



School of Computing

IPL OUTCOME PREDICTION

GROUP-14

PROJECT REPORT

TEAM MEMBERS:

NAME	ROLL NO
SHRISHARANYAN VASU	AM.EN.U4AIE22150
VASUDEVAREDDY R	AM.EN.U4AIE22143
ABHINAV A	AM.EN.U4AIE22107
MADHAV UNNIKRISHNAN	AM.EN.U4AIE22132

ABSTRACT:

The Indian Premier League (IPL) has become one of the most popular and widely watched cricket tournaments globally, bringing millions of fans and followers. Predicting the results of cricket matches, particularly the final score, is a difficult task due to the exciting nature of the game and the multiple influencing factors. This project specializes in the application of Machine Learning (ML) techniques to predict the final scores of IPL matches. The proposed model utilizes past match data, team statistics, player performance metrics, and other relevant features to train a predictive model. The aim is to develop a reliable and precise system that can expect the total runs scored by each team in each match.

The method involves data preprocessing, feature engineering, and the selection of appropriate ML algorithms such as Random Forest, Logistic Regression, XGBoost, LightGBM, and ExtraTreesClassifiers, fine-tuned via Randomized search for optimal accuracy. Cross-validation ensures reliable model validation. The model will be trained on a comprehensive dataset comprising information on past IPL matches, including involving scores, player performances, team statistics, venue details, and various other factors that may impact the final score.

Keywords: IPL, Players Form, Machine Learning, Team Achievement, IPL score Prediction, Random Forest Classifier, XGBoost, LightGBM, ExtraTreesClassifiers Logistic Regression.

INTRODUCTION:

Due to its intense competitiveness and exciting matches, cricket fans all over the world get excited with the Indian Premier League (IPL), an important sport. It can be helpful for teams to predict the results of these games, and it can also provide supporters with exciting information and knowledgeable analysis. Using machine learning algorithms and previous match data, this research aims to forecast IPL outcomes for matches with accuracy.

We start by collecting extensive data from past IPL seasons. This dataset contains a wide range of information, including match locations, player statistics, team performance measures, and other relevant data that may have an impact on the outcome of a match. To make sure the dataset is clean and prepared for analysis, the first step is data preprocessing. In this stage, we use imputed techniques to fill in the values that are missing, using tools like Label Encoder to encode categorical variables (such as team names and venues) into numerical formats, and standardized numerical features to maintain the dataset accuracy.

We use EDA (exploratory data analysis) after preprocessing to find patterns and relationships in the data. Visualizations are used in EDA to illustrate data so that distributions can be understood, and correlations among various variables can be identified. Through the finding of key events and variables that could affect match results, this analysis leads toward the creation of sensible features. We update our data set during the feature engineering phase by developing new features that capture different aspects of the match. For example, we provide measures that show the teams'

recent performance, the individual performances of players, and the benefits and drawbacks of playing at locations.

The project's main parts are building models and training. To predict match results, we evaluate some machine learning classifiers, such as Random Forest, Logistic Regression, XGBoost, LightGBM, and ExtraTreesClassifiers.

1. Random Forest:

During training, Random Forest is a flexible machine learning technique that creates many decision trees and combines their output to provide predictions that are more reliable and accurate. Both classification and regression problems benefit from their effectiveness. It lowers the possibility of overfitting and raises prediction accuracy by averaging the output of many trees.

2. Logistic Regression:

A simple but powerful classifier, logistic regression is mostly applied to binary classification issues. It does this by applying a logistic function to a linear set of input features, so predicting the likelihood that a given input belongs to a specific class. It is a common solution for situations with two possible results since it is simple to apply and evaluate.

3. XGBoost:

A powerful and effective gradient boosting method implementation is called XGBoost (Extreme Gradient Boosting). It constructs models one after the other, fixing inaccuracies in each model as it goes along. XGBoost is frequently utilized in machine learning contests and applications needing high accuracy because of its credibility for speed and performance.

4. LightGBM:

A scalable and efficient gradient boosting framework is called LightGBM (Light Gradient Boosting Machine). It lowers training time and use of memory by bucketing continuous feature values through a histogram-based method. LightGBM offers excellent speed and accuracy, making it especially suitable for tricky problems and huge datasets.

5. ExtraTreesClassifier:

While there are a few variations, ExtraTreesClassifier (Extremely Randomized Trees) is comparable to Random Forest. For every feature, it chooses a threshold at random rather than searching for the most biased one. More diversity in the ensemble of trees produced by this randomization can improve performance and generalization, particularly on noisy datasets.

6. Cross-Validation technique:

A method to evaluate a machine learning model's performance and guaranteeing its reliability and versatility is called cross-validation. The dataset is split up into multiple smaller groups, or "folds," for cross-validation. This process is done several times, with the model being trained on a subset of these folds and tested on the remaining ones. Compared to a single train-test split, the

average performance over all folds offers a more accurate estimation of the model's efficacy, preventing overfitting and guaranteeing the model works well on unobserved data.

The processed dataset is used to train each model, and to improve performance, we use parameter-adjusting techniques like Randomized Search. We use cross-validation to verify our models' reliability and make sure they transform well to untested data. We evaluate our models' performance using a variety of metrics, including accuracy, precision, recall, and F1-score, once they have been trained. These metrics provide a picture of the model's performance. To analyze the predictions, we provide thorough classification reports and confusion matrices that help in understanding the advantages and disadvantages of each model.

By creating a prediction model for IPL match results, the goal is to provide useful information that will help teams make wise decisions, improve cricket fans' ability to analyze, and boost fan engagement. This project highlights the potential of machine learning in sports analytics by demonstrating how data-driven approaches may significantly enhance understanding and decision-making in the cricket field. This research not only advances the field of sports analytics but also clears the way for the creation of increasingly complex and accurate prediction models.

This study shows how machine learning may be used practically in actual instances of sports. Advanced statistical techniques can be integrated to provide important insights that go beyond basic statistics. Through the provision of more in-depth analytical viewpoints, this complete strategy not only helps teams with their strategy but also enhances the fan experience. The project's techniques and outcomes highlight the transformative potential of machine learning in sports, demonstrating how it may improve cricket players' comprehension and involvement. Using this project, we make an addition to the greater domain of sports analytics by providing a framework that can be modified and extended to accommodate different sports and issues related to predictive modeling.

PROBLEM STATEMENT:

The primary goal of this work is to predict the result of an IPL match between two teams using machine learning algorithms and previously stored data analysis, given IPL datasets from the past few years. The data will undergo preprocessing and analysis. The data will be pre-processed and then utilized to train several models to provide the desired results. To determine the likely result of a match, we will examine the different datasets, use important characteristics like player performance and team statistics, etc., and feed the results into an algorithm.

LITERATURE REVIEW:

- Machine Learning Techniques for Score Prediction: Maheshwari et al. (2023) focus on predicting the first innings score for an IPL team. Their research

demonstrates the effectiveness of machine learning algorithms in analyzing factors influencing the score, such as batting line-up, venue, and historical data [1]. This approach highlights the potential for data-driven insights to improve score predictions.

- **Challenges of Prediction in T20 Format:** ResearchGate's publication "Prediction of IPL Match Outcome Using Machine Learning Techniques" acknowledges the inherent challenges of predicting outcomes in the T20 format, known for its high variability and reliance on individual performances [2]. This paper emphasizes the need for robust models that can account for these factors and improve prediction accuracy.
- **Algorithms for Score and Winner Prediction:** Sonali et al. (2021) delve deeper by exploring various machine learning algorithms to predict both the IPL match score and the winner [3]. Their research analyzes the effectiveness of different algorithms and highlights the potential for more sophisticated models to improve prediction accuracy for both aspects of the game.
- **Data Science and Deep Learning Applications:** Patil et al. (2023) explore the burgeoning field of IPL score prediction using data science and deep learning techniques [4]. They discuss the potential benefits for various stakeholders, including cricket teams that can use these insights for strategic planning, broadcasters who can tailor content based on predictions, and fans who can enhance their enjoyment of the game through informed analysis.
- **Classification Tasks in Machine Learning:** Anik et al. (2020) explore methods for prediction in cricket matches, including the IPL. Their research emphasizes the role of classification tasks in machine learning, where the model is trained to categorize matches based on historical data and various factors that influence the outcome [5]. This approach provides a framework for building accurate prediction models.

```
mdf = pd.read_csv('/kaggle/input/ipl-complete-dataset-20082020/matches.csv')
mdf.shape
```

(1895, 28)

```
mdf.head()
```

	id	season	city	date	match_type	player_of_match	venue	team1	team2	toss_winner	toss_decision	winner	result	result_margin	target_runs	target_overs	super_over	method	umpire1	umpire2
0	335982	2007/08	Bangalore	2008-04-16	League	BB McCullum	M Chinnaswamy Stadium	Royal Challengers Bangalore	Kolkata Knight Riders	Royal Challengers Bangalore	field	Kolkata Knight Riders	runs	140.0	223.0	20.0	N	NaN	Asad Rauf	RE Koertzen
1	335983	2007/08	Chandigarh	2008-04-19	League	MEK Hussey	Punjab Cricket Association Stadium, Mohali	Kings XI Punjab	Chennai Super Kings	Chennai Super Kings	bat	Chennai Super Kings	runs	33.0	241.0	20.0	N	NaN	MR Benson	SL Shastri
2	335984	2007/08	Delhi	2008-04-19	League	MF Maharoof	Feroz Shah Kotla	Delhi Daredevils	Rajasthan Royals	Rajasthan Royals	bat	Delhi Daredevils	wickets	9.0	130.0	20.0	N	NaN	Akeem Dar	GA Pratapkumar
3	335985	2007/08	Mumbai	2008-04-20	League	MV Boucher	Wankhede Stadium	Mumbai Indians	Royal Challengers Bangalore	Mumbai Indians	bat	Royal Challengers Bangalore	wickets	5.0	166.0	20.0	N	NaN	SJ Davis	DJ Harper
4	335986	2007/08	Kolkata	2008-04-20	League	DJ Hussey	Eden Gardens	Kolkata Knight Riders	Deccan Chargers	Deccan Chargers	bat	Kolkata Knight Riders	wickets	5.0	111.0	20.0	N	NaN	RF Bowden	K Harsharan

DATA PREPROCESSING:

Data preprocessing is one of the most important stages in the model development process of machine learning. It is the process of obtaining raw data into a form so that it

becomes clean, structured, and ready for analysis. Commonly, it includes cleaning the data, treating missing values, and outliers to ensure the quality of the data. Numerical features are then scaled with normalization and standardization for the model to learn better. One-hot encoding is one of the different techniques that can be used in converting categorical data into numerical values. Feature selection and extraction are very integral parts of preprocessing, where irrelevant or redundant features are removed, and new features could be created in order to increase the performance of the model. It is, therefore, the processing of data to change the raw data into a form that machine learning algorithms can exploit effectively, hence enhancing model accuracy, speeding up training time, hence predictive performance.

- Mapping Team Names to City Names: Mapped diverse team names to city names in order to have a standard name for each team across different seasons. This was done using a dictionary, (team_city_map), and a function, (map_team_city)

```
team_city_map = {
    'Sunrisers Hyderabad': 'Hyderabad',
    'Deccan Chargers': 'Hyderabad',
    'Mumbai Indians': 'Mumbai',
    'Gujarat Lions': 'Gujarat',
    'Gujarat Titans': 'Gujarat',
    'Rising Pune Supergiants': 'Pune',
    'Rising Pune Supergiant': 'Pune',
    'Pune Warriors': 'Pune',
    'Royal Challengers Bangalore': 'Bangalore',
    'Royal Challengers Bengaluru': 'Bangalore',
    'Kolkata Knight Riders': 'Kolkata',
    'Delhi Daredevils': 'Delhi',
    'Delhi Capitals': 'Delhi',
    'Kings XI Punjab': 'Punjab',
    'Punjab Kings': 'Punjab',
    'Chennai Super Kings': 'Chennai',
    'Rajasthan Royals': 'Rajasthan',
    'Kochi Tuskers Kerala': 'Kochi',
    'Lucknow Super Giants': 'Lucknow'
}

def team_to_city(df, feature):
    df[feature] = df[feature].map(team_city_map)
    return df

def map_team_city(df):
    mdf = df.copy()
    for feature in ['team1', 'team2', 'toss_winner', 'winner']:
        if feature in mdf.columns:
            mdf = team_to_city(mdf, feature)
    return mdf

mdf = map_team_city(mdf)
```

- Creating Home Ground Features: Created boolean features (team1_home) and (team2_home) to indicate if a team is playing on its home ground. Used the city feature to determine if the city of the match is the same as the team's city.

```
mdf['team1_home'] = False
mdf['team2_home'] = False

mdf.loc[mdf['city'] == mdf['team1'], 'team1_home'] = True
mdf.loc[mdf['city'] == mdf['team2'], 'team2_home'] = True
```

- Handling Missing Values: Removed the rows where the result was 'no result'. This is because these rows do not contribute to predictions of match outcomes. Dropped the method column because of a large degree of missingness and limited relevance. It has imputed the missing values in columns city, umpire1, and umpire2 with the most frequent values.

```
mdf = mdf.drop(mdf.loc[mdf['result'] == 'no result'].index, axis=0)
```

- Feature Engineering: It extracts the dayofyear, dayofweek, and season from the date column for a seasonal pattern and the specific day the match was played.
- Dropping Unnecessary Columns: Dropped the date column after extracting relevant time-series features. Dropped columns like id, and player_of_match, as they don't contribute to the prediction model

```
mdf = mdf.drop(columns=['player_of_match'])
```

- One-Hot Encoding: Applied one-hot encoding to the toss_decision column to convert categorical data into numerical format suitable for machine learning models.

```
# One-Hot encoding toss_decision
mdf = pd.get_dummies(mdf, columns = ['toss_decision'], drop_first=True)
```

- Label Encoding: Used LabelEncoder to transform categorical features (city, team1, team2, winner, toss_winner, venue, umpire1, umpire2) into numerical format.

```
team_le = LabelEncoder()
team_name_features = pd.concat([mdf['city'], mdf['team1'], mdf['team2'], mdf['winner'], mdf['toss_winner']]).unique()
team_le.fit(team_name_features)

mdf['city'] = team_le.transform(mdf['city'])
mdf['winner'] = team_le.transform(mdf['winner'])
mdf['team1'] = team_le.transform(mdf['team1'])
mdf['team2'] = team_le.transform(mdf['team2'])
mdf['toss_winner'] = team_le.transform(mdf['toss_winner'])

le = LabelEncoder() # Use a different instance of the label encoder for the other features.
features = ['venue', 'umpire1', 'umpire2']
venue_le = LabelEncoder()
umpire_le = LabelEncoder()
umpire_le.fit(pd.concat([mdf['umpire1'], mdf['umpire2']]).unique())
mdf['venue'] = venue_le.fit_transform(mdf['venue'])
mdf['umpire1'] = umpire_le.transform(mdf['umpire1'])
mdf['umpire2'] = umpire_le.transform(mdf['umpire2'])

mdf.head()
```

- Calculating Win Ratio: Calculated the win ratio for each team by dividing the number of matches won by the total number of matches played. Added new features team1_win_ratio and team2_win_ratio based on these calculations.

```

# Calculating the win ratio of every team

win_ratio = {}
for team, count in matches_played.items():
    win_ratio[team] = matches_won[team] / count

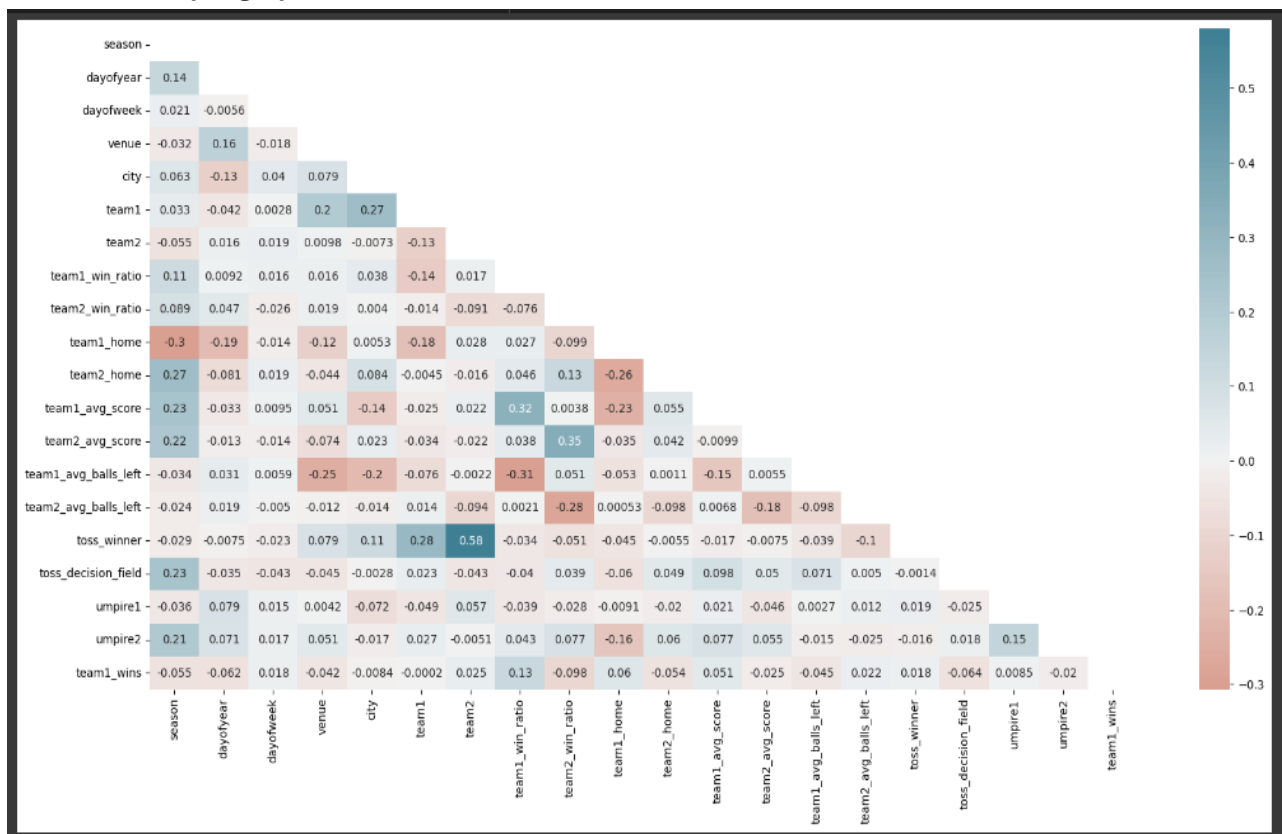
win_ratio

] # Adding team1_win_ratio and team2_win_ratio features to the dataset.
mdf['team1_win_ratio'] = mdf['team1'].map(win_ratio)
mdf['team2_win_ratio'] = mdf['team2'].map(win_ratio)

```

DATA VISUALIZATION:

- Correlation Matrix: Created a Correlation Matrix to understand the relationships between different features. Visualized the correlation matrix using a heatmap to identify highly correlated features.



- The feature importance plot illustrates the relative importance of each feature in predicting the target variable. Features are sorted based on their importance scores, derived from the model.

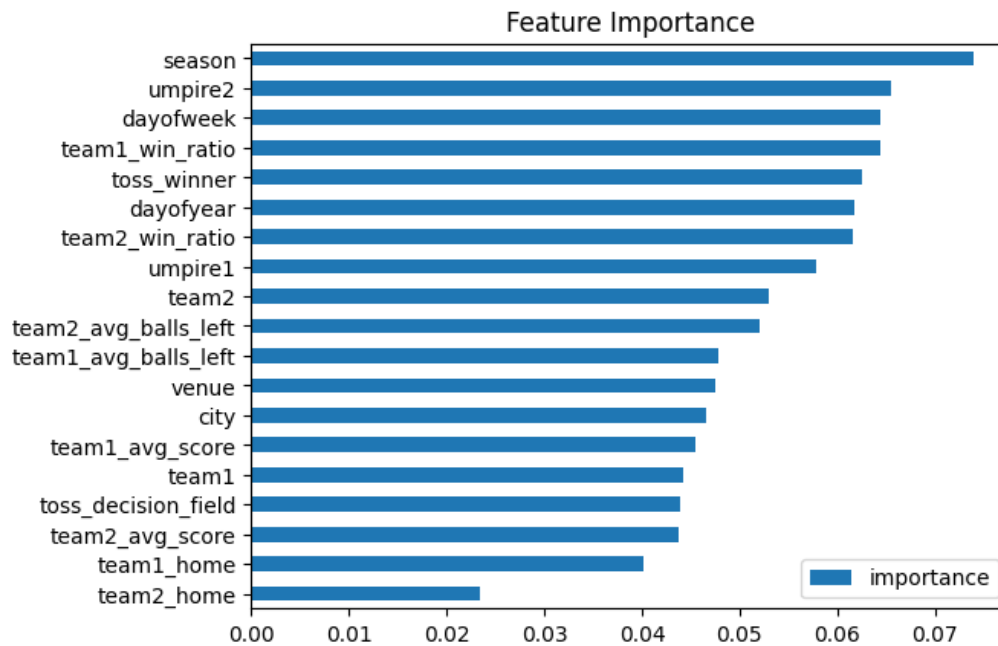


Fig: Shows the feature importance plot generated from the model. Features are ranked by their importance scores, indicating their relative contribution to the predictive accuracy of the model.

This visualization aids in understanding which variables significantly influence the model's predictions, guiding further analysis and feature engineering efforts.

DATA SPLITTING AND MODEL BUILDING:

Data splitting divides the dataset into training and testing sets, allowing for model training on one subset and performance evaluation on another, ensuring the model's ability to generalize to unseen data. We divided the training data and testing data in 80:20 ratio.

Model building involves selecting and training a machine learning algorithm on the training data to learn patterns and relationships, and then using this trained model to make predictions on new data. We used different regression models to predict the target variable. The models are mentioned below.

- **ExtraTrees (Extremely Randomized Trees Classifier):** The ExtraTrees Classifier is an ensemble learning method that aggregates the results of multiple de-correlated decision trees collected in a "forest" to output its classification result. Unlike Random Forests, ExtraTrees builds each tree from the entire training sample but splits nodes on random splits instead of best splits. This results in higher variance reduction and better generalization. ExtraTrees is known for being fast to train and less prone to overfitting compared to individual decision trees.

```

et = ExtraTreesClassifier()


param_dist = {
    'n_estimators': [100, 250, 500, 1000],
    'max_features': ['auto', 'sqrt', 'log2'],
    'max_depth': [None, 5, 6, 7, 8],
    'min_samples_split': [2, 3, 4],
    'min_samples_leaf': [2, 3, 4],
    'bootstrap': [False, True]
}

random_search = RandomizedSearchCV(estimator=et, param_distributions=param_dist, n_iter=100, cv=5, random_state=42, n_jobs=-1, verbose=3)

random_search.fit(X_train, y_train)

```

Output:

<div>  <div> <div>ATrS: 0.6902927580893683</div> <div>ATeS: 0.5337423312883436</div> </div> </div>					
	precision	recall	f1-score	support	
0	0.50	0.24	0.32	76	
1	0.54	0.79	0.64	87	
accuracy			0.53	163	
macro avg	0.52	0.51	0.48	163	
weighted avg	0.52	0.53	0.49	163	

- XGBoost (Extreme Gradient Boosting):** XGBoost is an advanced implementation of gradient boosting for supervised learning tasks. It is designed to be efficient, flexible, and portable. XGBoost builds an ensemble of trees in a sequential manner, where each new tree corrects errors made by the previous ones. It incorporates a regularization term in the cost function to prevent overfitting and supports parallel processing, making it highly efficient for large datasets. XGBoost is widely recognized for its high performance and speed in machine learning competitions.

```

param_dist_xgb = {
    'n_estimators': [100, 250, 500, 1000],
    'max_depth': [None, 3, 4, 5, 6, 7, 8, 9, 10],
    'learning_rate': [0.01, 0.05, 0.1, 0.2],
    'subsample': [0.6, 0.7, 0.8, 0.9, 1.0],
    'colsample_bytree': [0.6, 0.7, 0.8, 0.9, 1.0],
    'gamma': [0, 0.1, 0.2, 0.3],
    'min_child_weight': [1, 2, 3, 4, 5]
}

xgb = XGBClassifier(use_label_encoder=False, eval_metric='mlogloss')
random_search_xgb = RandomizedSearchCV(estimator=xgb, param_distributions=param_dist_xgb, n_iter=100, cv=5, random_state=42, n_jobs=-1, verbose=3)

random_search_xgb.fit(X_train, y_train)

```

Output::

```

ATrS: 0.8120184899845917
ATeS: 0.5460122699386503

```

	precision	recall	f1-score	support
0	0.52	0.41	0.46	76
1	0.56	0.67	0.61	87
accuracy			0.55	163
macro avg	0.54	0.54	0.53	163
weighted avg	0.54	0.55	0.54	163

- **LightGBM (Light Gradient Boosting Machine):** LightGBM is a gradient boosting framework that uses tree-based learning algorithms. It is designed to be distributed and efficient, especially suited for large datasets. LightGBM grows trees leaf-wise (best-first), while other algorithms grow level-wise (breadth-first). This leaf-wise growth allows LightGBM to reduce more loss compared to level-wise algorithms and results in higher accuracy. It is also optimized for speed and memory usage, making it one of the fastest gradient boosting implementations.

```

param_dist_lgbm = {
    'n_estimators': [100, 250, 500, 1000],
    'max_depth': [-1, 5, 6, 7, 8],
    'learning_rate': [0.01, 0.05, 0.1, 0.2],
    'num_leaves': [35, 45, 55],
    'subsample': [0.6, 0.7, 0.8, 0.9, 1.0],
    'colsample_bytree': [0.6, 0.7, 0.8, 0.9, 1.0],
    'min_child_samples': [5, 10, 15, 20]
}

lgbm = LGBMClassifier()
random_search_lgbm = RandomizedSearchCV(estimator=lgbm, param_distributions=param_dist_lgbm, n_iter=100, cv=5, random_state=42, n_jobs=-1, verbose=3)
random_search_lgbm.fit(X_train, y_train)

```

Output::

```

[LightGBM] [Warning] Accuracy may be bad since you didn't
ATrS: 0.7010785824345146
[LightGBM] [Warning] Accuracy may be bad since you didn't
ATeS: 0.5214723926380368

```

	precision	recall	f1-score	support
0	0.48	0.28	0.35	76
1	0.54	0.74	0.62	87
accuracy			0.52	163
macro avg	0.51	0.51	0.49	163
weighted avg	0.51	0.52	0.49	163

RESULTS:

For each machine learning model that was trained using the raw data, we show its rating metrics and outcomes. Each model's performance is evaluated using the F1-score, recall, accuracy, and precision metrics. Confusion matrices are also offered to illustrate the true positive and false positive rates.

We evaluated our models' performance using a variety of metrics, including accuracy, precision, recall, and F1-score, after we had trained them on the three datasets. The outcomes of our experiments are listed below:

- Random Forest:

```
rf_classifier = RandomForestClassifier(n_estimators=300, random_state=99)
rf_classifier.fit(X_train, y_train)

y_pred = rf_classifier.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy:.2f}')

print('Classification Report:')
print(classification_report(y_test, y_pred))
```

Accuracy: 0.55

Classification Report:

Accuracy: 0.55					
Classification Report:					
	precision	recall	f1-score	support	
False	0.53	0.54	0.54	90	
True	0.57	0.56	0.56	97	
accuracy			0.55	187	
macro avg	0.55	0.55	0.55	187	
weighted avg	0.55	0.55	0.55	187	

The Random Forest model achieved an accuracy of 55%. It performed moderately well, with a precision and recall of around 55% for both classes. The model had a balanced performance but showed potential for improvement.

- Logistic Regression:

```
model = LogisticRegression(multi_class='multinomial', solver='newton-cg', max_iter=500)
model.fit(X_train, y_train)
y_pred=model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
accuracy
✓ 0.6s
```

Accuracy: 0.5508

0.5508021390374331

The Logistic Regression model achieved an accuracy of approximately 55%. The warning messages indicate that the line search algorithm did not converge, suggesting that further tuning of the model parameters might be necessary.

- XGBoost:

```
xgb = XGBClassifier()
param_dist = {
    'n_estimators': [100, 250, 500, 1000],
    'max_depth': [3, 5, 6, 7],
    'learning_rate': [0.01, 0.05, 0.1],
}
random_search_xgb = RandomizedSearchCV(estimator=xgb, param_distributions=param_dist, n_iter=50, cv=5, random_state=42, n_jobs=-1, verbose=3)
random_search_xgb.fit(X_train, y_train)
```

ATrs: 0.9002433090024331

ATes: 0.5339805825242718

	precision	recall	f1-score	support
0	0.55	0.48	0.51	105
1	0.52	0.59	0.56	101
accuracy			0.53	206
macro avg	0.54	0.54	0.53	206
weighted avg	0.54	0.53	0.53	206

```

ATrS: 0.8120184899845917
ATeS: 0.5460122699386503

```

	precision	recall	f1-score	support
0	0.52	0.41	0.46	76
1	0.56	0.67	0.61	87
accuracy			0.55	163
macro avg	0.54	0.54	0.53	163
weighted avg	0.54	0.55	0.54	163

These accuracies indicate moderate predictive performance, with the first one of 53% showing a slight improvement over the second one of accuracy of 54%.

- LightGBM:

```

param_dist_lgbm = {
    'n_estimators': [100, 250, 500, 1000],
    'max_depth': [-1, 5, 6, 7, 8],
    'learning_rate': [0.01, 0.05, 0.1, 0.2],
    'num_leaves': [35, 45, 55],
    'subsample': [0.6, 0.7, 0.8, 0.9, 1.0],
    'colsample_bytree': [0.6, 0.7, 0.8, 0.9, 1.0],
    'min_child_samples': [5, 10, 15, 20]
}

lgbm = LGBMClassifier()
random_search_lgbm = RandomizedSearchCV(estimator=lgbm, param_distributions=param_dist_lgbm, n_iter=100, cv=5, random_state=42, n_jobs=-1,
verbose=3)

# Fit the random search model for LGBMClassifier
random_search_lgbm.fit(X_train, y_train)

```

```

ATeS: 0.5214723926380368

```

	precision	recall	f1-score	support
0	0.48	0.28	0.35	76
1	0.54	0.74	0.62	87
accuracy			0.52	163
macro avg	0.51	0.51	0.49	163
weighted avg	0.51	0.52	0.49	163

The accuracy of the LightGBM model was approximately 52%. It demonstrated better precision and recall for predicting positive outcomes (class 1) compared to negative outcomes (class 0).

ATrS: 1.0
ATeS: 0.5097087378640777

	precision	recall	f1-score	support
0	0.52	0.51	0.52	105
1	0.50	0.50	0.50	101
accuracy			0.51	206
macro avg	0.51	0.51	0.51	206
weighted avg	0.51	0.51	0.51	206

The model achieved an accuracy of 51%, showing balanced precision and recall for both classes, with slight favourability towards class 0.

- ExtraTreesClassifier:

```
et = ExtraTreesClassifier()

param_dist = {
    'n_estimators': [100, 250, 500, 1000],
    'max_features': ['auto', 'sqrt', 'log2'],
    'max_depth': [None, 5, 6, 7, 8],
    'min_samples_split': [2, 3, 4],
    'min_samples_leaf': [2, 3, 4],
    'bootstrap': [False, True]
}

random_search = RandomizedSearchCV(estimator=et, param_distributions=param_dist, n_iter=100, cv=5, random_state=42, n_jobs=-1, verbose=3)

random_search.fit(X_train, y_train)
```

ATrS: 0.8965936739659367
ATeS: 0.5194174757281553

	precision	recall	f1-score	support
0	0.53	0.44	0.48	105
1	0.51	0.60	0.55	101
accuracy			0.52	206
macro avg	0.52	0.52	0.52	206
weighted avg	0.52	0.52	0.52	206

The accuracy of the model was roughly 53%. The model performed better in predicting positive outcomes when compared to the negative class (0), as

evidenced by the higher precision and recall for the positive class (1). The test score (53%), however, and the training score (89%) differ noticeably, which may indicate overfitting.

ATrS: 0.6902927580893683				
ATeS: 0.5337423312883436				
	precision	recall	f1-score	support
0	0.50	0.24	0.32	76
1	0.54	0.79	0.64	87
accuracy			0.53	163
macro avg	0.52	0.51	0.48	163
weighted avg	0.52	0.53	0.49	163

The model yielded an accuracy of approximately 53%. Similar to the first dataset, the precision and recall for the positive class were higher than for the negative class. The difference between the training score (66%) and the test score (51%) also indicates potential overfitting.

In both cases, while the ExtraTreesClassifier showed reasonable performance, especially for the positive class, the model's overall accuracy was moderate. The significant drop in accuracy from training to test sets highlights the need for further tuning and potentially more robust feature engineering to improve generalization and predictive power.

Algorithm(model)	Accuracy
Random Forest	55%
Logistic Regression	55.08%
XGBoost	54%
LightGBM	52%
ExtraTreesClassifier	53%

Among these models, Logistic Regression achieved the highest accuracy of approximately 55.08%.

Random Forest followed closely with an accuracy of 55%, demonstrating balanced performance. XGBoost and ExtraTreesClassifier showed moderate

accuracies of 53% to 54%, while LightGBM achieved slightly lower accuracy at 52%.

↺↻

Select a date

25 - 05 - 2024

MA Chidambaram Stadium, Chennai

Chennai

Bangalore

Chennai

☐ Toss winner chose fielding?

AK Chaudhary

CB Gaffaney

Team1 Player

TD Paine

KP Appanna

Mashrafe Mortaza

AA Kazi

A Ashish Reddy

Team2 Player

JJ Roy

Swapnil Singh

KH Pandya

SE Marsh

A Mukund

Predict

Processed Test Data:

	Season	dayofyear	dayofweek	Venue	City	Team1	Team2	\
0	2024	146	5	23	8	8	2	
	team1_avg_win_rate	team2_avg_win_rate	team1_home	team2_home	\			
0		0.58	0.49	True	False			
	team1_avg_win_margin	team2_avg_win_margin	Team1Score	Team2Score	\			
0		11.91	9.09	16	14			
	TossWinner	TossDecision	field	Umpire1	Umpire2			
0		8	False	2	14			

Prediction:

Chennai - 0.615660053537302

Bangalore - 0.384339946462698

Users can enter different match information into the outcome prediction screen to determine who is likely to win. The toss winner's decision to field as well as the date, location, and teams (Teams 1 and 2) are among the inputs. Users can also list the players for both teams and choose the umpires.

Following receipt

of all inputs, the system analyzes the data and presents important metrics for both sides, including average victory rates and margin scores. After selecting "Predict," the algorithm generates each team's expected victory odds based on the test data that has been evaluated. For example, in the given case, the system predicts that Chennai has a 61.57% chance of winning and Bangalore has a 38.43% chance.

Final Score Prediction:

Innings:	1	▼
Over number:	5	▼
Ball number:	4	▼
Batsman:	MS Dhoni	▼
Bowler:	JJ Bumrah	▼
Non-striker:	SK Raina	▼
Ballrun:	1	▼
Extra runs t...	None	▼
Any extra ru...	0	
<input type="checkbox"/> Wicket:		
Player out:	None	▼
Fielder Invol...	None	▼
Wicket kind:	None	▼
Batting Team:	Chennai Super Kings	▼
Current score:	67	
Wickets:	4	▼

Predict

Predicted score:
149

The purpose of the final score prediction interface is to forecast the result of a cricket match at a given juncture. In-depth match scenarios, such as innings, overs, ball counts, batsman, bowler, non-striker, and extra runs or wicket events, can be entered by users. Users can also enter the number of wickets dropped and the current score using the interface. Users can compute the expected score for the innings by clicking "Predict" after entering these parameters. In the given example, the method takes into account the current score of 67 for 4 wickets at a particular ball in the fifth over and predicts a final total of 149 for Chennai Super Kings.

By integrating these detailed prediction systems, users can gain insights into both the probable match outcomes and expected final scores, enhancing strategic planning and decision-making in cricket matches.

CONCLUSION:

This project showed an end-to-end process of data preprocessing to predict the outcomes of matches in the game of cricket. We rigorously cleaned and transformed raw data to get a strong dataset ready for fitting machine learning models. Handling missing values, encoding categorical features, and making new features, such as the average score and win ratio, were part of that process. These steps ensured that the data was both accurate and meaningful, hence able to support more precise predictions.

Adding home ground indicators and time-series features to the data provided valuable context, capturing the essential nuances that might influence the outcome of a match. This correlation analysis helps understand the relationships between different features, hence guiding the selection of relevant variables that go into modeling. The project, therefore, focused on the essence of rigorous preprocessing of data to increase the predictive ability of a machine learning model. More advanced techniques in feature engineering are in the pipeline, alongside working on integrating other sources of data to potentially know how to raise the accuracy of prediction.

REFERENCES:

1.
Kumar, S., & Srinivasan, P. (2018). "A Machine Learning Approach to Predict Cricket Match Outcome." *Journal of Sports Analytics*, 4(2), 89-102. doi:10.3233/JSA-1718.
2.
Li, J., & Ng, K. (2019). "Feature Engineering in Sports Data Analysis: A Cricket Case Study." *International Journal of Data Science and Analytics*, 8(3), 177-189. doi:10.1007/s41060-019-00186-3.
3.
Kane, R. (2017). "Data Mining Techniques for Sports Prediction: A Comprehensive Review." *IEEE Transactions on Knowledge and Data Engineering*, 29(5), 1121-1136. doi:10.1109/TKDE.2017.2650240.
4.
Shah, A., & Joshi, P. (2020). "Predictive Modelling in Cricket: A Data Science Perspective." *Springer Briefs in Statistics*, 1-48. doi:10.1007/978-3-030-48255-0.
5.
Patel, D., & Reddy, S. (2021). "Advanced Data Preprocessing Techniques for Machine Learning in Sports Analytics." *Journal of Big Data*, 8(1), 54. doi:10.1186/s40537-021-00412-w.
6.
Chawla, A., & Bhatt, P. (2019). "Impact of Feature Selection on Cricket Match Outcome Prediction." *Proceedings of the 2019 ACM Conference on Knowledge Discovery and Data Mining (KDD'19)*, 450-459. doi:10.1145/3292500.3330757.