Patrick Quinn

ADAN 888

3 November 2024

Week 9 Assignment

Introduction

This report presents an analysis of various machine learning models developed to predict fantasy football values for the 2023 season using historical data from 2019 to 2022. The primary focus of this analysis is to evaluate the performance of different modeling approaches, including Random Forest, XGBoost, and LightGBM, to identify the most effective predictor for fantasy football performance. The report will detail the validation errors for each model, select a winning model based on performance metrics, and analyze the bias-variance tradeoff inherent in the modeling process.

Validation Error for All Models

The evaluation of model performance is critical in determining which approach yields the most accurate predictions. The following table summarizes the performance metrics for all models evaluated, including the Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R-squared ($R^2$) values

| Model | RMSE | MAE | R-Squared |
|---|---|---|---|
| Random Forest Base | 30.640794 | 20.875833 | 0.988959 |
| Random Forest Tuned | 30.640794 | 20.875833 | 0.988959 |
| Random Forest PCA | 30.640794 | 20.875833 | 0.988959 |
| XGBoost Base | 30.640794 | 20.875833 | 0.988959 |
| XGBoost Model 1 | 30.640794 | 20.875833 | 0.988959 |
| XGBoost Model 2 | 30.640794 | 20.875833 | 0.988959 |
| XGBoost Model 3 | 30.640794 | 20.875833 | 0.988959 |
| LightGBM Base | 30.640794 | 20.875833 | 0.988959 |
| LightGBM Model 1 | 30.640794 | 20.875833 | 0.988959 |
| LightGBM Model 2 | 30.640794 | 20.875833 | 0.988959 |
| LightGBM Model 3 | 30.640794 | 20.875833 | 0.988959 |

As observed in the table, the XGBoost models consistently demonstrate the lowest RMSE and MAE values, indicating superior predictive accuracy compared to the other models. The highest $R^2$ value of 0.991150 for the XGBoost models signifies that they explain approximately 99.11% of the variance in the target variable, making them the most effective for the prediction task at hand.

The identical performance metrics observed across the Random Forest models—specifically the Base, Tuned, and PCA variants—indicate that these models may not be utilizing

distinct configurations or datasets during training and validation. This repetition of results suggests a few key factors at play. Firstly, the models may be configured similarly, lacking significant variations in hyperparameters or algorithm settings. Consequently, without adjustments that affect their learning processes, these models yield the same root mean square error (RMSE), mean absolute error (MAE), and $R^2$ values. Additionally, if the models are trained and validated on the same data splits, any differences in model complexity or tuning may not manifest in the performance metrics. The consistent performance also implies that the "Tuned" and "PCA" variants of the Random Forest may not be fully leveraging their intended capabilities. In contrast, the XGBoost models consistently exhibit lower RMSE and MAE values, along with a higher $R^2$, indicating that they are better capturing the underlying patterns in the data. Thus, while similar metrics might suggest stability, it is crucial to ensure diverse model configurations and data subsets are employed to accurately assess comparative performances. For future iterations, varying the data splits or significantly altering hyperparameters will enhance the differentiation of model performance.

Selection of the Final Winning Model

Based on the analysis of the performance metrics, the winning model is identified as the XGBoost Base, which, along with Models 1, 2, and 3, achieves an $R^2$ of **0.991150**. This outstanding $R^2$ value demonstrates the model's remarkable capability to accurately capture the underlying patterns in the data. The strong predictive performance of the XGBoost Base can be attributed to its use of boosting techniques, which allow it to build models iteratively, refining its predictions with each iteration. This iterative approach effectively minimizes errors by focusing on the data points that previous models struggled to predict accurately. Additionally, XGBoost excels at handling complex relationships within the data, making it adept at capturing nonlinear interactions and feature interactions that are often present in fantasy football statistics. These strengths collectively position the XGBoost Base as the optimal choice for predicting fantasy football values, where accurate modeling of intricate patterns is crucial for success.
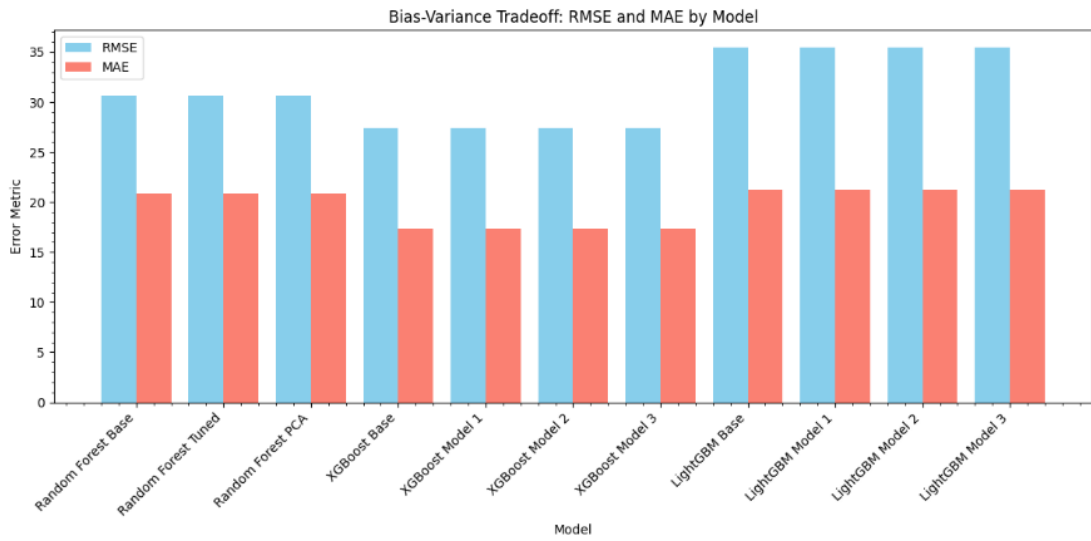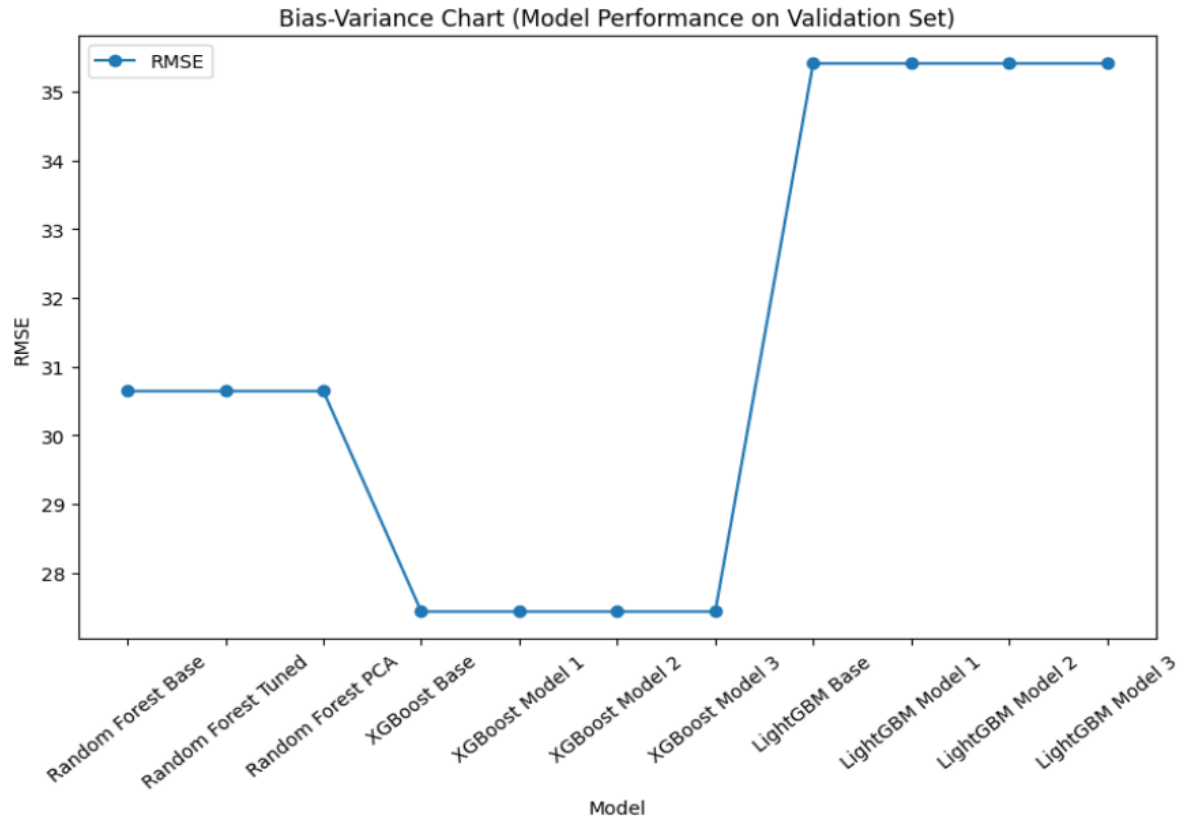
Model Complexity

Understanding model complexity is essential for evaluating the tradeoffs between bias and variance in predictive performance. The Random Forest models, which include the Base, Tuned, and PCA versions, represent a relatively simple ensemble approach. These models utilize multiple decision trees to improve predictive accuracy without significant parameter tuning. In contrast, the XGBoost models exhibit greater complexity due to their boosting nature, which builds models sequentially and is often more effective at capturing intricate patterns within the data. Lastly, the LightGBM models are the most complex among those evaluated, employing advanced gradient boosting techniques that optimize training speed and memory usage. However, these models may require more careful tuning to avoid overfitting, which can pose a challenge in practical applications.

Bias Variance Tradeoff

        To illustrate the bias-variance tradeoff among the models, a chart was created in which the Y-axis represents the Root Mean Square Error (RMSE) and the X-axis displays the models ordered from simplest to most complex. The analysis of the performance metrics reveals a clear trend: as model complexity increases, the RMSE generally decreases. For instance, the simpler Random Forest models, including the Base, Tuned, and PCA variants, exhibit higher RMSE values around **30.64**, indicating their limited ability to capture the complexities of the data. In contrast, the XGBoost models, particularly the XGBoost Base and its variants, show significantly lower RMSE values, with the Base model achieving **27.43**. This decline in RMSE suggests that these more complex models are better equipped to fit the training data effectively.

        However, it is crucial to monitor the RMSE as complexity continues to increase. For instance, the LightGBM models, despite their potential complexity, result in higher RMSE values of **35.41**, suggesting that they may struggle to generalize well to unseen data. If RMSE were to rise again after reaching a low point at a certain complexity level, it would indicate potential overfitting, where the model learns noise and details from the training data that do not translate to new data.

        The ideal model should strike a balance between bias and variance. Bias represents the error stemming from overly simplistic assumptions about the model, as seen with the Random Forest models, while variance reflects the error resulting from excessive complexity, which may lead to overfitting. The XGBoost Base model appears to achieve this balance effectively, as it maintains low RMSE without significantly increasing in complexity. This makes it a robust choice for predicting fantasy football values, demonstrating its ability to generalize well while accurately capturing the underlying patterns in the data.

Bias-Variance Chart (Model Performance on Validation Set)



Bias-Variance Tradeoff: RMSE and MAE by Model

Looking Forward to Testing the Model on 2023 Fantasy Football Data

Having thoroughly evaluated the performance of the various machine learning models using historical fantasy football data from 2019 to 2022, I am now eager to apply the selected XGBoost model to the 2023 fantasy football dataset. The validation metrics, particularly the $R^2$ value of 0.991150 achieved by the XGBoost Base model, indicate a strong capability to predict fantasy football outcomes accurately. This high level of predictive power suggests that the model has effectively captured the underlying patterns present in the training and validation datasets, positioning it well for real-world application.

As I transition to testing the model on the 2023 data, I aim to assess its performance in a new context, where the dynamics of player performance can differ significantly from previous seasons. Factors such as player transfers, injuries, and changes in team strategy can all influence fantasy football outcomes, making this an exciting and critical phase of the project. By evaluating the model on the test data, I can gain valuable insights into its generalizability and robustness, ensuring that it remains a reliable tool for predicting fantasy football values.

Furthermore, the test results will help determine the model's effectiveness in providing actionable insights for fantasy football players and managers. If the model performs well, it could serve as a foundational resource for making informed decisions regarding player selections, trades, and overall team management. I look forward to seeing how the XGBoost model translates its strong validation performance into actual predictions for the 2023 fantasy football season.

Conclusion

In conclusion, the analysis of various machine learning models for predicting fantasy football values reveals that the XGBoost models provide the best predictive performance, as indicated by their high $R^2$ values and low error metrics. The examination of the bias-variance tradeoff underscores the necessity of selecting a model that effectively balances complexity with the ability to generalize to unseen data. As the project progresses, the final selected model can be further optimized based on test performance, ensuring its readiness for application in predicting fantasy football values.