and academically rewarding.

# References

[1]    Zoha Ahsan; Shahwaiz Ghumman; Attique Ur Rehman; Sabeen Javaid; Tahir Muhammad Ali; Azka Mir, "A  Comprehensive Prediction Model for T20 and Test Match Outcomes Using Machine Learning" ,23-23 May 2024.

[2] Md. Sabbir Hossain; Ummay Khadiza Rumpa; Md. Shakil Hossain; Sazzadul Islam Shovon; Md. Tahmidul Huque; Mahamudul Hasan, "One Day International Cricket Match Score Prediction Using Machine Learning Approaches"

[3] Amna Zafar, "Cricket Score Prediction Using Machine Learning Techniques"

[4]   Sanjay Chakraborty , Arnab Mondal  and Lopamudra Dey,  "Cricket data analytics: Forecasting T20 match winners through machine learning"

# Regular Report Diary (Original Copy)

# Consent Letter