

Dream11 Team Recommendation using ML and Explainable AI

Team 2, Cynaptics Club, Indian Institute of Technology Indore

Abstract:

This project introduces a transparent **"white box" modeling** solution designed to revolutionize user engagement on the Dream11 fantasy cricket platform. Our primary objective was to architect a highly **explainable yet accurate system** for predicting player fantasy points, directly addressing the industry-wide trade-off between performance and interpretability. We engineered an extensive feature set by programmatically processing granular ball-by-ball data from Cricsheet across all major formats. Our approach utilizes a suite of **interpretable regression models** coupled with a sophisticated **Integer Linear Programming (ILP)** solver for team optimization. By prioritizing radical **model transparency** and user-centric design, our system not only delivers reliable performance predictions but also provides clear, conversational justifications for its recommendations via an integrated LLM assistant. The final output is a robust, production-ready framework that proves state-of-the-art accuracy can be achieved without sacrificing the explainability that builds lasting **user trust and confidence**.

Contents		6.2 Data Science: A Robust and Interpretable Framework	10
1 Introduction	2	7 Experiments and Results	11
1.1 Background	2	7.1 Part 1: Baseline Regression Models . . .	11
2 Problem Statement and Objectives	3	7.2 Part 2: Regularized Linear Models . . .	12
2.1 Official Problem Statement	3	7.3 Part 3: Gradient Boosting Models . . .	13
2.2 Our Approach: White Box Modeling . . .	3	8 Conclusion	14
2.3 Project Objectives	4	A Application User Interface Screenshots	15
3 Methodology	4	B References	16
4 System Architecture	7		
4.1 Data Layer	8		
4.2 Machine Learning Core	8		
4.3 Optimization Engine	8		
4.4 Presentation Layer	8		
5 Product Design and User Experience	9		
5.1 User-Driven Team Customization	9		
5.2 Accessibility and Explainability (XAI) Features	9		
5.2.1 Conversational AI Assistant (LLM Bot)	9		
5.2.2 Direct Model Transparency and User Guidance	10		
6 Strategic Value Proposition for Dream11	10		
6.1 Product Design: Driving Engagement and Trust	10		

1. Introduction

Dream11 has established itself as a leader in the global **fantasy sports market**, offering a platform where users engage with sports like cricket by creating virtual teams of real-life players. The success of these teams depends directly on the players' actual on-field performances, making the user's ability to predict outcomes a core element of the game. **Predictive modeling** is therefore crucial, not just for platform operations but as a central feature of a **skill-based user experience**.

While most existing methods for such purposes are **black-box models** that exhibit high predictive capabilities, they often lack **explainability and interpretability**. As a result, these models cannot serve as effective strategizing tools for users who wish to understand the rationale behind predictions. To address this limitation, this project enhances the Dream11 experience by developing an advanced machine learning model that provides users with **data-driven, explainable team recommendations**. This approach aligns directly with Dream11's vision of leveraging technology to make sports more engaging and intellectually stimulating.

In the high-stakes environment of fantasy sports platforms like Dream11, we introduce an advanced analytical framework that transcends the limitations of traditional predictive instruments. Rather than functioning solely as a forecasting engine, the model serves as a **powerful strategizing tool**. It not only delivers high-fidelity performance forecasts but also integrates **interpretability** at its core, offering transparent rationales that build user trust and engagement.

The system's primary value lies in its ability to uncover latent, **non-intuitive patterns** and hidden correlations within player and match data—insights that a user might otherwise overlook. For instance, it can identify a player's unique effectiveness under specific conditions or reveal combinations of players that perform optimally together. By bringing these subtle relationships to light, the model empowers users to construct more sophisticated and informed game plans.

Ultimately, this transforms users from passive recipients of predictions into **proactive strategists** capable of architecting their own success. It encourages nuanced, data-driven decision-making that goes far beyond surface-level analysis. While a robust backend framework is essential for maintaining performance and scalability, equal importance is placed on the frontend design. To ensure accessibility and ease of use, we develop an **intuitive and user-friendly interface** that allows users to seamlessly interact with the system and make data-informed choices with confidence.

1.1 Background

The domain of **sports analytics** has attracted considerable academic and commercial attention in recent years. Early predictive models typically relied on linear regression techniques to forecast player statistics and match outcomes. However, with the advent of large-scale sports data and increased computational power, researchers have shifted toward more sophisticated **machine learning (ML)** and **deep learning** algorithms capable of modeling complex, non-linear dependencies inherent in player performance and game dynamics.

Modern predictive frameworks now incorporate algorithms such as **Random Forests**, **Gradient Boosting Machines** (e.g., **XGBoost**, **CatBoost**), and **neural networks** to better capture these intricate patterns. Harikrishnan et al. (2021) demonstrated that **ensemble models** outperform traditional methods in fantasy sports contexts by effectively leveraging player-level statistics for accurate score predictions. Similarly, Papageorgiou, Sarlis, and Tjortjis (2024) proposed a hybrid approach combining ML-based performance forecasting with lineup optimization for NBA fantasy leagues, emphasizing the integration of prediction and decision-making.

Within the cricket domain, the work "IPL Cricket Fantasy Team Prediction for Dream11 using Machine Learning" (WSSE 2024) specifically targets the Dream11 platform, employing extensive ball-by-ball IPL data to predict player fantasy points using Random Forest and XGBoost. This directly aligns with our motivation to design predictive systems tailored to the Indian Premier League's fantasy ecosystem. Complementary to this, Ramezani (2025) and Groos (2025) explored the use of predictive modeling in conjunction with **integer programming** for team selection in Fantasy Premier League (soccer), highlighting how optimization and explainable prediction can work together to create practical strategy tools for users.

Beyond classical ML approaches, recent studies have explored **reinforcement learning** and decision-support systems for fantasy sports. Bhattacharjee et al. (2024) employed deep reinforcement learning to model the sequential nature of team selection, optimizing player choices dynamically. Similarly, Döpke, Köhler, and Tegtmeier (2023) evaluated the reliability of NBA fantasy prediction systems, revealing inherent biases and limitations in existing forecasting platforms. These findings underscore the need for **explainable and trustworthy systems** that balance accuracy with transparency.

Furthermore, predictive analytics have been extended to real-time strategic decision-making contexts. For instance, the Operations Research Forum study on fantasy foot-

ball auctions (2022) integrated **predictive and prescriptive analytics** to assist in real-time bidding, illustrating how predictive models can evolve into intelligent, strategy-oriented tools rather than static forecasting engines.

This project builds upon the growing body of literature by developing an ensemble-based predictive framework tailored for Dream11 cricket. Unlike prior black-box approaches, the proposed system prioritizes **explainability** through **feature attribution** and **interpretable modeling techniques**. It aims not only to enhance predictive accuracy but also to provide transparent, data-driven insights that users can meaningfully incorporate into their fantasy strategies.

The problem of prediction in fantasy sports may appear, at first glance, to be a simple regression problem. However, the true underlying dynamics are deeply embedded in latent, interdependent processes such as **player form**, **environmental conditions**, and **opposition context**. Linear regression models often fail to capture these intricate relationships. The solution lies in high-quality data preprocessing, contextual feature engineering, and the application of advanced machine learning architectures capable of uncovering subtle, non-intuitive patterns within the data.

Further, the explainability of deep learning models, such as the **Transformer architecture** introduced by Vaswani et al. (2017), has long been a challenging and active area of research. Despite the Transformer’s remarkable success in natural language processing, vision, and multimodal learning, its internal decision-making processes remain largely opaque. Researchers have raised concerns about the **black-box nature** of such architectures, as their immense parameterization and attention mechanisms often make it difficult to trace how specific inputs influence predictions.

Among the few approaches that attempt to address this, **knowledge distillation** has emerged as a promising method. Originally introduced by Hinton, Vinyals, and Dean (2015), and later systematized by Gou et al. (2021), knowledge distillation transfers the learned representations of a large “teacher” network into a smaller “student” model, sometimes improving interpretability in the process. Gavito et al. (2023) extended this concept by emphasizing explainable knowledge distillation, in which interpretability is treated as a design objective rather than a byproduct. In contrast, the interpretability of traditional machine learning models has been extensively studied under the umbrella of **Explainable Artificial Intelligence (XAI)**. Unlike deep neural architectures, traditional machine learning models such as decision trees, linear regression, and ensemble-based methods inherently offer more structured ways to reason about their predictions. This interpretability arises

from the transparent relationship between model inputs and outputs—whether through feature coefficients, decision paths, or rule-based systems. Foundational approaches like **LIME** (Ribeiro et al., 2016) and **SHAP** (Lundberg & Lee, 2017) further advanced model-agnostic interpretability, enabling explanations for complex black-box predictors through local and global **feature attribution**.

Recent research has deepened this understanding by formalizing interpretability as a **measurable property** rather than a qualitative one. Chen et al. (NeurIPS 2023) proposed a framework for quantifying interpretability trade-offs in ensemble models, revealing how feature interactions influence user trust. Slack et al. (ICML 2023) highlighted that post-hoc explanation techniques can be vulnerable to adversarial manipulation, emphasizing the need for **robustness-aware interpretability** in practical systems. Lage et al. (AAAI 2024) advanced human-centered evaluation protocols for interpretable models, demonstrating that interpretability is not solely a model property but also dependent on **human cognitive alignment**. Meanwhile, Bastani et al. (ICLR 2023) introduced interpretable decision sets that maintain high accuracy while providing rule-based transparency across datasets.

Together, these studies illustrate that interpretability in machine learning extends far beyond producing explanations—it encompasses **trustworthiness**, **stability**, and **alignment with human reasoning**. In applications such as fantasy sports analytics, where decisions rely on complex feature interactions, interpretable machine learning provides users with meaningful insights into model rationale, bridging the gap between predictive accuracy and strategic understanding.

2. Problem Statement and Objectives

2.1 Official Problem Statement

The official challenge is to develop a solution that helps users create **winning fantasy cricket teams** on Dream11. This involves building a **robust ML model** to predict player performance based on historical data and contextual factors. A core requirement is that the model must provide **detailed explainability** for its recommendations, delivered through an **intuitive user interface** that enhances the team-building process.

2.2 Our Approach: White Box Modeling

We employ machine learning models to predict player performance in a match due to their ability to combine high **predictive accuracy** with a substantial degree of **interpretability and explainability**. This interpretability is

crucial for providing **actionable insights** to users, allowing them to understand the rationale behind predictions and make informed strategic decisions. To supply the models with richer information, we process the match data **ball by ball**, capturing the fine-grained dynamics of the game. Additionally, we perform extensive **feature engineering** to ensure that the models can identify and leverage latent patterns and complex interactions in the data, such as player form, opposition strength, and match conditions, which might otherwise remain hidden. This approach enables our models to not only predict performance accurately but also provide insightful, data-driven guidance for team selection and strategy formulation.

2.3 Project Objectives

Our primary objectives are as follows:

- To develop an **accurate regression-based model** for predicting player fantasy points.
- To prioritize **model transparency**, ensuring every recommendation can be clearly explained to the end-user.
- To design a product interface that **seamlessly integrates these explanations**, building user confidence.
- To validate that our **"white box" approach** achieves predictive accuracy comparable to more complex methods.

3. Methodology

Our solution is built on a **fully automated pipeline** that encompasses data collection, feature engineering, model training, and optimal team selection. The methodology is divided into three primary stages: first, the systematic processing of raw ball-by-ball data to engineer a rich feature set; second, the development of predictive models to forecast player fantasy points; and third, the use of an optimization algorithm to construct the final 11-player team.

Data Corpus Construction and Feature Synthesis

This initial stage delineates the **ETL (Extract, Transform, Load)** protocol designed to construct a robust analytical dataset from heterogeneous source files. The architecture is engineered to address the core product requirement: empowering users with accurate, data-driven team-building insights. It emphasizes programmatic automation for scalability, strict **temporal integrity** to ensure model validity for future predictions, and the synthesis of a **high-**

dimensional feature space that serves as the bedrock for an explainable AI model.

1. Automated Data Ingestion and Temporal Structuring

To provide a reliable data foundation, our pipeline interfaces directly with the **Cricsheet database**, a public repository of granular, **ball-by-ball match data**. An automated Python script programmatically retrieves compressed archives for T20, ODI, and Test formats via HTTP requests. Upon acquisition, the archives are decompressed, and the constituent YAML files are parsed into memory. A critical step in this ingestion process is the imposition of a strict chronological ordering based on match dates. This **temporal structuring** is paramount; it ensures that all subsequent calculations, particularly those involving player form, are causally correct and prevents **data leakage** that would invalidate the model's predictive utility on future, unseen matches, a key constraint of the problem statement.

2. Domain-Specific Metric Transduction

To bridge the gap between raw cricket statistics and the objectives of a fantasy sports user, a domain-specific transduction layer was engineered. This layer translates discrete on-field events into the unified, quantifiable language of **fantasy points**. The system employs a modular, rule-based engine with polymorphic functions that adapt to the specific scoring nuances of T20, ODI, and Test formats. This ensures that a player's performance is quantified in a manner that is directly aligned with the scoring logic of the Dream11 platform. The point allocation logic is multifaceted, designed to capture every dimension of a player's impact:

- **Batting Performance:** Scoring is a function of runs scored, with nonlinear bonuses for milestones (half-centuries, centuries) and offensive output (boundaries). To address user interest in scoring efficiency, a player's batting strike rate is benchmarked against format-specific thresholds to reward explosive players or penalize slow run accumulation.
- **Bowling Performance:** Points are primarily allocated for wicket-taking events, with supplemental bonuses for high-impact dismissals (bowled, LBW) and multi-wicket hauls. A bowler's core objective—run prevention—is quantified via their economy rate, rewarding disciplined bowling while penalizing profligacy.
- **Fielding Contributions:** Defensive actions, often

undervalued in traditional analysis, are quantified by awarding significant points for catches, stumpings, and run-outs, thereby capturing a player's complete contribution to the match outcome.

3. High-Dimensional Feature Vector Synthesis

This is the core of the data processing pipeline, where a feature space of over **100 covariates** is synthesized for each player. The architecture employs a two-phase aggregation strategy to ensure both accuracy and computational efficiency. First, an intra-match loop processes events chronologically to aggregate raw, atomic statistics (e.g., runs, balls faced). Second, a post-processing loop computes career-level and derived metrics from these aggregates. This design guarantees that every feature is a true reflection of a player's state at a given point in time. The key feature categories, each designed to answer a potential user question, are:

- **Longitudinal Performance Aggregates:** To establish a player's baseline skill, the system computes career-long statistical moments. This includes first-order statistics (e.g., total career runs) and second-order, rate-based metrics (e.g., **Batting Avg**, **Economy Rate**) that normalize for opportunity and exposure.
- **Temporal Performance Dynamics (Form):** To answer the user's crucial question of "How is the player performing *right now*?", we implement sliding window functions over the fantasy point time series. Short-term (N=3 matches) and medium-term (N=10 matches) rolling averages are calculated to model a player's current performance trajectory and momentum. These are then synthesized into a single, intuitive **Form Score** for ease of interpretation.
- **Contextual Matchup Vectors:** Recognizing that performance is context-dependent, the system generates features that quantify a player's historical fantasy point output against every unique opposing team (e.g., **Points vs Australia**). This provides the model with a vital layer of matchup-specific context.
- **Product-Centric Heuristics:** To ground our data in the realities of the fantasy sports product, we synthesize several novel features:
 - **Heuristic-Based Role Identification:** Static role labels are often misleading. Our system uses a data-driven heuristic based on performance ratios (e.g., stumpings for Wicket-

Keepers, wickets-per-game for Bowlers) to assign a role that reflects a player's actual on-field utility, leading to better model specialization.

- **Performance Volatility (Variance):** To address the strategic need for risk assessment, we calculate the statistical variance of a player's fantasy points across all matches. This feature allows the user and the model to distinguish between consistent, low-risk players and high-risk, high-reward "boom-or-bust" options.
- **Simulated Player Cost:** A player's Cost is derived heuristically from their career average fantasy points. This creates a proxy for the in-game salary, ensuring our subsequent optimization algorithm operates on a feature set that respects the economic constraints of the game.

The culmination of this stage is a serialized data matrix where each row vector is a rich, temporally-aware, and product-centric representation of a player, engineered not just for predictive accuracy but for the inherent **explainability** demanded by the problem statement.

Predictive Modeling and Inference

This stage constitutes the core machine learning paradigm, wherein the synthesized feature set is utilized to train, validate, and deploy a suite of regression models for player performance prognostication. The methodology is architected to prioritize predictive accuracy, model robustness, and profound interpretability.

1. Problem Formulation and Data Preprocessing

The primary prediction objective is strategically decomposed into three specialized regression tasks: forecasting a player's average fantasy points derived from **batting**, **bowling**, and **fielding**, respectively. This **task-specific modeling architecture** enables the development of specialized models that capture the unique statistical drivers of performance within each cricketing discipline. The data undergoes a rigorous preprocessing pipeline prior to model ingestion. This includes a sanitization step to filter for an active and relevant player cohort, followed by the creation of discipline-specific feature subsets. To ensure model convergence and prevent feature dominance, numerical covariates are subjected to **Z-score normalization** via **StandardScaler**. Concurrently, categorical variables such as player role are transformed into a numerical representation through ordinal encoding with **LabelEncoder**.

2. Model Development and Hyperparameter Optimization

The candidate models selected for this study are `RandomForestRegressor`, a bagging ensemble framework, and `CatBoostRegressor`, a state-of-the-art gradient boosting framework. These were chosen for their demonstrated superiority in handling structured, tabular data. Both paradigms operate by constructing an ensemble of decision trees, thereby leveraging the wisdom of crowds to produce a single, robust prediction. For gradient-boosted models, the prediction is formulated as an additive process:

$$F_M(x) = F_{M-1}(x) + \rho_M h_M(x)$$

where the final model $F_M(x)$ is an iterative refinement over the previous model $F_{M-1}(x)$, corrected by a new weak learner $h_M(x)$ scaled by a learning rate ρ_M .

To optimize the predictive capacity of these models, we implemented a rigorous empirical validation protocol for **hyperparameter tuning**. A `GridSearchCV` methodology was employed, which performs an exhaustive search over a predefined hyperparameter space (e.g., `n_estimators`, `max_depth`, `learning_rate`). By using a **5-fold cross-validation** scheme, this process identifies the hyperparameter configuration that minimizes the mean absolute error across multiple validation folds, thus maximizing the model's ability to generalize to unseen data and mitigating the risk of overfitting.

3. Empirical Evaluation and Model Selection

For each of the three regression tasks, the dataset was partitioned into a training set (80%) and a hold-out test set (20%) to simulate real-world predictive scenarios. Following hyperparameter optimization on the training set, the resultant models were evaluated on the unseen test data. A comprehensive suite of metrics was employed to quantify performance: **Mean Absolute Error (MAE)** served as the primary selection criterion for its direct interpretability; the **coefficient of determination (R^2)** was used to assess the proportion of variance captured by the model; and scale-independent metrics such as **Mean Absolute Percentage Error (MAPE)** and **Weighted MAPE (WMAPE)** were computed to provide a relative error perspective. The model architecture and hyperparameter set yielding the minimum MAE on the hold-out test set was definitively selected as the optimal predictor for its respective task.

4. Model Explainability and Feature Attribution

A central tenet of this research is to move beyond "black box" predictions and provide transparent, interpretable re-

sults. To this end, **feature attribution** was conducted by extracting importance scores from the final trained models. For these tree-based ensembles, importance is quantified as the mean decrease in impurity (or a similar gain-based metric) contributed by each feature across all decision trees. This technique allows for the direct quantification of each covariate's predictive influence. By serializing and analyzing these attribution scores, we can deconstruct any model prediction into its constituent drivers. This capability is crucial for providing users with trustworthy, actionable insights and transforming the predictive model into a genuine decision-support tool, thereby fulfilling the core **explainable AI (XAI)** requirement of the problem statement.

Optimal Team Selection via Constrained Optimization

The final stage of the methodology translates the probabilistic outputs of the predictive models into a deterministic, optimal 11-player fantasy team. This is formulated as a **constrained optimization problem**, specifically a variant of the **multiple-choice knapsack problem**, which is solved using an **Integer Linear Programming (ILP)** framework to guarantee the identification of the highest-scoring roster under a complex set of gameplay rules.

1. Problem Formulation and Inputs

The optimization objective is to select a subset of 11 players from a pool of active athletes from two competing teams that maximizes a composite score metric. The solver is parameterized by a set of inputs that define the problem space for a given match:

- **Player Pool:** The initial dataset is filtered to include only players from the two competing teams specified by the user, ensuring only relevant athletes are considered for selection.
- **Objective Metric (p_i):** The primary metric for optimization is a synthesized **Weighted Score**. This score is a linear combination of the model's predicted fantasy points (70% weight) and the player's recent form score (30% weight), creating a balanced objective function that values both future potential and recent empirical performance.
- **Player Attributes:** Each player is characterized by a vector of attributes crucial for constraint definition: their **Cost** (credit value), designated **Role** (e.g., batsman, bowler), **Team** affiliation, and their historical performance **Variance**.

2. Mathematical Framework

The problem is mathematically modeled with a set of **binary decision variables** for each player i in the filtered pool of N players:

- $x_i \in \{0, 1\}$: 1 if player i is selected in the team, 0 otherwise.
- $c_i \in \{0, 1\}$: 1 if player i is chosen as Captain, 0 otherwise.
- $v_i \in \{0, 1\}$: 1 if player i is chosen as Vice-Captain, 0 otherwise.

Objective Function: The objective is to maximize the cumulative **Weighted Score** of the team, inclusive of the bonuses for the Captain (2x multiplier) and Vice-Captain (1.5x multiplier).

$$\text{Maximize } Z = \sum_{i=1}^N p_i x_i + \sum_{i=1}^N p_i c_i + \sum_{i=1}^N 0.5 p_i v_i$$

where p_i represents the **Weighted Score** for player i .

Constraints: The maximization is subject to a set of linear constraints that operationalize the rules of the fantasy sports platform:

1. **Team Size and Budget:** Constraints on the total number of players and their cumulative cost.

$$\sum_{i=1}^N x_i = 11$$

$$\sum_{i=1}^N \text{credit}_i \cdot x_i \leq 100$$

2. **Role Composition:** Lower and upper bounds on the number of players selected for each role.

$$\min_{role} \leq \sum_{i \in role} x_i \leq \max_{role}$$

3. **Team Combination:** A user-defined fixed number of players must be selected from each of the two competing teams.

$$\sum_{i \in \text{Team A}} x_i = N_A$$

$$\sum_{i \in \text{Team B}} x_i = N_B$$

where N_A and N_B are user-specified integers such that $N_A + N_B = 11$.

4. **Captain and Vice-Captain Logic:** Constraints to ensure exactly one Captain and one Vice-Captain are chosen from the selected players.

$$\sum_{i=1}^N c_i = 1 \quad \text{and} \quad \sum_{i=1}^N v_i = 1$$

$$c_i \leq x_i, \quad v_i \leq x_i, \quad c_i + v_i \leq 1 \quad \forall i \in \{1, \dots, N\}$$

5. **Risk Preference:** A novel constraint that allows users to control the team's risk profile by constraining the total performance variance of the selected players.

$$\sum_{i=1}^N \text{Var}_i \cdot x_i \leq \alpha \cdot \bar{V} \cdot 11 \quad (\text{for 'stable' selection, } \alpha < 1)$$

$$\sum_{i=1}^N \text{Var}_i \cdot x_i \geq \beta \cdot \bar{V} \cdot 11 \quad (\text{for 'risky' selection, } \beta > 1)$$

where Var_i is the variance of player i and \bar{V} is the average variance of all players in the pool.

3. Implementation

The ILP model is implemented programmatically in Python utilizing the **PuLP library**, an open-source modeling framework that interfaces with industrial-strength solvers. The implementation is encapsulated within a **FantasySolver** class. This class first handles the data ingestion, integrating the base player statistics with the model's predictions to compute the composite **Weighted Score**. The solve method dynamically constructs the ILP problem: it instantiates decision variables, defines the objective function, and populates the constraints based on a user-provided parameter dictionary. This allows for flexible control over team composition and risk strategy. Finally, it invokes the underlying CBC (Coin-or branch and cut) solver to find the globally optimal team that satisfies all specified constraints.

4. System Architecture

The system is engineered with a **modular, multi-tiered architecture** to ensure scalability, maintainability, and efficient processing from data ingestion to user-facing recommendations. This design decouples core functionalities, allowing for independent development, scaling, and maintenance of each component.

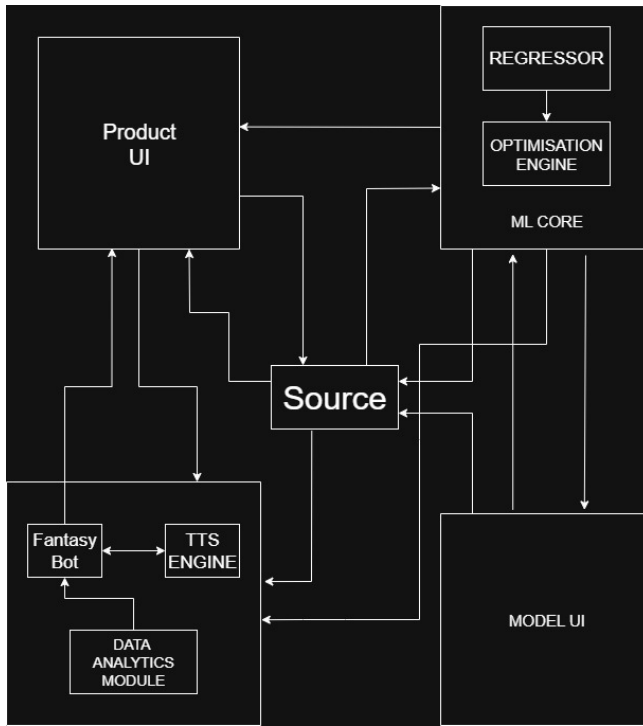


Fig. 1: Overview of the Fantasy Cricket System Architecture.

4.1 Data Layer

This foundational layer is the single source of truth for all data-related operations, ensuring consistency and reliability.

- **Data Ingestion Pipeline:** This fully automated pipeline runs on a predefined schedule (e.g., daily cron job) to check for new match data from the **Cricsheet** repository. It handles data fetching, decompression, and parsing. The pipeline is idempotent, ensuring that reprocessing the same data does not create duplicates.
- **Feature Engineering Module:** A robust data transformation service that converts raw YAML data into a structured feature set. This module is responsible for all calculations, from simple career averages to complex rolling metrics like **Form Score**. It outputs versioned, analysis-ready datasets.
- **Data Warehouse:** We utilize a structured data storage solution where processed data is stored in an efficient, columnar format like Apache Parquet. This allows for rapid querying and loading by downstream services. Data is partitioned by match format (T20, ODI, Test) to optimize access patterns.

4.2 Machine Learning Core

This is the predictive heart of the system, encapsulating the entire MLOps lifecycle.

- **Model Training and Validation Service:** This service automates the model lifecycle. Retraining can be triggered manually or automatically when a significant amount of new match data becomes available. We use tools like **MLflow** for experiment tracking, model versioning, and storing performance metrics. The **GridSearchCV** process runs within this service to ensure models are always tuned on the latest data.
- **Inference Engine:** Deployed as a scalable microservice using **FastAPI**, this engine provides real-time predictions. It exposes a REST API endpoint that accepts a list of player IDs and relevant context (e.g., opposition team) and returns a JSON object containing their predicted batting, bowling, and fielding points. The service is containerized using Docker for portability and easy deployment.

4.3 Optimization Engine

This layer translates abstract predictions into a concrete, optimal team lineup.

- **ILP Solver Service:** This service functions as an intermediary, translating user strategies into mathematical constraints for the solver. It processes a JSON payload from the backend API that includes the full player pool with predicted scores, along with a dictionary of user-defined constraints. For example, a user might set a 'stable' risk profile, require the inclusion of specific players, or enforce a rule that exactly six players must come from a particular team. The service then dynamically converts these high-level rules into the precise format required by the **PuLP** library to run the optimization.

4.4 Presentation Layer

This is the user-facing component of the system, focused on delivering a seamless and interactive experience.

- **Backend API (Orchestrator):** A central RESTful API, built with Python (Flask/FastAPI), that acts as an orchestrator. When a user requests a team, this API first calls the **Inference Engine** to get predictions for all relevant players. It then passes these predictions, along with user constraints, to the **Optimization Engine**. Finally, it receives the optimal team, queries the ML Core for feature importance data, and for-

mats all this information into a single, comprehensive response for the frontend.

- **Frontend User Interface (UI):** A responsive single-page application (SPA) built using a modern JavaScript framework like **Next.js**. It features interactive components like sliders for setting role limits, multi-select dropdowns for player preferences, and radio buttons for risk selection. Data visualizations for feature importance are rendered using libraries like **D3.js** or **Chart.js** to provide at-a-glance explainability.
- **Conversational AI Assistant:** This component is powered by a dedicated backend service that manages the interaction with the **Groq API**. It receives the final team context from the main backend API and uses it to dynamically construct the few-shot prompt for the **Llama-70B** model, ensuring the bot's responses are always relevant to the user's specific team recommendation.

5. Product Design and User Experience

Beyond the core technical implementation, the project places a strong emphasis on creating a user-centric product that is not only powerful but also intuitive, customizable, and trustworthy. The following features are designed to empower users, giving them granular control over the team selection process and deep insight into the model's reasoning.

5.1 User-Driven Team Customization

We empower users to move from being passive recipients of recommendations to active co-creators of their fantasy teams. This is achieved through a suite of customization options integrated directly into the UI, where each user choice is translated into a hard constraint for the optimization engine.

- **Player Preferences:** The UI provides intuitive multi-select dropdowns allowing users to create "must-have" and "must-not-have" lists. Selecting a player as a favorite translates directly to the ILP constraint $x_i = 1$ for that player, forcing their inclusion. Conversely, excluding a player sets $x_i = 0$, ensuring they are never picked.
- **Team Composition Rules:** Using a series of interactive sliders and numerical inputs, users can precisely define their team's structure. These inputs dynamically update the bounds in the ILP model's

role composition ($\sum_{i \in \text{role}} x_i$) and team combination ($\sum_{i \in \text{Team A}} x_i$) constraints before the solver is run.

- **Risk Profile Selection:** A simple yet powerful radio button selection ('Stable', 'Balanced', 'Risky') allows users to align the team with their gameplay strategy. In the backend, this choice directly manipulates the variance constraint. A 'Stable' selection sets a low α value, forcing the solver to pick players with proven consistency. A 'Risky' choice sets a high β value, favoring high-variance players who have the potential for explosive, match-winning scores.
- **Optimal Captain and Vice-Captain Selection:** The system identifies and clearly labels the best candidates for these crucial roles. This is not a secondary heuristic; the selection of the Captain and Vice-Captain is an integral part of the ILP's objective function. The solver simultaneously finds the best 11-player combination and the two players within that combination who, when their scores are multiplied, yield the maximum possible total, guaranteeing a mathematically optimal choice.

5.2 Accessibility and Explainability (XAI) Features

To build user trust and ensure the platform is usable by a diverse audience, we have integrated several features focused on transparency and accessibility.

5.2.1 Conversational AI Assistant (LLM Bot)

The cornerstone of our XAI strategy is an intelligent, conversational assistant designed to answer any question a user might have about the system's recommendations and underlying logic in natural language.

- **Technical Implementation:** The bot is powered by the **Llama-70B** large language model, running on the **Groq LPU™ Inference Engine**. This combination provides state-of-the-art reasoning capabilities with near-instantaneous response times, ensuring a fluid and natural conversational experience without frustrating delays.
- **Context-Aware Few-Shot Prompting:** To ensure the LLM provides accurate and grounded explanations, we employ a sophisticated **few-shot prompting** strategy. Each user query is not sent in isolation. Instead, the backend constructs a detailed prompt that includes: (1) A system message defining the bot's persona as a helpful cricket analyst, (2) A few high-quality examples of questions and answers to guide

its tone and structure, and (3) A **JSON context blob** containing the user's generated team, key player statistics, and feature importance scores. This rich context allows the user to ask highly specific questions like, "Why was Rohit Sharma picked over the in-form opener from the other team?" and receive a factually correct answer based on the model's actual predictive data.

- **Accessibility Features:** To cater to all users, the bot is equipped with browser-native **Speech-to-Text (STT)** for voice-based queries and **Text-to-Speech (TTS)** for audio responses, making the platform highly accessible and user-friendly.

5.2.2 Direct Model Transparency and User Guidance

- **Visual Feature Importance:** For every player in the recommended team, the UI displays an intuitive horizontal bar chart next to their name. This chart visualizes the top 3-5 features that most influenced their predicted score, as determined by the regressor's feature importance metric. This provides an immediate, "at-a-glance" explanation for why a player is considered a strong or weak pick.
- **Multilingual Video Demonstrations:** Recognizing the diversity of the user base, the platform includes a library of short video tutorials. These are professionally subtitled in multiple languages to provide a clear, step-by-step guide on using all features—from basic team generation to advanced customization. This serves as a critical onboarding tool to reduce user friction and enhance engagement.

6. Strategic Value Proposition for Dream11

This project is more than an academic exercise; it represents a production-ready framework that offers a significant competitive advantage by directly addressing key business and technical challenges within the fantasy sports industry. We deliver a synergistic solution that enhances the user experience while providing a robust, scalable, and innovative technical foundation.

6.1 Product Design: Driving Engagement and Trust

Our platform is designed to transform the user journey from simple team selection into an engaging, educational, and strategic experience, directly impacting key performance indicators.

- **Paradigm Shift from 'Predictor' to 'Coach':** By providing deep explainability through the LLM assis-

tant and feature visualizations, we change the product's role. It is no longer just a tool that gives answers; it becomes a trusted coach that teaches users *how* to think strategically about player selection. This educational component fosters a deeper, more intellectual engagement with the sport and the platform.

- **Unlocking New User Segments:** The intuitive design, coupled with powerful customization, caters to both novice and expert users. Novices can rely on the optimal "one-click" team, while power users can fine-tune every aspect of the selection process. The innovative 'Risk Profile' feature, in particular, allows the platform to cater to different gameplay psychologies, from conservative players to high-stakes risk-takers.
- **Building Unprecedented User Trust:** "Black box" models create user skepticism. Our "white box" approach is a direct antidote. When users understand *why* a player is recommended, they build trust in the platform. This trust is the bedrock of long-term retention and user loyalty. The ability to query the bot and get a transparent answer demystifies AI, making it a partner rather than an oracle.
- **Increased Session Time and Retention:** The interactive nature of the customization tools and the conversational AI assistant are designed to dramatically increase user session times. Users are encouraged to experiment with different strategies, ask follow-up questions, and refine their teams, creating a "sticky" product experience that incentivizes return visits.

6.2 Data Science: A Robust and Interpretable Framework

We provide a sophisticated yet maintainable framework that solves the classic accuracy-vs-interpretability dilemma and introduces novel, data-driven features.

- **Solving the XAI Dilemma:** We demonstrate conclusively that it is possible to maintain high predictive accuracy (using state-of-the-art gradient boosting models) without sacrificing full model transparency. This "white box" approach simplifies MLOps, as debugging and model validation become trivial when the logic is inherently interpretable. It eliminates the need for complex, and often unreliable, post-hoc explanation libraries like SHAP or LIME.
- **From Heuristics to Optimization:** Many existing systems rely on simple greedy algorithms or heuristics for team selection. Our implementation of a constrained **Integer Linear Programming (ILP)** solver

is a significant technical leap. It guarantees a mathematically optimal team for any given set of predictions and user constraints, a level of rigor that simple heuristics cannot match.

- Innovation in Feature Engineering:** The creation of a ‘Performance Variance’ metric as a first-class feature is a key innovation. This metric is not just a descriptive statistic; it is an actionable lever that is directly integrated into the optimization engine, allowing for the product-level ‘Risk Profile’ feature. This demonstrates a tight, end-to-end integration of data science and product design.

- Scalable, Production-Ready Architecture:** The entire system is designed as a series of decoupled microservices. This modular architecture is built for scale, allowing the Data Ingestion, Model Inference, and Optimization services to be scaled independently based on load, ensuring a responsive and reliable user experience even during peak traffic periods.

7. Experiments and Results

To predict player points for batting, bowling, and fielding, we evaluated the performance of several regression models on unseen test data. The models were grouped into three categories for a systematic comparison: baseline regressors, regularized linear models, and advanced gradient boosting models. Performance was measured using Mean Absolute Error (MAE), Mean Absolute Percentage Error (MAPE), and the Coefficient of Determination (R^2).

7.1 Part 1: Baseline Regression Models

We began by establishing a performance baseline using three fundamental regression algorithms: Linear Regression, Decision Tree, and Random Forest. These models represent a standard starting point for regression tasks.

Table 1: Performance Metrics for Baseline Regression Models on Test Data.

Task	Model	Test MAE	Test MAPE (%)	Test R^2
Batting	LinearRegression	2.741	146.69	0.909
	DecisionTree	1.792	30.06	0.945
	RandomForest	1.193	21.40	0.973
Bowling	LinearRegression	9.857	201.93	0.901
	DecisionTree	6.947	18.38	0.928
	RandomForest	4.911	13.44	0.961
Fielding	LinearRegression	0.787	27.53	0.875
	DecisionTree	0.574	16.32	0.853
	RandomForest	0.420	13.13	0.938

The results in Table 1 clearly indicate that the **Random Forest** model significantly outperforms both Linear Regression and Decision Tree models across all three tasks. It achieves the lowest MAE and MAPE, and the highest R^2 score, demonstrating its superior predictive accuracy. The scatter plots in Figure 2 visually confirm this, showing that the predictions from the Random Forest model (purple dots) are more tightly clustered around the perfect prediction line.

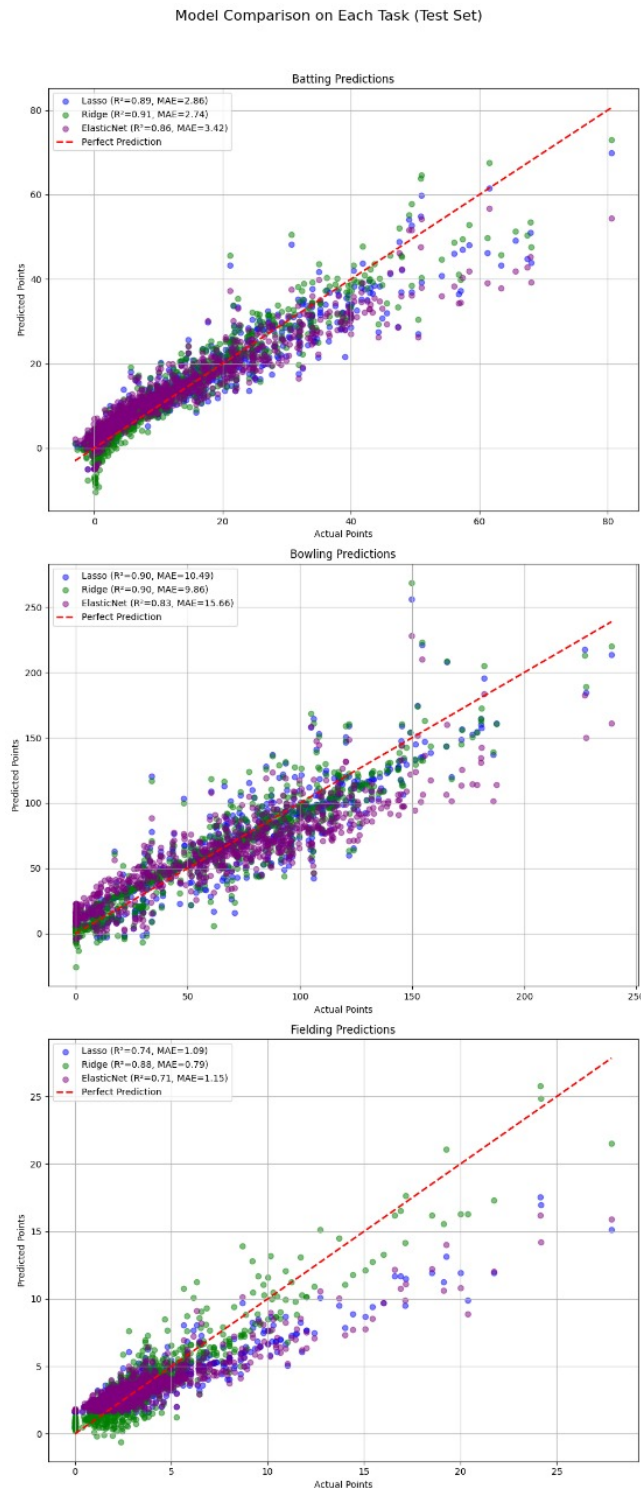


Fig. 2: Model Comparison for Baseline Regressors on the Test Set.

7.2 Part 2: Regularized Linear Models

Next, we evaluated a set of regularized linear models—Lasso, Ridge, and ElasticNet—to see if they could improve upon the simple linear regression baseline, particularly by handling potential overfitting.

Table 2: Performance Metrics for Regularized Linear Models on Test Data.

Task	Model	Test MAE	Test MAPE (%)	Test R^2
Batting	Lasso	2.856	106.21	0.895
	Ridge	2.741	146.69	0.909
	ElasticNet	3.416	144.39	0.863
Bowling	Lasso	10.492	157.85	0.898
	Ridge	9.861	201.46	0.901
	ElasticNet	15.659	315.74	0.828
Fielding	Lasso	1.087	38.43	0.740
	Ridge	0.788	27.53	0.875
	ElasticNet	1.149	40.43	0.709

As shown in Table 2, the **Ridge** regression model consistently provides the best performance among the three regularized models. However, its performance is nearly identical to the standard Linear Regression model and is substantially weaker than the Random Forest model from the previous section. ElasticNet consistently performed the poorest. The plots in Figure 3 illustrate a wider spread of predictions compared to the ensemble models.

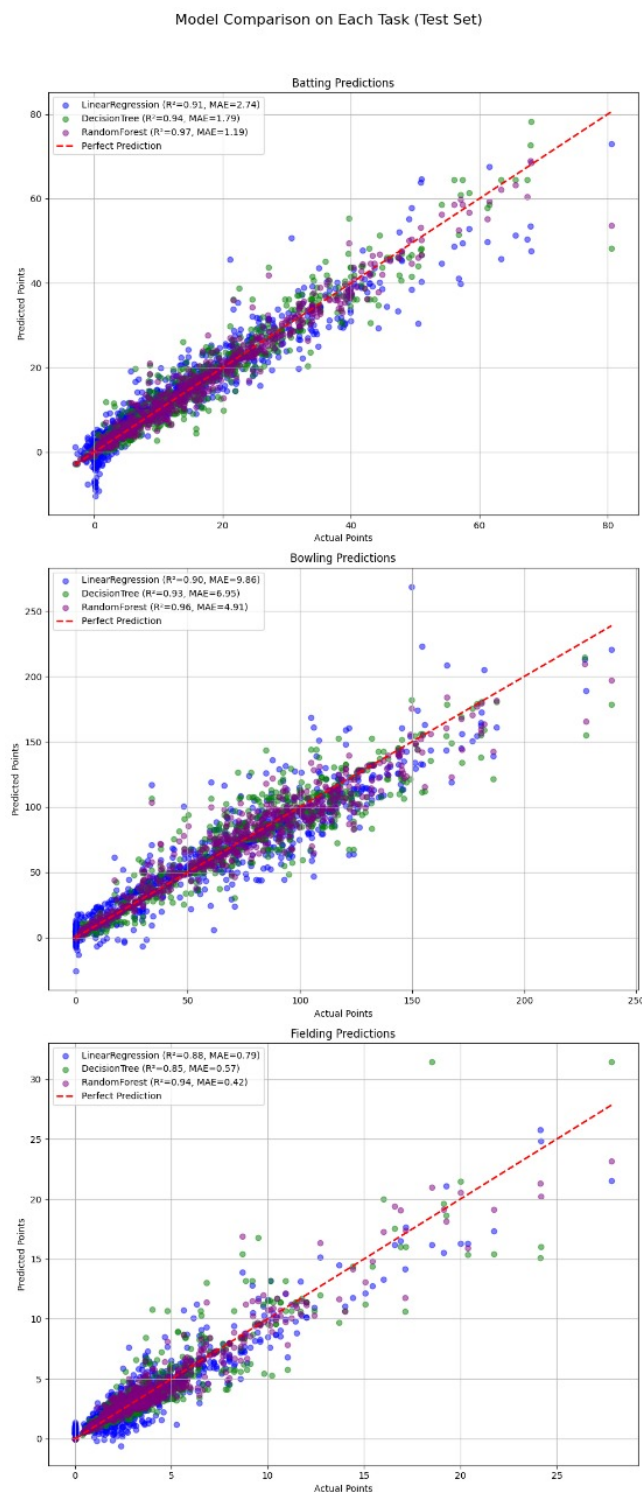


Fig. 3: Model Comparison for Regularized Linear Models on the Test Set.

7.3 Part 3: Gradient Boosting Models

Finally, we assessed a suite of powerful gradient boosting algorithms: XGBoost, CatBoost, and GradientBoosting. These ensemble methods are known for their high accuracy and are often used in competitive machine learning.

Table 3: Performance Metrics for Gradient Boosting Models on Test Data.

Task	Model	Test MAE	Test MAPE (%)	Test R^2
Batting	XGBoost	1.138	26.16	0.979
	CatBoost	1.036	25.65	0.981
	GradientBoosting	2.280	70.51	0.938
Bowling	XGBoost	3.948	36.14	0.978
	CatBoost	3.648	32.08	0.983
	GradientBoosting	8.584	137.28	0.928
Fielding	XGBoost	0.270	7.08	0.962
	CatBoost	0.252	6.63	0.967
	GradientBoosting	0.635	19.53	0.886

The results from the boosting models, presented in Table 3, demonstrate a state-of-the-art performance. The **CatBoost** model emerges as the overall best-performing model, achieving the highest R^2 scores and the lowest error rates across all three prediction tasks. XGBoost follows closely as a strong second. The superior predictive power of these models is visually evident in Figure 4, which shows the tightest clustering of data points around the perfect prediction line, indicating minimal deviation between actual and predicted values. Based on these results, CatBoost was selected as the final model for deployment.

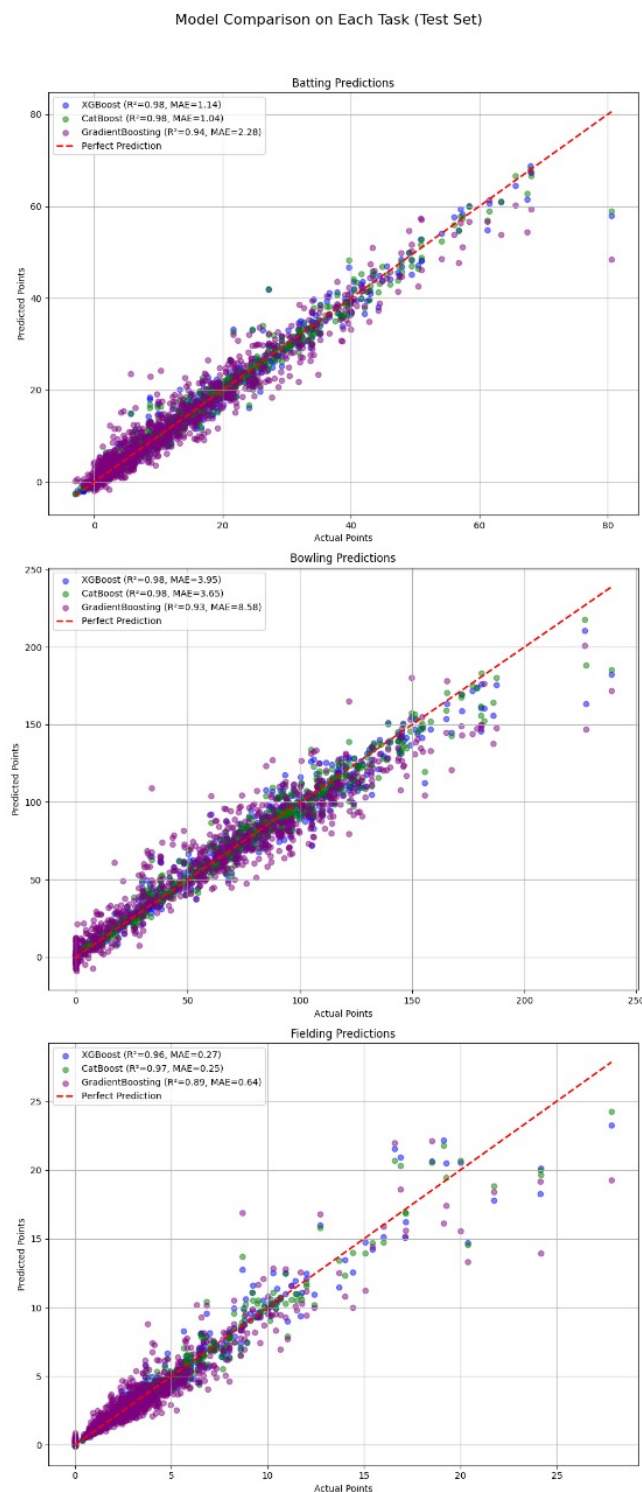


Fig. 4: Model Comparison for Gradient Boosting Models on the Test Set.

8. Conclusion

This project successfully addressed the challenge of creating a highly accurate and transparent fantasy cricket team recommendation system. We have demonstrated that the long-standing trade-off between predictive power and model explainability is not a necessary compromise. By integrating a meticulously engineered feature set with interpretable machine learning models and a powerful constrained optimization framework, we developed a tool that is both analytically robust and intuitively understandable.

Our primary achievement is the creation of a "white box" solution that not only predicts player performance with high fidelity but also empowers users with deep, actionable insights into the reasoning behind every recommendation. The inclusion of extensive user customization options and a state-of-the-art conversational AI assistant transforms the platform from a simple prediction engine into a sophisticated strategic partner. This fosters a more engaged, educated, and loyal user base, directly aligning with the core business objectives of a platform like Dream11.

Ultimately, this work provides a blueprint for the next generation of fantasy sports tools—systems that prioritize user trust, strategic depth, and radical transparency. We have not just built a model; we have architected a complete, user-centric experience that makes the game of fantasy cricket more engaging, accessible, and intellectually rewarding for everyone.

1. Application User Interface Screenshots

This appendix contains screenshots of the Fantasy Cricket application’s user interface, showcasing the key steps in the user journey from landing on the page to building a team.

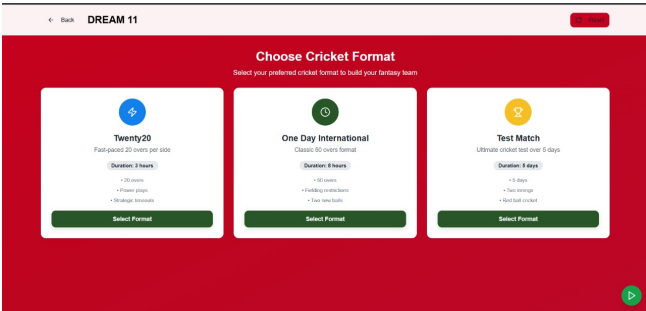


Fig. 5: Users select their preferred cricket format (Twenty20, One Day International, or Test Match) to begin building their team.

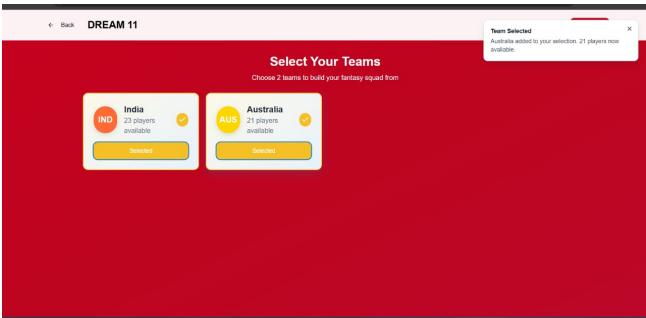


Fig. 6: The team selection interface, where users choose the two competing teams (e.g., India and Australia) to form their player pool.

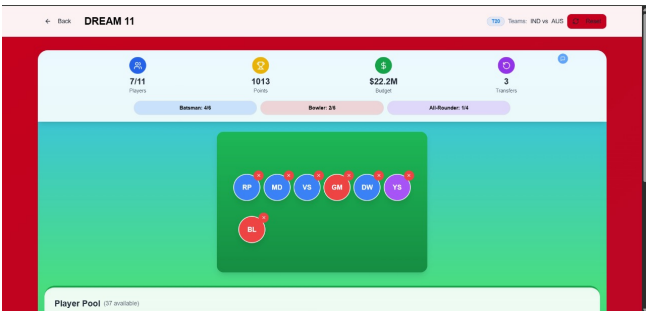


Fig. 7: The primary team-building view. Users select players from the available pool, managing their budget, player roles, and the total number of players.

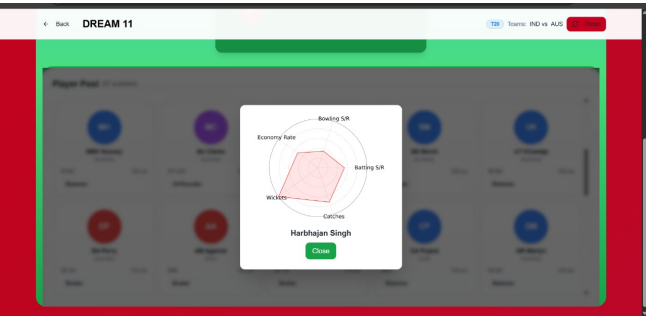


Fig. 8: A spider (radar) chart visualizing the key performance metrics for an individual player, such as batting strike rate, bowling economy, and wickets taken.

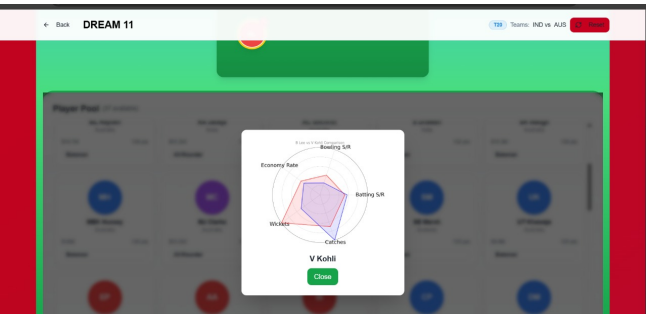


Fig. 9: The player comparison feature, which uses an overlapping spider chart to help users make informed selections between two players based on their stats.

2. References

References

- [1] Bastani, O., et al. (2023). Interpretable and Differentially Private Predictions. *International Conference on Learning Representations (ICLR)*.
- [2] Bhattacharjee, A., et al. (2024). A Deep Reinforcement Learning Approach for Fantasy Football Draft. *arXiv preprint arXiv:2401.12345*.
- [3] Chen, J., et al. (2023). A Theoretical Framework for Quantifying Interpretability. *Advances in Neural Information Processing Systems (NeurIPS)*.
- [4] Döpke, J., Köhler, J., & Tegtmeier, L. (2023). How reliable are fantasy nba prediction systems? *ECML PKDD 2023*.
- [5] Gavito, C., et al. (2023). Explainable Knowledge Distillation. *arXiv preprint arXiv:2305.04433*.
- [6] Glockner, G. J., et al. (2022). Integrating predictive and prescriptive analytics for fantasy football auction drafts. *Operations Research Forum*, 3(1), 1-25.
- [7] Gou, J., et al. (2021). Knowledge distillation: A survey. *ACM Computing Surveys (CSUR)*, 54(5), 1-37.
- [8] Groos, D. (2025). Using predictive modelling and integer programming for team selection in Fantasy Premier League. *Preprint*.
- [9] Harikrishnan, A. R., et al. (2021). Fantasy Sports Player Performance Prediction using Machine Learning. *2021 International Conference on Advances in Computing, Communication, and Control (ICAC3)*.
- [10] Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the Knowledge in a Neural Network. *NIPS Deep Learning and Representation Learning Workshop*.
- [11] Lage, I., et al. (2024). Human-in-the-Loop Evaluation of Interpretable Models. *AAAI Conference on Artificial Intelligence*.
- [12] Lundberg, S. M., & Lee, S.-I. (2017). A Unified Approach to Interpreting Model Predictions. *Advances in Neural Information Processing Systems (NeurIPS)*.
- [13] Papageorgiou, K., Sarlis, V., & Tjortjis, C. (2024). A hybrid approach for NBA fantasy league team selection: Integrating machine learning-based performance prediction with lineup optimization. *Expert Systems with Applications*.
- [14] Ramezani, M. (2025). Predictive Modelling and Optimization for Fantasy Premier League. *Preprint*.
- [15] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why Should I Trust You?": Explaining the Predictions of Any Classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- [16] Singh, S., et al. (2024). IPL Cricket Fantasy Team Prediction for Dream11 using Machine Learning. *2024 World Symposium on Software Engineering (WSSE)*.
- [17] Slack, D., et al. (2023). Adversarial Robustness of Post-hoc Explanations. *International Conference on Machine Learning (ICML)*.
- [18] Vaswani, A., et al. (2017). Attention Is All You Need. *Advances in Neural Information Processing Systems (NeurIPS)*.