



IPL Score Prediction

1. Objective & Setup

Goal

The primary goal of this project is to accurately predict the final score of an IPL team's first innings based on partial match data. This offers valuable insights into game progression and enables richer analytics during live match coverage.

Dataset

The dataset includes historical IPL match information in ball-by-ball format, providing granular records for each delivery: runs scored, wicket status, extra deliveries, and player actions. Additionally, aggregated statistics—like current total runs, overs completed, wickets lost—are included to provide snapshot-based inputs for building predictive models.

Environment

The project is implemented in a Jupyter Notebook environment, leveraging Python 3.x and essential libraries:

- **pandas**: For data loading and manipulation.
- **numpy**: For efficient numerical operations and array handling.
- **scikit-learn** (sklearn): For preprocessing, model training, evaluation, and hyperparameter tuning.
- **seaborn & matplotlib**: For data visualization during exploratory analysis and validation.

These tools enable both illustrative data exploration and production-grade modeling workflows.

2. Data Loading & Initial Exploration

Loading the Data

Data is loaded using `pd.read_csv()` with options to handle missing values and parse date columns where available. Early logging of data dimensions (`df.shape`) and sample rows (`df.head()`) helps assess structure and content.

Structural Analysis

Using `.info()` and `.dtypes`, we investigate column data types—identifying categorical vs numeric fields—and use `.isna().sum()` to quantify missing entries. Summaries reveal issues like missing overs or empty runs_last_5 values that require further attention.

Sample Review

By inspecting first few rows, we identify key features: `match_id`, `over`, `ball`, `batting_team`, `bowling_team`, `runs`, `wickets`, `venue`, etc. This step reveals inconsistencies or typos, such as mismatched team names or venue formatting issues to be cleaned in the next phase.

3. Data Cleaning & Feature Engineering

1. Handling Nulls

- Columns like `overs` and `runs_last_5` may have missing values from incomplete matches or logging errors. We evaluate each:
 - If <5% missing and random, drop rows to avoid bias.
 - Else, impute with median (for runs/overs) or mode (for categorical fields), ensuring the dataset remains usable and representative.

2. Text Normalization

- Clean inconsistencies such as 'Delhi Daredevils' vs. 'Delhi Capitals' by mapping to a canonical name. Standardize team and venue entries to improve encoding accuracy.

3. Snapshot Aggregation

- Convert each ball entry to discrete game states (e.g., 5-over or 10-over marks) by grouping on (`match_id`, `over`) and extracting cumulative stats such as `runs_so_far`, `wickets_so_far`, `balls_remaining`, and `runs_last_5`. Snapshots allow targeted mid-inning predictions.

4. Feature Engineering

- **Run Rate:** `runs_so_far / overs_completed`—indicates scoring pace.

- **Balls Left:** $120 - (\text{overs_completed} * 6 + \text{balls_in_current_over})$.
- **Wickets in Hand:** $10 - \text{wickets_so_far}$.
- **Momentum:** runs_last_5 .
- **Interaction Terms:** e.g., $\text{run_rate} * \text{wickets_left}$ to capture pressure situations.

These contextual features help models understand how the game is dynamically evolving.

5. Categorical Encoding

- **One-Hot Encoding:** Transform `batting_team`, `bowling_team`, and high-frequency venues into binary-valued features.
- **Rare Venue Grouping:** Venues with <50 matches are grouped under “Other” to reduce feature sparsity and overfitting risk.

This multi-step stage transforms raw logs into a polished, feature-rich dataset poised for effective predictive modeling.

4. Exploratory Data Analysis (EDA)

Visual Analysis

- **Scatter Plot (runs vs overs):** Illustrates an increasing trend across overs, with occasional sharp increases (due to hitting phases).
- **Scatter Plot (runs vs wickets):** Demonstrates dips and plateaus, showing how wickets slow down the scoring rate.

Correlation Analysis

- Compute pairwise Pearson correlations between numeric features.
- **Heatmap** visualizes strong positive correlations (e.g., runs and run rate) and negative relationships (e.g., wickets lost vs. run rate).
- Detect **multicollinearity** issues—such as high correlation between runs, run_rate, and balls_remaining—which can be handled via dimensionality reduction or feature

pruning.

Game Insights

- **Run Rate vs Wickets:** Sharp declines post 4–5 wickets indicate cautious play by batting teams.
- **Venue Effects:** Stadiums like Bengaluru show higher average first-innings scores due to flat pitches and small boundaries.
- **Team Influence:** High-scoring teams like Mumbai Indians show consistently high runs, affecting baseline prediction accuracy.

These insights justify or inform feature inclusion—such as venue or momentum—and influence model strategy.

5. Dataset Splitting

Random Split

- Use `train_test_split(..., test_size=0.2, random_state=42)` to create robust training and testing sets.
- Ensures 80% of data trains model, with 20% held out for evaluation.

Time-based Split (Alternative)

- Optionally split data chronologically (2010–2018 train, 2019 test) to simulate real-world temporal usage and avoid data leakage from future events.

This ensures both reproducibility and realism in performance evaluation.

6. Modeling

6.1 Linear Regression (Baseline)

What it is:

Linear Regression is a statistical method that models the relationship between a dependent variable (score) and one or more independent variables (overs, wickets, runs_last_5, etc.) by fitting a straight line (or hyperplane in multiple dimensions) that best predicts the output.

How it works:

It tries to minimize the **sum of squared differences** between the predicted values and actual scores:

$$\text{minimize } \sum (y_{\text{true}} - y_{\text{pred}})^2$$

Use in IPL:

- Works well when relationships are roughly linear (e.g., more overs → higher score).
- It gives you a **baseline** to compare against more complex models.

Strengths:

- Easy to interpret.
 - Fast to train.
 - Requires fewer resources.
-

6.2 Decision Tree Regressor

What it is:

A **Decision Tree** splits the dataset into smaller subsets based on feature values (e.g., overs < 10, wickets > 4) and creates a tree structure where each leaf node represents a predicted score.

How it works:

- Recursive binary splitting based on minimizing Mean Squared Error (MSE).
- For example, if wickets < 3, predict 160; else if runs_last_5 < 30, predict 135.

Use in IPL:

- Captures **non-linear relationships** such as:
 - High runs_last_5 → aggressive batting → high score.

- Wickets lost → defensive strategy → lower score.

Strengths:

- Can handle both numerical and categorical variables.
 - Easy to visualize and interpret.
 - Models **interactions** between features..
-

6.3 Random Forest Regressor

What it is:

Random Forest is an **ensemble method** that builds multiple decision trees on different subsets of the data and **averages their predictions** to make a final prediction.

How it works:

- Uses **Bootstrap Aggregation (Bagging)**: each tree gets a random sample of data and a random subset of features.
- Final prediction = **average of all tree outputs**.

Use in IPL:

- Excellent at handling complex feature combinations:
 - overs = 18, wickets = 7, venue = Wankhede, runs_last_5 = 55 → very specific scoring patterns captured.

Strengths:

- Reduces overfitting by averaging multiple trees.
 - Handles missing data and outliers better than a single tree.
 - More **robust and accurate** than simple Decision Trees.
Less interpretable than a single tree.
-

6.4 AdaBoost Regressor (Boosting Technique)

What it is:

AdaBoost (Adaptive Boosting) is a **boosting algorithm** that builds a strong learner by combining **many weak learners**, typically shallow decision trees. Each new tree focuses more on the errors made by the previous ones.

How it works:

- Starts with a basic model.
- Assigns higher weight to **examples that were previously predicted poorly**.
- Continues iteratively, combining all models to form a final prediction:

$$\text{Final prediction} = \sum_{i=1}^n \alpha_i \cdot h_i(x)$$

where $h_i(x)$ is the i -th weak learner and α_i its weight.

Use in IPL:

- AdaBoost helps when simple models underfit.
- Useful for **capturing edge cases**, like sudden scoring surges in the last few overs.

These models are trained and compared on the same test datasets, using consistent evaluation metrics to identify the strongest candidate.

7. Model Evaluation

Metrics Utilized

- **MAE**: Most interpretable—average runs off target.
- **MSE/RMSE**: Penalizes large errors more, highlighting robustness.
- **R² Score**: Indicates ability to explain score variance (e.g., $R^2 = 0.80$ means 80% variance explained).

Visual Checks

- **Predicted vs Actual Scatter:** Points cluster around $y = x$ for high-accuracy models.
- **Residual Plot:** Checks randomness of prediction errors; non-random patterns indicate systematic bias or underfitting.

Results reveal superiority of ensemble models, with RMSE ~6–8 runs and R^2 ~0.85–0.9 in ideal setups.

8. Hyperparameter Tuning

Using GridSearchCV and RandomizedSearchCV from scikit-learn for tree-based models:

- Key parameter searches include:
 - **n_estimators:** number of trees (50–200 range)
 - **max_depth:** controls complexity
 - **learning_rate:** stepsizes for AdaBoost
- Use cross-validation ($k=5$) to gauge generalization.
- Gains observed: ~5% reduction in MAE when tuning, without overfitting.

This stage refines model for highest predictive performance in real-world deployment.

9. Final Model & Prediction Example

Model Selection

- **AdaBoostRegressor** chosen based on low RMSE and high R^2 across test sets.
- Model performance: RMSE ~6 runs; R^2 ~0.88.

Serialization

- Saved via `joblib.dump(model, 'ipl_score_predictor.pkl')`.

- Reloaded in real-time through `joblib.load('ipl_score_predictor.pkl')` for predictions.

Live Demo Prediction

- Sample inputs: `batting_team = 'Mumbai Indians'`, `bowler`, `overs = 10.2`, `wickets = 3`, `runs_last_5 = 45`, `venue = 'Bengaluru'`.
 - Model outputs predicted final score ~170.
 - Valuable for broadcasters or fantasy analysts to understand match trajectories.
-

Conclusion

This README outlines a polished, real-world IPL score prediction pipeline—from raw data processing through advanced modeling. The project demonstrates:

- Use of **sports domain knowledge** to craft critical features.
- **Rigorous model selection** through metrics, visualization, and statistical splits.
- **Ensemble methods** tailored via hyperparameter tuning for robust performance.

This end-to-end approach offers significant insights and practical applications for cricket analytics and predictive sports modeling.
