

Journal of Computational Algorithms and Information Technology

Article

Predicting IPL Match Outcomes Using Machine Learning Techniques

Geetha A. J.

Central University of Andhra Pradesh, India.

* Correspondence: geethabilinele@gmail.com

Received: 12 May 2025; Accepted: 03 June 2025; Published: 21 June 2025.

Abstract: Cricket is amongst the most popular sports in the world. Indian Premier League, more commonly known as IPL, is the biggest domestic cricket league in the world. It generates a lot of revenue along with excitement among fans. Many bookers, bettors, and fans like to predict the outcome of a particular match which changes with every ball. This project studies and compares different Machine Learning techniques that can be applied to predict the outcome of a match. Features like team strength and individual strength of a player are also included along with conventional features like toss, home ground, weather and pitch conditions that are taken into account for predicting the result. Machine Learning algorithms such as Naive Bayes, Random Forest Classifier, Logistic Regression, XGBoost, AdaBoost, and Decision Tree are selected to determine the predictive model with highest accuracy.

Keywords: AdaBoost, Decision Tree, Indian Premier League, Machine Learning, Naïve Bayes, Logistic Regression, Random Forest Classifier, XGBoost

1. Introduction

1.1. Problem Statement

IPL is considered to be the most popular T20 league in the world and it creates a lot of excitement amongst the fans of the game. Many people like to predict the outcome of the game beforehand and like to formulate their own strategies and create their own fantasy teams. There are many systems already present that predict the outcome of an IPL match, but all the factors affecting a particular match are not taken into consideration. The main aim of this work is to create a prediction system with accuracy as high as possible, which will help users to create their fantasy teams with much ease.

1.2. Motivation

There is always a rush in our heads regarding the results of a cricket match. Friendly betting was one of the factors which motivated us for predicting the winners of IPL matches. As the law-obliged sports gambling industry is growing at a faster rate than ever, it can be useful for people to get the idea of the outcome before a match begins. So, we thought of creating a prediction system with accuracy as high as possible which will make it easy to create fantasy teams.

2. Literature Review

Many researchers have contributed towards predicting the results of cricket matches. Here, some of the prominent works done by researchers are discussed. In [1] the researchers investigated the match by considering multiple features. The authors have taken two cases into consideration:

1. Home ground advantage

2. Winning the toss

When the algorithms were applied to the given two cases, the results were as follows:

Case 1: Home Ground Advantage

1. Naïve Bayes – 57%
2. Model decision tree – 56%
3. Random forest – 54%
4. KNN – 52%

Case 2: Winning the Toss

1. KNN – 62%
2. Naïve Bayes – 52%

Their model provided better prediction results by using the second case of winning the toss and applying KNN algorithm to it. They made their model on the basis of past results of matches and have not included the player's data. In [2] the researchers applied different algorithms on their dataset and the results they obtained were as follows:

1. Naïve Bayes – 71%
2. KNN – 66%
3. Logistic Regression – 64%
4. Random Forest – 63%
5. Support Vector Machines – 59%

They did not use features like weather conditions, venue, and pitch conditions which play a vital role in a particular cricket match.

In [3] the researchers applied various Data Science techniques to foretell the result of an IPL match. The datasets were scraped and obtained from espnricinfo [4], ipl2020 [5], and Kaggle [6]. Batting, bowling, and winning strength for the past 10 years were analyzed and the strengths of all IPL teams were considered and the final prediction was made considering the past data from 2008-2017. Cumulative team strength was also considered for predicting match outcomes.

In [7] the analytic hierarchy process was used for getting priority order and weights of attributes that are considered for calculating batting and bowling average. Win rate strength was calculated by multiplying team strength with win rate which was then used to predict the outcome of a particular match. The results were as follows:

1. Naïve Bayes – $58.23\% \pm 5.5\%$
2. Adaboost – $60.03\% \pm 6.2\%$
3. Logistic Regression – $57.77\% \pm 5.8\%$
4. Support Vector Machines – $58.42\% \pm 5.69\%$
5. KNN – $53.47\% \pm 5.2\%$
6. XGBoost – $55.42\% \pm 5.9\%$
7. Extra Trees Classifier – $59.51\% \pm 5.9\%$
8. Random Forest Classifier – $60.04\% \pm 6.3\%$

Highest accuracy was provided by Random Forest and AdaBoost ($60.4\% \pm 6.3\%$ and $60.3\% \pm 6.2\%$ respectively), and lowest accuracy of $53.47\% \pm 5.2\%$ was provided by KNN. Even if you consider the best case result, the maximum accuracy will not go beyond 67% which is not that good a result. Though the methodology was good as compared to other papers, there were a few drawbacks related to team strength calculation and win rate calculation. Nothing was mentioned about the debutants of the competition. Win rate was figured out by fractionating the total number of matches won by the team to the total number of matches played by the team, but there were few factors that affected the match like weather conditions, home ground advantage, and toss winning advantage which were not considered.

The authors of [8] used various Machine Learning algorithms to predict results and analyze an IPL match. Two datasets were used, one having ball-to-ball information of the whole IPL till 2019 and the other one having information of all the matches. The authors made use of a vast number of factors like weather, venue, and pitch that affect a match were taken into consideration to make a training model. Batting and bowling performances of individual players were also considered in the model. In addition to the past performance factors the condition of the factors were also considered by them. The results were as follows:

Table 1. Comparison of Algorithms Used in Previous Papers

| Algorithms | Maximum Accuracy |
|--------------------------|------------------|
| Naïve Bayes | 81.63% |
| Random Forest Classifier | 83.67% |
| Logistic Regression | 95.91% |
| XGBoost | 55.42% ± 5.9% |
| AdaBoost | 60.4% ± 6.3% |
| Decision Tree | 87.5% |
| SVM | 83.67% |
| KNN | 66% |

2.1. Combined Study

On studying the given related papers thoroughly, we got some useful information, which is presented in Table I.

3. Proposed Methodology

3.1. Block Diagram

Fig. 1 shows the block diagram.

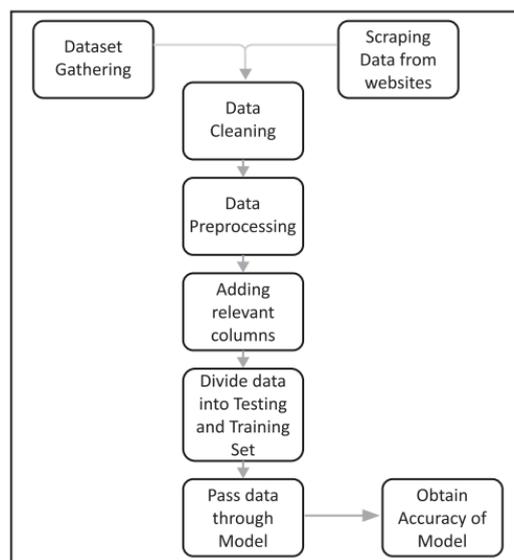


Figure 1. Block Diagram

3.2. Dataset

The dataset used for analysis and prediction was collected from Kaggle [6] and furthermore, data was scraped from espnricinfo [4] and ipl2020 [5] using the BeautifulSoup library of Python [9].

3.3. Data Preprocessing

For applying the Machine Learning algorithms, string data will be converted into numerical data as the algorithms work better with numerical values.

All the unnecessary attributes such as names of umpires, venue, date, player of the match, method, and eliminator were removed from the dataset to get more accurate results.

Further, all the match rows that were dismissed, drawn, or null were eradicated.

Newer needed features such as team strength and team points were added to the dataset. Team points were calculated on the basis of previous matches of each team.

3.4. Features

1) **Toss:** Toss is considered as the most important factor sometimes. After winning a toss, a team can select either to bat or bowl first, which when analyzed can provide an upper hand to the toss winning team over the other team even before the match starts. Toss analysis of our dataset tells us that 55.8% times toss winning team wins the match (Fig. 2). However, toss is not the only factor on which the outcome of a match is decided.

2) **Team Points:** Team points are used to represent the performance of a team by using the results of the team's

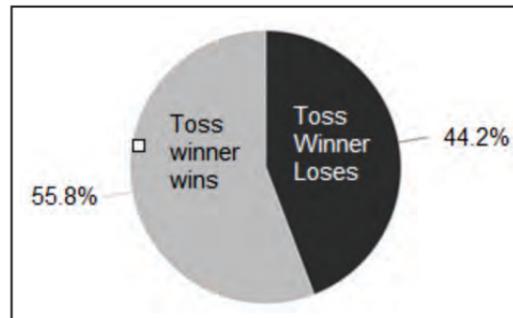


Figure 2. Toss Factor

previous matches. The previous 5 matches of each team were considered for calculating team points. A win was counted as 2 points and a loss was counted as 0 points as it is done officially in the IPL. Team points analysis of our dataset tells us that 59.3% times team having greater points wins the match (Fig. 3).

3) **Team Strength:** Team strength is used to represent the cumulative strength of all the players of a team. For calculating team strength we made use of the player-points data of every year available on IPL's official website. Each team's top 11 players were selected and their average points were taken to get team strength. Team strength analysis of our data-set tells us that 60.9% (Fig. 4) times team having greater strength wins the match.

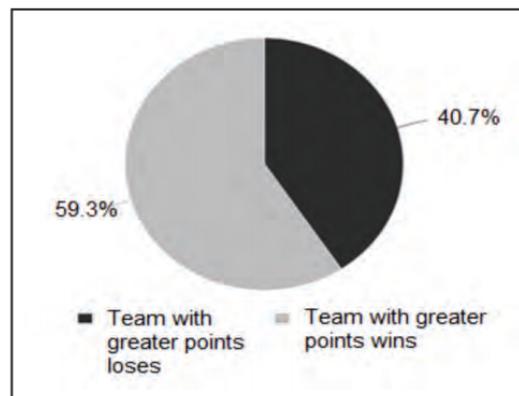


Figure 3. Team Points Factor

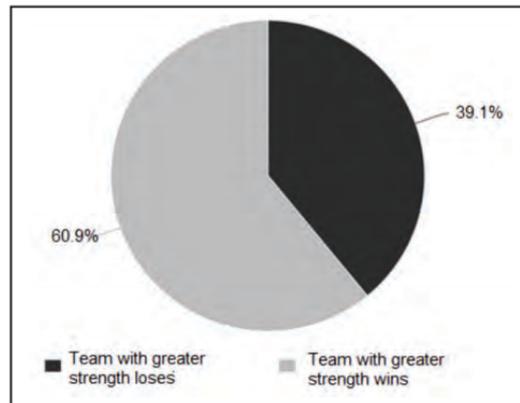


Figure 4. Team Strength Factor

3.5. Algorithms

1) AdaBoost: Ada Boost is short for Adaptive Boosting. It is an ensemble ML method. It learns from the mistakes of weak classifiers in every iteration and hence, converts the same to strong classifiers.

Parameters of AdaBoost algorithm used are:

1. n_estimators
2. learning_rate

2) Decision Tree: Decision Tree algorithm is used for regression as well as classification. It begins with root node and finishes with decision node called leaf. It has root, decision, and terminal nodes. It is very convenient in case of decision based examples.

Parameters of Decision Tree algorithm used are:

1. criterion
2. min_samples_split
3. min_samples_leaf
4. max_features
5. max_depth

3) Random Forest: The Random Forest classifier contains numerous separate decision trees. Each tree in the forest calculates a class prediction and the class that has highest votes will be given out as the prediction of the whole model. It operates like an ensemble. For higher accuracy the number of trees should be higher.

Parameters of Random Forest algorithm used are:

1. n_estimators
2. min_samples_split
3. min_samples_leaf
4. max_features
5. max_depth
6. bootstrap

4) XGBoost: XgBoost stands for Extreme Gradient Boosting that uses decision tree algorithm to predict small to medium structured or tabular data. Ensemble method is used which gives us predictive power of multiple iterations. Using Bagging or Boosting, highly similar behaviour is reduced, which in turn helps to increase the accuracy.

Parameters used for XGBoost algorithm are:

1. n_estimators
2. learning_rate
3. max_depth

5) Logistic Regression: This algorithm is used in cases where the target variable is categorical. It is predicted using independent variables. It is used in the case of classification and not regression problems.

Parameters of Logistic Regression that were used are:

1. penalty
2. tol
3. solver
4. max_iter
5. C

6) Naïve Bayes: This algorithm is based on the Bayes Theorem of conditional probabilities. Classification problems can be solved using this algorithm. In this, each feature is supposed to be independent of the other and it contributes to the result independently.

Parameter used for Naïve Bayes algorithm are:

1. var_smoothing

3.6. Hyper Parameter Tuning

Hyper-parameters of algorithms are the parameters the values of which need to be set before building the model and beginning the Machine Learning process. Their values cannot be obtained from training. Tuning a model for a certain problem can be done using Random Search or Grid Search. By doing this, the most efficient prediction of the learning model can be obtained.

a) RandomizedSearchCV: RandomizedSearchCV (Cross-Validation) is used to basically optimize the hyper-parameters. Random combinations are used to train the model. A sampling distribution is used for every hyper-parameter to do random search which allows us to administer the attempted combinations.

b) GridSearchCV: GridSearchCV (Cross-Validation) is used to optimize the hyper-parameters of a model. It uses traditional trial and error to get all the combinations. After this cross validation is done, it helps us to get the best accuracy. Cross Validation checks show how a model generalizes itself to an unconstrained dataset.

For tuning the algorithms, GridSearchCV was used in the case of AdaBoost and Naïve Bayes algorithm while the other algorithms were tuned using RandomizedSearchCV.

4. Results and Conclusion

The conclusion that was obtained after building the models of all the algorithms mentioned earlier is presented in Table II.

It can be seen in Fig. 5 that the highest accuracy achieved by the current investigation is 89.04% provided by XGBoost algorithm and the highest accuracy achieved was 95.91% provided by Logistic Regression algorithm for the previous investigation.

Table 2. Results and Comparison of Various ML Techniques

| Algorithm | Maximum Accuracy |
|--------------------------|------------------|
| Naïve Bayes | 59.36% |
| Random Forest Classifier | 84.726% |
| Logistic Regression | 56.19% |
| XGBoost | 89.049% |
| AdaBoost | 59.942% |
| Decision Tree | 73.017% |

factor because of which the accuracy of current investigation is differing from previous investigation. Other than that, parameters of a particular algorithm chosen in current investigation may be different from the parameters chosen in previous investigation for that algorithm. There are many different features selected by previous investigation as compared to current investigation which can also be the reason of difference in accuracy for same algorithms.

5. Future Scope

As the IPL seasons go on, the dataset can be updated to the current season to get more precise results and for fantasy team builders the best 11 players of the match can also be predicted based on the fantasy points they have previously obtained.

Author Contributions: All authors contributed equally to the writing of this paper. All authors read and approved the final manuscript.

Conflicts of Interest: “The authors declare no conflict of interest.”

References

- [1] K. Kapadia, H. Abdel-Jaber, F. Thabtah, and W. Hadi, “Sport analytics for cricket game results using machine learning: An experimental study,” *Appl. Comput. Inform.*, vol. ahead-of-print, 2019, doi: 10.1016/j.aci.2019.11.006.
- [2] P. K. Dubey, H. Suri, and S. Gupta, “Naïve Bayes algorithm based match winner prediction model for T20 Cricket,” in S. S. Dash, S. Das, B. K. Panigrahi (eds), *Intell. Comput. Appl., Advances Intell. Syst. Comput.*, vol. 1172, Springer, Singapore, 2021. doi: 10.1007/978-981-15-5564-4_38.
- [3] A. Tripathi, R. Islam, V. Khandor, and V. Murugan, “Prediction of IPL matches using Machine Learning while tackling ambiguity in results,” *Indian J. Sci. Technol.*, vol. 13, no. 38, pp. 4013–4035, 2020, doi: 10.17485/IJST/v13i38.1649.
- [4] Espnricinfo. [Online]. Available: <https://www.espnricinfo.com/>
- [5] Indian Premier League. [Online]. Available: <https://www.ipl2020.com/>
- [6] Kaggle. [Online]. Available: <https://www.kaggle.com/datasets>
- [7] Dream11. [Online]. Available: <https://www.dream11.com/>
- [8] H. Barot, A. Kothari, P. Bide, B. Ahir, and R. Kankaria, “Analysis and prediction for the Indian Premier League,” *Int. Conf. Emerg. Technol.*, 2020, pp. 1–7, doi: 10.1109/INCET49848.2020.9153972.
- [9] Beautiful Soup Python library. [Online]. Available: <https://pypi.org/project/beautifulsoup4/>



© 2025 by the authors; licensee BIESR, Lahore, Pakistan. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) license (<http://creativecommons.org/licenses/by/4.0/>).