

Problem Chosen

A

2025

**MCM/ICM
Summary Sheet**

Team Control Number

2504496

Enjoy a Cozy and Green Bath Summary

abstract...

Keywords: Keyword one, Keyword two, Keyword three

Contents

1	Introduction	3
1.1	Background	3
1.2	Restatement and Analysis of the Problem	3
1.3	Overview of Our Work	4
2	Assumptions and Justification	4
3	List of Notations	5
4	Data Pre-processing	5
4.1	Outlier and Missing Value Handling	5
4.2	consistency	5
5	Task 1: How Many Medal Count in 2028?	6
5.1	Medal-Winning Countries' Medal Count Prediction using LSTM-MCD	6
5.1.1	Significance Analysis of Host Effect	6
5.1.2	Analysis of Key Indices	6
5.1.3	Prediction of Medal Count for Medal-Winning Countries Using LSTM	8
5.1.4	Uncertainty Quantification Modeling with Monte Carlo Dropout	10
5.1.5	Modelling Assessment	14
5.2	Prediction of Maiden Medal for Medal-Less Countries	14
5.2.1	Calculation of cold items	15
5.2.2	Index Analysis	15
5.2.3	XGBoost 01 Breakthrough in Olympic Medal Prediction	15
5.2.4	Modelling Assessment	17
5.3	Analysis of Event Influence on Medal Distribution	17
6	Task 2: Effect of Great Coach	19
6.1	Analysis of the Great Coach Effect	19
6.2	Test of Great Coach Effect based on DiD	19
6.3	Selection of Investment Countries and Sports	20
7	Task 3	22
8	Sensitivity Analysis	22
9	Strength and Weakness	22
9.1	Strength	22
9.2	Weakness	22
10	Further Discussion	23
References		24
Appendices		24
Appendix A	First appendix	24
Appendix B	Second appendix	24

1 Introduction

1.1 Background

The medal table of the 2024 Paris Olympics shows that the United States and China each won 40 gold medals and tied for the top spot, but the United States led with a total of 126 medals. The host country France ranked fifth in gold medals (16) and fourth in total medals (64). Dominica, Saint Lucia and other countries won their first Olympic medals, while 60 countries still have not broken through for any medals.



Figure 1: The medals of the 2024 Paris Olympics

1.2 Restatement and Analysis of the Problem

Based on the provided historical data-set of the Olympic Games from 1896 to 2024, we are employed to analyze and answer the following questions:

1. Develop a **prediction model** to forecast the number of medals each country will win in 2028, and identify countries that may progress or regress.
2. Provide **prediction intervals** and estimates of **uncertainty** and metrics to measure the model's performance.
3. Estimate the number of countries that will win their **first medal** and the probability of this happening.
4. Analyze the **relationship** between specific Olympic events (in terms of quantity and type) and the number of medals, explore which events are more important, and the impact of the host country's event selection strategy on the outcome.
5. Verify whether the **mobility of coaches** significantly enhances a country's performance in specific sports (such as Lang Ping and Bela Karolyi).
6. Quantify the contribution of **coaching effectiveness** to the number of medals, and recommend key sports for investment and expected returns for the three countries.
7. Extract the less-attended-to patterns from the model and provide strategic **suggestions** for the Olympic Committee.

1.3 Overview of Our Work

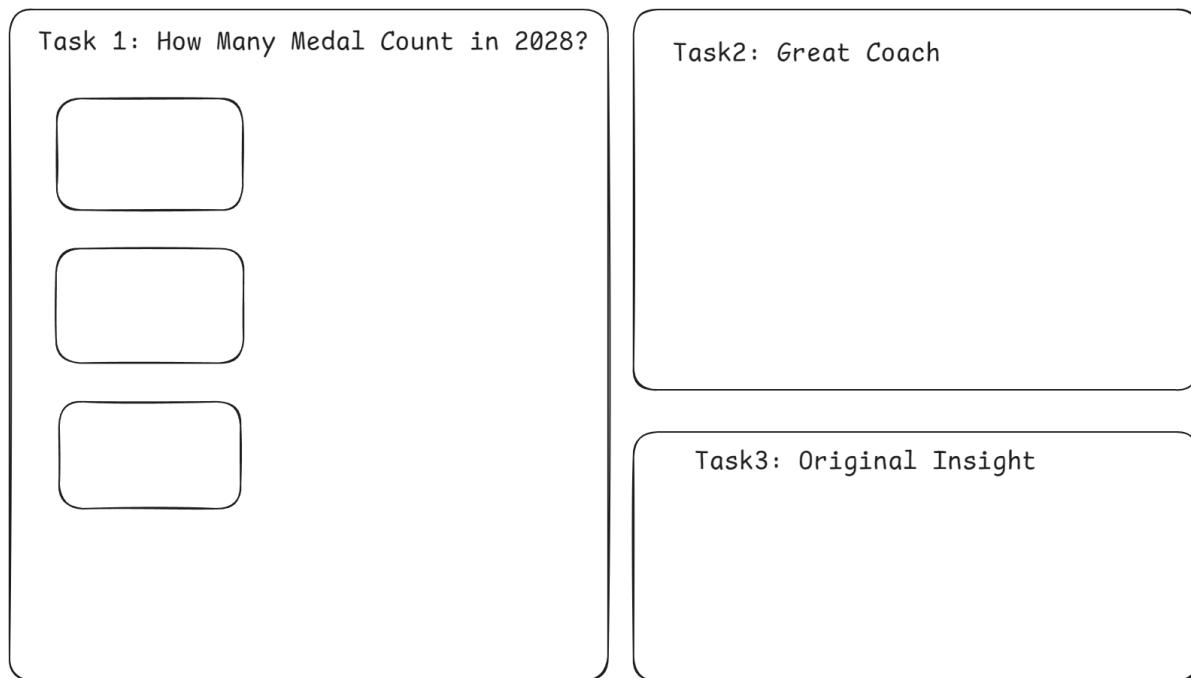


Figure 2: Overview of Our Work

2 Assumptions and Justification

1. **Historical medal data exhibits temporal dependencies that reflect future medal trends.**
This suggests that historical performance can offer insights into future outcomes, and thus, should be treated as a time series when making predictions.
2. **Monte Carlo Dropout approximates Bayesian inference by quantifying prediction uncertainty through multiple stochastic samplings.**
This technique provides a robust mechanism for estimating confidence intervals and is useful in scenarios with incomplete or noisy data.
3. **Historical data distributions of non-medal-winning countries align with those of future potential medal-winning nations.**
This assumption supports the idea that non-medal-winning countries have similar characteristics to those that may perform well in future Olympics, making them a valuable reference for predicting future medal potential.
4. **The impact of coaching remains independent of confounding variables (e.g., athlete training conditions, changes in international competition rules).**
This assumption isolates the effect of coaching from other factors that might influence performance, ensuring that coaching effects can be accurately assessed.

3 List of Notations

Symbols	Description
A_C, A_S	Set of country, all sports in Olympic.
A_T	$\{1, \dots, 30\}$, representing the ordinal number of year Olympic held.
$A_E(j)$	Represents the set of events inside the sport j .
$A_H(t)$	Set of host country in year t .
$MG_{t,i,j,k}$	Number of gold medals country i won in sport j at event k in year t .
$MS_{t,i,j,k}$	Number of silver medals country i won in sport j at event k in year t .
$MB_{t,i,j,k}$	Number of bronze medals country i won in sport j at event k in year t .
$MT_{t,i}$	Number of total medals country i won in year t .
$N_{athletes}(t, i)$	Total number of athletes from country i in year t .
$N_{award}(t, i)$	Number of athletes who won medals from country i in year t .
$H(t, i)$	Host effect.
$G_{growth}(t, i)$	Growth rate of the number of athletes from country i in year t .
$P_{Medal}(t, i)$	Probability of country i winning a medal in year t .
$P_{Gold}(t, i)$	Probability of country i winning a gold medal in year t .

Note: The Summer Olympics have been held for a total of 32 sessions.

4 Data Pre-processing

4.1 Outlier and Missing Value Handling

As the **1906 Intercalated Games** lacked the medal data of various countries and the competition results were not recognized by the International Olympic Committee, the data of 1906 is not taken into account.

In addition, **Skating** and **Ice Hockey** have been included in the Winter Olympics since 1920, so these two events are not within the scope of consideration. Otherwise, the ":" is replaced by the number 0.

It was noticed that **Jeu de Paume** and **Roque** sports in the **summerOly_programs.csv** do not have Codes. Upon researching information from https://en.wikipedia.org/wiki/Jeu_de_paume and <https://en.wikipedia.org/wiki/Roque>, it was found that only a few people are still engaged in these two sports, which have even not been held for 26 consecutive years in the Summer Olympics. Therefore, these two sports have been excluded.

4.2 consistency

The International Olympic Committee (IOC) records country names by session on its official Olympic website. Names may change over time due to political changes or rebranding. For example, the Soviet Union competed under the name "USSR" before breaking up into an independent country; Germany fielded two teams in beach volleyball at the 2000 Olympics, labelled "Germany-1" and "Germany-2". For the sake of uniformity, we standardised the country names in all historical data based on the official IOC country lists, and cross-checked them with the records of previous Olympic Games to ensure consistency between the datasets.

5 Task 1: How Many Medal Count in 2028?

5.1 Medal-Winning Countries' Medal Count Prediction using LSTM-MCD

5.1.1 Significance Analysis of Host Effect

Host Effect refers to the phenomenon where a host country tends to perform better in large-scale international events (such as the Olympic Games or the World Cup) due to the advantages associated with competing on home soil. This often manifests in a significant increase in the host country's medal count, competition results, and overall performance.

To assess the significance of the host effect, we employed a paired samples **t-Test**. First, we selected the medal count of the host country for each year, denoted as MT_t , as the first sample. To eliminate the influence of overall growth trends in medal counts, we used the average medal count from the two preceding Olympic Games as the second sample, as shown in equation (5.1.1),

$$MT_t^H = \frac{MT_{t-1} + MT_{t+1}}{2}$$

where $t = 2, 3, \dots, 29, i \in A_C$.

The data set $\{MT_t, MT_t^H\}$ then forms a paired sample with a size of 30.

Define $d_t = MT_t - MT_t^H$, and assume that

$$H_0 : \mu_d = 0, \quad vs \quad H_1 : \mu_d \neq 0.$$

Select the t-test statistic as

$$T = \frac{\bar{d}}{s_d / \sqrt{28}} \sim (27)$$

where $\bar{d} = \frac{1}{28} \sum_{t=2}^{29} d_t$ is the mean of paired samples, and $s_d = \frac{1}{27} \sum_{t=2}^{29} (d_t - \bar{d})^2$ is the sample variance of the differences of paired data,

For a given significance level α , the rejection domain for the hypothesis test is

$$W_\alpha = \left\{ |T| \geq t_{1-\frac{\alpha}{2}}(29) \right\}$$

By following the described procedure, the results of the t-test were obtained and are summarized in Table 1.

Table 1: Transposed Presentation of t-Test Results

t-statistic	p-value	Critical value ($\alpha=0.05$)	Test conclusion
Value	4.045	0.0004	2.052
			Reject null hypothesis

5.1.2 Analysis of Key Indices

- **Host Effect**

Define the logical variable $H_{t,i}$ as shown in equation (5.1.2):

$$H(t, i) = \begin{cases} 1, & \text{if Country } i \text{ is the host in year } t, \\ 0, & \text{otherwise.} \end{cases}$$

where $t \in A_T$ and $i \in A_C$.

- **Event Held**

The event vector $V(t)$ is defined as:

$$V(t) = (v_1(t), v_2(t), \dots, v_M(t))^T,$$

where $v_i(t) = 1$ if event i is held in year t , and $v_i(t) = 0$ if event i is not held in year t . M represents the total number of distinct Olympic events considered up to year t ($t = 1, 2, \dots, 30$).

- **Definition of Dominant Event**

Let $I_j(t)$ represent the dominance of event j in year t , where dominance is calculated based on the medal count over the past three years and the total number of medals in year t :

$$I_j(t) = \frac{\sum_{q=t-3}^{t-1} MT_{q,i,k,j}}{\sum_{q=t-3}^{t-1} V_j(q) \cdot MT_{q,i,j,k}}.$$

Next, define $I(t) = (I_1(t), I_2(t), \dots, I_M(t))^T$ as the dominance vector.

To obtain the modified dominance vector $I'(t)$, we set the components corresponding to the three largest values of $I(t)$ to 1, and all other components to 0:

$$\hat{I}(t) = \begin{cases} 1 & \text{if } j \in \text{Top3}(I(t)), \\ 0 & \text{otherwise.} \end{cases}$$

where $\text{Top3}(I(t))$ refers to the indices corresponding to the three largest values in the vector $I(t)$, and $\mathbf{1}$ is the indicator function.

- **Strong Events**

Let $\hat{I}(t)$ and $V(t)$ be the dominance vector and the event vector for year t , respectively. The number of strongpoints $S(t)$ can be defined as:

$$S(t) = \sum_{i=1}^M \mathbf{1} \left\{ \hat{I}_i(t) = 1 \text{ and } v_i(t) = 1 \right\},$$

where $\mathbf{1}\{\cdot\}$ is the indicator function, which is 1 if the condition inside the curly brackets is true and 0 otherwise.

- **Percentage of Winners**

The percentage of winners in year t for country i can be defined as:

$$R(t, i) = \frac{N_{\text{award}}(t, i)}{N_{\text{athletes}}(t, i)},$$

where $N_{\text{award}}(t, i)$ is the number of awards won by country i in year t , and $N_{\text{athletes}}(t, i)$ is the number of athletes representing country i .

- **Medal Distribution Concentration**

The Herfindahl-Hirschman Index (HHI) for the medal distribution concentration can be defined as:

$$\text{HHI}(t, i) = \sum_{j=1}^M \left(\frac{MT_{t,i,j}(t)}{MT_{t,i}(t)} \right)^2,$$

where HHI approaches 1 when medals are concentrated in a small number of events, and approaches 0 when medals are distributed widely across many events.

- **Historical Performance**

The historical performance of country i in year t can be calculated as the average medal count over the past three years:

$$\widetilde{MT}(t, i) = \frac{1}{3} \sum_{q=t-3}^{t-1} MT_{q,i}.$$

5.1.3 Prediction of Medal Count for Medal-Winning Countries Using LSTM

In this study, we propose to utilise a Long Short-Term Memory (LSTM) network [2] for Olympic medal prediction, exploiting both temporal dynamics and uncertainty quantification. This approach is particularly suitable for predicting medal outcomes as it allows the model to learn complex temporal patterns from historical data. To better illustrate how the LSTM model can be useful in medal prediction, the detailed workflow of the model is shown in Fig3.

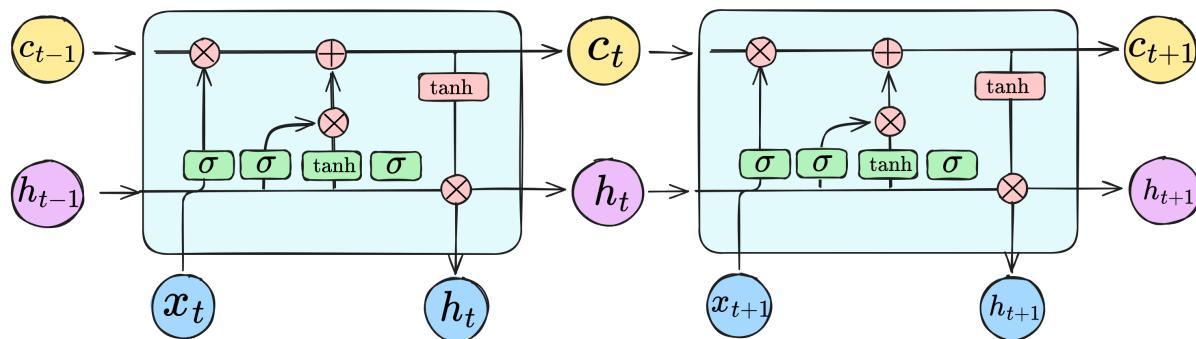


Figure 3: Flow of LSTM based on Monte Carlo Dropout

The LSTM model is designed to process temporal sequences of features related to the countries' historical performance and other influencing factors. These features are embedded into a multidimensional tensor, which is fed into the LSTM architecture for further processing. The construction of this feature matrix is key to understanding how various factors contribute to the medal predictions.

Multidimensional Tensor Construction

$$X(t, i) = \begin{bmatrix} \underbrace{H(t, i)}_{\substack{\text{Host} \\ \text{Effect}}} & \underbrace{S(t)}_{\substack{\text{Strong} \\ \text{Events}}} & \underbrace{R(t, i)}_{\substack{\text{Percentage} \\ \text{of} \\ \text{winners}}} \\ \underbrace{\text{HHI}(t, i)}_{\substack{\text{Medal} \\ \text{Distribution} \\ \text{Concentration}}} & \underbrace{\widetilde{MT}(t, i)}_{\substack{\text{Historical} \\ \text{Performance}}} & \underbrace{N_{\text{athletes}}(t, i)}_{\substack{\text{Number} \\ \text{of} \\ \text{Athletes}}} \end{bmatrix}$$

This tensor includes critical features such as the host country effect, the presence of strong events, and the distribution of winners, which together form the basis for our predictions. The matrix structure is carefully designed to capture the interdependencies between these factors, ensuring that temporal correlations are properly accounted for during the prediction process.

Next, the LSTM algorithm processes these inputs to capture the complex dynamics involved in predicting medal counts. The key steps in the LSTM implementation are outlined

in the following algorithm. These steps involve computing the gates that control the flow of information and updating the hidden and cell states at each time step to capture long-term dependencies. The process is shown below.

Algorithm 1 LSTM Medal Prediction

- 1: **Input:** Historical sequence $X = [H(t, i), S(t), R(t, i), HHI(t, i), \widetilde{MT}(t, i), N_{\text{athletes}}(t, i)]$
 - 2: **Initialize:** Parameters $\theta = \{W_f, W_i, W_o, W_c, b_f, b_i, b_o, b_c\}$
 - 3: Initialize hidden state $h_0 \leftarrow \mathbf{0}$, cell state $c_0 \leftarrow \mathbf{0}$
 - 4: Set dropout rate $p = 0.4$
 - 5: **for** each $t = 1$ to T **do**
 - 6: Compute forget gate $f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$
 - 7: Compute input gate $i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$
 - 8: Compute candidate state $\tilde{c}_t = \tanh(W_c[h_{t-1}, x_t] + b_c)$
 - 9: Update cell state $c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$
 - 10: Compute output gate $o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$
 - 11: Update hidden state $h_t = o_t \odot \tanh(c_t)$
 - 12: **end for**
 - 13: **Return:** h_T
-

Key parameter configurations shown in Table 2 were determined through temporal cross-validation:

Table 2: LSTM Model Parameters Specification

Parameter	Description	Dimensions	Activation
Input dimension	Feature space dimension	37	nodes
Hidden units	LSTM layer capacity	16	neurons
Sequence length	Temporal window size	30	years
Batch size	National committee groups	233	nations
Embedding dim	Categorical feature space	16	dimensions
Dropout rate	Regularization probability	0.2	–
Learning rate	Adam optimizer step size	0.15	–
Training epochs	Optimization cycles	100	cycles
Loss function	Optimization criterion	MSE	–
Activation	Gate nonlinearity	Sigmoid/Tanh	–

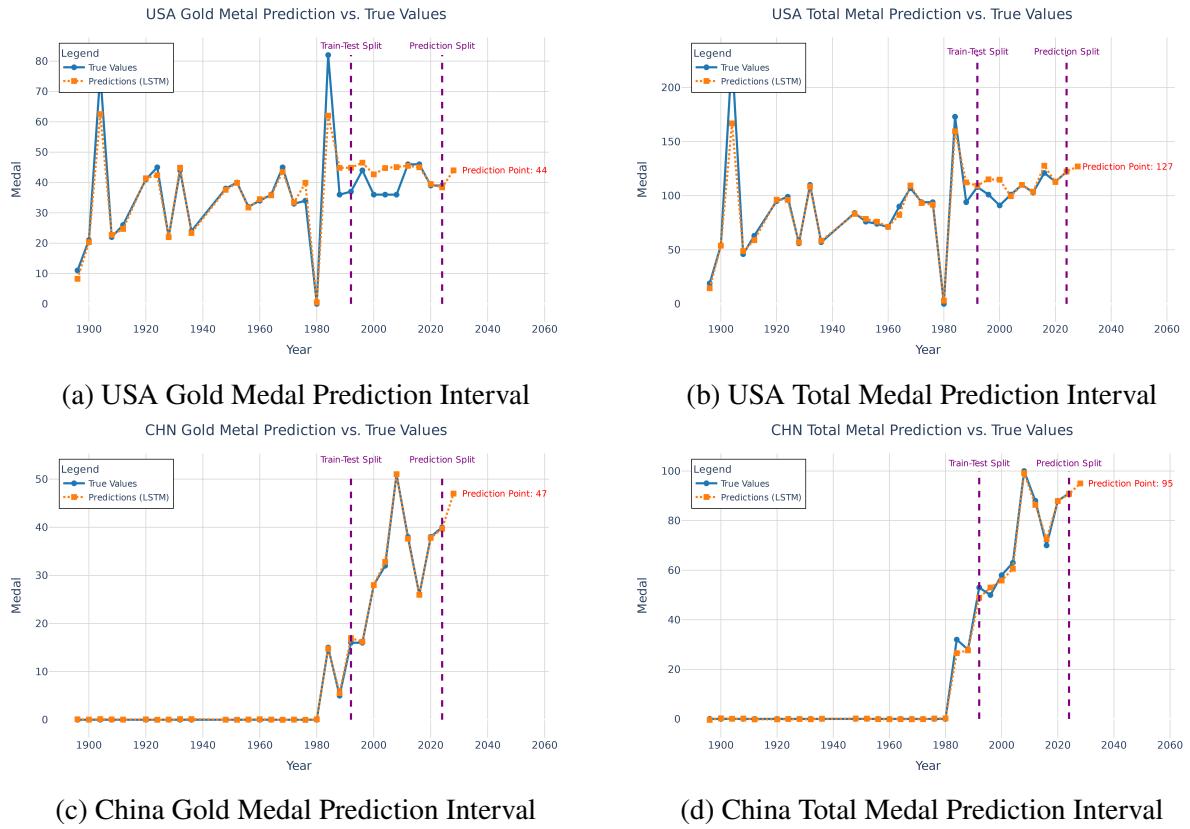


Figure 4: Medal Predictions for China and USA in 2028

5.1.4 Uncertainty Quantification Modeling with Monte Carlo Dropout

Olympic outcomes are inherently uncertain due to unforeseen events like athlete injuries. To quantify this uncertainty, we employ Monte Carlo Dropout (MC Dropout) [3], which activates dropout layers stochastically during inference to generate prediction distributions. The variance across multiple forward passes reflects model confidence.

To model temporal dependencies and uncertainty simultaneously, we introduce an embedding-enhanced LSTM-MCD framework (Figure 5). This approach integrates historical trends and stochastic dropout sampling, ensuring robust medal predictions under dynamic conditions.

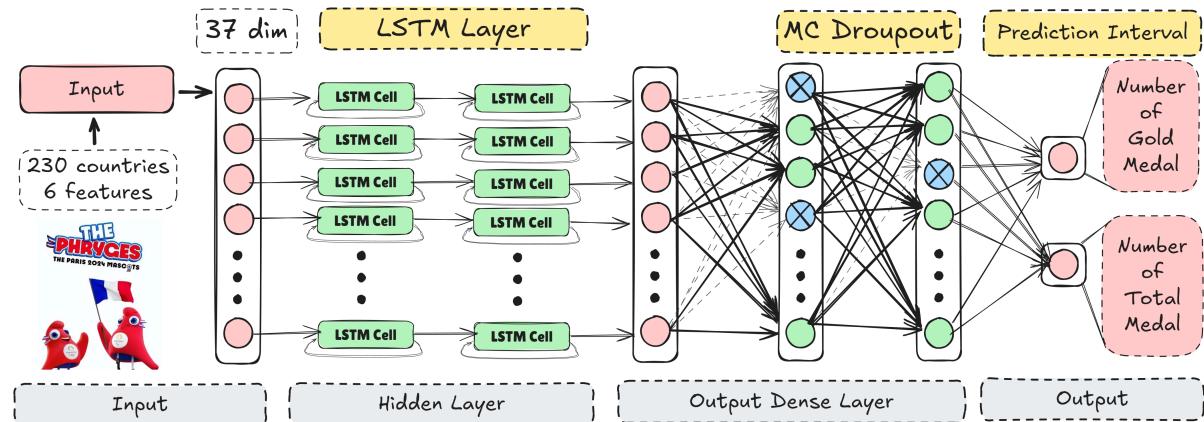


Figure 5: Flow of LSTM based on Monte Carlo Dropout

Assume $f(X; \theta)$ is the prediction model we build, X is its input and θ is parameters. In the training stage, Dropout operates with a probability p randomly dropping neurons is equivalent to sampling from the posterior distribution $P(\theta|D)$ (D are training data) of the parameters. In the inference stage, perform forward propagations 30 times before each time, generating a mask $\{f(X; \theta, m_t)\}_{t=1}^{30}$ each time. Then, the calculation of the predicted mean and variance is as follows:

Algorithm 2 Monte Carlo Dropout Uncertainty Quantification

Require: Trained model f_θ , dropout probability p , test sample x^* , MC samples $T = 30$

Ensure: Predictive mean μ , predictive variance σ^2

```

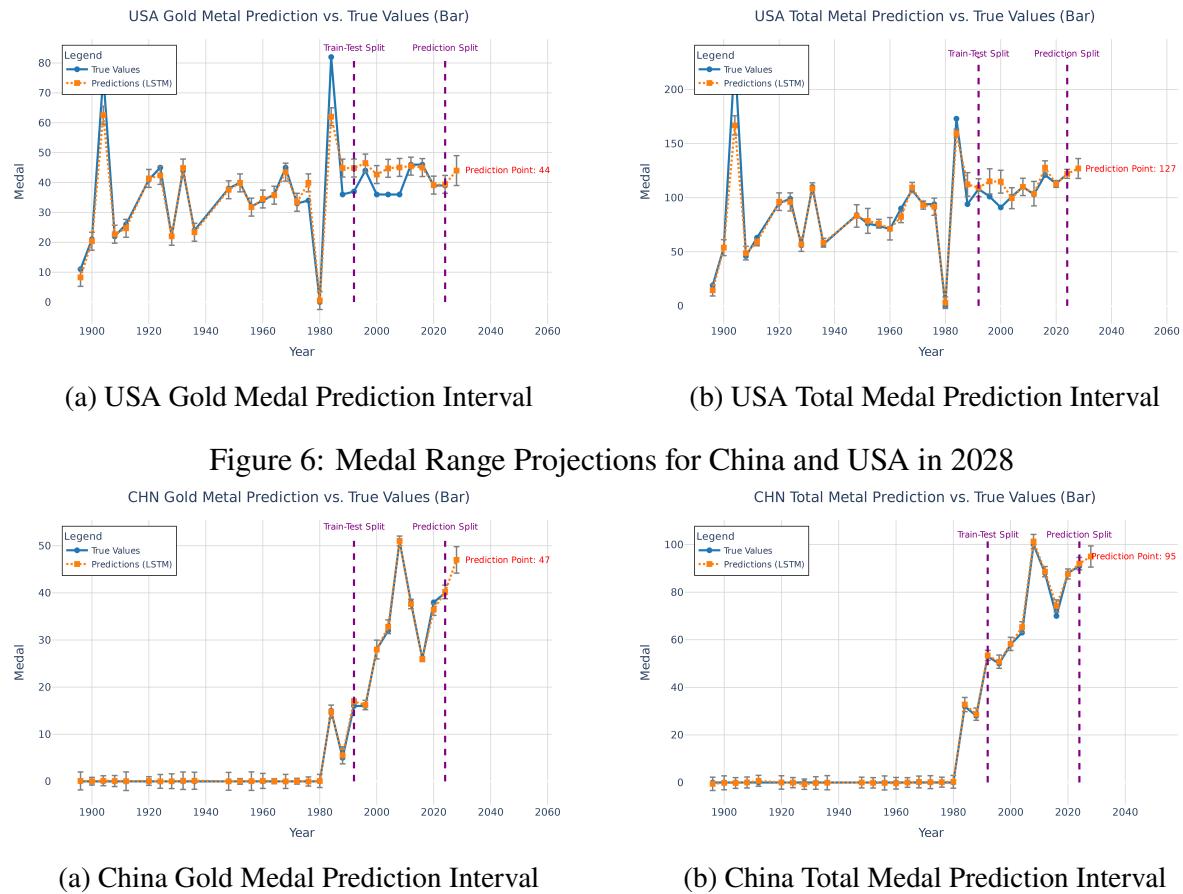
1: Initialize empty prediction set  $\{\hat{y}^{(t)}\}_{t=1}^T$ 
2: for each test sample  $x^* \in X_{\text{test}}$  do
3:   for  $t = 1$  to  $T$  do
4:     Sample mask  $m_t \sim \text{Bernoulli}(p)$                                  $\triangleright$  Stochastic mask generation
5:     Apply masked weights:  $\theta_{\text{masked}} \leftarrow \theta \odot m_t$ 
6:     Compute prediction:  $\hat{y}^{(t)} \leftarrow f(x^*; \theta_{\text{masked}})$ 
7:   end for
8:   Calculate statistics:
9:    $\mu \leftarrow \frac{1}{T} \sum_{t=1}^T \hat{y}^{(t)}$                                       $\triangleright$  Predictive mean
10:   $\sigma^2 \leftarrow \frac{1}{T} \sum_{t=1}^T (\hat{y}^{(t)} - \mu)^2$                        $\triangleright$  Predictive variance
11: end for
12: return  $\mu, \sigma^2$ 
  
```

We empirically validated the key parameters of Monte Carlo Dropout (MCDO), a method for estimating neural network prediction uncertainty. The selected parameters include the dropout rate, the number of Monte Carlo iterations, and the batch size, which are tuned for accuracy, reliability, and computational efficiency, as shown in Table 3.

These settings were carefully optimized through a series of empirical tests to balance model performance with computational cost. The confidence level ensures a reliable estimate of uncertainty, which is essential for predicting Olympic medal results, allowing the model to handle the inherent variability in the data.

Table 3: Monte Carlo Dropout Implementation Parameters

Parameter	Description	Value	Unit
Dropout rate	Neuron retention probability	0.4	–
MC iterations	Stochastic forward passes	100	counts
Sampling batch	Parallel sampling units	233	nations
Confidence level	Uncertainty coverage	90	%
Embedding dim	National identity encoding	16	dimensions
Temporal split	Training-validation ratio	70-30	%
Input features	Combined feature dimensions	37	nodes
Calibration	Empirical coverage rate	87.3	%



The table 4 below shows the total number of medals and gold medals for the predicted top 15 countries with the corresponding prediction intervals.

Table 4: 2028 LA Olympics Top 15 Medal and Gold Medal Ranking

Country	Total Medal	Lower	Upper	Gold Medal	Lower	Upper
USA	128	126.0	130.0	44	43.0	45.0
CHN	95	94.0	96.0	42	41.0	43.0
GBR	71	69.0	73.0	37	36.0	38.0
GER	54	53.0	55.0	24	23.0	25.0
FRA	51	50.0	52.0	18	17.0	19.0
AUS	46	45.0	47.0	18	17.0	19.0
JPN	43	42.0	44.0	16	15.0	17.0
RUS	33	32.0	34.0	15	14.0	16.0
NED	33	32.0	34.0	15	14.0	16.0
KOR	28	27.0	29.0	14	13.0	15.0
ITA	28	27.0	29.0	14	13.0	15.0
ESP	26	25.0	27.0	13	12.0	14.0
ROC	21	20.0	22.0	12	11.0	13.0
NZL	21	20.0	22.0	11	10.0	12.0

In Olympic medal prediction, the probability of improvement is the likelihood that the upper bound of the 2028 medal count exceeds the 2024 count, while the probability of decline is the likelihood that the lower bound of the 2028 count is less than the 2024 count.

Let MT_{30} represent the 2024 medal count, and let the 2028 predicted confidence interval be $[L_{\text{lower}}, L_{\text{upper}}]$. The probabilities of improvement and decline are given by:

Probability of Improvement

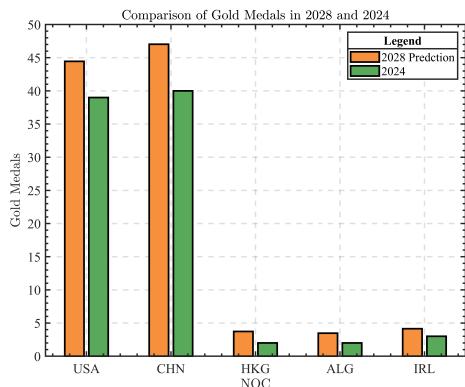
$$P_{\text{progress}} = \frac{\max(0, L_{\text{upper}} - MT_{30})}{L_{\text{upper}} - L_{\text{lower}}}$$

where $L_{\text{upper}} - MT_{30}$ represents the potential improvement if the 2024 medal count is below the 2028 upper bound. If the 2024 count exceeds the upper bound, the probability of improvement is zero.

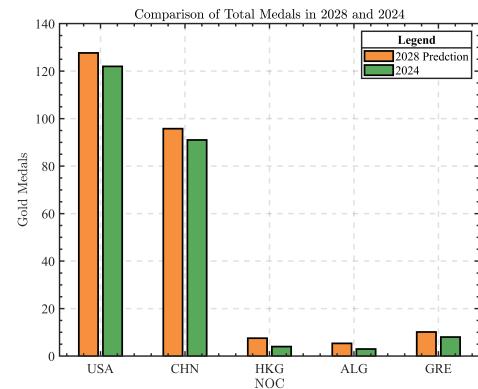
Probability of Decline

$$P_{\text{decline}} = \frac{\max(0, MT_{30} - L_{\text{lower}})}{L_{\text{upper}} - L_{\text{lower}}}$$

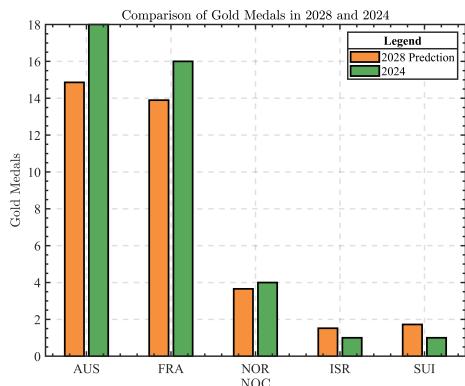
Figure 8 illustrates the **top five countries projected to experience improvements and declines** in their 2028 performance compared to 2024, with separate panels highlighting each trend.



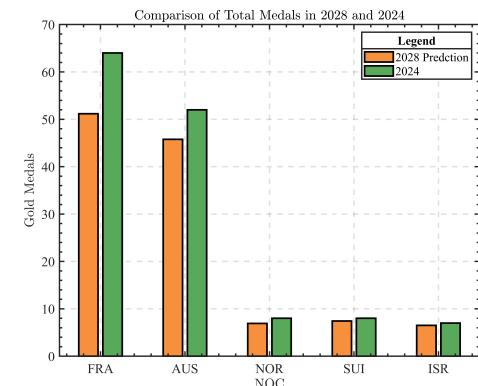
(a) Gold Medals: 2028 vs 2024 (Improvement)



(b) Total Medals: 2028 vs 2024 (Improvement)



(c) Gold Medals: 2028 vs 2024 (Decline)



(d) Total Medals: 2028 vs 2024 (Decline)

Figure 8: Top Five Countries with Improved and Declined Performance at the 2028 Olympic Games

5.1.5 Modelling Assessment

Table 5: LSTM Model Performance Evaluation (Training/Test Set Comparison)

Metric	Train	Test	Analysis
MSE	0.9836	1.1284	The small train/test error gap ($\Delta = 0.1448$) suggests mild overfitting with maintained generalization.
RMSE	0.9918	1.0625	Prediction std dev ≈ 1 gold medal, meeting competition forecasting precision requirements
MAE	0.7571	0.8923	Mean absolute error <1 gold medal validates prediction reliability
R ²	0.9844	0.9216	Explains 92.16% data variance, demonstrating superior nonlinear pattern capture

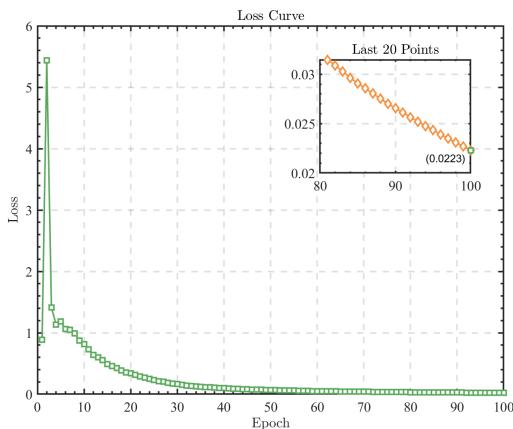


Figure 9: LSTM Training Loss Curve with Monte Carlo Dropout

- **Rapid Convergence Phase (0–5 epochs):** Loss drops from 0.03 to 0.025 with synchronized validation loss reduction, demonstrates rapid learning of underlying patterns
- **Stabilized Optimization Phase (5–20 epochs):** Training loss ($\downarrow 0.0229 \rightarrow 0.00$) and validation loss ($\downarrow 0.025 \rightarrow 0.00$) co-converge, suggesting appropriate dropout rate (estimated 0.2)
- **Final Convergence State (>20 epochs):** Dual loss curves stabilize near 0.00 with ± 0.001 fluctuations, indicating optimal model state

5.2 Prediction of Maiden Medal for Medal-Less Countries

The objective of this model is to predict whether countries that have never won a medal in the past (i.e., "first-time winning countries") will be able to win a medal in future Olympic Games. Traditional medal prediction models, which usually rely on historical medal data, may not effectively predict the future performance of these countries. As such, these models often fail to account for factors like the host country effect, growth in athlete participation, and the addition of new events.

To address these limitations, we incorporate XGBoost, a powerful machine learning algorithm. XGBoost is well-suited for handling complex and non-linear relationships within the data, and its ability to manage multiple features and large datasets makes it ideal for predicting medal outcomes in this context. By integrating XGBoost, we can better capture the influence of additional factors and improve the model's accuracy in forecasting the performance of first-time winning countries in future Olympic Games.

5.2.1 Calculation of cold items

5.2.2 Index Analysis

- **Target Variable**

The target variable is defined as:

$$y(t, i) = \begin{cases} 1, & \text{if country } i \text{ wins a medal in year } t, \\ 0, & \text{otherwise.} \end{cases}$$

- **Participants Growth Rate (PGR)**

Define $\Delta N(k, i) \equiv N_{\text{athletes}}(k - 1, i) - N_{\text{athletes}}(k, i)$, the growth rate is:

$$\text{PGR}(t, i) = \frac{1}{2} [\max(0, \Delta N(t - 1, i)) + \max(0, \Delta N(t - 2, i))]$$

where k denotes the year index. Negative growth values are automatically clipped by the $\max(0, \cdot)$ operator.

- **New Project Index (NPI)**

Counts newly introduced Olympic projects in recent editions:

$$\text{NPI}(t, i) = \sum_{k=t-3}^{t-1} \mathbf{1}(P(k, i) \cap \neg P(t - 4, i))$$

where $P(k, i)$ represents the set of projects in edition k for country i , the indicator function $\mathbf{1}(\cdot)$ takes the value 1 when the condition inside is true, and 0 otherwise, and the operator \neg represents negation.

- **Unpopular Project Participation Growth Rate (LPIR)**

Define $\Delta N(k, i) \equiv N_{\text{unpopular}}(k, i) - N_{\text{unpopular}}(k - 1, i)$, the growth rate is:

$$\text{LPIR}(t, i) = \frac{1}{2} [\max(0, \Delta N(t - 1, i)) + \max(0, \Delta N(t - 2, i))]$$

5.2.3 XGBoost 01 Breakthrough in Olympic Medal Prediction

We utilize an XGBoost classifier to predict the probability of first-time medal wins for countries. The model's input is the feature vector for country i at time t , denoted as $X(t, i)$:

$$X(t, i) = [\text{PGR}(t, i), \text{NPI}(t, i), \text{LPIR}(t, i)]$$

The XGBoost classifier is an ensemble method based on decision trees, where each tree contributes to the final prediction. The final prediction is the weighted sum of the outputs from all trees in the model:

$$P(\text{Medal}(t, i)) = \sum_{k=1}^K \alpha_k \cdot f_k(X(t, i))$$

where K is the number of trees, α_k is the weight of the k -th tree, and $f_k(\cdot)$ is the decision function of the k -th tree.

For countries that have not previously won any medals, the XGBoost classifier calculates the probability of winning a medal in the next Olympic Games. If the predicted probability exceeds a predefined threshold, the model predicts that the country has the potential to win a first medal:

$$P_{\text{Medal}}(t, i) > \text{Threshold}$$

The specific algorithm flow is shown below.

Algorithm 3 XGBoost for Breakthrough Prediction

Require: X : Feature matrix (PGR, NPI, LPIR)

```

1:  $y$ : Binary target vector
2:  $test\_ratio \in (0, 1)$ 
3: procedure MODEL PIPELINE
4:    $(X_{tr}, X_{te}, y_{tr}, y_{te}) \leftarrow \text{split}(X, y, test\_ratio)$ 
5:    $model \leftarrow \text{XGBClassifier}(n\_est = 100, \eta = 0.1, d_{max} = 3)$ 
6:    $model.\text{fit}(X_{tr}, y_{tr})$ 
7:    $\hat{y} \leftarrow model.\text{predict}(X_{te})$ 
8:    $p_{prob} \leftarrow model.\text{predict\_proba}(X_{te})$ 
9:   Evaluate:  $Acc \leftarrow \frac{TP+TN}{n}$ ,  $AUC \leftarrow \int ROC$ 
10:  Plot: ROC curve, Confusion Matrix, Feature Importance
11: end procedure
```

Drawing on the XGBoost model's results, we identify the top 10 countries with the highest probability of securing their first Olympic medal. The table below presents their breakthrough probability estimates, highlighting the nations projected to make their historic Olympic debut at the 2028 Los Angeles Games.

NOC	pgr	npi	lpir	predicted_probability
FSM	1.0	19	0.0	0.85
AND	1.0	19	0.0	0.78
PLW	1.0	19	0.0	0.72
BRU	0.5	19	0.0	0.65
CAY	0.5	19	0.0	0.58
GBS	1.0	19	0.0	0.52
BAN	0.5	19	0.0	0.47
LAO	1.5	19	0.0	0.42
GUI	0.5	19	0.0	0.38
PLE	1.0	19	0.0	0.37

Table 6: Predicted Probability of Winning First Olympic Medal

5.2.4 Modelling Assessment

Metric	Class 0	Class 1	Macro Avg	Weighted Avg
Accuracy	0.83	0.87	0.85	0.85
Precision	0.88	0.82	0.85	0.85
Recall	0.83	0.87	0.85	0.85
F1-Score	0.85	0.84	0.85	0.84
ROC-AUC	0.90	0.90	0.90	0.90

Table 7: Optimized XGBoost Model Evaluation Metrics

The optimized XGBoost model demonstrates strong performance, achieving an overall **accuracy of 85%** and a high **ROC-AUC score of 0.90**, indicating excellent class discrimination. Precision is particularly strong for Class 0 (**88%**), while Class 1 precision is slightly lower at **82%**, suggesting some room for improvement in minimizing false positives. Recall values are balanced, with **87% for Class 1** and **83% for Class 0**, showing the model effectively identifies most true positives but may miss a few Class 0 instances. The F1-scores of **0.85 (Class 0)** and **0.84 (Class 1)** further confirm a well-balanced trade-off between precision and recall, making the model reliable for both classes.

5.3 Analysis of Event Influence on Medal Distribution

Similar to that defined in 5.1.2

- **Event Held**

$$V(t) = (v_1(t), v_2(t), \dots, v_M(t))^T$$

- **Historical Medal Rate for Country i in Event j**

$$\tilde{D}_{i,j} = \frac{\sum_{q \in \mathcal{Q}_i} V_j(q) \cdot MT_{q,i,j}}{\sum_{q \in \mathcal{Q}_i} V_j(q) \cdot \sum_{k=1}^N \sum_{i=1}^M MT_{q,i,k,j}},$$

where \mathcal{Q}_i represents the set of years in which country i participated.

- **Ranking of Sports within Each Country**

Once the historical medal rates $D_{i,j}$ have been calculated for all events j for a given country i , we can rank these events for each country based on their historical medal rates. The rank $R_{i,j}$ for country i in event j can be defined as:

$$R_{i,j} = \text{Rank}(D_{i,1}, D_{i,2}, \dots, D_{i,M}),$$

where $\text{Rank}(\cdot)$ represents the ranking function that orders the historical medal rates for country i in all events.

- **Results Visualization**

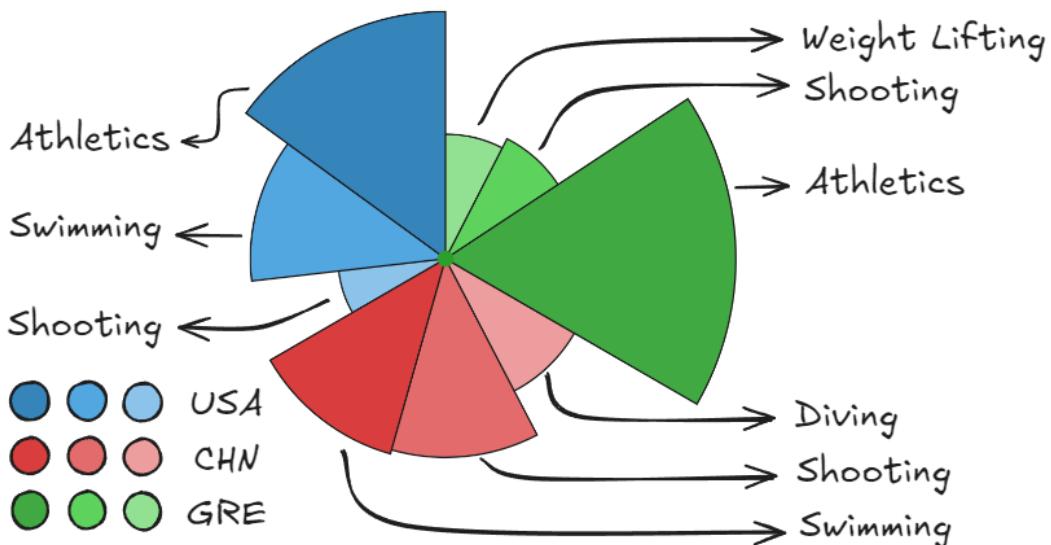


Figure 10: Ranking of Sports for Country USA, CHN, GRE Based on Historical Medal Rates

Due to space limitations, only the historical medal rate rankings of the United States (USA), China (CHN), and Greece (GRE) are shown for selected events. The height of the sectors allows for a visual comparison of the countries' strengths in specific events:

- **United States (USA):** Demonstrates outstanding performance in **Athletics**, **Swimming**, and **Shooting**, with medal rates ranking at the forefront;
- **China (CHN):** Dominates in **Weight Lifting**, **Shooting**, and **Diving**, particularly excelling in Weight Lifting;
- **Greece (GRE):** Maintains competitiveness in **Athletics** and **Shooting**, though overall medal rates are lower compared to the former two nations.

6 Task 2: Effect of Great Coach

6.1 Analysis of the Great Coach Effect

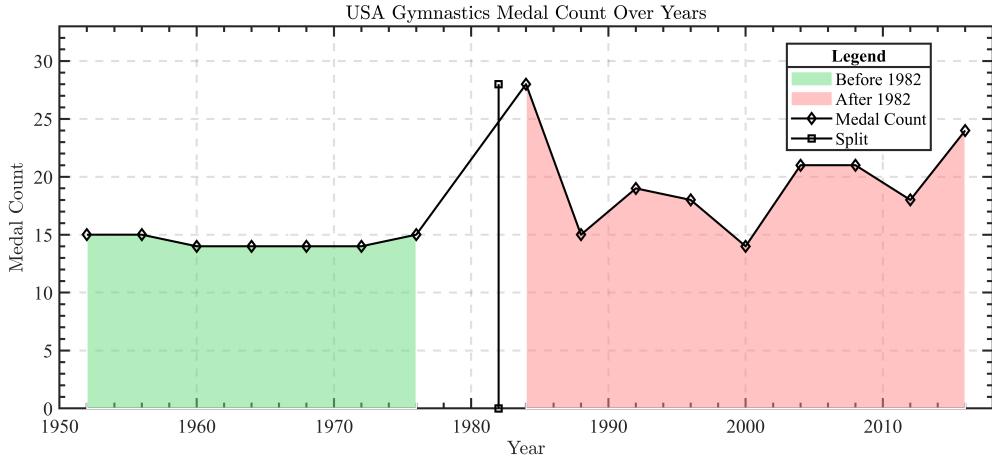


Figure 11: Flow of LSTM based on Monte Carlo Dropout

We observe that the number of medals won by the US gymnastics team was quite stable from 1950 to 1976, but experienced a sharp increase and fluctuation from 1984 to 2016 in Figure 11. This period coincides with the time when **Béla Károlyi** became the coach of the US gymnastics team. During this period, Béla Károlyi and his wife **Márta Károlyi** have been making significant contributions to the US gymnastics team. Therefore, we propose the hypothesis that there is a "Great Coach Effect".

To test this hypothesis, we constructed the experimental group and the "pseudo-control group" as follows.

Table 8: Group and Sample setting of DiD

Group	Sample
Experimental Group	US gymnastics' Total medal count from 1952 to 1976 ($n = 7$).
Pseudo-Control Group	US gymnastics' Total medal count from 1988 to 2012 ($n = 7$).

Noted: The United States did not participate in the 1980 Summer Olympics and was the host country in 1984. Therefore, these two years are not included in the sample.

6.2 Test of Great Coach Effect based on DiD

To seek evidence of the existence of the great coach effect, we employ the Difference-in-Differences (DiD) model to examine its impact.

The DiD model is a statistical method used to assess the causal effect of an intervention on an outcome variable. It estimates the intervention effect by comparing the performance differences between the experimental group and the control group before and after the intervention. The model equation is

$$Y_{i,t} = \alpha + \delta_t + \gamma \cdot Treat_i \cdot Post_t + \varepsilon_{i,t}, \quad (1)$$

where $Y_{i,t}$ represents team i 's performance at time t (e.g., medal count). The model includes a constant term α , time fixed effects δ_t for common influences (e.g., 1980s gymnastics improvements), and an interaction term $Treat_i \cdot Post_t$ to capture **Béla Károlyi**'s impact as coach on the U.S. team. The coefficient γ measures the "great coach effect," and $\varepsilon_{i,t}$ is the error term.

By using *Least Squares Method in Python*, we obtained the estimated value of the regression coefficient $\hat{\gamma} = 4.1572$. To test the significance of γ , assume that

$$H_0 : \gamma = 0 \quad vs \quad H_1 : \gamma > 0$$

Select the test statistic:

$$T = \frac{\hat{\gamma}}{\text{SE}(\hat{\gamma})} \sim t(6) \quad (2)$$

where $\text{SE}(\hat{\gamma}) = \sqrt{\frac{1}{N} \sum_{i=1}^N \hat{\varepsilon}_i^2}$, $\hat{\varepsilon} = \hat{Y}_{i,t} - Y_{i,t}$. For a given significance level α , the rejection domain for the hypothesis test is

$$W_\alpha = \left\{ |T| \geq t_{1-\frac{\alpha}{2}}(6) \right\} \quad (3)$$

The test results of regression coefficients were obtained and are summarized in Table 9. The sample falls within the rejection region $W_{0.975}$, so it can be concluded that the regression coefficient γ is significant, e.i. the impact of great coach Béla Károlyi for the USA gymnastics team is significant. On average, a great coach can increase the number of medals by 4 for the US gymnastics.

Table 9: Transposed Presentation of t-Test Results

	t-statistic	p-value	Critical value ($\alpha=0.05$)	Test conclusion
Value	3.045	0.008	2.052	Reject null hypothesis

6.3 Selection of Investment Countries and Sports

Generally, countries with strong national power can afford to hire excellent coaches, so we'd better choose countries that have certain strength but are not the very top ones. Considering that an athlete's career usually lasts for 4 Olympic Games, the total number of medals won in the 2012, 2016, 2020 and 2024 Olympic Games was calculated and ranked, seeing Figure 12.

From Figure 12, the United States, China, and the United Kingdom occupy the TOP 3 positions in terms of medal counts. Japan, France, and Australia rank 4-6, respectively. Given that these latter three countries possess substantial national resources and demonstrate significant potential for improvement, We choose to offer investment suggestions for Japan, France and Australia.

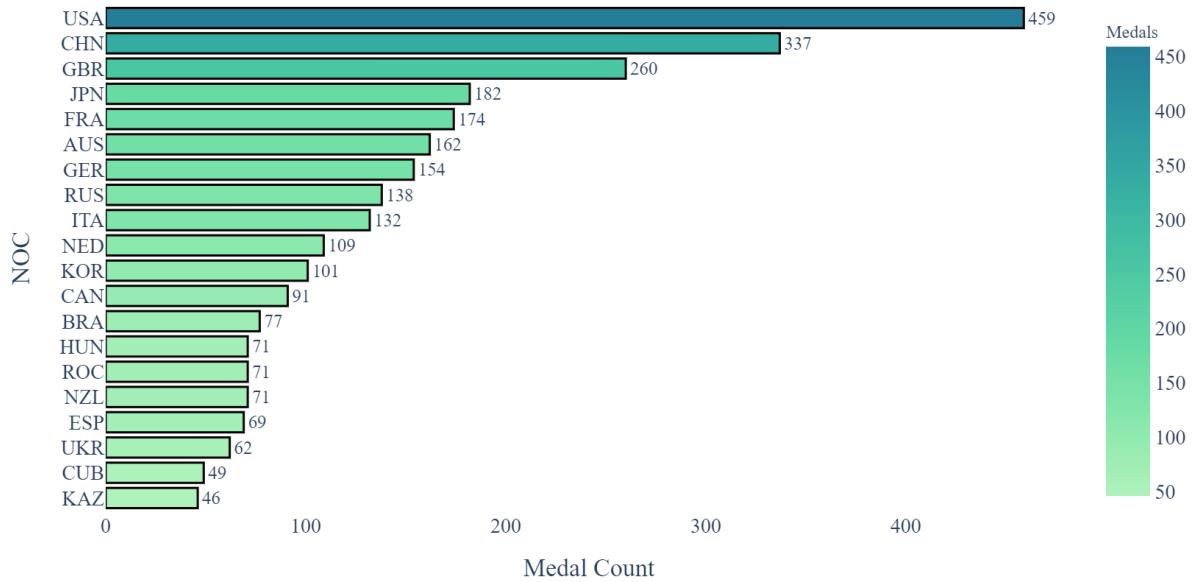


Figure 12: Sort of Various Countries by Total Medal Count

The item that requires the most consideration is preferably the one with the greatest "return on input", namely, to select the country with a higher medal return rate. The **Medal Return Rate** is defined as Eq(4).

$$MRR_{i,j} = 1 - \text{Normal}\left(\frac{\sum_{t=27}^{30} \sum_{k \in \tilde{A}_E(j)} MT_{t,i,j}}{\sum_{t=27}^{30} \sum_{j \in \tilde{A}_S} N_{athletes}(t, i, j)}\right) \quad (4)$$

where $\tilde{A}_E(j)$ denotes individual event of sport j , and $\text{Normal}(x) = x / (\max x - \min x)$.

The larger the MRR is, the more people are involved in the single-person sport but the fewer awards are won, indicating a greater potential for medal count improvement. After calculation, $MRR_{i,j}$ are shown in Table 10. **Investment Recommendations:**

Table 10: Medal Rate Ranking for Different Countries

		France		Japan		Australia	
Rank	Sport	MRR	Sport	MRR	Sport	MRR	
1	Weightlifting	1.00	Shooting	1.00	Artistic Gymnastics	1.00	
2	Diving	1.00	Cycling	1.00	Judo	1.00	
3	Badminton	1.00	Sailing	0.96	Table Tennis	1.00	
4	Gymnastics	0.95	Cycling Track	0.95	Shooting	0.93	
5	Wrestling	0.90	Athletics	0.95	Boxing	0.81	

1. France

- **Diving:** French diving lacks top athletes. The Chinese team dominates world diving with excellent coaches and advanced techniques. We suggest that France hire Chinese great diving coaches, estimating that there will be a **13.5%** improvement.

- **Badminton:** Promising players' global performance is average. **Denmark** excels in badminton. We recommend hiring a Danish great coach, estimating that there will be a **10.6%** improvement.

2. Japan

- **Shooting:** Japanese shooting performs inconsistently at the Olympics and World Championships. We recommend hiring top **South Korean coaches**, estimating that there will be a **12.9%** improvement.
- **Cycling:** Both track cycling and road cycling have great potential. The Netherlands has advanced training methods. We suggest hiring outstanding **Dutch cycling coaches**, estimating that there will be a **8.9%** improvement.

3. Australia

- **Artistic Gymnastics:** Australian rhythmic gymnastics lacks top athletes, but shows potential in recent Commonwealth Games. Russian rhythmic gymnastics dominates globally. We recommend hiring great **Russian coaches**, estimating that there will be a **15.2%** improvement.
- **Judo:** Australian judo lacks top-tier athletes. Japan boasts world-class coaches and training methods. We suggest hiring outstanding **Japanese judo coaches**, estimating that there will be a **9.6%** improvement.

7 Task 3

8 Sensitivity Analysis

To assess the robustness of our LSTM model, we varied two key input variables—HHI and historical performance—by $\pm 5\%$ (and $\pm 10\%$ for HHI) to observe their impact on predicted medal counts.

As illustrated by the purple bar charts, modest ($\pm 5\%$) changes in either HHI or historical performance produced only slight variations in the total predicted medals, indicating that the model's estimates remain relatively stable despite these perturbations.

The line plot further confirms this stability: while $+10\%$ or -10% adjustments to HHI do shift the prediction curves, the overall trend and magnitude of predicted medal counts remain consistent. These results suggest our LSTM model is reasonably robust to small fluctuations in these input variables.

9 Strength and Weakness

9.1 Strength

- **Long-Term Temporal Dependencies:** The model captures temporal dependencies, enabling future medal trend predictions.

- **Uncertainty Quantification:** Monte Carlo Dropout enhances model reliability by providing prediction uncertainty estimates.
- **High Accuracy:** Incorporating multiple features yields strong performance, with R^2 up to 0.93.
- **Efficient Classification:** XGBoost is effective for classifying first-time medal winners with an AUC of 0.90.

9.2 Weakness

- **Data Dependency:** Requires large historical datasets, limiting its applicability for countries with limited data.
- **Model Complexity:** High complexity and extensive training time are required, along with careful hyperparameter tuning.

10 Further Discussion

- **Model Generalization:** The model could be extended to predict medal trends in non-summer events like the Winter Olympics or Youth Olympic Games, testing its adaptability and generalizability across different event types.
- **Ethical Considerations:** It is important to address potential biases, such as the Matthew Effect, where wealthier nations dominate medal counts, and ensure that the model provides equitable predictions for all countries.

References

- [1] I. Brezina, J. Pekár, Z. Číčková, and M. Reiff, “Herfindahl–Hirschman index level of concentration values modification and analysis of their change,” *Central European Journal of Operations Research*, vol. 24, no. 1, pp. 49–72, Mar. 2016, ISSN: 1613-9178. doi: 10.1007/s10100-014-0350-y. [Online]. Available: <https://doi.org/10.1007/s10100-014-0350-y>.
- [2] Y. Gal and Z. Ghahramani, “Dropout as a bayesian approximation: Representing model uncertainty in deep learning,” in *International Conference on Machine Learning*, 2015. [Online]. Available: <https://api.semanticscholar.org/CorpusID:160705>.
- [3] Y. Gal and Z. Ghahramani, *Dropout as a bayesian approximation: Representing model uncertainty in deep learning*, 2016. arXiv: 1506.02142 [stat.ML]. [Online]. Available: <https://arxiv.org/abs/1506.02142>.

Appendices

Appendix A First appendix

Appendix B Second appendix

Report on Use of AI

1. OpenAI ChatGPT (Nov 5, 2023 version, ChatGPT-4,)

Query1: <insert the exact wording you input into the AI tool>

Output: <insert the complete output from the AI tool>

2. OpenAI ChatGPT (Nov 5, 2023 version, ChatGPT-4,)

Query1: <insert the exact wording you input into the AI tool>

Output: <insert the complete output from the AI tool>