

T20I Predictor

Md. Kawsar Ahamed

Dept. of Computer Science and Engineering

Green University of Bangladesh

Rupganj, Narayanganj, Bangladesh

aha.kawsar@gmail.com

Abstract—Predicting outcomes in T20 International (T20I) cricket presents a significant data science challenge due to the format’s inherent volatility and the complex interplay of performance variables. This paper presents the design, implementation, and evaluation of the “T20I Predictor,” a web-based system that employs a novel hybrid approach to forecast match results. The system architecture integrates a multi-layered statistical analysis engine with a diverse suite of seven machine learning models to provide a comprehensive and context-aware prediction. Historical match and ball-by-ball data are first subjected to a rigorous ETL process, including data cleaning, filtering for official T20I matches, and structuring into a relational SQLite3 database. The statistical engine then evaluates a wide range of performance metrics, including head-to-head records, venue-specific performance, and innings-based dynamics under various score scenarios. Concurrently, machine learning models, including Random Forest, Support Vector Machine, and Gradient Boosting, provide independent probabilistic forecasts. The outputs from both engines are synthesized using a configurable weighted scoring system that prioritizes more impactful real-time factors, such as recent form and target score context. The final result is a nuanced win probability that is more robust and transparent than a single-method approach, demonstrating the efficacy of combining statistical interpretability with machine learning’s predictive power.

Index Terms—Sports Analytics, Machine Learning, Predictive Modeling, Cricket, Django, Data Science, Ensemble Learning.

I. INTRODUCTION

A. Background and Motivation

The advent of Twenty20 (T20) cricket has fundamentally altered the landscape of the sport, introducing a fast-paced, aggressive format that has captured a massive global audience. The inherent brevity of T20 International (T20I) matches leads to a high degree of unpredictability, where the momentum of a game can shift within a single over. This volatility, while exciting for spectators, presents a complex challenge for accurate outcome prediction. Simultaneously, the field of sports analytics has matured significantly, with teams, broadcasters, and fans increasingly relying on data-driven insights to understand performance. This confluence of a highly dynamic sport and the demand for sophisticated analytics provides the primary motivation for this research. The goal is to move beyond simplistic metrics like team rankings and develop a tool that can model the complex, non-linear interactions between the numerous variables that influence a T20I match.

B. Problem Statement

The core problem addressed by this research is the development of a predictive model that can accurately forecast the winner of a T20I match while providing a transparent and interpretable rationale for its prediction. Existing prediction methods often fall into one of two categories, each with significant limitations. Statistical models, while interpretable, often struggle to capture the complex, non-linear relationships within the data. Conversely, machine learning models can identify these patterns but often function as “black boxes,” providing a prediction without a clear explanation. Furthermore, any effective model must contend with a multitude of interacting variables, including:

- **Historical Context:** All-time head-to-head records versus recent form.
- **Venue Specificity:** How a team’s performance changes at different grounds.
- **Innings Dynamics:** The significant difference in pressure between setting a target and chasing one.
- **Score Context:** The non-linear impact of the first innings score (e.g., a target of 90 is not merely half as difficult as a target of 180).
- **Player Composition:** The collective strength of the 11 players selected for the match.

The challenge, therefore, is to design a system that can process all these factors, weigh their relative importance, and synthesize them into a single, reliable, and explainable forecast.

C. Proposed Solution and Contributions

To address this problem, we propose the “T20I Predictor,” a hybrid system that synergistically combines a statistical analysis engine with a multi-model machine learning engine. The primary contribution of this work is the design and evaluation of this hybrid architecture. Unlike single-method approaches, our system generates two parallel streams of analysis: one based on human-interpretable statistical factors and another based on the consensus of seven distinct machine learning algorithms.

The novelty of our approach lies in the final synthesis stage, which is handled by a configurable weighted scoring system. This module assigns a priority score to each statistical factor, allowing the model to intelligently prioritize more impactful information (e.g., giving “Recent Head-to-Head Form” a higher weight than “All-Time Win Ratio”).

This creates a transparent "glass-box" model, where the final prediction is directly traceable to the underlying data, addressing the interpretability limitations of pure machine learning approaches. The entire system is implemented as a fully functional and interactive Django web application, making the complex analytical process accessible to a non-technical audience.

D. Paper Structure

The remainder of this paper is organized as follows. Section II details the system's architecture, including the data preprocessing pipeline, the design of the statistical and machine learning engines, and the logic of the weighted scoring system. Section III presents the performance evaluation of the system, discussing the accuracy of the individual machine learning models and the logical coherence of the final hybrid predictions. Finally, Section IV concludes the paper, summarizing the key findings, acknowledging the project's limitations, and outlining potential directions for future research.

II. SYSTEM ARCHITECTURE AND IMPLEMENTATION

The T20I Predictor is architected as a modular, three-tier web application to ensure a clear separation of concerns between data, logic, and presentation. This section details the design of each component, from initial data processing to the final synthesis of the prediction.

A. Data Preprocessing and Management

The foundation of the project is the "Men's T20I Cricket Complete Dataset" sourced from Kaggle. The initial and most critical phase of implementation was the development of a custom data import script which performs a comprehensive ETL (Extract, Transform, Load) process.

1) *Extraction*: The script begins by extracting data from the two source CSV files, `matchwise_data.csv` and `deliverywise_data.csv`, using the Pandas library.

2) *Transformation*: A series of crucial transformation steps are applied to clean and structure the raw data. First, the dataset is filtered to include only official T20I matches by validating that both competing teams are present in a predefined list of recognized international teams, thereby excluding all domestic and franchise league matches. Second, venue names are standardized using a mapping dictionary to merge common duplicates (e.g., "Shere Bangla National Stadium, Mirpur" is mapped to "Shere Bangla National Stadium"). Finally, player-to-team associations are established by analyzing the `deliverywise_data` to determine the national team each player has most frequently represented.

3) *Loading*: The transformed data is loaded into a lightweight SQLite3 database managed by Django's Object-Relational Mapper (ORM). The database schema, defined in `predictor/models.py`, consists of five primary models: Team, Player, Venue, Match, and Delivery. This relational structure provides an efficient and queryable foundation for all subsequent analytical tasks.

B. The Dual-Analysis Engine

The core of the application is its dual-analysis engine, which generates two parallel streams of predictive insight: one based on interpretable statistics and another on algorithmic pattern recognition.

1) *Statistical Analysis Engine*: This engine is implemented as a Python package (`predictor/analysis/`) containing distinct modules for different analytical contexts. It performs complex database queries to calculate a wide range of performance metrics. Key modules include `team_analysis.py` for evaluating head-to-head records and recent form, and `venue_analysis.py` for assessing team performance at specific grounds. A significant feature is the adaptive target score analysis within `innings_analysis.py`, which provides a context-aware evaluation based on whether a first innings total is categorized as very low (<90), low (90-130), average (131-180), high (181-220), or very high (≥ 220), ensuring a more realistic assessment for extreme scores.

2) *Machine Learning Engine*: To provide an independent forecast, a suite of seven machine learning models was implemented using the Scikit-learn library: Random Forest, K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Gradient Boosting, Logistic Regression, Decision Tree, and Gaussian Naive Bayes. A custom training script (`predictor/ml/trainer.py`) prepares the data by converting categorical features (teams, venue) into numerical format using label encoding. To address data imbalance and ensure all teams can be predicted, the encoders are fitted on the complete list of teams from the database, not just those present in the training set. The models are trained on an 80/20 train-test split of the historical data and serialized to disk using Joblib. For predictions, the models' raw probability outputs are intelligently normalized to ensure the forecast is always constrained to one of the two competing teams.

C. Weighted Prediction Synthesis

The final prediction is generated by the "Weighted Scorer" (`predictor/analysis/weighted_scorer.py`), a module that synthesizes the outputs of the statistical engine. It employs a dictionary of predefined weights, assigning a numerical importance to each statistical factor. This allows the system to prioritize more impactful information; for example, "Recent Head-to-Head" form is assigned a weight of 25, whereas the less critical "All-Time Win Ratio" is assigned a weight of 5. The system aggregates these weighted scores for each team and applies a smoothing algorithm, adding a baseline score to each team to prevent unrealistic 100% win probabilities. This results in a final, nuanced prediction that is directly traceable to the underlying statistical evidence.

III. PERFORMANCE EVALUATION

A. Experimental Setup

The development and performance evaluation of the T20I Predictor were conducted in a standardized local environment. The system was run on a Windows 11 operating system, with Python 3.12 serving as the core programming language. All

project dependencies, including Django, Pandas, and Scikit-learn, were managed within an isolated virtual environment (venv). The database backend was SQLite3, and the application was served locally using Django’s built-in development server. The evaluation procedure involved two phases: a quantitative assessment of the machine learning models’ accuracy and a qualitative assessment of the final hybrid prediction’s logical coherence.

B. Model Accuracy Evaluation

The quantitative evaluation focused on the predictive accuracy of the seven machine learning models. After the training process was completed on 80% of the dataset, the models were tasked with predicting the outcomes of the 20% hold-out test set. The accuracy was calculated as the percentage of correctly predicted winners. The results are summarized in Table I.

TABLE I
MACHINE LEARNING MODEL ACCURACY

Model	Accuracy on Test Set (%)
Support Vector Machine	68.83
Gradient Boosting	67.21
Random Forest	65.59
Logistic Regression	65.18
Decision Tree	63.97
K-Nearest Neighbors	62.35
Gaussian Naive Bayes	61.13

The models achieved a commendable level of predictive accuracy, with most falling within the 65-75% range. This result is significant as it aligns with academic findings on the inherent stochasticity of T20 cricket, where a perfect prediction is impossible. The variation in accuracy across different algorithms validates the project’s multi-model approach, as it provides a more reliable consensus forecast by mitigating the biases of any single model.

C. Qualitative System Evaluation

The qualitative evaluation assessed the logical coherence of the final hybrid prediction through a series of test cases. In scenarios involving historical matches, the statistical engine consistently identified the correct historical advantages. For example, in a simulated match between a top-ranked team and an underdog with a strong record at a specific venue, the statistical factors correctly reflected this nuance, and the weighted scorer appropriately adjusted the final prediction. The adaptive target score analysis performed as expected, correctly identifying very low scores (≤ 90) as a massive advantage for the chasing team and very high scores (≥ 220) as a near-certain win for the defending team, providing a common-sense fallback even when direct historical data was sparse. The final prediction, which combines the weighted statistical score with the ML consensus, was found to be a robust and transparent measure of a team’s chances, successfully meeting the project’s core design objectives.

IV. CONCLUSION

A. Summary of Contributions

This paper has detailed the design, implementation, and evaluation of the T20I Predictor, a hybrid system that successfully combines statistical analysis and machine learning to forecast cricket match outcomes. The project’s primary contribution is the development of a “glass-box” predictive model that provides not only a probabilistic forecast but also a transparent, data-driven narrative explaining the rationale behind it. The synergy between the human-interpretable statistical factors and the pattern-based machine learning models creates a more powerful and trustworthy prediction tool than either method could achieve in isolation. The implementation of the configurable weighted scoring system and the adaptive target score analysis represents a significant step towards creating a more intelligent and context-aware sports analytics tool.

B. Limitations and Future Work

Despite its robust design, the system has several limitations. Its accuracy is fundamentally dependent on the quality and completeness of the historical dataset. The system is also unable to account for real-time, qualitative variables that are crucial to a cricket match, such as player injuries, sudden changes in weather, or the specific condition of the pitch on match day. Future work could focus on several key enhancements. The feature engineering for the machine learning models could be expanded to incorporate more granular player statistics, such as performance against specific bowling types or in different phases of the game. Integrating a live data API would be a significant upgrade, allowing for real-time updates to player form and potentially enabling in-match win probability calculations. Finally, the frontend could be enhanced with more advanced data visualizations to better illustrate performance trends and statistical comparisons.

REFERENCES

REFERENCES

- [1] N. Muruganantha, “Men’s T20I cricket complete dataset (2005-2025),” *Kaggle*, 2024. [Online]. Available: <https://www.kaggle.com/datasets/nishanthmuruganantha/mens-t20i-cricket-complete-dataset>
- [2] D. R. Greenfeld and A. R. Greenfeld, *Two Scoops of Django 3.x: Best Practices for the Django Web Framework*. Two Scoops Press, 2020.
- [3] F. Pedregosa *et al.*, “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [4] A. C. Müller and S. Guido, *Introduction to Machine Learning with Python: A Guide for Data Scientists*. O’Reilly Media, 2016.
- [5] D. K. J. S. Sankalpa and K. D. K. Wanniarachchi, “A data analytic approach to predict the winner in a game of cricket,” in *2017 IEEE International Conference on Industrial and Information Systems (ICIIS)*, 2017, pp. 1–6.
- [6] W. McKinney, *Python for Data Analysis*, 2nd ed. O’Reilly Media, 2017.