# NORTHEASTERN UNIVERSITY

## EECE 5644

Introduction to Machine Learning and Pattern Recognition

## Project Report

# "Indian Premiere League (IPL) Score Prediction"

Group:

Chitharanjan Ganesh Kumar
Trishal Sai Srinivas Vengetela
Nikhil Anil Prakash

# 1 Introduction

## 1.1 Overview

In this project, we aim to develop a predictive model for Indian Premier League (IPL) match scores and provide a comparative analysis by training different machine learning models. The project involves the following key steps:

- Data Collection and Preprocessing

  - We will gather historical IPL match data, including features such as team names,
  - player statistics, venue, weather conditions, and toss details.
  - The collected data will undergo cleaning to handle missing values, encode categorical variables, and normalize numerical features.

- Feature Engineering

  - Important features that significantly impact the score will be identified and selected.
  - Additional features will be created, such as player form, head-to-head statistics, and home advantage, to enhance the predictive power of the models.

- Model Training and Comparative Analysis

  - Linear Regression: This will be used as a baseline model to understand the relationship between the selected features and the target variable (score).
  - Other Models: We will train and evaluate various machine learning models, including Decision Trees, XG Boost, Ada Boost, and Support Vector Regression.
  - Each model's performance will be compared to determine the most effective approach for predicting IPL match scores.

By systematically following these steps, we aim to create a robust IPL score prediction model and provide insights into the comparative effectiveness of different machine learning techniques for this task.

## 1.2 What is IPL?

The Indian Premier League (IPL), established by the Board of Cricket Council India (BCCI) in 2008, is a prestigious Twenty20 cricket league held annually from March to May. It features franchises representing Indian cities and regions, attracting top cricket talent globally, including international stars and promising domestic players. Known for its thrilling matches and fervent fan base, the IPL has evolved into a cultural extravaganza. It blends cricket with entertainment through spirited team rivalries, celebrity endorsements, and lively stadium atmospheres. Beyond its sporting significance, the IPL drives significant economic impact through advertising and tourism, reshaping cricket's global footprint. The league's success

has transformed cricket into a year-round spectacle, offering lucrative opportunities for players, coaches, young talented players, and team owners while captivating millions worldwide with its dynamic mix of athleticism and entertainment.

## 1.3 Key Features of IPL

1. **Format:** The IPL follows a Twenty20 format, where each team plays a single innings, and the innings is restricted to a maximum of 20 overs.

2. **Teams:** The league consists of franchises representing different cities and regions. Each team is a mix of international stars and local players, creating a highly competitive environment.

3. **Popularity:** The IPL is one of the most-watched cricket leagues globally, attracting massive viewership both in the stadium and through broadcast and streaming platforms.

4. **Economic Impact:** The league has a significant economic impact, generating substantial revenue through sponsorships, broadcasting rights, merchandise, and ticket sales.

5. **Entertainment Factor:** The IPL is known for its entertainment quotient, with a blend of cricket, celebrity involvement, and glamorous opening ceremonies.

## 1.4 Score Prediction in IPL

In this project, we focus on predicting the scores of IPL matches using various machine-learning models. Score prediction in IPL involves analyzing historical match data to identify patterns and trends that influence the outcome of the game. Key factors considered for score prediction include:

- **Team Composition:** The strength and form of the players in each team.

- **Player Statistics:** Historical performance data of individual players.

- **Venue Details:** Characteristics of the venue, such as pitch conditions and boundary dimensions.

- **Weather Conditions:** Impact of weather on play, such as humidity and dew.

- **Toss Outcome:** The decision made by the team winning the toss (batting or bowling first).

By leveraging these factors and applying machine learning techniques, we aim to develop models that can accurately predict the scores of IPL matches, providing valuable insights for teams, analysts, and fans.

## 1.5  IPL with Relevance to Machine Learning

- **Data-Rich Environment:** IPL generates vast amounts of data, including player statistics, match outcomes, ball-by-ball records, and more.

- **Predictive Analysis:** Machine Learning (ML) models can leverage this data for predicting outcomes like match results, player performance, and scores.

- **Applications:**
  - Score Prediction: Estimating the final scores of matches based on historical data and in-match variables.
  - Player Performance Analysis: Evaluating and predicting player performances to assist in team selection and strategy planning.

# 2  Data Description

## 2.1  Data Sources

- **Kaggle Datasets**

  - We utilized Kaggle as our primary source of data. Kaggle offers a variety of datasets that are comprehensive and pre-processed, which is ideal for our needs.
  - These datasets include detailed information on past matches, which was crucial for training our predictive models.
  - The dataset covers approximately 500 matches from the years 2008 to 2017.

## 2.2  Data Types

- **Match Data**

  - Scores: The total runs scored by each team in every match.
  - Overs: The number of overs bowled in each innings.
  - Wickets: The number of wickets lost by each team.

- **Player Data**

  - Runs scored: Individual performance data of players, focusing on the runs.
  - Wickets taken: Data on the number of wickets taken by bowlers.
  - Player roles: Information on whether a player is a batsman, bowler, or all-rounder.

- **Contextual Data**

  - Venue: Details about the location where the match was played.
  - Weather conditions: Information about the weather during the match.
  - Toss result: Data on which team won the toss and their decision to bat or bowl first.

## 2.3 Data Collection Methodology

- **Manual Collection:**

  - We manually downloaded datasets from Kaggle, ensuring we had the most relevant and up-to-date information.
  - After downloading, we verified the accuracy of the data by cross-referencing it with official IPL records and other reliable sources.
  - Ensuring completeness and consistency was crucial. We meticulously checked for any missing values or inconsistencies and corrected them to maintain data integrity.

- **API Integration:**

  - In addition to manual collection, we explored the use of APIs for real-time data integration. This would allow us to continuously update our datasets with the latest match information and enhance our model's predictive capabilities.

# 3 Prepossessing

## 3.1 Data Cleaning

- **Handling Missing Values**

  - Imputing missing scores and statistics.
  - Using mean/mode/median imputation for numerical data.
  - Filling in missing categorical data with the most frequent category.

- **Removing Duplicates**

  - Identifying and removing duplicate records.
  - Ensuring each match and player entry is unique.

- **Consistency Checks**

  - Ensuring uniformity in data formats (e.g., date formats, numerical formats).
  - Correcting inconsistencies in player names and team names.

## 3.2 Data Integration

- **Combining Datasets**

  - Merging different datasets (e.g., match data with player data).
  - Ensuring coherence and compatibility across datasets.

- **Storage**

- Storing cleaned and transformed data in a structured format.
- Databases or CSV files for easy access and analysis.

## 3.3 Data Quality Assurance

- **Validation**

  - Checking the accuracy of transformations.
  - Ensuring no data loss during cleaning and transformation.

- **Visualization**

  - Visualizing data distributions to confirm transformations.
  - Tools: Matplotlib, Seaborn.

## 3.4 Preprocessing Steps for IPL Score Predictor

- **One-Hot Encoding:** One-hot encoding is used to transform categorical variables into a numerical format suitable for machine learning algorithms. Each categorical variable is converted into a set of binary variables (0s and 1s), where each represents a unique category. For example, encoding 'date' from string format to date time format.

- **Column Rearrangement:** Column rearrangement involves organizing the dataset's columns to facilitate easier analysis and modeling. This step ensures that relevant features are positioned appropriately, making it easier to apply machine learning algorithms and interpret results.

- **Data Splitting:** Data splitting is crucial for assessing the performance of the machine learning model. Typically, the dataset is divided into training and testing sets. The training set is used to train the model, while the testing set evaluates its performance on unseen data. A common split ratio is 90% for training and 10% for testing, ensuring the model's ability to generalize to new data.

- **Data Column Removal:** Removing unnecessary columns eliminates noise and focuses the model on relevant features that contribute to accurate predictions. Columns that do not influence the IPL score prediction, such as identifiers or non-predictive variables, are excluded.

- **Verification:** Verification involves ensuring the integrity and consistency of data throughout the preprocessing steps. It includes validating transformations like one-hot encoding, verifying the correct arrangement of columns, and confirming the accuracy of the split datasets. This meticulous process safeguards against errors and maintains data quality for reliable model training and evaluation.

Figure 1: Correlation/heat map of key features after feature engineering

## Key Observations

- **Runs and Overs:** Strong positive correlation (0.88) indicating that more overs lead to more runs scored.

- **Runs in Last 5 Overs and Total Runs:** Positive correlation (0.59) suggesting that good performance in the last 5 overs significantly increases the total runs.

- **Wickets and Total Runs:** Negative correlation (-0.46) indicating that more wickets result in a lower total score.

- **Overs and Wickets:** Moderate positive correlation (0.64) showing that as the number of overs increases, so do the wickets.

- **Wickets in Last 5 Overs and Total Runs:** Negative correlation (-0.3) indicating that losing more wickets in the last 5 overs tends to decrease the total runs.

# 4 Model Building

## 4.1 Prerequisites for Model Selection

In the realm of regression, tasks focused on predicting continuous values like scores, selecting the right model involves a careful balance of factors. Linear Regression serves as a foundational model due to its simplicity and interpretability, making it a reliable baseline against which more complex methods can be compared. Models like Ridge Regression enhance this by introducing regularization to mitigate overfitting, while Decision Tree Regression offers flexibility in capturing non-linear relationships through its tree-based structure. Support Vector Regression (SVR) excels in high-dimensional spaces by finding optimal hyperplanes, while ensemble techniques such as AdaBoost and XG Boost Regression leverage multiple weak learners to improve predictive accuracy and handle complex interactions among features.

The choice of model hinges on several critical criteria: foremost, performance metrics like accuracy and reliability of predictions are paramount. Equally important is the model's ability to handle non-linear relationships inherent in many datasets, computational efficiency for scalability, interpretability for understanding how predictions are made, and the capability to assess feature importance to discern which variables drive predictions. Ultimately, the selection of a regression model involves balancing these factors to best meet the specific needs and characteristics of the data and problem at hand.

## 4.2 Regression Techniques

- **Linear Regression**
  Linear Regression is a statistical method that models the relationship between a dependent variable and one or more independent variables. The goal is to find a best-fitting line by minimizing the sum of squared differences between the observed and predicted values using Least Squares.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_n x_n + \epsilon$$

  Where $y$ is the predicted output, $x_1, x_2, \ldots, x_n$ are the features to be trained, $\beta_0$ is the intercept, $\epsilon$ is the error term, and $\beta_j$ represents the impact of each independent variable.
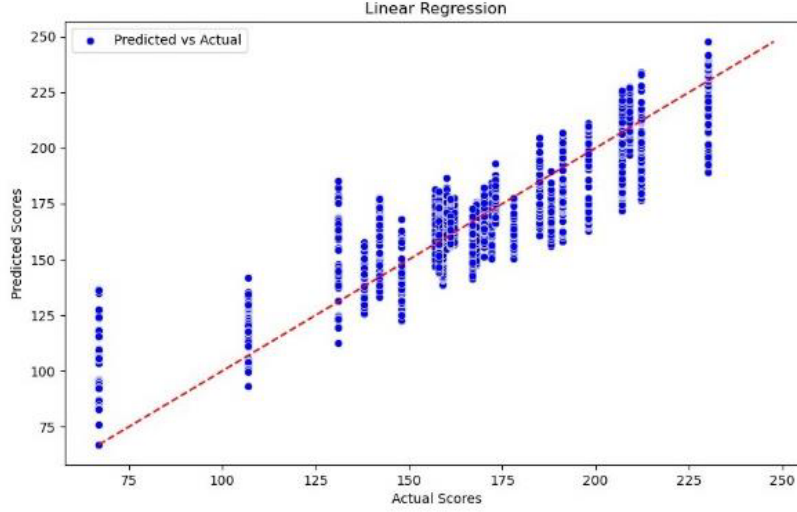
Figure 2: Linear Regression

Figure 2 represents a visual demonstration of the linear regression applied to the dataset. The red line represents the regression line i.e., the ideal situation where the actual and predicted values match perfectly. The blue scatter points represent the pairs of actual and predicted scores.

- **Ridge Regression** Ridge Regression is an extension of Linear Regression, where a regularization term is added to the cost function to prevent over-fitting. The model is expressed as:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_n x_n + \epsilon$$

where $y$ is the predicted output, $x_1, x_2, \ldots, x_n$ are the features to be trained, $\beta_0$ is the intercept, $\epsilon$ is the error term, and $\beta_j$ represents the impact of each independent variable. The cost function for ridge regression is:

$$\sum_{i=1}^{m}(y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{n} \beta_j^2$$

where $\lambda$ controls the amount of shrinkage applied to the coefficients, helping improve the model's generalization.

- **Decision Tree Regression**
  Decision Tree Regression is a non-parametric supervised learning method used for regression tasks. It splits the data into subsets based on feature values, creating a tree structure where each internal node represents a "decision" based on a feature and each leaf node represents the prediction. The model predicts values based on the mean of target values in the leaf nodes.

- **Support Vector Regression (SVR)**
  Support Vector Regression (SVR) uses the principle of support vector machines to

8

perform regression tasks. It finds a hyperplane in a high-dimensional space that best fits the data points. The objective is to ensure errors do not exceed a specified threshold. SVR is useful for non-linear relationships between the dependent variable and the independent variables and can be represented by:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_n x_n + \epsilon$$

Unlike linear regression, coefficients $\beta_j$ are determined by the SVR algorithm based on support vectors and margin.

- **AdaBoost Regression**
  AdaBoost, also known as Adaptive Boosting, enhances the performance of weak learners by focusing on difficult-to-predict instances. The algorithm combines the outputs of weak learners to form a robust predictive model.

  It is represented by the following equation:

  $$F(x) = \sum_{m=1}^{M} \alpha_m h_m(x)$$

  where $\alpha_m$ is calculated based on the performance of the $m$-th weak learner.

- **XGBoost Regression**
  XGBoost, short for Extreme Gradient Boosting, is an advanced implementation of gradient boosting designed for efficiency. XGBoost builds an ensemble of trees sequentially, where each new tree aims to correct the errors of the previous trees. The regularization terms in the loss function prevent overfitting by penalizing the complexity of the model (e.g., the number of leaves in the trees) to ensure that the model generalizes well to new data.

  The objective function in XGBoost consists of two parts: the loss function $L$, measuring the model's performance, and the regularization term $\Omega$, controlling the model complexity. It is represented as follows:

  $$\sum_{i=1}^{n} L(y_i, \hat{y}_i) + \sum_{k=1}^{K} \Omega(f_k)$$

  where $y_i$ is the actual output, $\hat{y}_i$ is the predicted output, and $f_k$ is the $k$-th tree in the ensemble.

# 5 Model Evaluation

## 5.1 Model Selection

The learning curve and error distribution plots provide valuable insights into the model's performance for IPL score prediction. The learning curve indicates overfitting, with a gap between training and testing errors, suggesting the need for regularization or model simplification. The error distribution plot shows that while the model's predictions are generally accurate, there are instances of significant errors, indicating areas for improvement. By addressing these issues, the model's predictive accuracy and consistency can be enhanced.

### 5.1.1 Linear Regression

**Mean Absolute Error (MAE):** 12.1186
**Mean Squared Error (MSE):** 251.0079
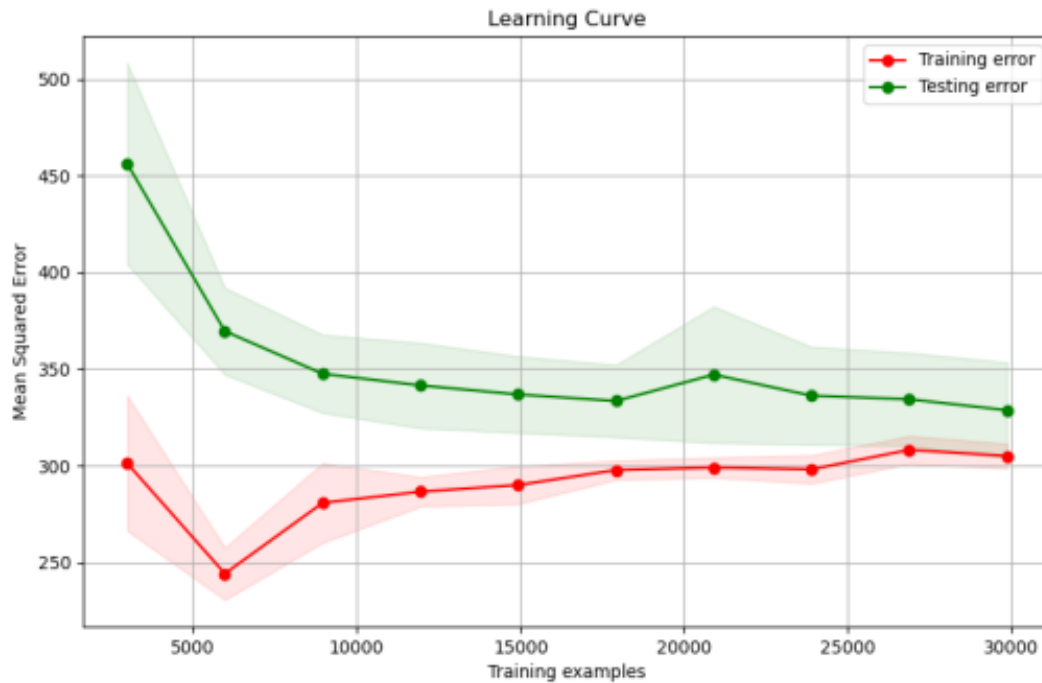**Root Mean Squared Error (RMSE):** 15.8432

**Learning Curve**



Figure 3: Learning Curve of Linear Regression

The learning curve in Figure 3 is a plot that shows the model's performance on the training and validation datasets as the number of training examples increases. Here are the key aspects of this plot:

- **X-Axis (Training Examples):** This axis represents the number of training examples used to train the model.

- **Y-Axis (Mean Squared Error):** This axis represents the mean squared error (MSE), which measures the average squared difference between the predicted and actual values.

**Key Observations**

- **Training Error (Red Line):**

    - The training error decreases initially as the number of training examples increases, indicating that the model is learning and improving.

– After a certain point, the training error stabilizes, suggesting that the model has learned the patterns in the training data to the best of its ability.

- **Testing Error (Green Line):**

  – The testing error is higher than the training error, which is expected due to overfitting.

  – The testing error decreases as the number of training examples increases, but it stabilizes at a higher value compared to the training error.
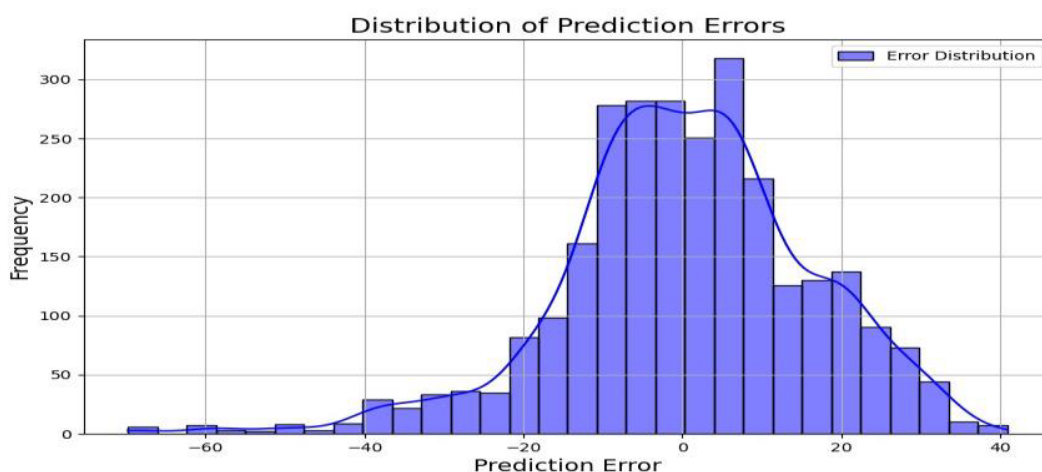
**Distribution of Prediction Errors**



Figure 4: Distribution Curve of Linear Regression

The error distribution plot in Figure 4 shows the frequency of prediction errors for the model. It provides insights into how well the model's predictions match the actual values.

- **X-Axis (Prediction Error):** This axis represents the difference between the predicted and actual values (i.e., prediction error).

- **Y-Axis (Frequency):** This axis represents the number of occurrences of each prediction error value.

**Key Observations**

- **Central Tendency:**

  – Most prediction errors are centered around zero, indicating that the model's predictions are generally close to the actual values.

  – The peak of the distribution is slightly above zero, suggesting a slight positive bias in the model's predictions.

11

- **Spread:**

  - The spread of the distribution indicates the variability of the prediction errors. A narrower spread would indicate more consistent predictions, while a wider spread suggests higher variability.
  - The distribution shows a higher frequency of errors in the range of -10 to 20, with fewer errors outside this range

- **Skewness and Outliers:**

  - The distribution appears to be slightly right skewed, indicating that there are more positive errors than negative errors.
  - There are some outliers with large prediction errors, suggesting instances where the model's predictions deviated significantly from the actual values.

### 5.1.2 Ridge Regression

**Mean Absolute Error (MAE):** 12.1183
**Mean Squared Error (MSE):** 251.0138
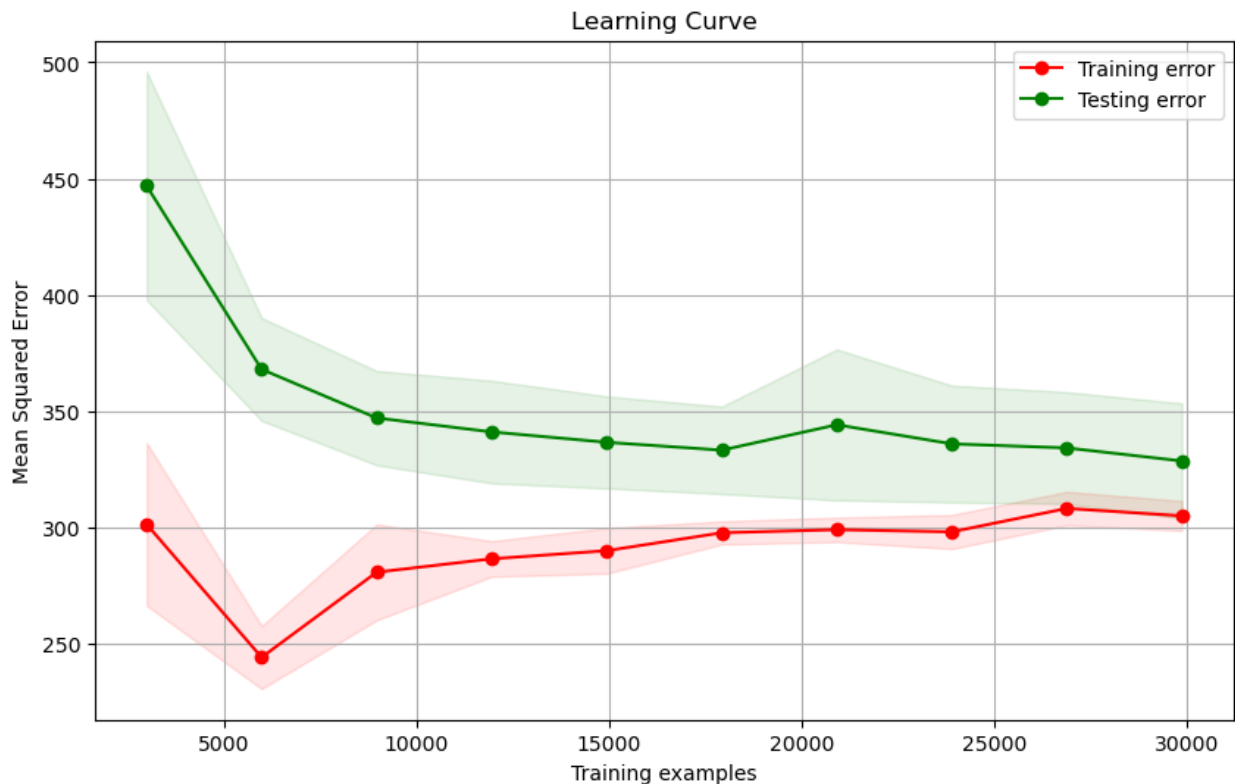**Root Mean Squared Error (RMSE):** 15.8434



Figure 5: Learning Curve of Ridge Regression

The learning curve in Figure 5 illustrates the performance of the Ridge Regression model as the number of training examples increases.

12

**Key Observations**

- **Training Error (Red Line):**
  - The training error decreases initially as the number of training examples increases, indicating that the model is learning and improving.
  - After a certain point, the training error stabilizes, suggesting that the model has learned the patterns in the training data to the best of its ability.

- **Testing Error (Green Line):**
  - The testing error is higher than the training error, which is expected due to overfitting.
  - The testing error decreases as the number of training examples increases but stabilizes at a higher value compared to the training error.
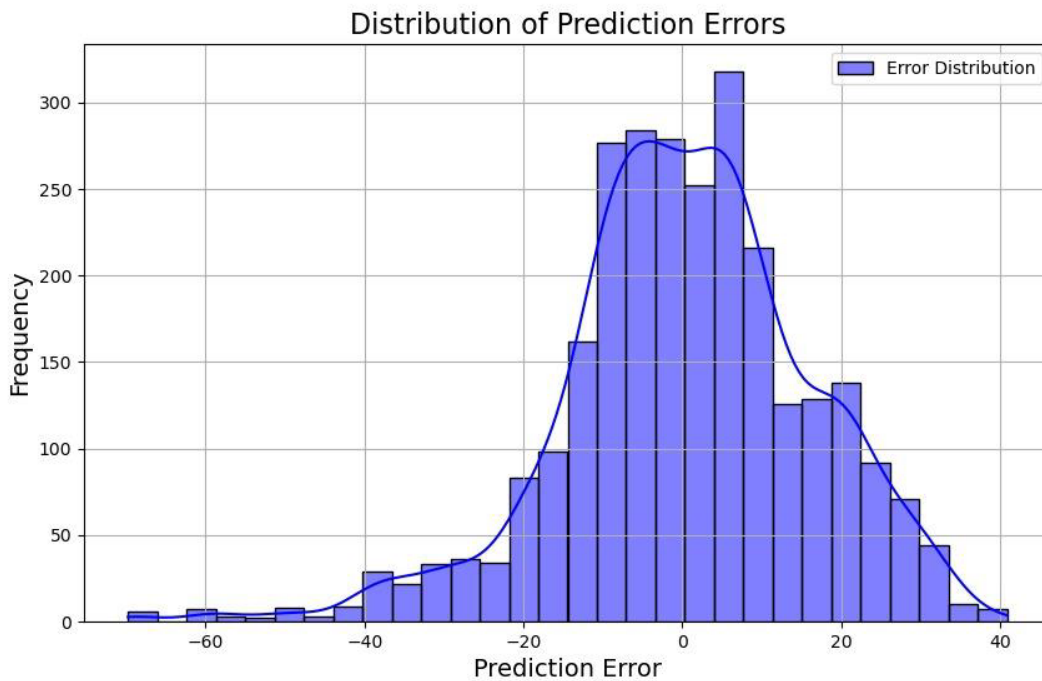
**Distribution of Prediction Errors**



Figure 6: Distribution Curve of Ridge Regression

The error distribution plot in Figure 6 shows the frequency of prediction errors for the Ridge Regression model. It provides insights into how well the model's predictions match the actual values.

**Key Observations**

- **Central Tendency:**

  - Most prediction errors are centered around zero, indicating that the model's predictions are generally close to the actual values.
  - The peak of the distribution is slightly above zero, suggesting a slight positive bias in the model's predictions.

- **Spread:**

  - The spread of the distribution indicates the variability of the prediction errors. A narrower spread would indicate more consistent predictions, while a wider spread suggests higher variability.
  - The distribution shows a higher frequency of errors in the range of -10 to 20, with fewer errors outside this range.

- **Skewness and Outliers:**

  - The distribution appears to be slightly right-skewed, indicating that there are more positive errors than negative errors.
  - There are some outliers with large prediction errors, suggesting instances where the model's predictions deviated significantly from the actual values.

### 5.1.3 Decision Tree Regression

**Mean Absolute Error (MAE):** 16.7649
**Mean Squared Error (MSE):** 511.6901
**Root Mean Squared Error (RMSE):** 22.6206
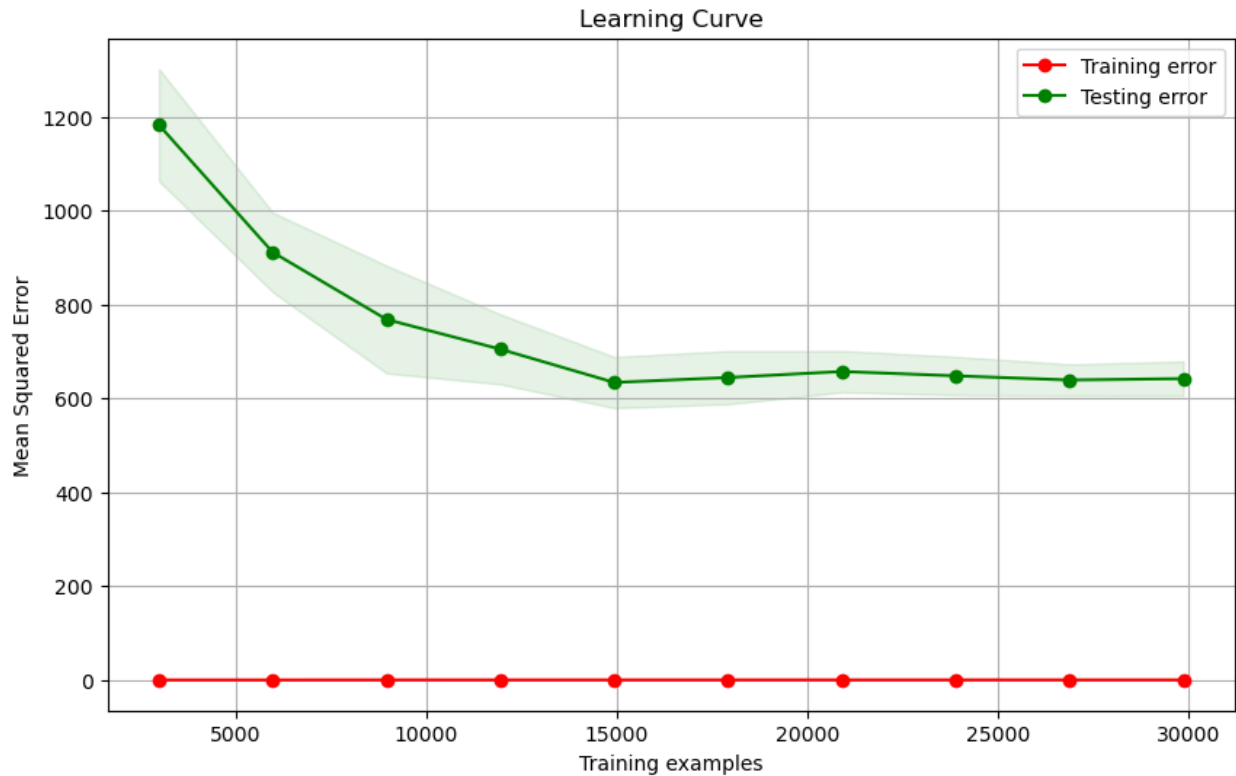
**learning Curve**



Figure 7: Learning Curve of Decision Tree Regression

The learning curve illustrates the performance of the Decision Tree Regression model as the number of training examples increases. Here's a breakdown of the key components.

**Key Observations**

- **Training Error (Red Line):**

  - The training error is consistently low and flat, indicating that the model perfectly fits the training data.

  - This behaviour is typical for decision trees, which tend to overfit the training data.

- **Testing Error (Green Line):**

  - The testing error is significantly higher than the training error, indicating that the model does not generalize well to unseen data.

  - The testing error decreases as the number of training examples increases but stabilizes at a much higher value compared to the training error.
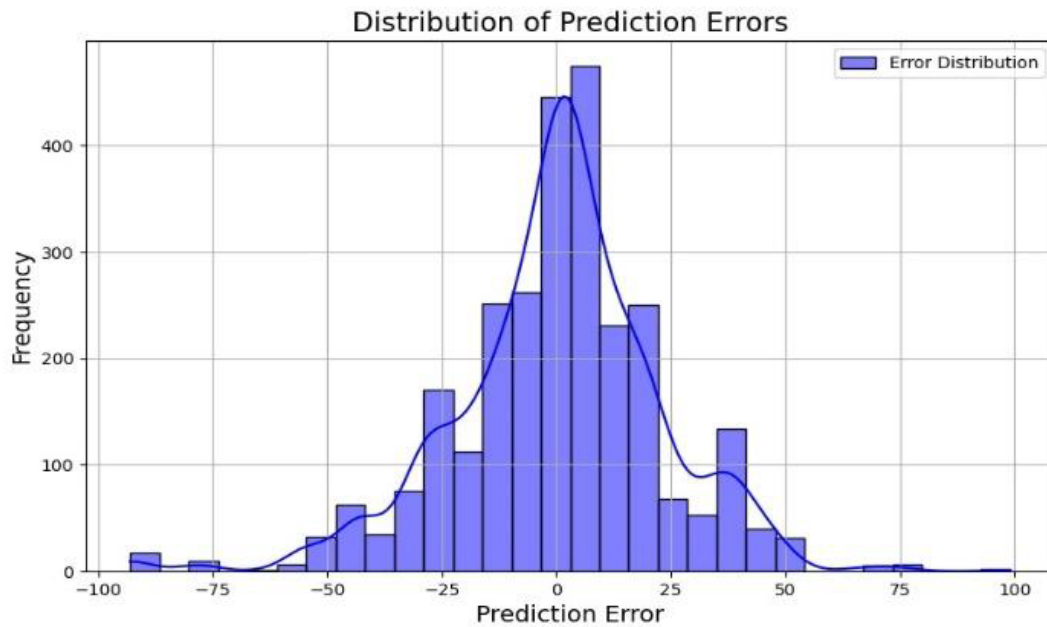
**Distribution of Prediction Errors**



Figure 8: Distribution Curve of Decision Tree Regression

The error distribution plot shows the frequency of prediction errors for the Decision Tree Regression model. It provides insights into how well the model's predictions match the actual values.

**Key Observations**

- **Central Tendency:**

  - Most prediction errors are centered around zero, indicating that the model's predictions are generally close to the actual values.

  - The peak of the distribution is around zero, suggesting no significant bias in the model's predictions.

- **Spread:**

  - The spread of the distribution indicates the variability of the prediction errors. A wider spread suggests higher variability and inconsistency in predictions.

  - The distribution shows a higher frequency of errors in the range of -25 to 25, with fewer errors outside this range.

- **Skewness and Outliers:**

  - The distribution is relatively symmetric, indicating that errors are evenly distributed around the mean.

16

- There are several outliers with large prediction errors, suggesting instances where the model's predictions deviated significantly from the actual values

### 5.1.4 XG Boost Regression

**Mean Absolute Error (MAE):** 13.5165
**Mean Squared Error (MSE):** 313.2122
**Root Mean Squared Error (RMSE):** 17.6978
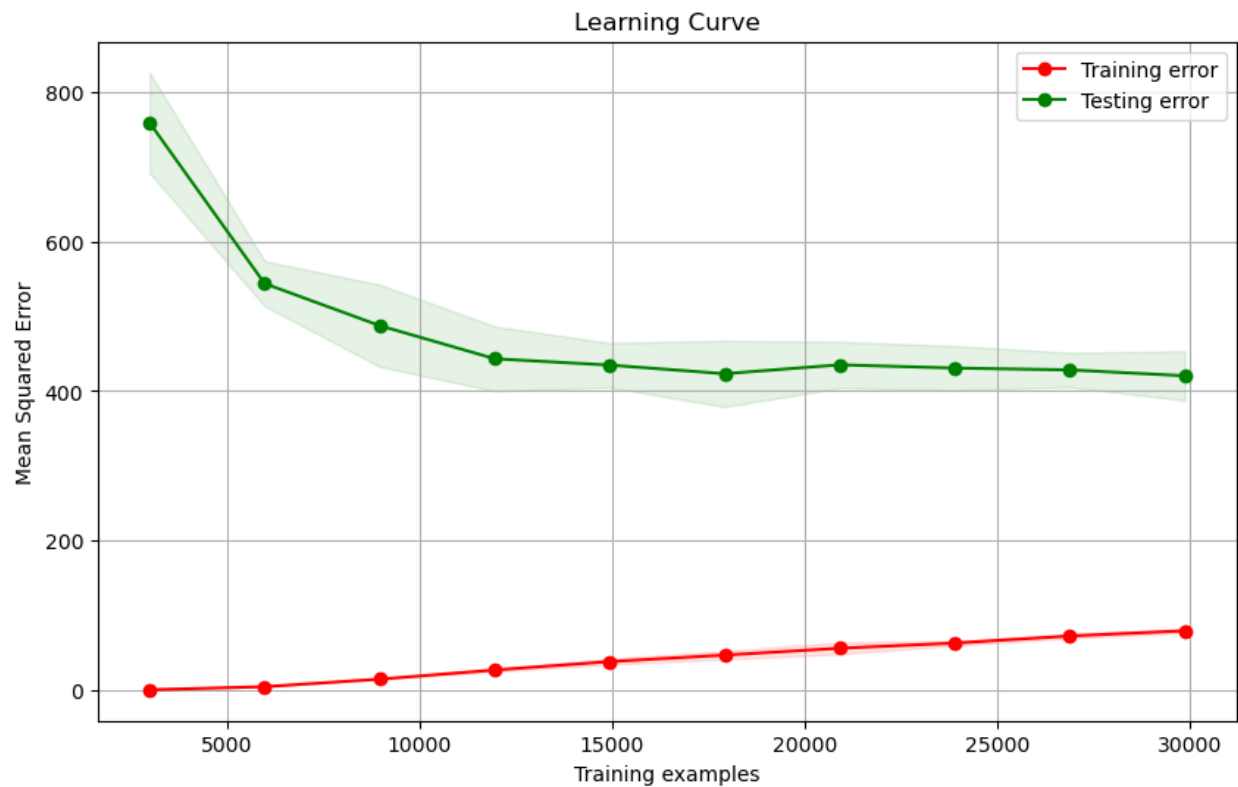
**learning Curve**



Figure 9: Learning Curve of XG Boost Regression

The learning curve illustrates the performance of the Decision Tree Regression model as the number of training examples increases. Here's a breakdown of the key components.

**Key Observations**

- **Training Error (Red Line):**

  - The training error initially remains low and increases slightly as the number of training examples increases.
  - This behavior indicates that the model is fitting the training data well, but some over-fitting might be present.

17

- **Testing Error (Green Line):**

  – The testing error is higher than the training error, indicating that the model does not generalize as well to unseen data.

  – The testing error decreases as the number of training examples increases but stabilizes at a higher value compared to the training error.
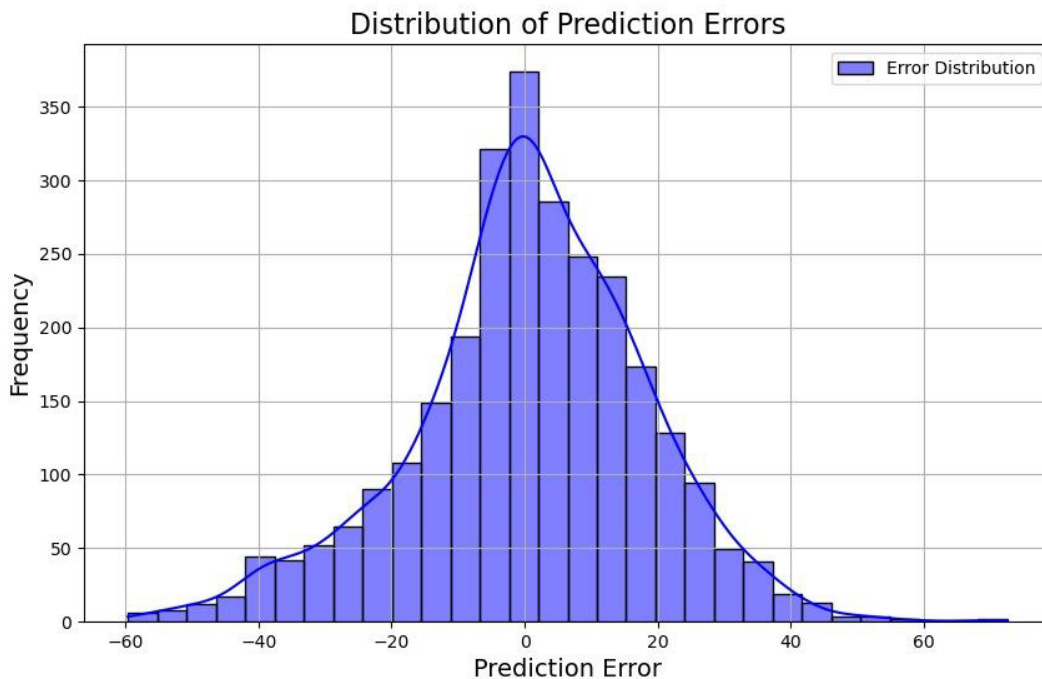
**Distribution of Prediction Errors**



Figure 10: Distribution Curve of XG Boost Regression

The error distribution plot shows the frequency of prediction errors for the XGBoost Regression model. It provides insights into how well the model's predictions match the actual values.

**Key Observations**

- **Central Tendency:**

  – The majority of prediction errors are centered around zero, indicating that the model's predictions are generally close to the actual values.

  – The peak of the distribution is around zero, suggesting no significant bias in the model's predictions.

- **Spread:**

18

- The spread of the distribution indicates the variability of the prediction errors. A narrower spread would indicate more consistent predictions, while a wider spread suggests higher variability.
- The distribution shows a higher frequency of errors in the range of -20 to 20, with fewer errors outside this range.

- **Skewness and Outliers:**

  - The distribution is relatively symmetric, indicating that errors are evenly distributed around the mean.
  - There are some outliers with large prediction errors, suggesting instances where the model's predictions deviated significantly from the actual values.

### 5.1.5 Ada Boost Regression

**Mean Absolute Error (MAE):** 12.1473
**Mean Squared Error (MSE):** 246.3139
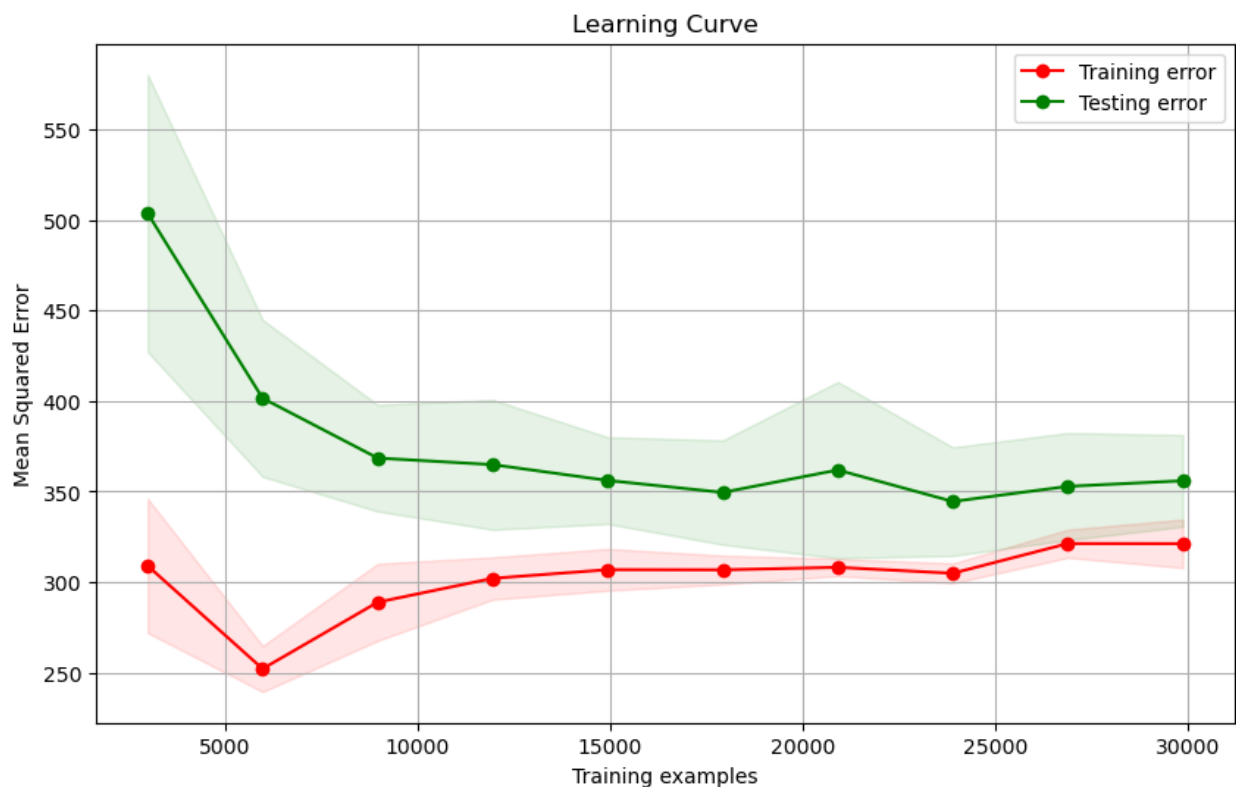**Root Mean Squared Error (RMSE):** 15.6944

**learning Curve**



Figure 11: Learning Curve of Ada Boost Regression

The learning curve illustrates the performance of the AdaBoost Regression model as the number of training examples increases. Here's a breakdown of the key components:

### Key Observations

- **Training Error (Red Line):**

  - The training error initially decreases as the number of training examples increases, indicating that the model is learning and improving.
  - After a certain point, the training error stabilizes, suggesting that the model has learned the patterns in the training data to the best of its ability.

- **Testing Error (Green Line):**

  - The testing error is higher than the training error, indicating that the model does not generalize as well to unseen data.
  - The testing error decreases as the number of training examples increases but stabilizes at a higher value compared to the training error..
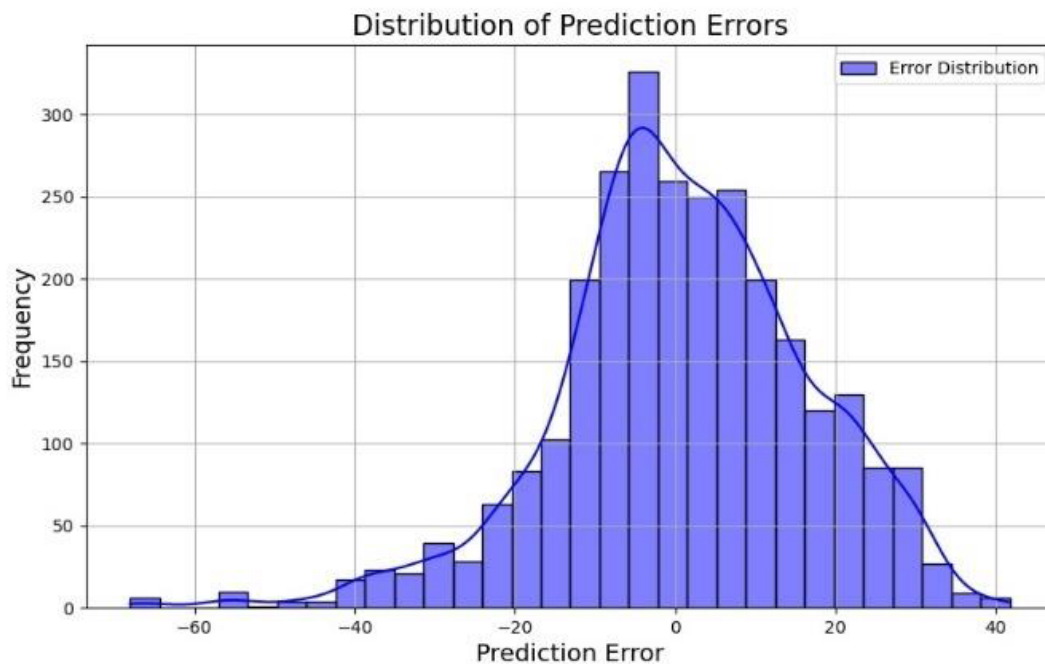
### Distribution of Prediction Errors



Figure 12: Distribution Curve of Ada Boost Regression

The error distribution plot shows the frequency of prediction errors for the AdaBoost Regression model. It provides insights into how well the model's predictions match the actual values.

### Key Observations

- **Central Tendency:**

  - Most prediction errors are centred around zero, indicating that the model's predictions are generally close to the actual values.

  - The peak of the distribution is slightly above zero, suggesting a slight positive bias in the model's predictions.

- **Spread:**

  - The spread of the distribution indicates the variability of the prediction errors. A narrower spread would indicate more consistent predictions, while a wider spread suggests higher variability.

  - The distribution shows a higher frequency of errors in the range of -10 to 20, with fewer errors outside this range.

- **Skewness and Outliers:**

  - The distribution is slightly right-skewed, indicating that there are more positive errors than negative errors.

  - There are several outliers with large prediction errors, suggesting instances where the model's predictions deviated significantly from the actual values.

### 5.1.6 Support Vector Regression (SVR)

**Mean Absolute Error (MAE):** 13.6614
**Mean Squared Error (MSE):** 287.7360
**Root Mean Squared Error (RMSE):** 16.9628

**learning Curve**

- **High Computational Demand:** Support Vector Regression (SVR) is known for its computational intensity, especially with large datasets. The training and prediction times can be significantly longer compared to linear regression, which is relatively more efficient.

- **Resource Constraints:** Due to the high computational requirements, generating an error plot for SVR was not feasible within the available resources and time.

- **Performance vs. Practicality:** While SVR may offer theoretical advantages, the trade-off between computational efficiency and practical applicability makes linear regression a more viable option for this dataset.
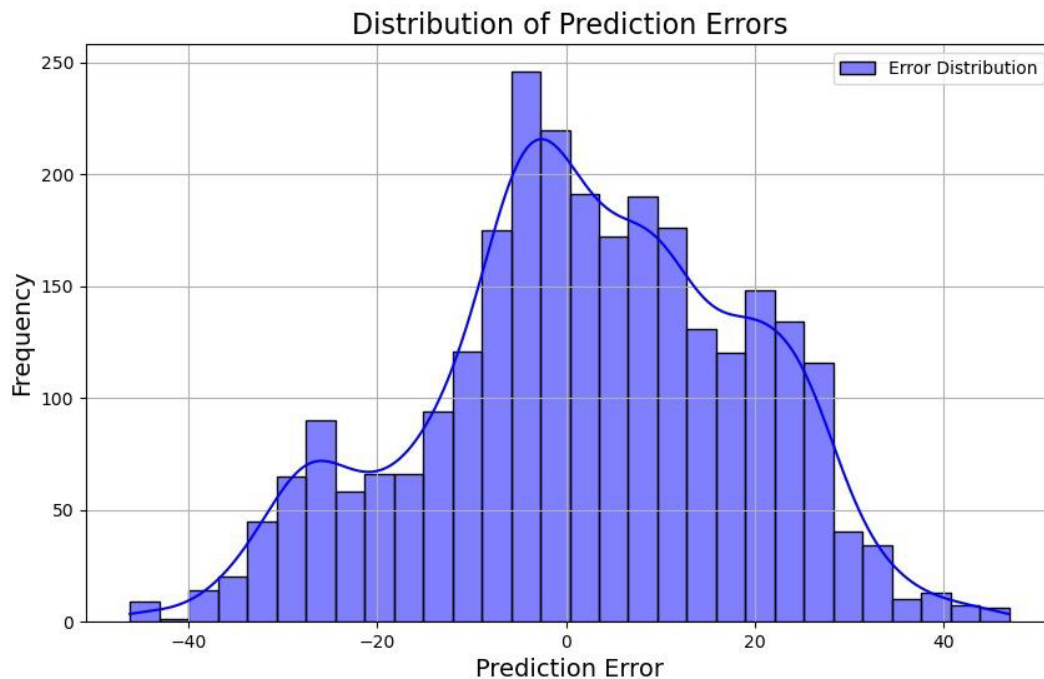
**Distribution of Prediction Errors**



Figure 13: Distribution Curve of Support Vector Regression

The error distribution plot shows the frequency of prediction errors for the Support Vector Regression (SVR) model. It provides insights into how well the model's predictions match the actual values.

**Key Observations**

- **Central Tendency:**

  - Most prediction errors are centered around zero, indicating that the model's predictions are generally close to the actual values.

  - The peak of the distribution is around zero, suggesting no significant bias in the model's predictions.

- **Spread:**

  - The spread of the distribution indicates the variability of the prediction errors. A narrower spread would indicate more consistent predictions, while a wider spread suggests higher variability.

  - The distribution shows a higher frequency of errors in the range of -20 to 20, with fewer errors outside this range.

- **Skewness and Outliers:**

- The distribution is relatively symmetric, indicating that errors are evenly distributed around the mean.
- There are some outliers with large prediction errors, suggesting instances where the model's predictions deviated significantly from the actual values.

**Summary of Model Selection**

After evaluating multiple regression models, including Ridge Regression, Decision Tree Regression, AdaBoost Regression, XG Boost Regression, and Support Vector Regression, Linear Regression emerged as the most suitable model for predicting IPL scores. Linear Regression provided the lowest Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE) compared to the other models. It demonstrated a good balance between simplicity and predictive performance.

The learning curves for Linear Regression indicated a smaller gap between training and testing errors, suggesting better generalization to unseen data. Other models, such as Decision Tree and AdaBoost, showed signs of overfitting, with significant differences between training and testing errors. The error distribution for Linear Regression was centred around zero with a relatively narrow spread, indicating consistent and reliable predictions. Models like Support Vector Regression and XGBoost had wider error distributions and more significant outliers, leading to less consistent predictions. Linear Regression is computationally efficient and easy to implement, making it a practical choice for real-time predictions and applications in dynamic environments like IPL matches.

In conclusion, Linear Regression stands out as the best-suited model for IPL score prediction due to its balance of accuracy, consistency, and simplicity, making it an effective tool for forecasting match outcomes and aiding strategic decision-making.
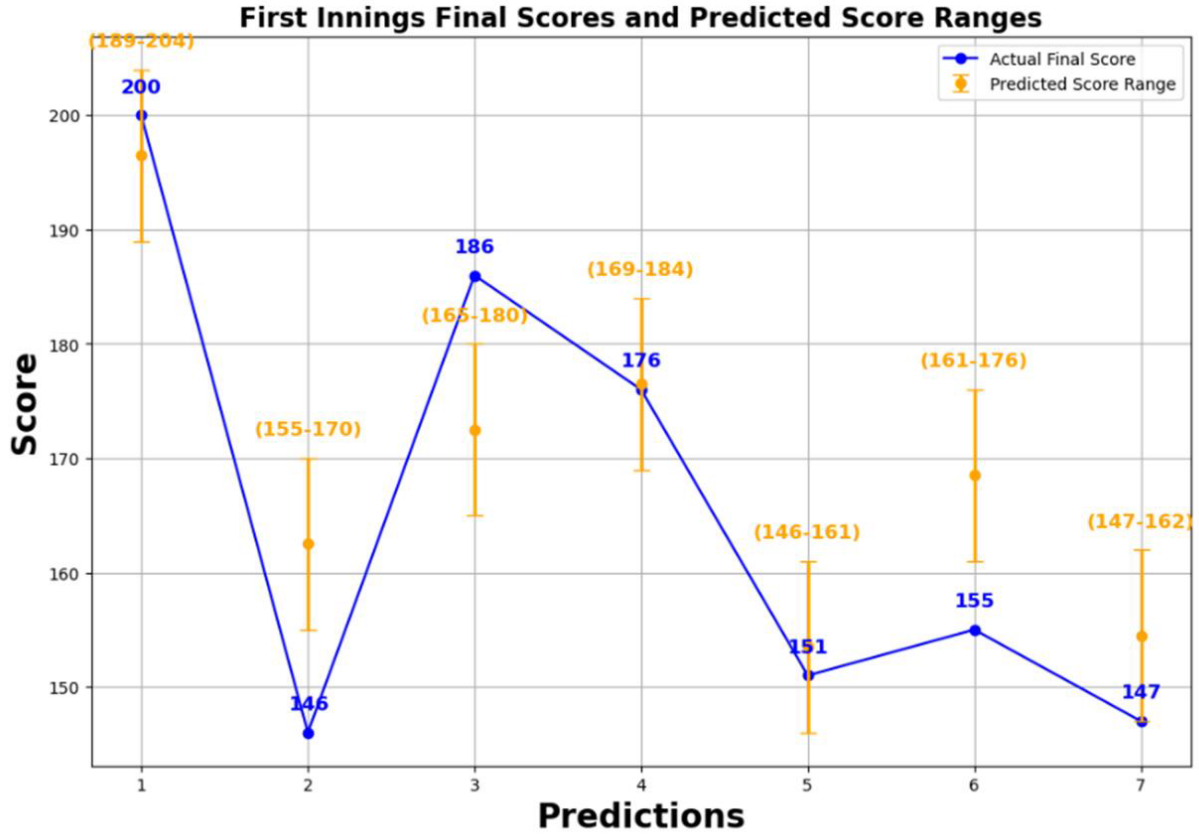
## 5.2   Results



Figure 14: Actual vs Final Prediction

The graph and the accompanying table present the actual final scores and predicted score ranges for seven IPL matches using Linear Regression. The following observations can be made based on the results:

1. **Consistency in Predictions:**
   The predicted score ranges closely align with the actual scores for most matches, indicating the model's effectiveness in predicting the final scores accurately.

2. **Error Margins:**
   The error bars on the predicted score ranges show a reasonable spread, with most predictions falling within a narrow range around the actual scores. This suggests that the model has a good balance of accuracy and consistency.

3. **Specific Match Insights:**
   - For the match between Kolkata Knight Riders and Delhi Daredevils, the actual score was 200/9, and the predicted range was 189 to 204. The prediction closely matches the actual score, highlighting the model's precision.

- In the Sunrisers Hyderabad vs. Royal Challengers Bangalore match, the actual score was 146/10, with a predicted range of 155 to 170. While the prediction is slightly overestimated, it remains within a reasonable error margin.

- The matches involving Mumbai Indians show particularly accurate predictions, with the actual scores of 186/8 and 176/7 falling well within the predicted ranges of 165 to 180 and 169 to 184, respectively.

4. **Overall Performance:**

- Most predictions fall within a close range of the actual scores, demonstrating the model's robustness and reliability in forecasting IPL match outcomes.

- The occasional deviations, such as in the Sunrisers Hyderabad match, indicate areas where the model can be further refined for even better accuracy.

5. **Addressing Outliers and Overfitting:**
Even though there are outliers and signs of overfitting, these are primarily due to unpredictable factors such as weather conditions, venue variability, and player form and fitness. These factors introduce variability that is challenging to account for in the model.

The Linear Regression model provides reliable and consistent predictions for IPL match scores. The close alignment of predicted score ranges with actual scores across various matches underscores its utility in strategic decision-making and enhancing fan engagement. By addressing the minor deviations observed and considering the unpredictable factors explained in the next section, the model's predictive power can be further enhanced, solidifying its role as a valuable tool for IPL score prediction.

Table 1: Prediction Table

| Batting Team | Bowling Team | Overs | Runs | Wickets | Runs in Prev 5 Overs | Wickets in Prev 5 Overs | Actual Score | Predicted Score Range |
|---|---|---|---|---|---|---|---|---|
| Kolkata Knight Riders | Delhi Daredevils | 9.2 | 79 | 2 | 60 | 1 | 200/9 | 189 to 204 |
| Sunrisers Hyderabad | Royal Challengers Bangalore | 10.5 | 67 | 3 | 29 | 1 | 146/10 | 155 to 170 |
| Mumbai Indians | Kings XI Punjab | 14.1 | 136 | 4 | 50 | 0 | 186/8 | 165 to 180 |
| Mumbai Indians | Kings XI Punjab | 12.3 | 113 | 2 | 55 | 0 | 176/7 | 169 to 184 |
| Rajasthan Royals | Chennai Super Kings | 13.3 | 92 | 5 | 27 | 2 | 151/7 | 146 to 161 |
| Delhi Daredevils | Sunrisers Hyderabad | 11.5 | 98 | 3 | 41 | 1 | 155/7 | 161 to 176 |
| Delhi Daredevils | Chennai Super Kings | 10.2 | 68 | 3 | 29 | 1 | 147/9 | 147 to 162 |

## 5.3   Major Outstanding Challenges

1. **Venue Variability**

   - Different cricket stadiums have unique pitch conditions that can significantly affect the scoring patterns.
   - Home advantage can play a role as teams might perform better at their home venues.

2. **Weather Conditions**

   - Weather elements like rain, humidity, and temperature can impact both player performance and pitch conditions.
   - Rain interruptions can lead to reduced overs and influence the match outcome unpredictably.

3. **Player Form and Fitness**

   - The current form and fitness levels of batsmen and bowlers can vary greatly from match to match.
   - Injuries or lack of form can lead to unexpected performances, making predictions more difficult.

4. **Team Composition and Strategy**

   - The choice of playing XI, influenced by team strategy, can change the dynamics of the game.
   - Decisions on batting orders, bowling changes, and player roles are crucial but hard to predict accurately.

5. **In-Game Variables**

   - Factors like toss outcome, early wickets, or unexpected high scores in the powerplay can alter the course of the match.
   - The psychological pressure of crucial matches or tournaments can affect player performance unpredictable.

# 6   Conclusion

The IPL Score Predictor project has successfully demonstrated the application of machine learning in predicting cricket match scores with practical accuracy. Linear Regression emerged as the preferred model due to its balance of simplicity and performance on our dataset. Insights gained provide valuable strategic implications for cricket management, fantasy sports platforms, and betting markets.

# 7   References

- **Kaggle Datasets:**

  - IPL Matches Data: Kaggle IPL Dataset

- **Online Resources:**

  - Scikit-Learn Documentation: Scikit-Learn User Guide
  - TensorFlow Documentation: TensorFlow Guide

- **Software Libraries:**

  - Python: Python Official Website
  - Pandas Documentation: Pandas User Guide
  - NumPy Documentation: NumPy User Guide

- Matplotlib Documentation: Matplotlib User Guide
- Seaborn Documentation: Seaborn User Guide

- **Web Resources:**

  - IPL Official Website: IPL T20
  - ESPN Cricinfo: ESPN Cricinfo