| Problem Chosen | 2025 | Team Control Number |
|:---:|:---:|:---:|
| C | MCM/ICM<br>Summary Sheet | 2502268 |

# Decoding Olympic Glory: A Dual-Modal Computing Research

## Summary

Accurate prediction of the medal distribution for the next edition is crucial for the development of national sports strategies. External factors such as the quality of the athletes, the setting of the event and the host effect will significantly affect the medal distribution when the program is not yet decided. This study aims to build model to accurately predict the distribution of gold and total medals for the 2028 Olympic Games by combining historical Olympic medal data and external factors.

For Model I, in order to the distribution of gold and total medals for the 2028 Olympic Games, this study uses an innovative **Dual-Modal Model** which combines **ARIMA** for time-series historical data and **XGBoost** for non-time-dependent external features. This research method considers the prediction intervals at different confidence level to determine the predicted intervals. Next, quantitative indicators for predicting whether each country will progress or not are established by comparing the difference between the predicted number of gold medals for each country in 2028 and in 2024. Afterwards, to predict whether countries that never won gold medals will win their first gold medal in the next Olympics, we changed XGBoost to a **Binary Classification Model** and estimated the probability of each country to win. A quantitative analysis of the features importance of theXGBoost model demonstrates that new events and **Host Effects** contribute the most to gold medal predictions. Further, **Spearman Correlation Analysis** was used to demonstrate the dominant events for each country. Finally, a **T-test** is used to test that there is a significant difference between when a country is a host and when it is not, give the significant impact of the host effect on medal predictions.

For Model II, after utilizing **Pair-T Test** of Olympic women's volleyball data to determine the existence of the great coach effect, the p-value is less than 0.05, so the great coach effect exists. Next, quantifying the great coach effect uses **Difference-in-difference Regression Model** with 255 Olympic samples related to this effect, and the experimental group and the control group is 1:1.5, and the result shows that the contribution rate of the great coach effect is **0.9574**, indicating a significant great coach effect. Finally, taking **USA, Italy and France**, use the **AHP-TOPSIS** model to assess the most invest of programs in great coaches. The answer of USA is **Volleyball, Cycling Road**, Italy is **3x3 Basketball, Baseball**, and France is **Fencing, Hockey**.

For Problem 3, by analyzing the data, certain countries have obvious advantages in specific events. Monopoly countries should focus on young athletes' training and event experience; countries with dominant programs should avoid relying on a single sport and improve medal diversity. Countries that are not competitive enough for gold medals need to improve their performance through targeted training and the introduction of high-level coaches.

**Keywords**: **Dual-Modal Model, ARIMA, XGBoost, Binary Classification Model, Host Effects, Spearman Correlation Analysis, T-test, Difference-in-Difference Regression Model, AHP-TOPSIS**

# Content

# 1 Introduction

## 1.1 Problem Background

With the development of the Olympic Games and the participation of more and more countries and regions, the Olympic Games is not only a platform for sports competition but also reflects the comprehensive achievements of countries in sports and other fields. The Olympic medal table has always been an important criterion for measuring national sports achievements. Although the number of medals is affected by various factors, such as the quality of athletes and the programs they participate in, countries such as the United States and China have maintained their leading position all year round. In addition, some smaller countries or regions may also excel in specific sports. However, external factors such as host effect and coach training also affect the distribution of medals. For example, host countries usually win more medals due to their familiarity with the venues and dominant events. In addition, changes in event settings and programs also affect the number of medals, making predictions more complex.

Therefore, accurate prediction of future Olympic medal distribution, especially the number of gold medals, is significant to countries' sports strategies and athletes' training directions. This study aims to construct a prediction model by analyzing previous Olympic medal data to help countries adjust their sports strategies better.

## 1.2 Restatement of the Problem

Considering the background information and restricted conditions identified in the problem statement, we need to solve the following problems:

◉ Problem 1

Constructing a predictive model for predicting future Olympic medal distributions, estimating the uncertainty and accuracy of the forecasts, and measuring the performance of the model. We need to forecast each country's gold and total medals for the 2028 Los Angeles Olympic Games and provide forecast intervals. At the same time, analyze which countries will likely improve or decline in performance.

◉ Problem 2

Examining the Olympic data for evidence of a "great coach effect" and assessing the impact of this effect on medal counts. Specifically, we will look at whether coaching turnover, particularly the involvement of high-profile coaches, significantly changes the number of medals won by certain countries in specific sports. After quantifying this effect, we will select three countries and analyze whether they should improve their medal counts in specific sports by bringing in world-class coaches.

◉ Problem 3

Exploring other key factors affecting the number of Olympic gold medals, such as the direction of national investment and the level of athletes. Based on the analysis's results, it is recommended to increase investment in key areas and programs, introduce high-quality coaches, and optimize the layout of sports training by using the country's advantageous resources. We can help countries make more scientific and accurate Olympic decisions through specific action plans.

## 1.3 Literature Review

The impact of various factors on the Olympic Games has been a topic of academic interest. Baimbridge examined uncertainty in the outcomes of international sporting events, including the Summer Olympics, and identified factors such as levels of participation, boycotts, and general trends as important determinants(1998). Recent studies have also delved into the predictors of medal counts for countries at the PyeongChang and Tokyo Olympics, considering demographic, economic, geographic, and religious demographics when estimating medal counts (Li et al., 2022). Additionally, the impact of hosting the Olympic Games on sports performance has been a topic of interest, and

researchers are exploring the home-field advantage effect and its impact on medal outcomes in host countries (Csurilla et al., 2023). Overall, the pursuit of Olympic medals continues to be a focus of attention worldwide, reflecting not only athletic achievement but also national pride, identity, and the socio-cultural dynamics that shape sporting success. Various factors influence the Olympic Games' success and require careful consideration and planning to achieve positive outcomes.
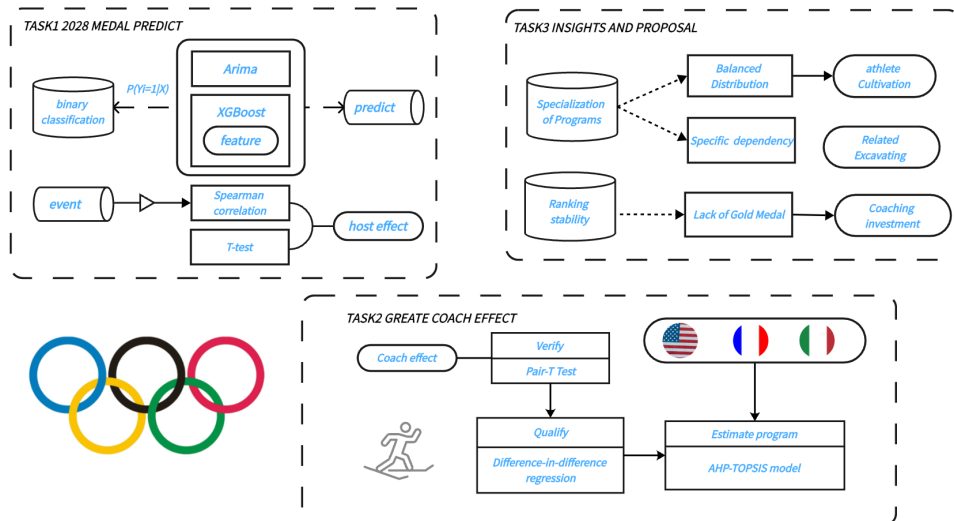
## 1.4 Our work



*Figure 1.4 Research work flow*

# 2 Assumptions and Justifications

Through the complete analysis of the problem, in order to simplify our model, we make the following reasonable assumptions and justify their use.

1. **Medals have a significant temporal correlation.** The performance of each Olympic Games not only reflects the training standards of that year but may also be influenced by the history of medal wins and losses. In addition, countries' investment in sports, such as training facilities and research, is usually a long-term process.

2. **Medals could be affected by the country hosting the event.** The country hosting the Olympic Games usually has a higher home-field advantage because of its unique psychological and environmental adaptation advantages; the host can also add its advantageous events to the competition program.

3. **Medals will be affected by the number of athletes involved.** Sending more athletes to the Olympics means that the country can compete in more events, thus increasing the probability of winning medals.

4. **International cooperation and foreign aid can affect medals.** Bringing in world-class coaches or exchanging sports resources with powerful countries can help a country boost its performance in specific sports.

# 3 Preparations

## 3.1 Data preprocessing

We used multiple data files provided by the title containing information related to Olympic medal distribution, athlete performance, program settings, and other relevant factors. The following is a detailed description of each data file.

*Table 3.1  Datasets information*

| Dataset | information |
| --- | --- |

| summerOly_athletes | Name, Sex, Team, NOC, Year, City, Sport, Event, Medal |
|---|---|
| summerOly_hosts | Year, Host_place |
| summerOly_medal_counts | Rank, NOC, Gold, Silver, Bronze, Total, Year |
| summerOly_programs | Sport, Discipline, Code, Sports Governing Body, Number of events in discipline for that year |

### 3.1.1 Missing Value Processing

a. Countries such as Djibouti have only one edition in the given medal data and are not competing for the first time in 2024. Since such countries do not have changing data and will not be host countries in the future, we can leave them unchanged in the projections.

b. Missing years and codes in some of the Olympiad sports program datasets. For the missing year data, we fill in the data by extracting the mean value of the data from two neighboring sessions. In the case of missing codes, since they did not affect our data analysis, we directly filled the corresponding codes to avoid any impact on subsequent data screening.

### 3.1.2 Data Merge

We have integrated the medals and athletes tables to improve the clarity of the data and facilitate subsequent operations. We adopt a merged processing approach for different teams in a country, e.g., Germany-1 and Germany-2. This operation ensures the consistency and integrity of the data.

### 3.1.3 Feature Extraction

Considering that some countries or regions may have been absent from previous Olympic Games or will not be competing in the future due to historical political factors, prior data from these countries or regions no longer affects subsequent gold medal predictions. Therefore, we can consider these data as anomalies or missing and ignore them in the model construction to ensure that the model accurately captures valid trends and patterns. For example, Germany split into East and West Germany after World War II until its merger in 1990. Considering that Germany has been represented during the re-unification period, we can merge the data for this period into Germany. In addition, some countries are no longer eligible to participate due to historical and political issues. For example, after the split of the former Soviet Union, Russia was banned, and athletes participated in a neutral capacity.

## 3.2 Notations

| symbol | description |
|---|---|
| $\phi$ | auto-regressive coefficient |
| $\theta$ | moving average coefficient |
| $M$ | medal count |
| $G$ | gold medal count |
| $H$ | host country indicator |
| $E$ | event |
| $\beta$ | regression coefficient |
| $\epsilon_t$ | random error |
| RSS | residual sum of squares |
| $C_{ij}$ | The importance ratio of factor i relative to factor j |
| $\widehat{\beta}$ | Optimal regression coefficients |

# 4 Model Ⅰ: Predicting Olympic Medal Counts Based on a Dual-Modal Time Series and Machine Learning Model

## 4.1 The basic concept of the Dual-Modal Model

(1) Objective: To predict the number of gold and total medals per country

(2) Methods: A time series model (ARIMA) and a machine learning model (XGBoost) will be used to capture the time trends in medal counts and the effects of external features. Then, the time series predictions from the ARIMA model were combined with the residual corrections from the XGBoost model to generate predictions for the number of gold and total medals per country in the 2028 Olympic Games in Los Angeles.

(3) ARIMA (Auto-Regressive Integrated Moving Average): By smoothing the data and fitting auto-regressive and moving average components, the ARIMA model can effectively predict future trends in the time series.

(4) XGBoost (Extreme Gradient Boosting): XGBoost minimizes the prediction error by constructing the decision tree step-by-step. It supports feature importance analysis and can handle a wide range of external feature types.

## 4.2 Model building

In order to accurately predict the number of Olympic medals for each country in 2028, we designed a Dual-Modal model that combines a time series analysis method (ARIMA) with a machine learning model (e.g., XGBoost).

(1) Data assumptions

i. Time correlation assumption: Medal count data have a significant time correlation, and predictions of future results can be made based on historical results.

ii. Residuals non-linearity assumption: ARIMA model can not fully consider all factors, the residuals are partially independent of the time series, and may contain non-linear features or other complex relationships that can be modeled by machine learning models.

iii. External influence assumption: the number of medals may be significantly influenced by external factors (e.g. host country effect, number of athletes, etc.) which cannot be fully modeled by ARIMA.

### 4.2.1 ARIMA Modeling

(1) Model formula

The ARIMA model characterizes the dynamics of the time series through three parameters $(p, d, q)$:

$p$: the order of the auto-regressive term, which indicates the linear relationship between the current value and the past $p$ values.

$d$: the number of differences, which is used to transform the non-stationary series into a stationary series.

$q$: the order of the moving average term, which represents the relationship between the current value and the past $q$ prediction error.

The mathematical formula is:

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \cdots + \theta_q \epsilon_{t-q}$$

Among them:

q: the current value of the time series

c: the constant term

$\phi_i$: Auto-regressive coefficient, which indicates the effect of past values on current values
$\theta_i$: The moving average coefficient, which indicates the effect of past prediction errors on the current value
$\epsilon_t$: noise term

    (2) Model fitting

We use ARIMA model to fit the gold medal time series for each country to obtain the gold medal prediction formula:

$$\hat{M}_{i,t+1}^{\text{ARIMA}} = c + \phi_1 M_{i,t} + \phi_2 M_{i,t-1} + \cdots + \theta_q \epsilon_{t-q}$$

Among them:

$\widehat{M_{i,t+1}^{ARIMA}}$: Country $i$'s predicted medal count for the $t+1$ year

$c$: the constant term

$\phi_p$: auto-regressive coefficient indicating the effect of past $p$ medal totals on the current predicted value.

$M_{i,t}$: country $i$'s medal count in the year $t$

$\theta_q$ : moving average coefficient indicating the effect of $q$ past forecast errors on the current forecast

$\epsilon_t$: noise term

the predicted total number of medals is obtained in the same way, then we have:

$$\widehat{M_{i,t+1}^{ARIMA}} \text{and} \widehat{G_{i,t+1}^{ARIMA}}$$

### 4.2.2 XGBoost modeling

    To address the effect of non-time-dependent external features on the predicted values, we further used XGBoost to correct the residuals of the ARIMA predicted values.

    (1) residual correction

i. target variables: the residuals of the ARIMA model

$$\epsilon_{i,t} = M_{i,t} - \hat{M}_{i,t}^{\text{ARIMA}}$$

ii. external features：

✤  Host country effect:

$$H_{i,t} = \begin{cases} 1, & \text{if country } i \text{ is the host of the } t\text{th Olympic Games} \\ 0, & \text{otherwise} \end{cases}$$

    where $H_{i,\,t}$ indicates whether country $i$ is the host country in year $t$ .

✤  Number of athletes:

$$A_{i,t} = \text{Number of athletes from country } i \text{ participating in the } t\text{th Olympic Games.}$$

✤  Event Characteristics:

$\cdot E_t$: the total number of events in the $t$ th Olympic Games.

$\cdot E_{new,t}$ : the number of new event added at the $t$ th Games.

$\cdot E_{type,t}$: distribution of the number of different types of event at the $t$ th Games.

$\cdot E_{adv,i,t}$ : distribution of country $i$'s medals in different events types at the $t$ th Olympic Games.

Then we have all input feature:

$X_{i,t} = \{M_{i,t}, M_{i,t-1}, M_{i,t-2}, \ldots, G_{i,t}, G_{i,t-1}, G_{i,t-2}, \ldots, H_{i,t}, A_{i,t}, E_t, E_{\text{new},t}, E_{\text{type},t}, E_{\text{adv},i,t}, \hat{M}_{i,t+1}^{\text{ARIMA}}\}$

iii. final predict outcome:

$$M_{i,t+1} = \hat{M}_{i,t+1}^{\text{ARIMA}} + f(X_{i,t})$$

## 4.3 Solution of model

After completing the model building, we use this Dual-Modal model to solve the three sub problems of Question 1

(1) Sub problem 1 :

i.  We need to predict the distribution of medals at the 2028 Olympic Games in Los Angeles, which is already solved since we get $\widehat{M}_{i,t+1}^{ARIMA}$ and $\widehat{G}_{i,t+1}^{ARIMA}$
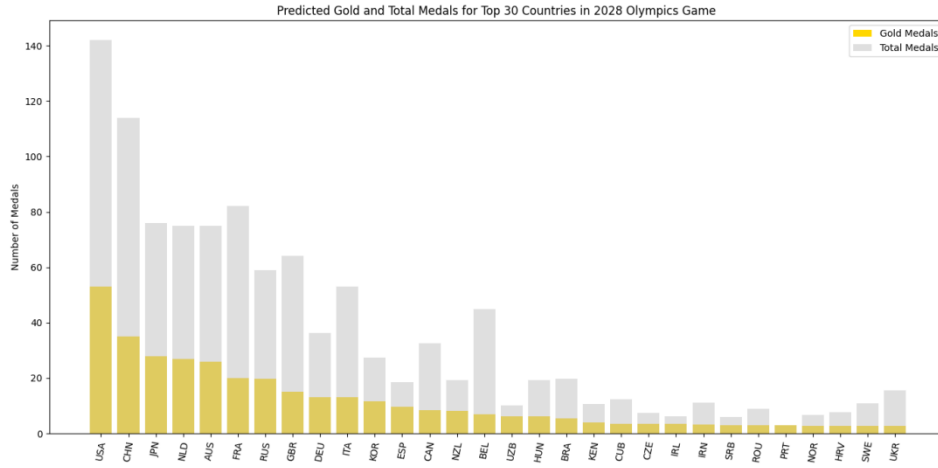


*Figure 4.3.1 predicted Gold Medals for 2028 Olympics Game*

Firstly, we note that the USA is predicted to win 40 in 2028, a distant 18 gold medals away from second place China. In the 2024 games, the US and China are tied in gold medals, whereas in the next Olympics in 2028, held in Los Angeles, USA, there is a significant difference, which may be closely related to the fact that the US is the host country, which we will explore further in the next questions. Secondly, the first-placed USA has 129 medals ahead of the second-placed China, which won 88 medals. This also shows the possibility of a host effect, which we'll look at further in the next question.

ii.  Point out prediction intervals for all results, we also can get from the Dual-Modal model:

$$\text{predicted Interval} = M_{i,t+1} \pm Z_{\alpha/2} \cdot \sqrt{\text{Var}(\hat{M}_{i,t+1}^{\text{ARIMA}}) + \text{Var}(\epsilon_{i,t+1})}$$

where:

$M_{i,t+1}$: country $i$'s medal count in the year $t+1$

$Z_{\alpha/2}$: Critical value of the standard normal distribution

$Var(\widehat{M}_{i,t+1}^{ARIMA})$: Variance of the prediction.

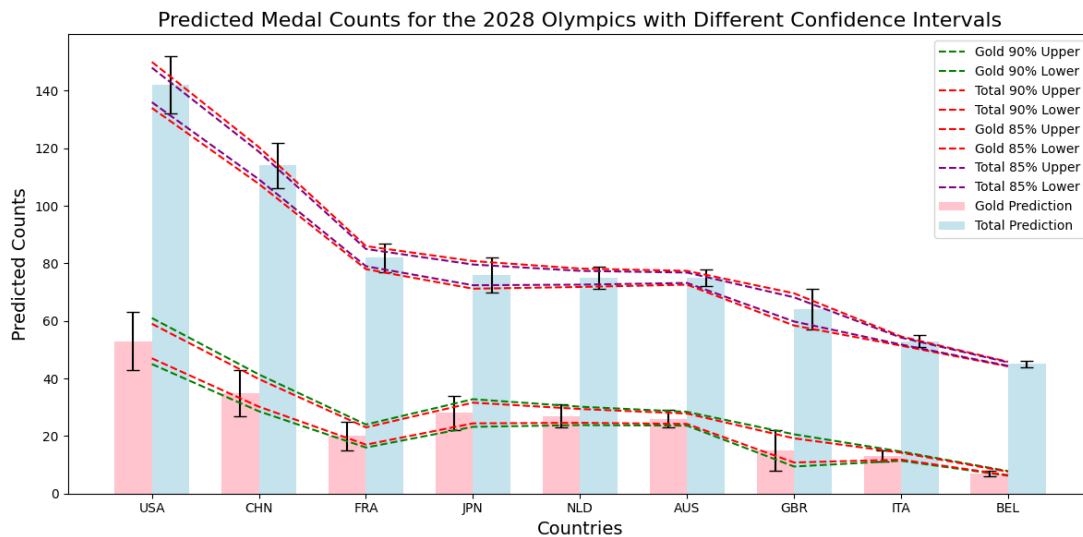$Var(\epsilon_{i,t+1})$: residual variance

*Figure 4.3.2 predicted Total Medals for 2028 Olympics Game*

First, the United States is predicted to win 53 gold medals and 142 total medals, placing it at the top of the list for the 2028 Olympics. And at a 95% confidence interval, the predict interval is (43,63) for gold medals and (132,152) for total medals. At 90% confidence interval, the predict interval of gold medal is (45,61) and the floating interval of total medal is (136,150). The predict medals count intervals at different confidence levels are considered.

iii. Analyse the possibilities for countries to progress and regress.

Here, we calculate national progress or regression based on the predicted values:

$$\Delta G_i = G_{i,2028} - G_{i,2024}$$

$\cdot if\ \Delta G_i > 0,\ then\ the\ country\ i\ progress$

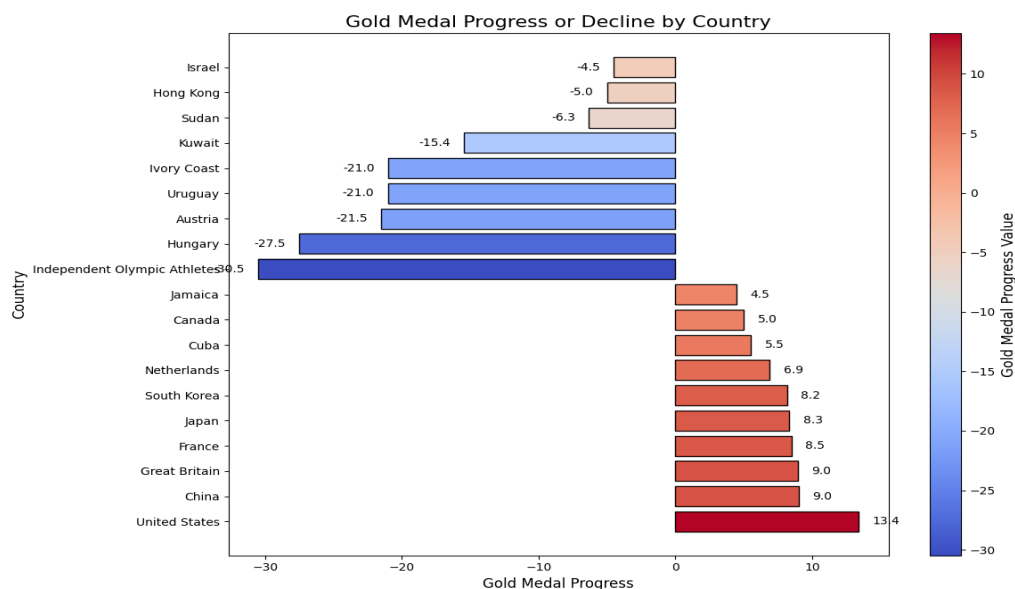$\cdot if\ \Delta G_i \leq 0,\ then\ the\ country\ i\ regress$



*Figure 4.3.3  predicted Gold Medal Progress or Decline of Country for 2028 Olympics Game*

This graph visualizes the top 10 countries that will progress and the top 10 countries that will regress in the predicted outcome. Firstly, USA is in the first place with an improvement value of 13.4, which may be related to the host effect. Secondly, Independent Olympic Athletes regressed the most with -30.5, which may be related to the instability of Independent Olympic Athletes.

(2) Sub problem 2 :

i.    give the number of country who will earn their first medal in the next Olympics

Next, we will look at countries that have not yet won an Olympic medal and predict the likelihood that they will win a medal for the first time in 2028. Since there is no time series data for countries that have not yet won a medal, it is not possible to use the ARIMA model. Therefore, we adapt the XGboost model as a classification model, with the model formula:

$$P(Y_i = 1|X_i) = f(X_i),$$

where $Y_i = 1$ indicates that country get a medal, $Y_i = 0$ indicates that country won't get a medal. Here we set the threshold to 0.5, if $P(Y_i = 1|X_i) > 0.5$, then the country will win a medal at the next Olympic Games and vice versa. Then we get the The total number of countries predicted to win medals is:

$$N_{\text{new}} = \sum_{i \in \text{NoMedal}} \mathbb{I}(P(Y_i = 1|X_i) > 0.5)$$

ii.    Give the possibilities of this estimate

The probability that at least k countries will win a medal:

$$P(\text{At least } k \text{ countries win medals}) = \sum_{j=k}^{N} \binom{N}{j} p^j (1-p)^{N-j}$$

Based on the calculations, we get that the following 8 countries are likely to win their first gold medal at the 2028 Olympics, with an average probability greater than 0.5.

*Table 4.3.4  country which probably will get their first gold medal in 2028 Olympics Game*

| country | probability | country | probability |
|---|---|---|---|
| Azerbaijan | 0.61 | Kiribati | 0.55 |
| Bangladesh | 0.57 | Bolivia | 0.52 |
| Yugoslavia | 0.56 | Serbia | 0.51 |
| Algeria | 0.56 | Chile | 0.51 |

(3) Sub problem 3 :

i.    Give the relationship between the events and how many medals countries earn.

Competition event related features were quantified in terms of their contribution to the number of medals by feature importance analysis of the XGBoost model. The feature importance formula is:

$$G_k = \frac{\sum_{j \in \text{Splits of } k} \text{Gain}_j}{\sum_{m \in \text{All Features}} \sum_{n \in \text{Splits of } m} \text{Gain}_n},$$

Where:

$G_k$: The gain share of feature $k$

$Gain_j$: Gain value of feature $k$ on split node j

$\sum_{j \in Splits\ of\ k} Gain_j$: The sum of the gains of feature $k$ over all split nodes.
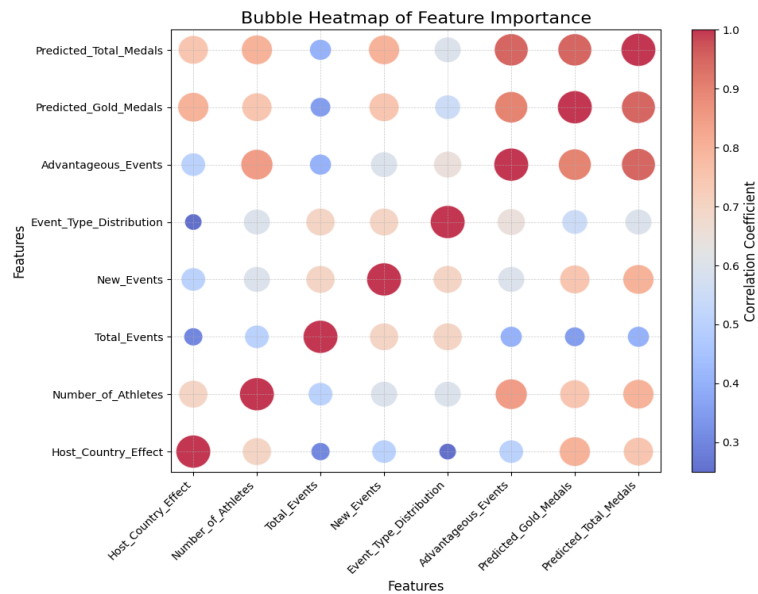
*Figure 4.3.5  Bubble Heatmap of Different Feature Importance for Predicted Country Medals*

Firstly, new events, advantageous events, and host country effects show high feature importance, and these factors contribute a lot to the number of gold medals predicted and the number of total medals predicted. Secondly, total_events and event_type_distribution exhibit very small feature importance and contribute very little to the number of gold medals and total medals predicted.

ii.   Give sports are most important for various countries and the reason

In order to analyse the relationship between each country and a particular event, i.e. to identify several dominant events in a country, we performed a Spearman correlation analysis.



*Figure 4.3.6  Spearman correlation of events with country*

(1) China (CHN)

From the analysis of the results, it can be seen that China's dominant events are Badminton, Artistic Gymnastics and Table Tennis, which show correlation values of 0.94, 0.96 and 0.93 respectively. These values indicate that an increase in the number of such events will directly affect China's performance and ranking.

(2) Australia (AUS)

From the analysis of the results, it can be seen that Australia's dominant events are swimming,

badminton and weightlifting, which show correlation values of 0.85, 0.81 and 0.78 respectively. These values indicate that an increase in the number of such events will directly affect Australia's performance and ranking.

(3) United States (USA)

From the analysis of the results, it can be seen that USA's dominant events are Athletics, Swimming and Basketball, which show correlation values of 0.81,0,79,0,71 respectively. These values indicate that an increase in the number of such events will have a direct impact on the performance and ranking of the USA.

iii.  Explain how host country impact results by choosing events

In order to explore the impact of the addition of new event in the host country on the distribution of medals, an independent samples t-test was used to analyse the differences in medal wins between a country when it is a host and a non-host country.

Here, we set assumption：

$H_0$：$\mu_1 = \mu_2$ ( $Host\ country\ status\ has\ no\ significant\ impact\ on\ medal\ wins$)

$H_1$：$\mu_1 \neq \mu_2$ ( $Host\ country\ status\ has\ a\ significant\ impact\ on\ medal\ wins$)

Where:

$\mu_1$ indicates the average number of medals when the country is the host country.

$\mu_2$ Indicates the average number of medals won by the same country when it was a non-host country.

Then use the t-statistic formula:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

where:

$\bar{x}_1, \bar{x}_2$ indicates the average number of medals for host and non-host samples

$s_1^2, s_2^2$ indicates the variances of the two sample groups

$n_1, n_2$ indicates the capacities of the two sample groups

Tested at a significance level of α = 0.05, the p-value is significantly less than the significance level of 0.05, and we reject the original hypothesis, suggesting that there is indeed a significant effect of host country status on medal winning. The data show that on average, a country wins 13.2 more medals when it is the host country than when it is not, and this increase accounts for about 34.7% of its regular medal count.

# 5 Model Ⅱ：Evaluation of Great Coach Effects Based on Difference-in-difference Regression and AHP-TOPSIS Model

**Objective:** to verify whether the great coach effect exists, to assess the contribution of the great coach effect, and to evaluate the return on investment.

**Methods:** Firstly, the paired t-test was used to preliminarily judge the changes between the validation group and the control group before and after the intervention with or without excellent coaches so as to preliminarily judge whether the coaching effect exists or not. Then, the double difference method was used to further quantify the specific contribution of the coaching effect and

assess the impact of the great coach effect; finally, the AHP-TOPSIS analysis model was used to decide which sports programs should be invested in "great coaches".

# 5.1 Verifying "great coach" effect

### 5.1.1 The concept of Pair-T Test

The paired t-test is a statistical method used to compare the difference in means between two related sets of samples. The method is often used to make comparisons between measurements of the same set of samples under different conditions, or two measurements on the same individuals under different treatments or experimental conditions. Here, we verify the existence of a coaching effect by measuring the difference in means between several groups of game data before the hiring of a quality coach and game data after the hiring of a quality coach.

### 5.1.2 Hypothesis

We set the hypothesis：

$H_0$： No significant change of the performance before and after the coach's tenure.

$H_\alpha$： Exists significant change of the performance before and after the coach's tenure.

### 5.1.3 Data Collection

We first filtered out the data with significant changes and looked for any effects of coaching changes in these data. Through our analysis, we found that the coaching change did have a significant impact on performance, as evidenced by an increase in the athletes' performance. To quantify this impact, we assigned weights to the results based on the number of medals won: 0 points for no medals won, and 1, 2 and 3 points for bronze, silver and gold medals won, respectively.

After searching the material, we found that Lang Ping served as head coach of the Chinese women's volleyball team in 1996, as head coach of the U.S. women's volleyball team in 2008, and was rehired by China in 2016. Her coaching resume has changed more significantly, so the data is representative. The data is consolidated below:

*Table 5.1.3  Chinese and US Volleyball data in Olympic*

| Coach | Country | Year of Appointment | Scores before | Scores after | Difference |
|---|---|---|---|---|---|
| Lang Ping | China | 1996 | 0 | 2 | 2 |
| Lang Ping | The United States | 2008 | 0 | 2 | 2 |
| Lang Ping | China | 2016 | 0 | 3 | 3 |
| Béla Károlyi | The United States | 1984 | 0 | 15 | 15 |

### 5.1.4 Pair-T Test

Calculating the Mean of Difference:

$$\overline{d} = \frac{2 + 2 + 3}{3} \approx 2.33$$

Calculating the Standard Deviation of Difference

$$S_d = \sqrt{\frac{0.33^2 + 0.33^2 + 0.77^2}{3 - 1}} \approx 0.6367$$

According to the T-test formula

$$t = \frac{\bar{d}}{S_d/\sqrt{n}}$$

We could get the result:

$$t \approx 7.0, \text{ and } p \approx 0.0198$$

Since the p-value is significantly less than 0.05, we can reject that coaching turnover has a significant effect on performance improvement.

## 5.2 Quantifying "great coach" effect

### 5.2.1 The concept of difference-in-difference regression
The difference-in-difference Regression method is a statistical experimental design method for policy effect-type studies. By comparing the changes in the treatment and control groups before and after the policy intervention through treatment, the data of the experimental groups are basically consistent, satisfying the 'parallel trend assumption' and the 'individual treatment stability assumption', and the net effect of the intervention is derived by eliminating the time trend and the differences between groups. In solving the present problem, the net effect of the intervention in the presence of an excellent coach was calculated by excluding the change in the number of medals due to external factors and the difference in the base level before the intervention.

### 5.2.2 Data reprocessing
Before building the multiple difference regression model, we focus on coaches with cross-country coaching experience and identify and extract the countries and programs these coaches have coached in the dataset given in the title. Through the changes in the performance of these countries and programs, we can construct the treatment group, (country-program combinations employing international master coaches) and the control group (country-program combinations not employing international master coaches) to build the following multiple difference regression model. For example, Lang Ping coaches women's volleyball, the Caroi's coach gymnastics, Wolfgang Mayer racewalking coaches, and Mihaly Berra coaches gymnastics.
In the end, using 225 total samples, the treatment and control groups were for example about 1:1.5.
**Quantification of medal scores:**
   0 points: no medals
   1 point: bronze medal awarded
   2 points: silver medal won
   3 points: gold medal won

### 5.2.3 Difference-in-difference Regression Model Building
**(1) Prerequisite assumptions**
There are no other significant events during the intervention period such as changes in policy affecting medal scores.

**(2) Parallel trend test**
The performance trends of the treatment group (teams that changed coaches) and the control group (teams that did not change coaches) were close to parallel.

**(3) Model formula**
The regression model of the double-difference approach explains the causal effect of the data through multiple beta regression coefficients jointly:

$\beta_0$: the baseline level of the control group before the intervention

$\beta_1$: inherent differences between the treatment and control groups before the intervention

$\beta_2$: the generalized effect of time trends on medal scores

$\beta_3$: The net effect of the intervention, the "great coaching effect"

**Two groups were included:**
Experimental group: the country hired a "great coach"
Control group: the country did not hire a "great coach"

**The formula is:**
$$Y_{it} = \beta_0 + \beta_1 \cdot Treatment_i + \beta_2 \cdot Post_t + \beta_3 \cdot (Treatment_i \times Post_t) + \epsilon_{it}$$

Among:
**Dependent variable:**
$Y_{it}$: Indicates the country's/program's medal score for the year, scored based on the number of medals won.
**Independent variable:**
$Treatment_i$: Indicates the indicator variable for the treatment group. A value of 1 means that the country hired a "great coach" (experimental group) and a value of 0 means that no coach was hired (control group).
$Post_t$: Indicator variable for time. A value of 1 indicates that it was in the post-intervention period, i.e., after the coach was hired, and a value of 0 indicates that it was in the pre-intervention period, i.e., when no coach was hired.
$Treatment_i \times Post_i$: Indicates an interaction term between treatment group and time, capturing the additional post-intervention change in the treatment group relative to the control group, reflecting the net effect of the Great Coach Effect.
**Constant term:**
$\beta_0$: The expected value of the dependent variable when all independent variables are zero, corresponding to the control group's average medal score before the intervention.
**Coefficient:**
$\beta_1$: indicates the mean difference between the treatment and control groups before the intervention.
$\beta_2$: indicates the generalized change in the country or program before and after the intervention due to time trends. Trend in medal scores over time for all countries or programs without the intervention.
$\beta_3$: Net contribution of the 'great coach effect', the impact of hiring a 'great coach' on medal scores.
**Error term:** $\epsilon_t$: Random error
**(4) Perform descriptive statistics on the data set:**

*Table 5.2.3 Datasets descriptive statistics*

| Treatment | Post | count | mean | std |
|---|---|---|---|---|
| 0 | 0 | 90 | 1.515 | 0.262 |
|  | 1 | 45 | 1.756 | 0.302 |
| 1 | 0 | 60 | 1.374 | 0.300 |
|  | 1 | 30 | 2.860 | 0.195 |

### 5.2.4 Model Training
The model was trained using python regression analysis tool, ordinary least squares OLS, to perform double difference analysis.

(1)First, the goal of training the model is to find the best linear relationship between the independent and dependent variables by minimizing the error. We use the ordinary least squares method with its core idea:

$$\widehat{\beta} = \text{argmin}_\beta \sum_{i=1}^{n} (y_i - X_i \beta)^2$$

Among:

$y_i$: The actual value of the dependent variable.

$X_{it}$: The matrix of independent variables.

β: The regression coefficients

$\widehat{\beta}$: Optimal regression coefficients, obtained by minimizing the sum of squared residuals.

(2)Pass the most important interaction term independent variable Treatment_Post to the model, who captures the net effect of the coaching effect.

(3)Define the model after adding the constant term intercept, which represents the baseline effect, the value of the dependent variable when all independent variables are zero. Without the constant term, the model may be forced through the origin, resulting in a poor fit.

### 5.2.4 Model Fitting
Ordinary Least Squares Model The OLS model calculates regression coefficients by minimizing the residual sum of squares (RSS):

$$\text{RSS} = \sum_{i=1}^{n} (y_i - X_i \beta)^2$$

Among:

$y_i$: The actual value of the dependent variable.

$X_{it}$: The matrix of independent variables.

β: The regression coefficients
The optimal regression coefficient is obtained by solving the following equation:

$$\widehat{\beta} = (X^T X)^{-1} X^T y$$

Among:

$X^T$: The transpose of the independent variable matrix.

$(X^T X)^{-1}$: Inverse matrix of the independent variable matrix.

$X^T y$: Covariance of the independent and dependent variables.

### 5.2.5 Use of robust standard errors
In real data, heteroskedasticity (i.e., the variance of the residuals is not constant) may exist. To cope with this in our model, robust standard errors (cov_type='HC1') are specified in the .fit() method to adjust the standard errors of the coefficients to make them more robust.

## 5.3 Model results analysis

### 5.3.1 Overall model performance

*R-squared (R² ): 0.891*

The model explains 89.1% of the variation in the dependent variable (score), indicating a good model fit.

*Adj.R-squared (Adjusted R² ): 0.890*

The adjusted $R^2$ is very close to the $R^2$ , which further indicates that the model has a strong explanatory power.

*F-statistic and Prob(F-statistic):*

The F-value is 427.2 and the p-value is 0.0011591, which is a small p-value indicating that the model is significant overall.

### 5.3.2 Coefficient results and analysis

*Table 5.3.2  DID Regression model result*

|                | coef    | std err | z      | P>\|z\|    | [0.025 | 0.975] |
|----------------|---------|---------|--------|----------|--------|--------|
| const          | 1.3868  | 0.014   | 93.511 | 0.000512 | 1.359  | 1.415  |
| Treatment      | -0.0143 | 0.024   | -0.597 | 0.550000 | -0.061 | 0.033  |
| Post           | 0.3001  | 0.025   | 11.793 | 0.000327 | 0.250  | 0.350  |
| Treatment_Post | 0.9574  | 0.046   | 20.610 | 0.00013  | 0.866  | 1.048  |

**Constant term:** the level of significance is extremely high (p-value < 0.05), indicating that the constant term is significant. The mean score of the control group before the addition of the great coach was 1.3868.

**Experimental group:** p-value of 0.550 is not significant (p-value > 0.05), indicating that the difference in scores between the treatment and control groups before the addition of the great coach is not significant. It is consistent with the parallel trend hypothesis, which means that the performance trend of the two groups before the great coach joined is close to each other.

**After the presence of the great coach:** significant (p-value < 0.05), indicating that there is a significant increase in the scores of the control group after the great coach joins. The mean score of the control group improved by 0.3001 points compared to the score before the great coach joined. The improvement in the scores of the control group may be due to external factors.
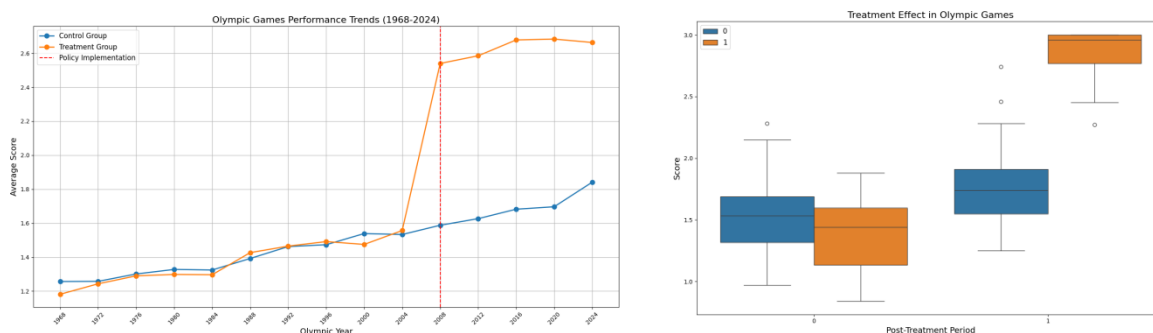


*Figure 5.3.2  Different set change in great coach effect*

**Interaction term:** significant (p-value < 0.05), indicating that the addition of the great coach had a significant effect on the treatment group.Indicates that the scores of the treatment group improved by an additional 0.9574 points over the control group after the addition of the great coach. The scores of the treatment group improved significantly and this improvement can be attributed to the effect of the policy.

## 5.4 Evaluation of Great Coach Program Selection Based on AHP-TOPSIS Modeling

### 5.4.1 Data Collection

Three characteristics, namely, the cumulative number of medals won, the total number of gold medals and the number of competitions, are selected to assess whether each program needs to invest in "great coach".We will choose the United States, Italy and France as examples for analysis.

### 5.4.2 AHP-TOPSIS modeling

The AHP-TOPSIS method is able to rank a limited number of evaluation objects according to their proximity to the idealized goal, and take into account a variety of factors affecting the national medals, so that it can give a more comprehensive weighting of the indicators.

**(1) The AHP-TOPSIS hierarchy** is constructed as follows:
**Goal layer:** select the best investment project
**Guideline layer:** Total medal count, Total gold count, Number of entries
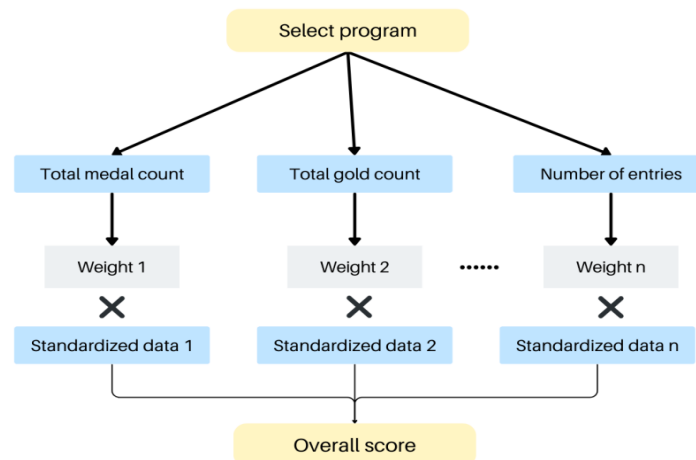**Program layer:** Standardized data



*Figure 5.4.2  AHP-TOPSIS hierarchy*

**(2) AHP modeling:**

According to the constructed evaluation system, the indicators at the same level are compared to determine the relative importance of the batches of indicators at the same level, and the comparative judgments at different levels are constructed by using Thomas Sethi's "1-9 scale method":

*Table 5.4.2  1-9 scale*

| Scale | Meaning |
|---|---|
| 1 | Two elements are equally important compared to each other |
| 3 | The former is slightly more important than the latter |
| 5 | The former is significantly more important than the latter |
| 7 | The former is more important than the latter. |
| 9 | Riding is absolutely more important than the latter. |
| 2,4,6,8 | Intermediate value of adjacency judgment |

| 1/2~1/9 | If element i scores $C_{ij}$ compared to element j, then element j is compared to element i by $1/C_{ji}$ |

**Matrix Satisfaction:** $C_{ij}>0$, $C_{ij}=1$, $C_{ij}=1/C_{ji}$ (i,j=1,2,3,...,n)

Assume $A_m$ is related to the factors $C_1, C_2, C_3, \ldots, C_n$ in the next level. The judgment matrix is shown above. $C_{ij}$ means the relative importance of $C_i$ compared to $C_j$ with respect to $A_m$.

**Calculate the weights of the indicators:**

Calculate the product of the elements of each row of the judgment matrix C

$$M_i = \prod_{j=1}^{n} aij(i = 1,2, \ldots, n)$$

$$\varpi m = \sqrt[n]{M_i}$$

vector normalization process, $\varpi = (\varpi_1, \varpi_2, \varpi_3, \ldots, \varpi_n)$, using formula:

$$\varpi_i = \frac{\varpi_i}{\sum_{j=1}^{n} \varpi_j}$$

Deriving weighting factors $\overline{\varpi}_i$. Then consistency testing:

$$\lambda_{max} = \frac{1}{n}\sum_{j=1}^{n} \frac{(B\omega)_j}{\omega_j}$$

$$CI = \frac{\lambda_{max} - n}{n - 1}$$

If the value of CI is large, the consistency of the matrix is worse.

Evaluating the random consistency index RI:

Matrices of different orders are evaluated with different consistency indicators RI, as represented in the figure below:

*Table 5.4.2  RI value in different matrix order*

| Matrix order | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| RI | 0 | 0 | 0.52 | 0.87 | 1.01 | 1.22 | 1.29 | 1.39 | 1.42 | 1.47 |

In order to determine whether the matrix is in satisfactory agreement, it is necessary to compare the CI with the RI:

$$CR = \frac{CI}{RI}$$

CR is the test coefficient.

CR=0: the judgment matrix has good consistency.

CR<0.1: the consistency is good

CR≥0.1: the judgment matrix consistency is not good, should modify the value of the matrix until CR<0.1 consistency, need to compare CI with RI.

**(3) TOPSIS model building:**

$$Z = (Z_{ij})_{n \times m} \ , \ Z_{ij} = w_j \times x_{ij}$$

Determine the positive and negative ideal solutions:

$$Z_j^+ = max(x_{ij}) \ , \ Z_j^- = min(x_{ij})$$

Calculate the combined appraisal value of each program:

$$S_i = \frac{d_i^-}{d_i^- + d_i^+}$$

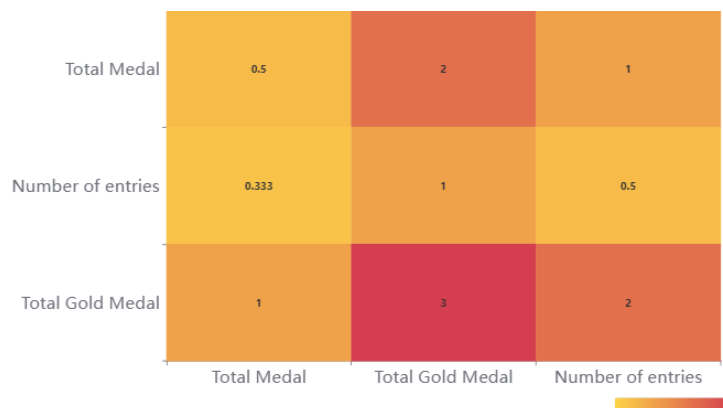**5.4.3 AHP-TOPSIS model solving:** The results of the judgment matrix heat map are as follows.



*Figure 5.4.3  judgment matrix*

Obtained from the above judgment matrix:

Consistency indicator CI=0.0047                    Consistency ratio CR=0.0078

Since CR<0.1, the consistency of this judgment matrix is acceptable. The weights of each indicator are obtained as shown in the table below:

*Table 5.4.3  Indicator weight*

| Indicator | Weight |
|---|---|
| Total medal count | 0.5412 |
| Total gold count | 0.1485 |
| Number of entries | 0.3103 |

After obtaining the weights for each indicator, the scores for each sport in each country and the three countries and sports most in need of investing in great coaches were calculated, with the following results:

**For USA: Volleyball, Cycling Road**

The USA has been doing well in women's volleyball sports, winning medals most of the time, but men's volleyball needs to be strengthened by bringing in great coaches to bring advanced techniques and training methods to USA volleyball players. The U.S. has a good foundation in cycling road racing, but there is still a gap compared to traditional cycling powerhouses, and could invest in bringing in coaches from cycling powerhouses like the Netherlands or Belgium.

**For Italy: 3x3 Basketball, Baseball**

Italy's performance in baseball still falls short of the world's powerhouses, and the introduction of foreign coaches from traditional baseball powerhouses, such as the United States or Japan, could be effective in improving the technical level of the players. In traditional basketball, there is still a lot of room for improvement, and foreign coaches from 3x3 basketball powerhouses such as Serbia and Latvia can be brought in to improve the effectiveness of decision-making in the game.
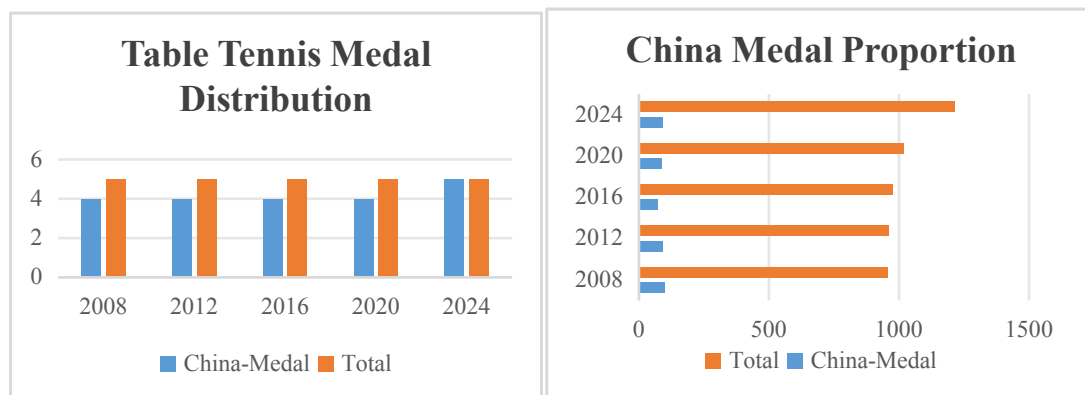
**For France: Fencing, Hockey**

France could invest in the introduction of experts from fencing powerhouses such as Italy or Hungary, which could provide French athletes with more diversified tactical ideas and technical improvements. France is relatively weak in sports such as ice hockey and field field hockey, and hiring coaches from Canada and Russia could help the French team improve their skills.

# 6 Problem Ⅲ: Medal Distribution and Strategy Optimization

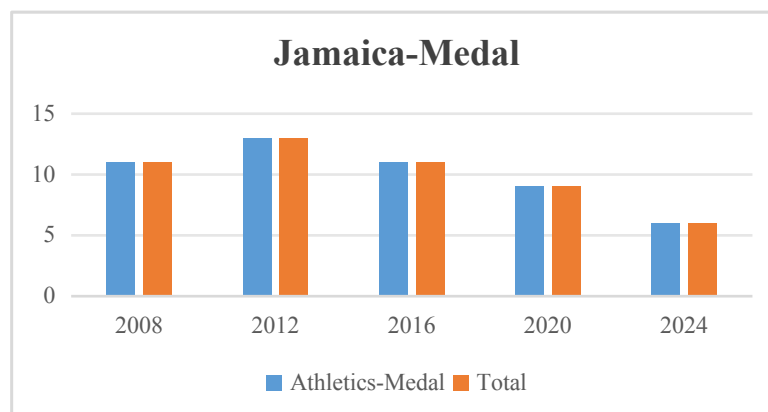## 6.1 Specialization and Centralization of Programs

By analyzing the data obtained from the model, it can be seen that certain countries win far more gold medals in specific sports than in others. For example, some countries may have won a large number of medals in traditional sports such as swimming and table tennis, and fewer medals in other sports such as shooting and judo. This imbalance suggests that these countries' dominant sports

account for a large portion of their medal totals. China is a striking example of this, as we analyzed medal data from only five Olympic Games and found that its table tennis program has been a monopoly in the world.



*Figure 6.1.1 China model value*

We can clearly see that China not only has essentially a monopoly on all medals in table tennis, but also consistently ranks favorably in the world's total medal count. For these fully developed monopolistic powers, there are unique historical advantages in specific sports which can provide a solid foundation for total medals. Therefore, the National Olympic Organizing Committee should focus on the training and early selection of young athletes, and cultivate young athletes with international competitiveness by strengthening the accumulation of experience in international competitions. Thus, it will pass on the favorable historical position and help the country continue to compete for medals in subsequent competitions.



*Figure 6.1.2 Jamaica medal value*

However, for some countries, there are instances where competitiveness is based solely on dominant events. Jamaica has a significant advantage in track and field sprinting events, especially in the 100 meters for both men and women. A big part of the reason for this is the fact that the country has a very good runner - Usain Bolt. According to the data analysis, Jamaica's medals are won through the sprint events, and as the country's athletic Bolt grows older and no longer has the advantage in the sprint events, the country's total number of gold medals has a downward trend. Therefore, for this kind of country only in a particular project dominance, the Olympic Committee needs to be targeted to optimize the allocation of resources, to avoid over-reliance on the cultivation of a particular project, in order to prevent the failure of a project as expected to play and lead to sharp fluctuations in the number of medals. In addition, the country needs to capitalize on its own genetic strengths to increase medal diversity by investing more in other related sports. Emerging sports could also be an important area of competition in future Olympics, and although many of

these sports do not yet hold a significant place in the traditional sports systems of countries that have a monopoly on such speciality sports, new sports could provide new breakthrough opportunities for such countries as the Olympic program continues to evolve.

## 6.2 Stabilization of Country Rankings

By compiling the data, it was found that some countries do not always rely on gold medals for their rankings, but rather secure their total medal count through undetermined silver and bronze medal performances. Germany, for example, is a country that excels in a number of Olympic sports, especially in areas such as shooting, weightlifting, swimming, gymnastics, track and field, and winter sports. Although it does not usually win as many gold medals as the United States or China, it is able to maintain a strong presence in the Olympic medal table with a more consistent number of silver and bronze medals.

For those countries that have been able to stabilize their competitive position through lower medals, we suggest that these countries can further improve their gold medal competitiveness through more targeted training, psychological coaching, and international exchanges, while maintaining a stable medal output. In addition, we have confirmed the existence of the "coaching effect" from our previous analysis. Therefore, these countries are best suited to invest in quality coaches, as they are already capable of winning medals, indicating that the average training quality of their athletes is high, and good coaches may be able to lead them to greater breakthroughs at this time.

# 7 Discussion

## 7.1 Strengths
For Model I, we combine the traditional time series model ARIMA with the machine learning model XGBoost. This feature enables the model to provide a more comprehensive analysis in the face of the complex distribution of Olympic medals, avoiding simple speculation relying only on historical data, and improving the accuracy and reliability of the prediction.
For model II, we combine differential regression and AHP-TOPSIS model with paired t-test, which can quantify the contribution of the "Great Coach Effect", and clearly show the importance of each indicator and its influence on the final decision through hierarchical analysis.
For Problem 3, our model not only takes into account countries that are already performing well, but also helps other countries that have yet to make their mark in certain sports to find opportunities to improve. After quantifying the impact of each factor in a comprehensive manner, the potential of future changes and emerging sports is taken into account, based on unique and innovative recommendations.

## 7.2 Possible improvements

Our model is based on the given datasets, which limits the influence of other possible factors to some extent, so in order for the model to have more accurate predictions, more dimensions can be added to the model in the future, such as adding data on national investment in sports, the degree of improvement of national athletes' training facilities, and so on. Secondly, more possible influencing factors can be considered, such as adding socio-economic indicators GDP, the size of the sports industry, and considering factors such as the influence of climate and geographical factors on the development of sports programs. These improvements will help to increase the accuracy and usefulness of the model.

# References

Gergely Csurilla; Imre Fertő;  *"The Less Obvious Effect of Hosting The Olympics on Sporting Performance"*,  *SCIENTIFIC REPORTS,  2023*

Feifei Li; Will G Hopkins; Patrycja Lipinska;  *"Population, Economic and Geographic Predictors of Nations' Medal Tallies at The Pyeongchang and Tokyo Olympics and Paralympics",  FRONTIERS IN SPORTS AND ACTIVE LIVING,  2022.*

Mark Baimbridge;  *"Outcome Uncertainty in Sporting Competition: The Olympic Games 1896-1996",  APPLIED ECONOMICS LETTERS,  1998.*