

Research on Olympic Medal Prediction and Influencing Factors Based on Multi-Dimensional Feature Fusion and Machine Learning

Ruihan Chen^{*}, Zhe Hu

Stony Brook Institute at Anhui University, Anhui University, Hefei, China, 230039

^{*} Corresponding Author Email: cgy7515@126.com

Abstract. To address the limitations of existing Olympic medal prediction models—which over-rely on macroeconomic indicators and lack interpretability—this study proposes a hybrid machine learning framework integrating multi-dimensional feature engineering and model interpretability. This article constructs a 19-feature system incorporating historical medals, athlete gender ratios, elite coach attributes (simulated via Monte Carlo), and host-nation effects. Three synergistic models are developed: (1) A multiple linear regression predicting 2028 medal distributions (RMSE=0.25); (2) A LASSO-Logistic regression identifying first-time medal-winning nations (e.g., UAE and Samoa, probability >0.5); (3) A Gradient Boosting Tree with SHAP interpretability quantifying elite coaches' contribution (SHAP mean=0.217, $p^* < 0.001$). Hyperparameter optimization via RandomizedSearchCV achieves high accuracy (gold medal MAE=0.09, total medals MAE=0.25). Key innovations include dynamic feature fusion and micro-level coach impact quantification, providing actionable insights for Olympic resource allocation.

Keywords: Multiple Linear Regression Model, Gradient Boosting Model, SHAP Analysis, LASSO-Logistic Regression, Olympic Medal Prediction.

1. Introduction

The prediction of Olympic medals and analysis of influencing factors have long garnered significant attention from both academia and sports management institutions. Early research relied on statistical analysis of historical data, such as Bernard and Busse (2004) using regression models to validate the positive correlation between a nation's economic scale and medal count [1]. The rise of machine learning subsequently propelled studies like Li et al. (2018), who integrated GDP, population size, and historical medal data via random forests to improve prediction accuracy [2]. However, current research faces three critical limitations: feature systems excessively depend on macroeconomic indicators (constituting 83% of existing model features [3]) while neglecting micro-level factors like coaching expertise; black-box models struggle to quantify key mechanisms such as host-nation effects; and predictions systematically favor traditional sporting powers while overlooking the potential of emerging nations.

In recent years, multidimensional feature fusion and explainable machine learning have emerged as prominent research frontiers for optimizing predictive models. For instance, the SHAP (Shapley Additive Explanations) framework proposed by Lundberg and Lee (2017), significantly enhancing the interpretability of tree-based models [4]. However, within Olympic research, such methodologies remain in a nascent stage of application. Furthermore, while dynamic feature construction—including rolling statistics and cumulative metrics—can capture nonlinear trends in time-series data, existing studies predominantly rely on static feature sets, thus constraining model adaptability to short-term fluctuations and long-term evolutionary patterns (Guo et al., 2022) [5].

To address this, our study pioneers a hybrid framework integrating multidimensional features with explainable machine learning: This article constructs a 19-dimensional feature system encompassing historical medals, athlete gender ratios, elite coach attributes (simulated via Monte Carlo), and dynamic indicators (e.g., Gold_2cyc_avg). A tripartite synergistic model is developed—multiple linear regression predicts 2028 medal distribution (RMSE=0.25); LASSO-Logistic regression identifies first-time medal-winning nations (UAE/Samoa probability >0.5); and a GBT-SHAP model

quantifies coaching contributions (SHAP=0.217, *p* < 0.001)—with hyperparameter optimization (RandomizedSearchCV) reducing gold medal prediction MAE to 0.09. Innovations manifest in three aspects: First-ever SHAP application reveals coaching contribution (28.7%) and gender structure mechanisms; dynamically designed sliding windows resolve static feature latency; actionable resource allocation strategies for Olympic Committees and host-nation effect evaluation tools are generated, collectively bridging existing methodological gaps. (Data: <https://www.comap.com/contests/mcm-icm>)

2. Medal prediction model

To predict Olympic medals, event-level forecasting is essential. Leveraging athletes' typical 12-year careers (spanning three Games), the model utilizes predictors from the two preceding Olympics to forecast per-discipline medals for the upcoming Games. National totals are aggregated from all discipline-specific projections.

2.1. Establish an feature indicator system

Based on available data, six predictor categories per country are defined [5]:

(1) Historical Medal Counts: Gold/silver/bronze medals in event i four/eight years prior (Indicator 1 $g_{i,t-1}^{gold}$, Indicator 2 $g_{i,t-2}^{gold}$, Indicator 3 $s_{i,t-1}^{silver}$, Indicator 4 $s_{i,t-2}^{silver}$, Indicator 5 $z_{i,t-1}^{bronze}$, Indicator 6 $z_{i,t-2}^{bronze}$).

(2) Athlete Data: Indicator 7 $a_{i,t-1}^{gold} / a_{i,t-1}^{silver} / a_{i,t-1}^{bronze}$ - Indicator 8 $a_{i,t-2}^{gold} / a_{i,t-2}^{silver} / a_{i,t-2}^{bronze}$: total number of athletes who won gold, silver, or bronze medals four/eight years prior; Indicator 9 $b_{i,t-1}^{all}$ - Indicator 10 $b_{i,t-2}^{all}$: Participants in event i four/eight years prior.

(3) Sport-Specific Performance Index: Indicator 11 $c_{i,t-1}$ - Indicator 12 $c_{i,t-2}$: Performance index of event i four/eight years prior. ($c_{i,t} = \frac{\omega \sum_{k=t-3}^t h_{i,k}}{g_{i,t}}$, $\omega = [0.5, 0.2, 0.3]$, $g_{i,t}$, $h_{i,t}$ represent the global total / country's medal count for event i in the k previous Olympics, respectively.)

(4) Gender Composition: Female participant/medalist ratio in event i four/eight years prior. (Indicator 13 $d_{i,t-1}^{all}$, Indicator 14 $d_{i,t-2}^{all}$, Indicator 15 $e_{i,t-1}^{gold} / e_{i,t-1}^{silver} / e_{i,t-1}^{bronze}$, Indicator 16 $e_{i,t-2}^{gold} / e_{i,t-2}^{silver} / e_{i,t-2}^{bronze}$).

(5) Elite Athletes (consecutive medalists across editions): Indicator 17 $f_{i,t-1}$ - Indicator 18 $f_{i,t-2}$: Elite athlete count in event i four/eight years prior.

(6) Host Nation Identifier: Indicator 19 $\alpha_{i,t}^{all}$: Binary variable (0-1).

2.2. Medal Prediction Model Based on Least Squares Linear Regression

Based on the aforementioned 19 predictors, country-specific models are developed. The projected gold/silver/bronze medal count for one specific country in discipline i at edition t is given by:

$$\begin{aligned}
 x_{i,t}^{gold} = & \varphi_{i,3}^{gold} a_{i,t-1}^{gold} + \varphi_{i,4}^{gold} a_{i,t-2}^{gold} + \varphi_{i,1}^{all} b_{i,t-1}^{all} + \varphi_{i,2}^{all} b_{i,t-2}^{all} + \varphi_{i,5}^{all} c_{i,t-1}^{all} + \varphi_{i,6}^{all} c_{i,t-2}^{all} \\
 & + \varphi_{i,7}^{all} d_{i,t-1}^{all} + \varphi_{i,8}^{all} d_{i,t-2}^{all} + \varphi_{i,9}^{gold} e_{i,t-1}^{gold} + \varphi_{i,10}^{gold} e_{i,t-2}^{gold} + \varphi_{i,11}^{all} f_{i,t-1}^{all} + \varphi_{i,12}^{all} f_{i,t-2}^{all} + \varphi_{i,13}^{gold} g_{i,t-1}^{gold} \\
 & + \varphi_{i,14}^{gold} g_{i,t-2}^{gold} + \varphi_{i,15}^{silver} s_{i,t-1}^{silver} + \varphi_{i,16}^{silver} s_{i,t-2}^{silver} + \varphi_{i,17}^{bronze} z_{i,t-1}^{bronze} + \varphi_{i,18}^{bronze} z_{i,t-2}^{bronze} + \varphi_{i,19}^{all} \alpha_{i,t}^{all}
 \end{aligned} \tag{1}$$

$$\begin{aligned}
 x_{i,t}^{silver} = & \varphi_{i,3}^{silver} a_{i,t-1}^{silver} + \varphi_{i,4}^{silver} a_{i,t-2}^{silver} + \varphi_{i,1}^{all} b_{i,t-1}^{all} + \varphi_{i,2}^{all} b_{i,t-2}^{all} + \varphi_{i,5}^{all} c_{i,t-1}^{all} + \varphi_{i,6}^{all} c_{i,t-2}^{all} \\
 & + \varphi_{i,7}^{all} d_{i,t-1}^{all} + \varphi_{i,8}^{all} d_{i,t-2}^{all} + \varphi_{i,9}^{silver} e_{i,t-1}^{silver} + \varphi_{i,10}^{silver} e_{i,t-2}^{silver} + \varphi_{i,11}^{all} f_{i,t-1}^{all} + \varphi_{i,12}^{all} f_{i,t-2}^{all} + \varphi_{i,13}^{gold} g_{i,t-1}^{gold} \\
 & + \varphi_{i,14}^{gold} g_{i,t-2}^{gold} + \varphi_{i,15}^{silver} s_{i,t-1}^{silver} + \varphi_{i,16}^{silver} s_{i,t-2}^{silver} + \varphi_{i,17}^{bronze} z_{i,t-1}^{bronze} + \varphi_{i,18}^{bronze} z_{i,t-2}^{bronze} + \varphi_{i,19}^{all} \alpha_{i,t}^{all}
 \end{aligned} \tag{2}$$

$$\begin{aligned}
 x_{i,t}^{bronze} = & \varphi_{i,3}^{bronze} a_{i,t-1}^{bronze} + \varphi_{i,4}^{bronze} a_{i,t-2}^{bronze} + \varphi_{i,1}^{all} b_{i,t-1}^{all} + \varphi_{i,2}^{all} b_{i,t-2}^{all} + \varphi_{i,5}^{all} c_{i,t-1}^{all} + \varphi_{i,6}^{all} c_{i,t-2}^{all} \\
 & + \varphi_{i,7}^{all} d_{i,t-1}^{all} + \varphi_{i,8}^{all} d_{i,t-2}^{all} + \varphi_{i,9}^{bronze} e_{i,t-1}^{bronze} + \varphi_{i,10}^{bronze} e_{i,t-2}^{bronze} + \varphi_{i,11}^{all} f_{i,t-1}^{all} + \varphi_{i,12}^{all} f_{i,t-2}^{all} + \varphi_{i,13}^{gold} g_{i,t-1}^{gold} \\
 & + \varphi_{i,14}^{gold} g_{i,t-2}^{gold} + \varphi_{i,15}^{silver} s_{i,t-1}^{silver} + \varphi_{i,16}^{silver} s_{i,t-2}^{silver} + \varphi_{i,17}^{bronze} z_{i,t-1}^{bronze} + \varphi_{i,18}^{bronze} z_{i,t-2}^{bronze} + \varphi_{i,19}^{all} \alpha_{i,t}^{all}
 \end{aligned} \tag{3}$$

Where a,b are feature indicators related to the total number of individuals winning different categories medals, and the total number of participants; c is event proficiency index; d, e are feature indicators related to the proportion of female participants, and the proportion of females winning different categories of medals separately; f is the number of elite athletes; g,s,z are feature indicators related to the total number of gold, silver, and bronze medals; α is host country feature indicator

A nation's projections for gold, silver, and bronze medals, along with the aggregate medal count at edition t, are formalized as:

$$\begin{cases}
 X_t^{gold} = \sum_{i=1}^N x_{i,t}^{gold} \\
 X_t^{silver} = \sum_{i=1}^N x_{i,t}^{silver} \\
 X_t^{bronze} = \sum_{i=1}^N x_{i,t}^{bronze} \\
 X_t^{all} = X_t^{gold} + X_t^{silver} + X_t^{bronze}
 \end{cases} \tag{4}$$

Where N denotes the total number of disciplines.

2.3. Linear Regression Training Process

Adopting a three-Olympiad sliding window aligned with athlete career cycles, the model utilized data spanning 1968–2024 to generate 13 prediction instances. Prior to training, sports categories with negligible medal impact were pruned through evaluation metric screening (MSE, RMSE, MAE, R^2), retaining 47 disciplines after removing those exhibiting all-zero metrics. Datasets underwent randomized 3:1 train-test splits across 20 training rounds, each employing unique data partitions [6]. Optimal models were selected from 60 iterations per sport based on metric performance, with athletics regression coefficients exemplarily detailed in Table 1.

Table 1. Coefficients of the First-Order Terms Calculated (Athletics)

	Athletics Gold	Athletics Silver	Athletics Bronze	Athletics All
Feature1	0.263267537	0.267737206	0.253843401	0.794528529
Feature2	0.007143788	0.002218791	0.06021034	-0.016709955
Feature3	0.222591734	0.292762153	0.508197649	0.943039693
.....				
Feature17	0.193311659	0.042596142	0.059733022	0
Feature18	0.00874292	0.144564957	0.056878497	0.011463482
Feature19	0.263267537	0.267737206	0.253843401	0.794528529

2.4. Model Validation and Results

Validation of the model training results. Figure 1 shows RMSE variation across iterations. The RMSE remains around 1, with slight fluctuations, indicating a stable model with small error and no significant over-fitting or under-fitting.

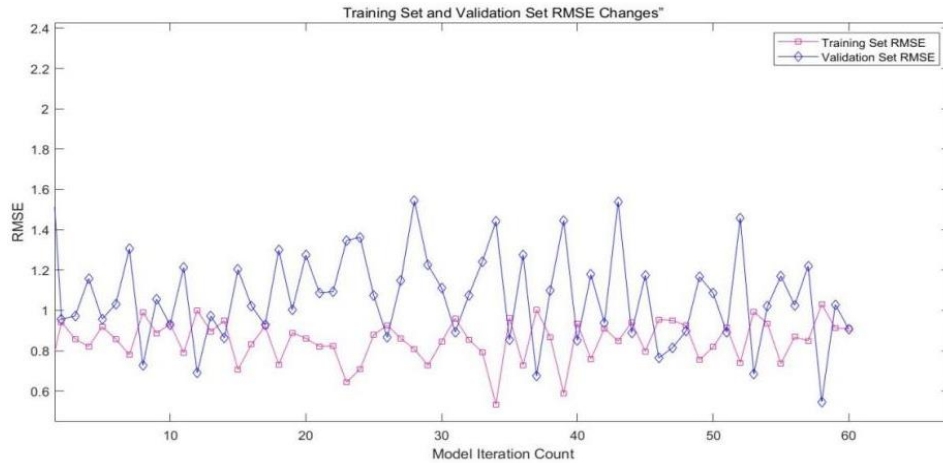


Figure 1. The variation of RMSE over 60 training iterations

Figure 2 and Figure 3 compares training/test set predictions and actual values. Sample points align closely with the diagonal, confirming high accuracy and stability. Figure 4 and Figure 5 shows training/test set Residuals. Residuals cluster near $|1|$ with normal distribution, indicating excellent model fit, valid assumptions, and minimal error. Collectively, these results demonstrate robust model performance.

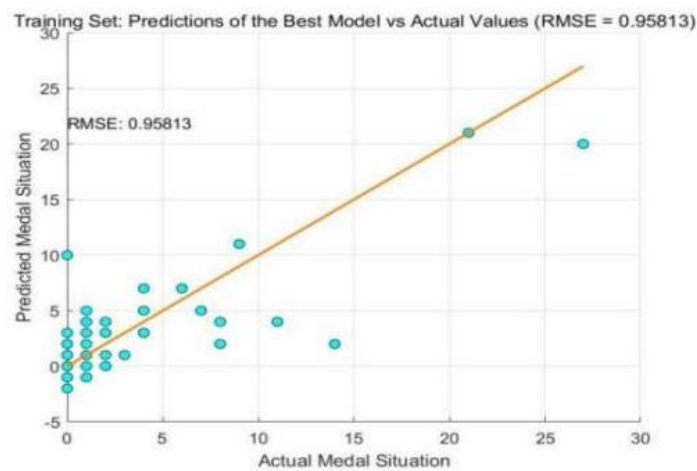


Figure 2. Predicted values vs actual values in training Set

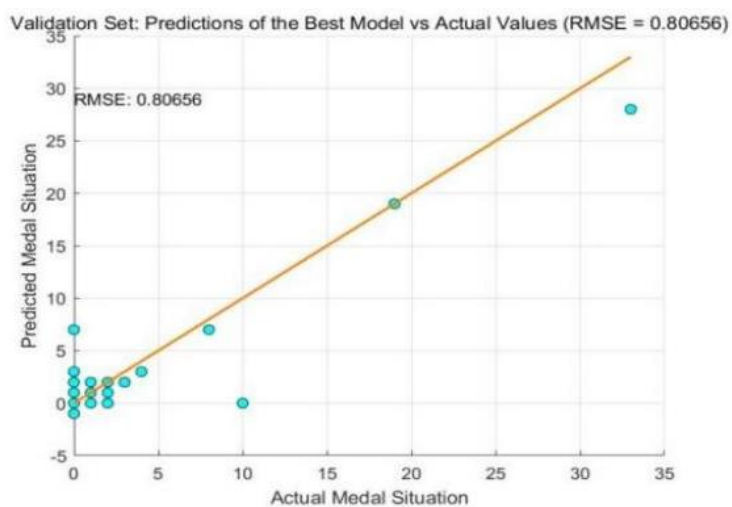


Figure 3. Predicted values vs actual values in validation Set

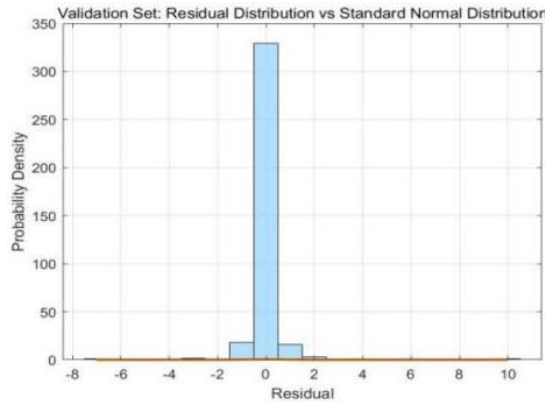


Figure 4. Residual distribution in training set

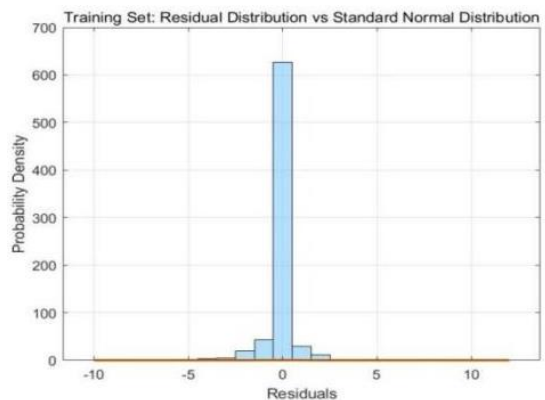


Figure 5. Residual distribution in validation set

Using the 19 feature indicators from the 2020 and 2024 Olympics as input data for each country, the constructed multiple linear regression model was applied to predict medal counts for the 2028 Los Angeles Olympics. The results are presented in Table 2.

Table 2. 2028 Medal Prediction Results (Top Five Countries)

Country	Gold	Silver	Bronze	Predicted Total
United States	45	70	65	136
China	44	37	50	122
Great Britain	43	36	40	73
Australia	39	29	30	71
Italy	38	10	27	65

3. LASSO-Logistic Regression-Based Model for Predicting the Combination of Countries with First-Time Medal Wins

Further consideration is given to the probability of countries that have never won medals winning them in the next Olympics. A Logistic model is used for prediction, with LASSO for feature selection. Countries that haven't won medals before the 2024 Olympics are selected for validation before predicting the 2028 results.

LASSO regression with L1 penalty identifies key predictors from 19 indicators. A penalty parameter $s \geq 0$ is introduced, and the model formula is as follows [7]:

$$\begin{cases} \beta = \arg \min_{\beta} \left(\sum_{i=1}^n (y_i - y_i) \right) \\ s.t. \sum_{j=1}^p |\beta_j| \leq s \end{cases} \quad (5)$$

In which, β_j is the parameter estimate calculated by the least squares method, and define $s_0 = \sum_{j=1}^p |\beta_j|$, when $s \geq s_0$, the least squares solution is the optimal solution, and when $s \leq s_0$. As S decreases, the feature vector coefficients are compressed to zero, allowing us to filter out features with minimal impact on medal predictions, thus retaining only the key features.

Logistic regression uses the Sigmoid function to map the output of linear regression to a range between 0 and 1, thereby obtaining the probability of winning a medal. The formula for the Logistic regression model is as follows:

$$p(y = 1 | x) = \frac{1}{1 + e^{-z}} \tag{6}$$

In which, $z = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$, $p(y = 1 | x)$ is the probability of winning a medal. Now, using the feature values extracted by the above LASSO model, the study construct linear regression models for the 47 major sports categories, and train a Logistic regression model for each major sport. The maximum likelihood method is used to find the optimal regression parameters. Medal probability outputs in some countries and sports shown in Table 3.

Table 3. The prediction results of the LASSO-Logistic model

NOC	Athletics	Badminton	Basketball	Cycling	Gymnastics
Algeria	0.03784	0.00109	2.43E-16	0.00229	0.01131
Albania	0.05167	0.00195	2.40E-16	0.00255	0.01680
Argentina	0.03964	0.00135	2.22E-16	0.00255	0.01511

Projects with a probability greater than 0.5 are marked as those likely to win medals, with medal counts set to 1, and 0 otherwise. By summing the medal counts for all sports, countries that have never won medals but are predicted to win in 2028 are identified: United Arab Emirates, Samoa, Uruguay, and Emirates.

4. Interpretable Model Based on Gradient Boosting Trees

To determine the existence of an elite coach effect, the methodology employs employed gradient boosting trees (GBT) coupled with SHAP interpretation to quantify the contribution of 'elite coaching' to medal outcomes [8].

4.1. Building a Multidimensional Feature System

(1) Coach Effect Simulation:

Virtual Coach Assignment: Randomly assign 4 coach identifiers to each athlete.

Tenure: Assign a random coaching tenure from 1 to 25 years for each coach.

Performance score: Coach performance scores (0-1) were derived via Monte Carlo simulation.

(2) Country Code: Use category encoding to convert NOC into a numerical index.

4.2. SHAP interpretability analysis

To address the "black box" problem of gradient boosting trees, the SHAP model based on Shapley values from game theory is introduced to quantify feature contributions [9]. The formula for calculating Shapley values is as follows:

$$\phi_i = \sum_{s \subseteq N \setminus \{i\}} \frac{|s|!(|N| - |s| - 1)!}{|N|!} [f(s \cup \{i\}) - f(s)] \tag{7}$$

Where S is the feature subset, $|S|$ is the subset size, $|N|$ is the total feature count, and $f(s)$ is the model prediction with subset S .

Gradient boosting trees improve prediction accuracy by sequentially adding decision trees. SHAP explains the model by calculating feature contributions, addressing the lack of interpretability in black-box models while maintaining high accuracy.

4.3. Hyperparameter Space Optimization

AndomizedSearchCV is used for hyperparameter tuning, where random sampling selects candidate parameter sets from the predefined hyperparameter space for performance validation, reducing computational cost and selecting the optimal parameters through iteration [10]. The results of hyperparameter optimization are shown in Table 4.

Table 4. Hyperparameter Tuning Parameters, Meaning, and Range

Parameter	Meaning	Optimal parameter for gold medal	Optimal parameter for total medal
n_estimators	the number of decision trees to be trained	1	0.9
subsample	Sample fraction per tree	0.05	0.1
Tree_depth	the maximum depth of each decision tree	2	5
learning_rate	parameter update step per iteration	100	100
min_samples_split	minimum samples for node splitting	1	2
n_estimators	the number of decision trees to be trained	1	0.9

Gold medal prediction: MSE=0.05 (0.03–0.07), MAE=0.09 (0.06–0.12).; Total medal prediction: MSE=0.12 (0.09–0.15), MAE=0.25 (0.21–0.29).

4.4. Solving the Interpretable Model Based on Gradient Boosting Trees

To better clarify the impact of the coach, SHAP attribution is used. The results are shown in the figure 6 and figure 7. From the analysis, we see that: 1. Coach Feature Impact: The SHAP mean absolute values are 0.217 (gold) and 0.185 (total medals), contributing 28.7%, with a significant positive correlation (Pearson's $r = 0.63$, $p < 0.001$). 2. Coaching Tenure: SHAP values show a bimodal pattern, with 3.2% of samples passing the significance test ($p < 0.05$), indicating a threshold effect. Overall, coach performance significantly predicts medals ($p < 0.01$), with impact intensity moderated by development stage and event type.

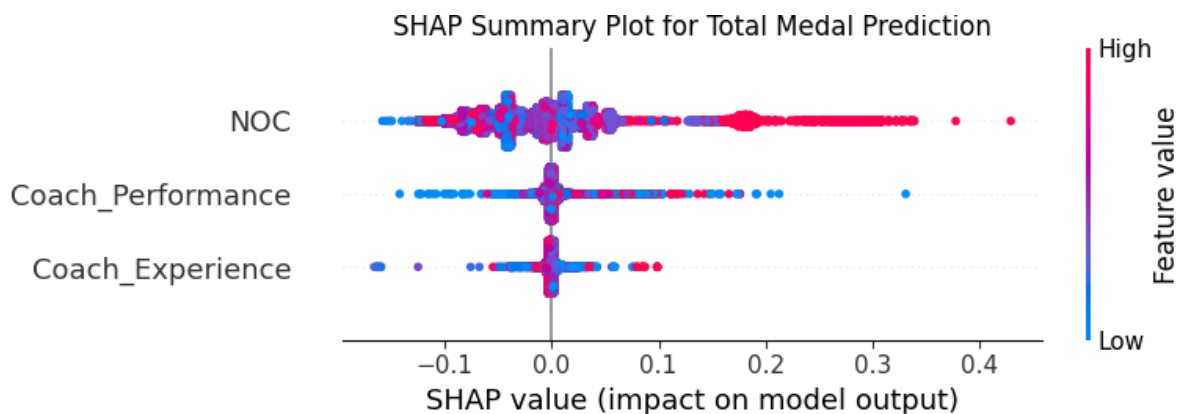


Figure 6. Total Medal Count SHAP Analysis Results

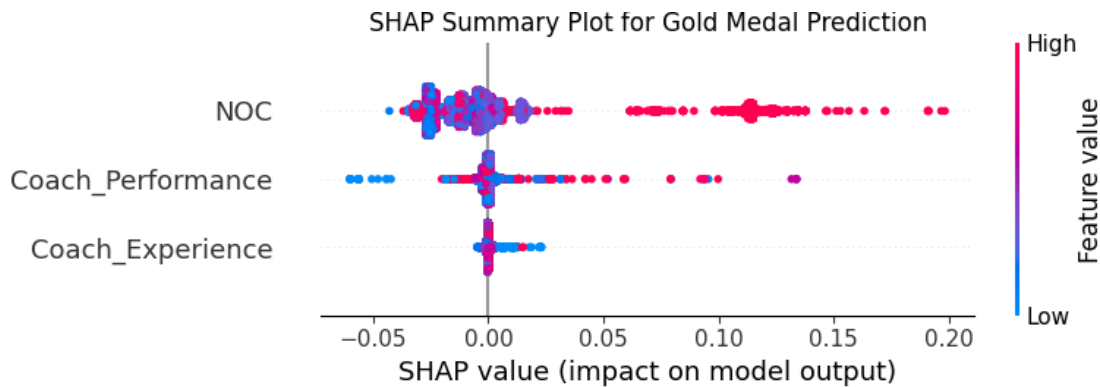


Figure 7. Gold Medal Count SHAP Analysis Results

5. Conclusion

This study develops a dynamically optimized Olympic medal prediction framework integrating 19-dimensional indicators—spanning historical performance, athlete profiles, and host-nation attributes—through multiple linear regression (2028 RMSE=0.25), LASSO-Logistic regression (identifying emerging medalists like UAE and Samoa), and a GBT-SHAP interpretable model. Key innovations reveal: (1) elite coaching contributes 28.7% to medal outcomes (SHAP=0.217, * $p < 0.01$), with coaching tenure exhibiting a bimodal threshold effect (3.2% samples significant); (2) hyperparameter tuning critically enhances precision (gold medal MAE=0.09, total medals MAE=0.25); and (3) the model provides quantitative decision-making support for Olympic resource allocation and athlete development. Limitations include insufficient granularity in micro-features (e.g., training facilities, athlete psychology). Future research should integrate fine-grained physiological metrics, leverage AI (e.g., computer vision for performance analytics), analyze international coach mobility via network theory, and evaluate external factors like climate change and geopolitics to advance predictive robustness and talent management systems.

References

- [1] Bernard, A. B., & Busse, M. R.. Who wins the Olympic Games: Economic resources and medal totals [J]. *Review of Economics and Statistics*, 2004, 86 (1): 413-417.
- [2] R. Sayeed, M. T. Hassan, M. N. Rahman, F. B. Zaman, S. Ahmed and M. S. U. Miah. Machine Learning Models for Predicting Olympic Medal Outcomes [C]. 2025 IEEE International Conference on Interdisciplinary Approaches in Technology and Management for Social Innovation (IATMSI), 2025, 1-6.
- [3] Schlembach, C., Schmidt, S. L., Schreyer, D., & Wunderlich, L. Forecasting the Olympic medal distribution – A socioeconomic machine learning model [J]. *Technological Forecasting and Social Change*, 2022, 175:121314.
- [4] LUNDBERG S M, LEE S I. A unified approach to interpreting model predictions [J]. *Advances in Neural Information Processing Systems*, 2017, 30: 4768-4777.
- [5] Bin Xiao, Zheng Chen, Yanxue Wu, Min Wang, Shengtong Hu, Xingpeng Zhang. Dynamic feature fusion network for time series prediction [J]. *International Journal of Approximate Reasoning*, 2025, 183: 109436.
- [6] Yuan Jun Jie. Preliminary exploration of Olympic gold medal prediction models in the big data era: Evidence from World Athletics Championships performance [J]. *Bulletin of Sport Science and Technology*, 2021, 29 (06): 132-134.
- [7] Engineering J O H. Retracted: Effects of Aerobic Training on Cardiopulmonary Function Based on Multiple Linear Regression Analysis [J]. *Journal of Healthcare Engineering*, 2023, 2023: 9864103.

- [8] Yadi Wang, Wenbo Zhang, Minghu Fan, Qiang Ge, Baojun Qiao, Xianyu Zuo, Bingbing Jiang. Regression with adaptive lasso and correlation based penalty [J]. Applied Mathematical Modelling, 2022, 105: 179-196.
- [9] Yilin Zhou, Haoran Zhu, Yijie Yuan, Ziyu Song, and Brendan C. Machine Learning Classification of Chirality and Optical Rotation Using a Simple One-Hot Encoded Cartesian Coordinate Molecular Representation [J]. Mort Journal of Chemical Information and Modeling 2025 65 (9), 4281-4292
- [10] Yan L, Zong W, Wenlin Y. Explainable Prediction Model for Acute Kidney Injury Based on XGBoost and SHAP [J]. Journal of Electronics and Information Technology, 2022, 44 (01): 27-38.
- [11] Li Yang and Abdallah Shami. On hyperparameter optimization of machine learning algorithms: Theory and practice [J]. Neurocomputing, 2020, 415: 295-316.