

Momentum: the invisible determinants of sports matches

Summary

The 2023 Wimbledon men's singles final was up and down, with dominant players constantly switching. It is considered to be the strength or force gained by motion or by a series of events, known as momentum. This paper aims to develop a quantitative model of player performance and momentum based on the provided dataset, to investigate the correlation between them, and to develop a predictive model to predict match swings.

For Problem 1: After data pre-processing, we build a player performance evaluation model to quantify the performance of players and describe the match flow. Firstly, we consider that the points can reflect the performance of the players, and calculate the winning rate of the server as **67.31%**. Then, we construct the **exponential function** and calculate that the weight of serve points is **0.4** and the weight of return points is **0.6**. Then, based on the quantitative expression of weighted players' performances, we describe the match flow in terms of the difference between the performances of the players on both sides. Finally, we apply the model to three matches to visualize the match flow.

For Problem 2: We establish a player momentum evaluation model to quantify momentum and analyze the correlation between performance and momentum. Firstly, based on data analysis, we determine six momentum evaluation indicators. Then, based on the **CRITIC method**, we identify the weights of the indicators as shown in the table2, to get the quantitative expression of momentum. Finally, based on **linear regression** and **Pearson's coefficient** method, we find that there is a strong linear correlation between player's momentum and match swings.

For Problem 3: We develop a match swings prediction model to predict match swings and get the most relevant factors. First, we find the difference between the six momentum evaluation indicators of the two sides of the match to get the six influencing factors. Then, based on the **Random Forest** algorithm, we build a prediction model that takes the influencing factors as feature inputs and the match swings points as target inputs, and derives the most relevant features as the **volley points** and **break points**. Finally, we propose match strategies under different momentum swings.

For Problem 4: We test and generalise the match swings prediction model. First, we select a test set, figure the **AUC** and combine it with a **confusion matrix** to evaluate the predictive performance of the model. We calculate **AUC = 0.76** and **F1 Score = 0.675**, which reflects the good predictive performance of our model. We then **generalise** the model to other types of matches and find that the model has high generalisability.

For Problem 5: We summarise the results in a memo to provide coaches with information relevant to momentum.

Finally, we conduct a sensitivity analysis of the model to investigate the effect of changes in the parameters of the model variables on the results.

Keywords: Performance; Momentum; Exponential Function; Criteria Importance Through Intercriteria Correlation; Pearson Coefficient; Random Forest

Contents

1	Introduction	3
1.1	Problem Background	3
1.2	Restatement of the Problem	3
1.3	Literature Review	4
1.4	Our Work	4
2	Assumptions and Explanations	5
3	Notations	6
4	Data Pre-processing	6
4.1	Missing Values Processing	6
4.2	Abnormal Values Processing	7
5	Player Performance Evaluation Model Based on Exponential Function	8
5.1	Quantification of Player Performance	8
5.2	Description of the Match Flow	9
5.3	Application and Visualization of the Model	10
6	Player Momentum Evaluation Model Based on CRITIC Method	12
6.1	Identification and Pre-processing of Evaluation Indicators	12
6.2	Determination of Indicator Weights Based on the CRITIC Method	13
6.3	Quantitative Results of Momentum	15
6.4	Correlation Analysis of Performance and Momentum	15
7	Match Swings Prediction Model Based on Random Forest	17
7.1	Analysis of Match Swings	17
7.2	Establishment of the Model Based on the RF	17
7.3	Determination of the Most Relevant Factors	18
7.4	Match Strategies Based on Momentum Swings	18
8	Testing and Extension of Match Swing Prediction Model	19
8.1	Testing of the Model	19
8.2	Identification of Other Impact Indicators	20
8.3	Extension of the Model	21
9	Sensitivity Analysis	21
10	Evaluation of Strengths and Weaknesses	22
10.1	Strengths	22
10.2	Weaknesses	22
	References	25

1 Introduction

1.1 Problem Background

"It's a dream come true for me. It's great to win. " Alcaraz said on court after the match.[1] In the 2023 Wimbledon Gentlemen's final, 20-year-old Spanish star Carlos Alcaraz defeated Novak Djokovic, who is one of the greatest Grand Slam players of all time, to become the third-youngest men's winner at Wimbledon in the Open era.



Figure 1: Alcaraz wins first Wimbledon title[2]

The match had its ups and downs. For the first four sets, the position of dominance shifted between the two players. In the tie-break, the exciting tennis ended in favor of Alcaraz, who took control three times throughout the set.

The constant transformation of the dominant player is considered to be the result of the power contrasts that athletes acquire throughout the match, also called **momentum**. However, momentum is difficult to measure and its sources and effects are unclear. This is also the problem that needs to be addressed by our model.

1.2 Restatement of the Problem

Through an in-depth analysis and study of the context of the problem, combined with the specific constraints given, the restatement of the problem can be formulated as follows:

- Develop a model to describe and visualize the flow of the match when points occur. Identify the performance of players at a specific point during the match.
- Based on the model, analyze whether momentum affects the match.
- Build a model to predicts the momentum swings during a match. and figure out the most relevant factors.

- Test the model with data from other matches. Analyze the performance and generalization capabilities of the prediction model.
- Summarize the results and prepare a one- to two-page memo to provide to the coach, elaborating the role of momentum in the match.

1.3 Literature Review

The research on this problem is mainly divided into three parts: the quantitative analysis of momentum, and the capture and prediction of the match flow. This section focuses on the models that have been proposed.

- Firstly, as far as the research methods of momentum are concerned, there are both qualitative and quantitative research methods, with the difference that the qualitative research model considered by Vallerand et al. is dominant[3][4][5], while Cooper's quantitative research methodology is only complementary to the qualitative model, which is currently underdeveloped[6].
- Secondly, the description of the match flow requires an evaluation class model. Jin Jiasheng proposes three methods that are commonly used for evaluating sports events and makes corresponding comments on the advantages and disadvantages of each method[7].
- Finally, for the prediction of match results, the vast majority of authors use machine learning algorithms[8]. Huang Yi used BP neural networks to predict the outcome of football matches[9], and Ban Yue et al. used random forests(RF)[10].
- The advantages and disadvantages of the different models built by different research methods can be visualized as shown below:

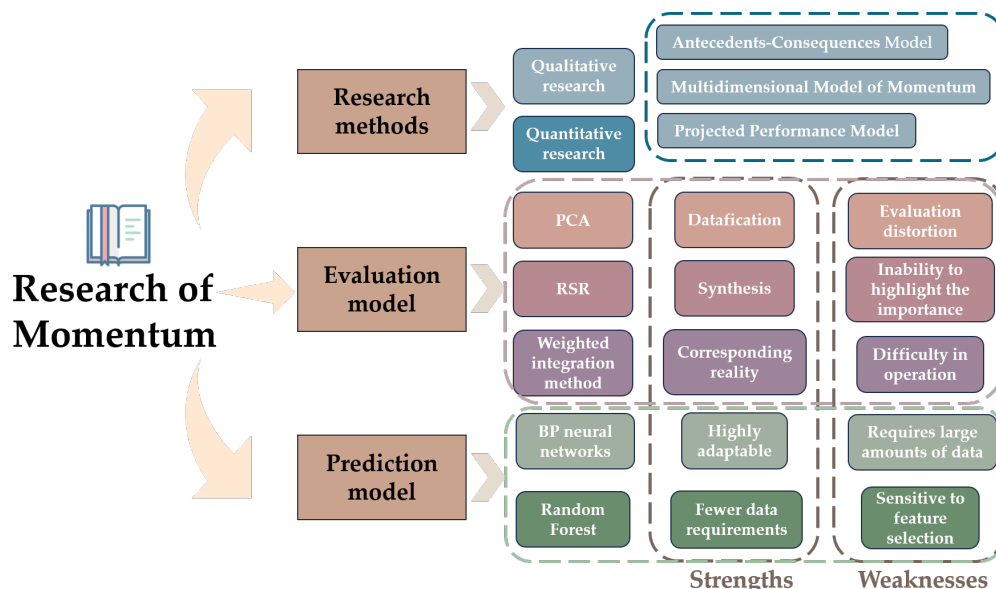


Figure 2: Literature Review Framework

1.4 Our Work

The problem requires us to Our work consists mainly of the following:

- 1) After data pre-processing, we build a player performance evaluation model to quantify the performance of the players and describe the match flow;
- 2) We develop a player momentum evaluation model to quantify a player's momentum and analyze the correlation between performance and momentum, thus reflecting the relationship between match swings and momentum;
- 3) We build a match swing prediction model to predict match swings and get the most relevant factors;
- 4) We test and generalise the match swings prediction model;
- 5) We summarise the results in a memo to provide coaches with information relevant to momentum.

To avoid cumbersome descriptions and to visualize our workflow, the flowchart is shown in figure below:

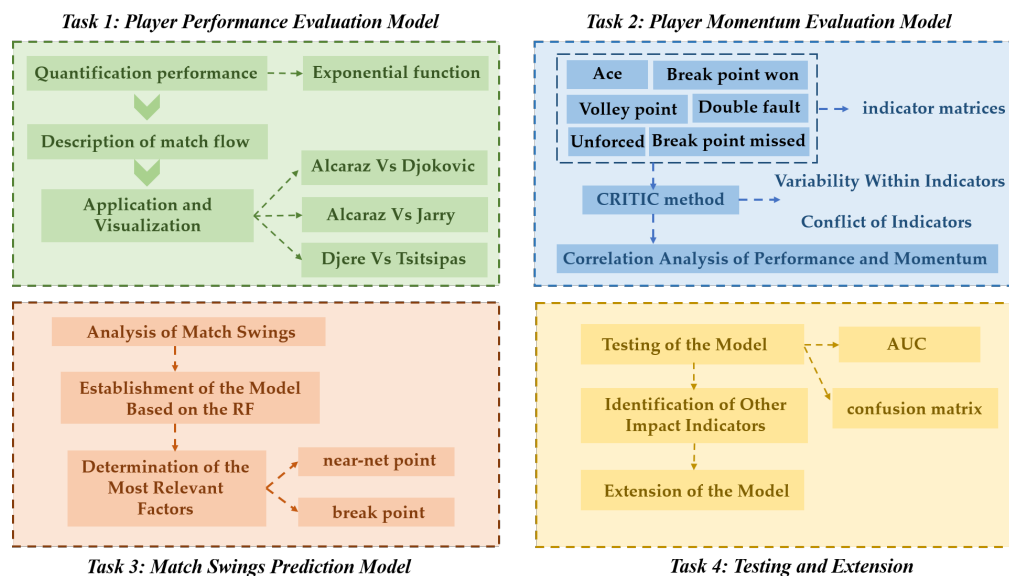


Figure 3: Flow Chart of Our Work

2 Assumptions and Explanations

Given that real problems always include many complicated factors, we first need to make reasonable assumptions to simplify the model, and each assumption is followed closely by the corresponding explanation:

Assumption 1: Players' performance and momentum are analyzed and quantified on a set basis.

Explanation: Since each set of matches is relatively independent and of moderate length, using the set as the unit of measurement to simplify and rationalize the problems.

Assumption 2: The momentum of both players does not interfere with each other.

Explanation: A player's momentum is determined only by the player's own behavior. A swing in one player's momentum does not directly lead to a swing in another player's momentum.

Assumption 3: The data obtained are accurate and the pre-processing of the data is reasonable.

Explanation: The data used to calculate and test the models are obtained from official websites with a high level of credibility. The pre-processing methods of the data are scientific and rigorous.

Assumption 4: Players compete fairly during the competition.

Explanation: The assumption that there are no speculative events or race flukes allows for a rationalized and realistic framing of the problem.

In addition, other assumptions have been made to simplify the analysis of individual parts. These assumptions will be addressed later in the text at the appropriate places.

3 Notations

Some of the key mathematical notations used in this paper are listed in the Table 1

Table 1: Notations

Symbol	Description
ξ	Probability of the server scoring
X	The number of times a player serves points
Y	The number of times a player returns points
W_1	The weight of the serve point
W_2	The weight of the return point
x_{ij}	Value of the j th point of the i th evaluation indicator for player 1
y_{ij}	Value of the j th point of the i th evaluation indicator for player 2
x_i	Sum of the i th indicators of player 1
y_i	Sum of the i th indicators of player 2
S	Variability of momentum evaluation indicators
R	Conflict of momentum evaluation indicators
W	Weight of momentum evaluation indicators

*There are some variables that are not listed here and will be discussed in detail in each section.

4 Data Pre-processing

4.1 Missing Values Processing

By observing the dataset, there are missing values in the speed_mph, serve_width, serve_depth, and return_depth columns, which are denoted by NA in the table. For this, the missing values are processed as follows:

1. **For the missing data in speed_mph:** Since these are numerical missing values, they can be approximated by the mean value method, which uses the average of the two adjacent speed values as an approximation of the missing value at that place.
2. **For the missing data in serve_width, serve_depth, and return_depth:** Since they are non-numeric missing values, it is not possible to obtain their original values by approximation. However, these missing values have no effect on the research.

4.2 Abnormal Values Processing

Using Python to analyze the dataset, three types of abnormal values are identified. The processing performed is shown below:

1. Elapsed time anomalies

In the elapsed_time column, there are data anomalies in rows 586 to 636. Therefore, 24 hours are subtracted from each elapsed time as processed data.

2. Average distance run anomalies

It is found that some of the data in columns p1_distance_run, p2_distance_run, and rally_count do not match the actual situation and are considered abnormal values. Calculating the associated mean and variance, and setting the sum of the mean and triple standard deviation as the threshold, the following figure was plotted:

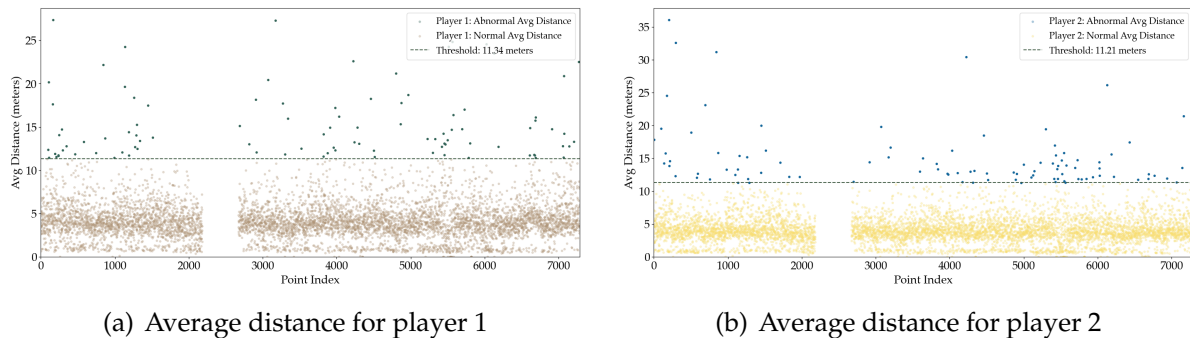


Figure 4: Results for average distance run

where the horizontal dotted line represents the threshold of the distance run. Scattered points below the dotted line represent normal values, while Scattered points above the dotted line represent abnormal values. In addition, the missing parts in both figures are caused by missing data in the dataset.

For the abnormal values, the average distance run when the number of shots is the same is used instead.

3. Distance run difference anomalies

When analyzing the data in the columns p1_distance_run and p2_distance_run, it is observed that the difference between the distance run by the two players in some cases is too large and impractical, therefore they are considered to be anomalies.

Accordingly, the results of the processing of the two columns of data are shown below:

Therefore, it is considered that when the distance run difference between two players exceeds 17.03 meters, the corresponding data is considered an anomaly.

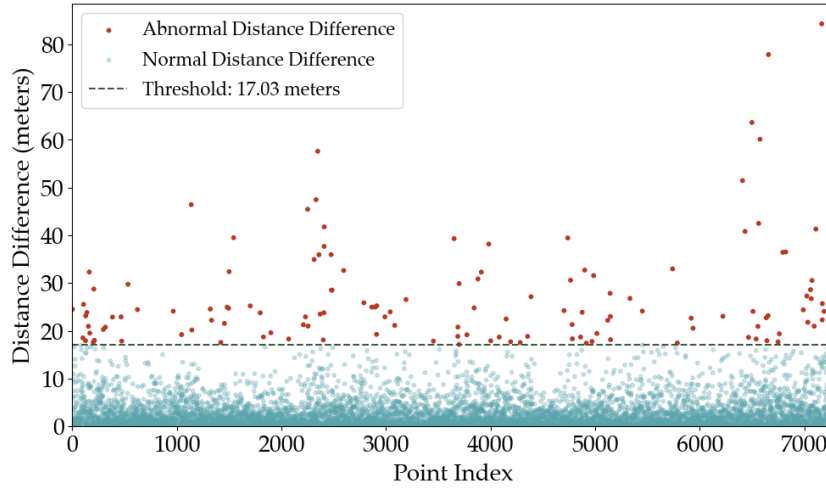


Figure 5: Results for distance run difference

However, since the actual situation at the time of the match is not known, these seemingly abnormal data may fall within the correct range. Therefore, these data are retained to avoid the loss of useful information due to the over-processing of data.

5 Player Performance Evaluation Model Based on Exponential Function

In this section, a player performance evaluation model is built to identify a player's performance at a given time of the match. Based on this model, the match flow can be described and visualized.

5.1 Quantification of Player Performance

The points in a match are the most direct reflection of a player's performance. Therefore, by obtaining the points of a match, the performance of a player can be quantified. Based on this, the player performance evaluation model can be established.

The quantified player performance is defined as the performance value. By Assumption, the points of the current set are used to reflect the performance of the player at a given time. Therefore, the performance value is determined by the player's points from the start of the current set until the given time. When a set ends, the player's performance value is reset to 0. The definition of a given time is shown below:

In tennis, the server has a greater likelihood of winning the point or the match. Analyzing the dataset, it can be seen that it is since the server can determine the speed, angle, and landing position of the serve. Therefore, the likelihood of serve points is higher and therefore the difficulty is lower. Considering that different point scoring difficulties have different responsiveness to player performance, the server points should be given a lower weighting compared to the return points.

To reflect the negative correlation between the probability of scoring and the weights, an exponential function is constructed to obtain the weights of the serve points. The

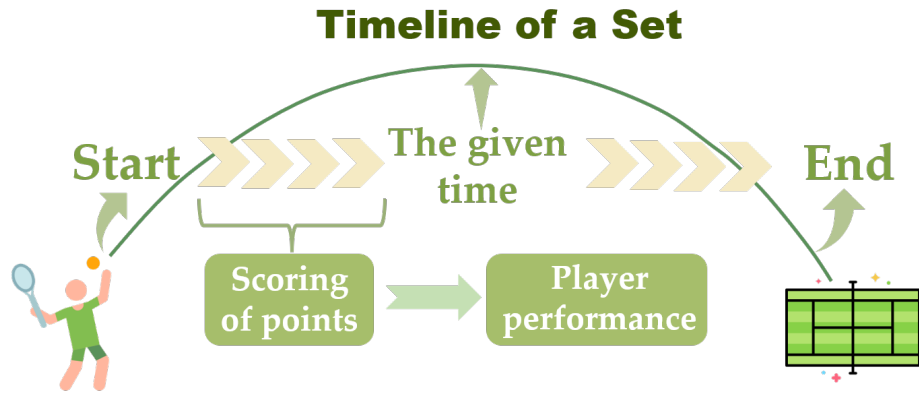


Figure 6: Definition of a given time

independent variable p represents the probability of the server winning the point, and the dependent variable $f(\xi)$ represents the weights corresponding to each probability. In addition, since serve points and return points should be weighted equally when the probability of serve points is 50 percent, the constant term $\frac{1}{2} - e^{-1/2}$ is added. The constructed function is shown below:

$$f(\xi) = e^{-\xi} + \frac{1}{2} - e^{-1/2} \quad (1)$$

Analyzing the entire dataset, the probability of serve points is obtained as 67.31%. By calculation, it is determined that the weight of the serve points is 0.40 and the weight of return points is 0.60. Thus, the performance value of a player at a given time can be quantitatively expressed as below:

$$G(X, Y) = W_1 \cdot X + W_2 \cdot Y \quad (2)$$

Where X represents the number of serve points, Y represents the number of return points, and $G(X, Y)$ represents the total performance value of the player from the beginning of the set to the given time. $W_1 = 0.4$, $W_2 = 0.6$.

5.2 Description of the Match Flow

The match flow can be described by the difference between the performance values of the two players (player1 - player2). To be specific, the visualization of the match flow is realized by using the points as the independent variable, the difference of the performance values as the dependent variable, and the sets as the separator to create a coordinate system.

By analyzing the results of the visualization, the following information can be obtained:

1. Which player scores points can be determined by the trend of the graph line. When the graph line goes up, player 1 scores points. When the graph line goes down, player 2 scores points.
2. The winners and losers of the game can be determined by the trend of the end of the graph line. When the end of the graph line is trending upwards, player 1 wins the game. When the end of the line is trending down, player 2 wins the game.

3. Which player is currently performing better can be determined by which side of $y = 0$ the point lies. When the point lies above $y = 0$, player 1 performs better. When the point lies below $y = 0$, player 2 performs better. The distance of the point from $y = 0$ reflects how much better the corresponding player performed.

5.3 Application and Visualization of the Model

To test the effectiveness of the model, three matches with the different number of sets are selected for the analysis and visualization of the match flow. Visualization results and analysis of the match flow are shown below.

- **Alcaraz Vs Djokovic (5 sets)**

According to Figure 7, it can be seen that in set 1, player 2 performs better and his performance is overwhelming in the late stages of the match. In sets 2 and 3, player 1 performs better and his performance is overwhelming in the middle and late stages of set 3. In sets 4 and 5, the points on Figure 7 fluctuate around $y = 0$ in the early stages, and it is difficult to tell which of the two performs better. But in the late stages of set 4, player 2 performs increasingly better, and his performance is overwhelming. In the late stages of the set 5, player 1 performed better.

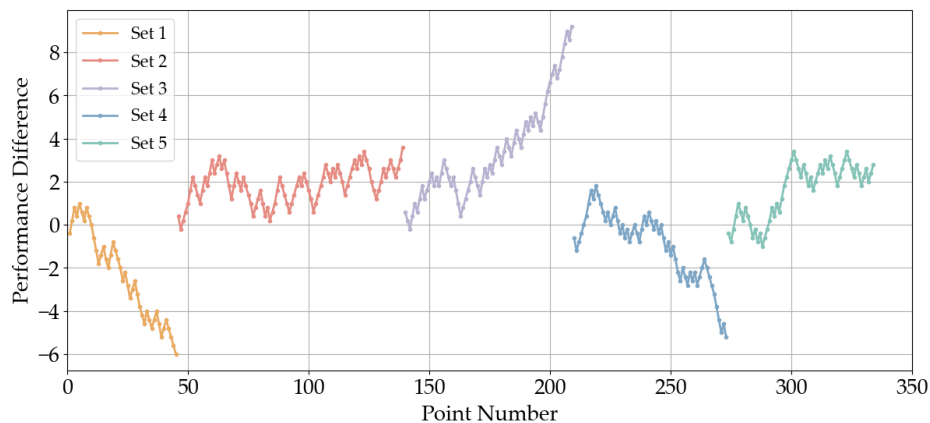


Figure 7: Alcaraz Vs Djokovic

- **Alcaraz Vs Jarry (4 sets)**

According to Figure 8, in set 1, the points on the figure all fluctuate around $y = 0$ and it is not easy to tell which of the two performs better. In set 2, player 2's performance is overwhelming at the beginning but decreases later in the match. In set 3, player 1 is better, and his performance is overwhelming later in the match. In set 4, the match was swung in the opposite direction, with player 2 playing better at the start, but player 1 overtaking in the middle and late stages of the match.

- **Djere Vs Tsitsipas (3 sets)**

According to Figure 9, the match flow is similar in sets 1 and 3, with about the same degree of performance. However, it is worth noting that in set 2, according to our quantitative results, player 1 performs better at the end of the set 2, but player 2 wins this set.

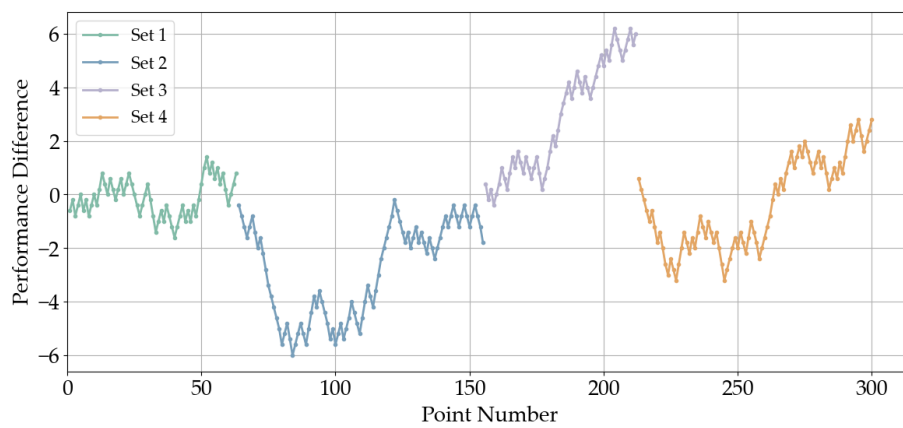


Figure 8: Alcaraz Vs Jarry

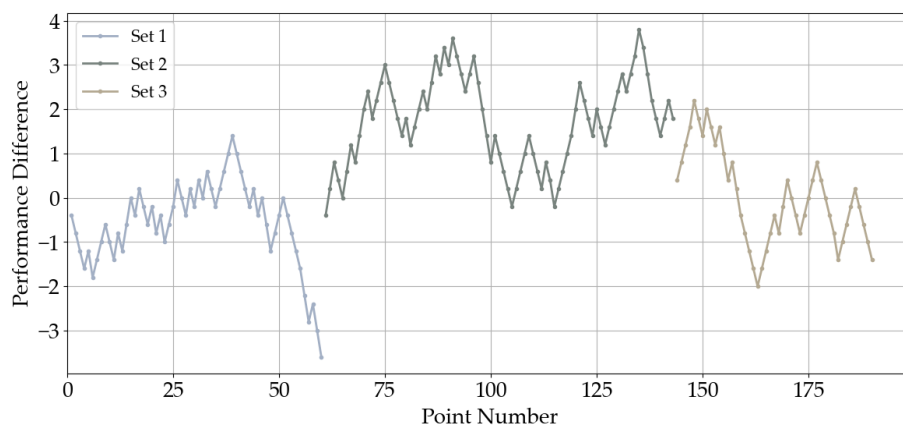
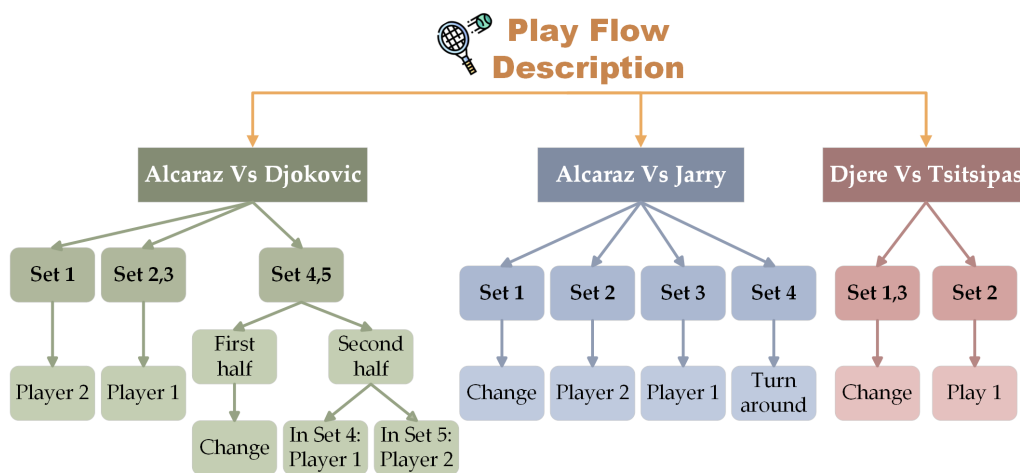


Figure 9: Djere Vs Tsitsipas

By analyzing and visualizing the three match flows as described above, who performed better overall in each set by each group of players is summarized in the figure below:



6 Player Momentum Evaluation Model Based on CRITIC Method

Momentum is defined as the strength or force gained by motion or by a series of events, which depends on a variety of factors. In this section, through quantifying the various influencing factors, corresponding weights are assigned to the different factors to quantify momentum so that the player momentum evaluation model can be established.

6.1 Identification and Pre-processing of Evaluation Indicators

First, through the review of literature and information, combined with the analysis of data, six momentum evaluation indicators are identified, which are divided into two main categories: positive and negative indicators. The specific meanings of the indicators are as follows:

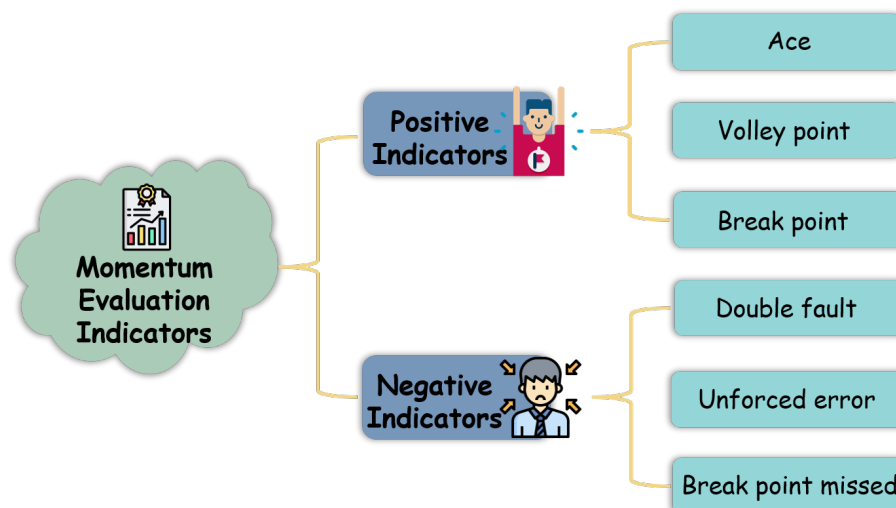


Figure 11: Momentum evaluation indicators

1. Ace

In tennis, a player who serves is far more likely to win points. When a player hits an untouchable winning serve, not only does it increase self-confidence, but it may also reduce the number of errors on the next serve to some extent.

2. Volley point

Volley point means hitting the ball before it has hit the ground, which is an aggressive tactic. Points scored near the net can increase a player's confidence.

3. Break point won

Break points are a win for the returner in unfavorable conditions and will undoubtedly increase self-confidence tremendously. Moreover, the next set is the returner's serve, which is favorable to the returner. Therefore, the momentum of the returner will increase.

4. Double fault

Double fault can be seen as a low-level error that can affect a player's mental stability to a certain extent.

5. Unforced error

Unforced errors are those that occur under normal circumstances when a player is faced with a ball that is within the player's ability to return. It can frustrate a player's self-confidence and bring about frustration.

6. Break point missed

Break point missed means that the returner loses the opportunity to win the match in one shot, and the next round is still unfavorable for the returner. As a result, it may cause the returner to feel a sense of regret and frustration.

Then, the indicator matrices corresponding to player 1 momentum and player 2 momentum are calculated respectively. For player 1, the above six indicators are defined sequentially as $x_{1j}, x_{2j}, x_{3j}, x_{4j}, x_{5j}, x_{6j}$, $j \in (1, 7284)$; for player 2, the above six indicators are defined sequentially as $y_{1j}, y_{2j}, y_{3j}, y_{4j}, y_{5j}, y_{6j}$, $j \in (1, 7284)$. With $n=7284$, the indicator matrix is generated as follows:

$$\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{16} \\ x_{21} & x_{22} & \cdots & x_{26} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{n6} \end{bmatrix}_{n \times 6} \quad (3)$$

$$\begin{bmatrix} y_{11} & y_{12} & \cdots & y_{16} \\ y_{21} & y_{22} & \cdots & y_{26} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \cdots & y_{n6} \end{bmatrix}_{n \times 6} \quad (4)$$

Next, the data for each indicator are positively normalized to eliminate differences in meaning, measurement, and magnitude between the indicators. The formulas for positive and negative indicator normalization are shown below, respectively:

$$x'_{ij} = \frac{x_{ij} - \min_{j=1, \dots, n} x_{ij}}{\max_{j=1, \dots, n} x_{ij} - \min_{j=1, \dots, n} x_{ij}} \quad (5)$$

$$x'_{ij} = \frac{\max_{j=1, \dots, n} x_{ij} - x_{ij}}{\max_{j=1, \dots, n} x_{ij} - \min_{j=1, \dots, n} x_{ij}} \quad (6)$$

In addition, since the methodology and steps for calculating the weights are the same for both indicator matrices, only the process of solving the player 1 indicator matrix is shown in the following.

6.2 Determination of Indicator Weights Based on the CRITIC Method

The CRITIC (Criteria Importance Through Intercriteria Correlation) method is a relatively well-established objective weighting method. The weight of an indicator is derived by comprehensively measuring the variability within the indicator and the conflict between indicators. Since it is impossible to subjectively judge the importance of the impact of different factors on momentum, and since there is correlation and variability among indicators, the CRITIC method is chosen to determine the weights of each momentum evaluation indicator.

Variability Within Indicators

The variability within the evaluation indicators is expressed in terms of the standard deviation of S_j , which indicates the size of the difference between the values of the various evaluation objects for the same indicator and is calculated as follows:

$$\begin{aligned}\bar{x}_j &= \frac{1}{n} \sum_{i=1}^n x_{ij}^2 \\ S_j &= \sqrt{\frac{\sum_{j=1}^n (x'_{ij} - \bar{x}_j)^2}{n-1}}\end{aligned}\quad (7)$$

where \bar{x}_j represents the mean value of the j th indicator and S_j represents the standard deviation of the j th indicator. The larger the standard deviation, the larger the difference in values of the indicators, the greater the variability, and the more unique information the indicator contains. Therefore a larger weight should be assigned to the indicator.

Conflict of Indicators

Conflict between evaluation indicators is based on correlation, which is quantified as the correlation coefficient of the different evaluation indicators R_j , which is calculated as follows:

$$R_j = \sum_{i \neq j} (1 - r_{ij}) \quad (8)$$

where r_{ij} represents the correlation coefficient between evaluation indicators i and j . The larger the correlation coefficient, the stronger the correlation between the indicators, the stronger the conflict, and the more duplicated information the indicator contains. Therefore a smaller weight should be assigned to the indicator.

Determination of Objective Weights

From the standard deviation and correlation coefficient of the indicators, the objective weight W_j of the j th indicator is calculated as:

$$W_j = \frac{S_j \cdot R_j}{\sum_{j=1}^6 (S_j \cdot R_j)} \quad (9)$$

The calculated weights of the different evaluation indicators for player 1 and player 2 respectively are shown in the table below:

	Serve point	Near-net point	Break point	Double fault	Unforced error	Break point lost
Player 1	0.168	0.170	0.165	0.163	0.173	0.166
Player 2	0.167	0.171	0.164	0.163	0.174	0.161

Table 2: Weights of different evaluation indicators

6.3 Quantitative Results of Momentum

Combined with the objective weights of the different evaluation indicators, the quantitative expression of player momentum is obtained as follows:

$$f_{p_1}(x_1, x_2, x_3, x_4, x_5, x_6) = 0.168x_1 + 0.170x_2 + 0.165x_3 + 0.163x_4 + 0.173x_5 + 0.166x_6 \quad (10)$$

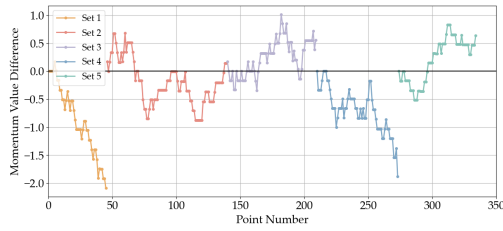
$$f_{p_2}(y_1, y_2, y_3, y_4, y_5, y_6) = 0.167y_1 + 0.171y_2 + 0.164y_3 + 0.163y_4 + 0.174y_5 + 0.161y_6 \quad (11)$$

where $x_k = \sum_{\text{start of a set}}^{\text{a given time}} x_{i_k}$, $y_k = \sum_{\text{start of a set}}^{\text{a given time}} y_{i_k}$.

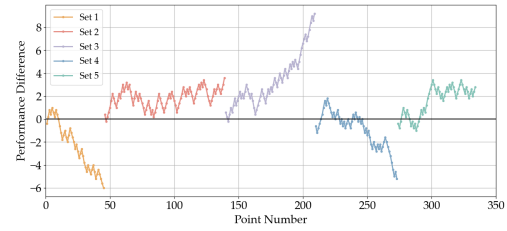
6.4 Correlation Analysis of Performance and Momentum

Performance in play can be reflected by changes in player performance. In the first question, the change in player performance has been derived from the player performance evaluation model. Therefore, the link between swings in play and momentum can be obtained by correlating performance in play and momentum of players.

Qualitative Analysis Based on Image Comparison



(a) Trends in momentum



(b) Trends in performance

Figure 12: Comparison of trends in momentum and performance

The match between Alcaraz and Djokovic is selected for analysis. Figure 12(a) depicts the trend of the difference between the two players' momentum during the match, represented by $f_{p_1} - f_{p_2}$. The figure reflects the relative change in player momentum at any given time throughout the match, and which player has the higher momentum. Figure 12(b) shows the trend of the player performance in the match, which is obtained from the first question of this paper.

By qualitatively analyzing the two figures, it is found that generally, there is a great similarity in the trend of the graph lines in the two figures. Hence, the following conclusions are obtained:

1. When the difference between the two players' momentum is increasing, the difference between the players' performances also increases.
2. When the difference in momentum decreases, the difference in the player's performance also decreases or even reverses.
3. Some of the graph lines in Figure 12(a) are horizontal, while those in Figure 12(b) fluctuate in moments. It indicates that the change in momentum is phasic and abrupt, while the player's performance has been dynamic and gradual.

Quantitative Analysis Based on Linear Regression and Pearson Correlation Coefficient

A linear regression method is developed to regress and analyze the changes in momentum and performance of the players. The results of the linear regression are shown in the figure below:

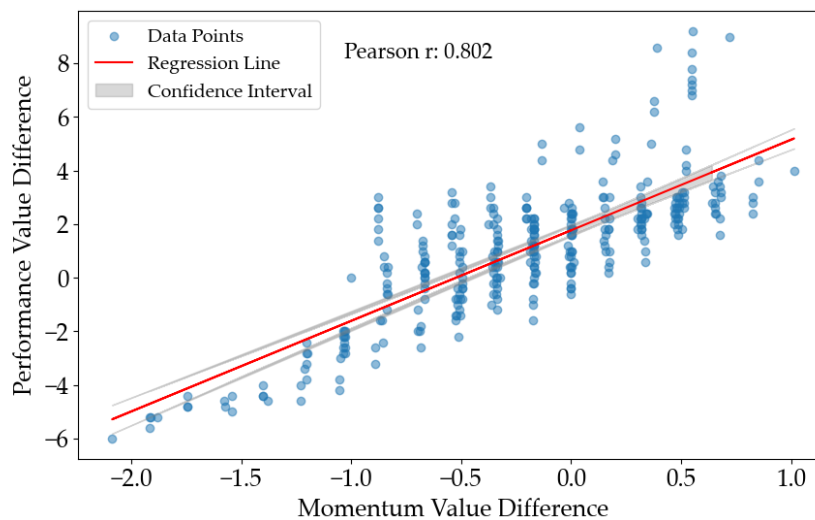


Figure 13: Results of the linear regression

Firstly, by observing this scatter plot, it is observed that most of the data points fall around the straight line of the linear regression results. Next, the Pearson correlation coefficient is calculated to be 0.802, and p-value is much less than 0.05, indicating that there is a strong correlation between the difference in momentum and the difference in performance of the players in this match.

Then, the Pearson correlation coefficients of the momentum and performance of the players in each match are calculated respectively, and the results are shown in the table below:

Match id (2023-wimbledon-)	1303	1305	1311	1406	1502
Pearson correlation coefficient	0.747	0.858	0.876	0.889	0.961

Table 3: Calculation of Pearson correlation coefficient

As can be seen from the data in the table, the Pearson correlation coefficients in each match are close to 1, indicating that there is a correlation between players' momentum and performance. Then, the p-values of each match are calculated respectively, which are much less than 0.05, indicating statistical significance. Therefore, it is reasonable to believe that there is a strong positive correlation between player momentum and performance in general.

Combining qualitative and quantitative analyses, it can be concluded that the coach's view is considered incorrect. It is believed that momentum indeed plays a role in the match and affects the match flow by influencing the players' performance. In addition, swings in play and runs of success by one player are not random but are the result of changes in momentum triggered by a variety of indicators such as Ace and break points, etc.

7 Match Swings Prediction Model Based on Random Forest

In this section, the match swing prediction model is built to predict the match swings. After analyzing the first two questions, it is determined that a player's momentum is related to performance. Therefore, momentum is used to predict the player's win or loss in the following set.

7.1 Analysis of Match Swings

First, match swings are defined. In this model, match swings are examined in terms of sets. For example, if player 1 is victorious in one set and player 2 is victorious in the next set, the match is considered to swing.

Then, factors affecting the swing of the match are identified. From Problem 2, ace, volley points, break points, double fault, unforced errors, and break points missed are related to a player's momentum. In addition, momentum has a strong positive correlation with performance. Therefore, the player's performance is measured by the indicators affecting the momentum, and thus the match swings is predicted. Calculate the difference between player 1 and player 2 corresponding to these indicators (player1-player2), and add up the differences of these indicators from the first set of a given match to the end of a set to get the total difference, which can indicate which player has performed better on that indicator up to that set.

Thus, the identified factors affecting match swings are: Ace difference, Volley point difference, Break point difference, Double fault difference, Unforced error difference, Break point missed difference.

7.2 Establishment of the Model Based on the RF

Random Forest (RF) is a statistical learning theory method proposed by Breimad in 2001[11], which is a Bagging algorithm with decision trees as estimators.

Firstly, multiple data subsets are randomly extracted from the original dataset, and then the data subsets are used to build the corresponding decision tree to form a random forest. Some features are randomly selected as inputs each time, the number of votes for each category is obtained as output, with the category that gets the most votes as the predicted result.

In this model, except for the three matches that will be predicted later, nine matches are randomly selected as the training set, and the six factors affecting the match swings in this training set are inputted into the random forest as features while the match swings are inputted as targets.

Since the size of the random forest is 100, the above process is repeated 100 times to generate 100 decision trees. The majority classification results of these 100 decision trees are used as the final prediction.

7.3 Determination of the Most Relevant Factors

The importance of each feature is viewed using this random forest model to determine which factors seem most related. By ranking the importance of the features, a bar ranking figure is generated as shown below:

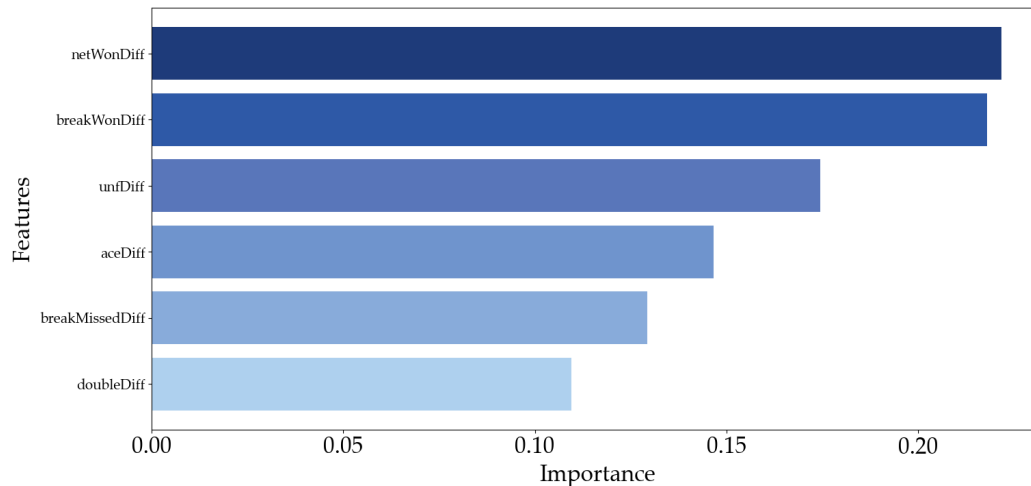


Figure 14: Ranking of the importance of indicators

From the images, it is clear that the near-net point difference and the break point difference are the most relevant factors, both of which have an importance of more than 20% for the match swings prediction model.

7.4 Match Strategies Based on Momentum Swings

Through the previous analyses, it is clear that the player's performance in a match is related to momentum swings. To be more specific, when the momentum increases, the performance as an advantage player is more overwhelming, and the performance as a disadvantage player gradually becomes better or even reverses. By reviewing the relevant literature, based on the differential in past match "momentum" swings, match situations can be classified into the following five categories, with countermeasures and recommendations provided for each situation.

1. Complete disadvantage

When the opponent's momentum is overwhelming, the disadvantaged player is advised to adopt a prudent strategy. Instead of rushing to catch up with the points, the disadvantaged player should calm down, think deeply, recall plans prepared in advance, and draw on experience in turning around unfavorable situations.

2. Turn disadvantage

When the momentum of the opponent is gradually moving in the direction of being stronger than oneself, the disadvantaged player is advised to adopt a competitive strategy. The disadvantaged player should concentrate on increasing energy and competing aggressively to reverse the disadvantage in a short time.

3. Considerable momentum

When both players have roughly the same momentum, it is recommended that players use a winning with stability strategy. Players should balance error reduction

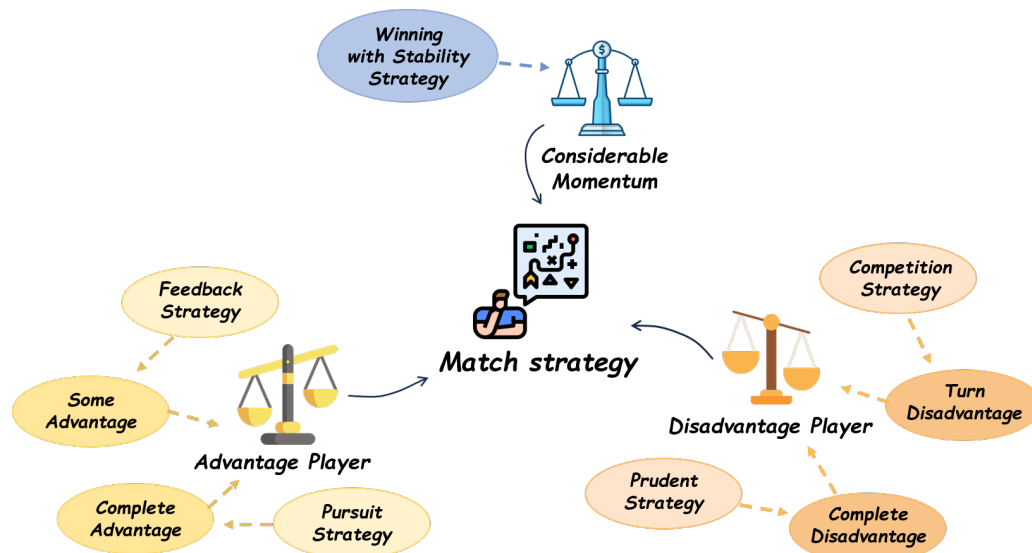


Figure 15: Ranking of the importance of indicators

with aggressive competition, and use their preparation patterns on serve and return to compete for control of the match.

4. Some advantage

When a player has better momentum in comparison to the opponent, the player should use a feedback strategy. The advantage player should think about the reasons why the advantage is established and apply the experience to the rest of the match. At the same time, it is important not to be overconfident in the current situation to prevent making unforced errors.

5. Complete advantage

When a player is in complete control of the match and the momentum is overwhelming, the advantage player is advised to adopt a pursuit strategy. The advantage player should aggressively adjust his tactics and keep compressing the opponent's room until winning the match.

8 Testing and Extension of Match Swing Prediction Model

8.1 Testing of the Model

Three matches visualizing the match flow in Problem 1 are selected as the test set with match_id are '2023-wimbledon-1301', '2023-wimbledon-1701', '2023-wimbledon-1308'. The test is performed using the Random Forest based match swing prediction model. The preliminary prediction accuracy is 71.3%, which indicates that the prediction model we developed has a good prediction ability on the given test set.

In order to analyse the performance of the model more comprehensively, ROC curves are plotted and further evaluated using AUC figures combined with confusion matrices.

Using a combination of the two methods provides insight into the model's overall classification ability (via AUC) and the model's predictive performance on different categories (via the confusion matrix).

Taking the above analyses into account, the results for the test set are as follows:

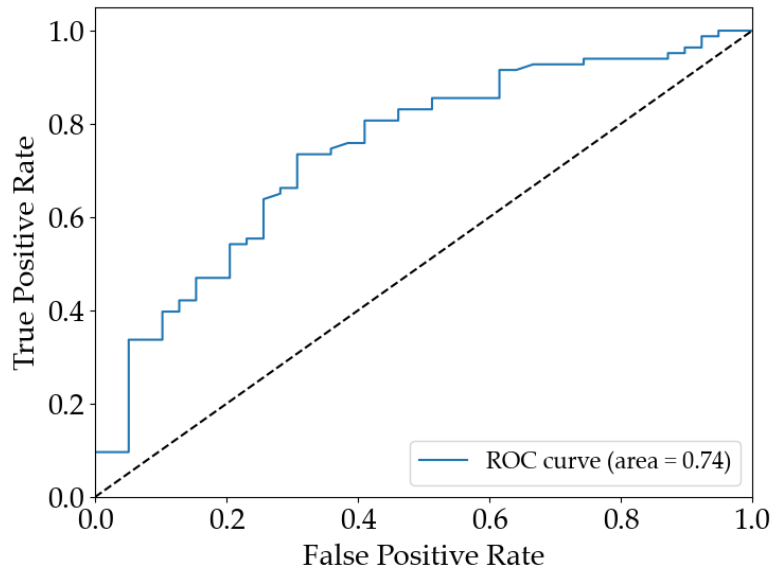


Figure 16: AUC figure for the test set

Firstly, the AUC method is used: in combination with Figure 16, the calculation of $AUC = 0.76$ indicates that the model has a good classification ability and is able to differentiate between positive and negative class samples to a large extent.

Then, the confusion matrix method is used: the Precision and Recall are calculated as:

$$P = \frac{TP}{TP + FP} \quad (12)$$

$$R = \frac{TP}{TP + FN} \quad (13)$$

By calculation, $[[TP, FN], [FP, TN]] = [[25, 14], [10, 73]]$.

Accordingly, the F1 Score is calculated with the following formula:

$$F1 = \frac{2(P \cdot R)}{P + R} \quad (14)$$

Calculating the $F1 = 0.675$, it shows that our model performs relatively well in terms of the balance between precision and recall.

Combining the results of AUC and $F1$ calculation, it reflects the good performance of our model.

8.2 Identification of Other Impact Indicators

There are other indicators that affect a player's momentum, such as consecutive points. When a player scores consecutive points, he is stimulated and inspired, and

his momentum increases. The opponent will feel pressure and frustration, and his momentum will decrease accordingly.

In addition, the hitting style, serving speed, and other indicators can reflect the player's competitive state and also can reflect the player's momentum. At the same time, some off-court factors such as the weather and the reaction of the spectators will also have an impact on a player's momentum.

8.3 Extension of the Model

To test the generalization ability of the model, data from the 2023 Wimbledon women's singles tournament and the 2020 Tokyo Olympics men's table tennis singles tournament are obtained and also evaluated using an AUC figure combined with a confusion matrix.

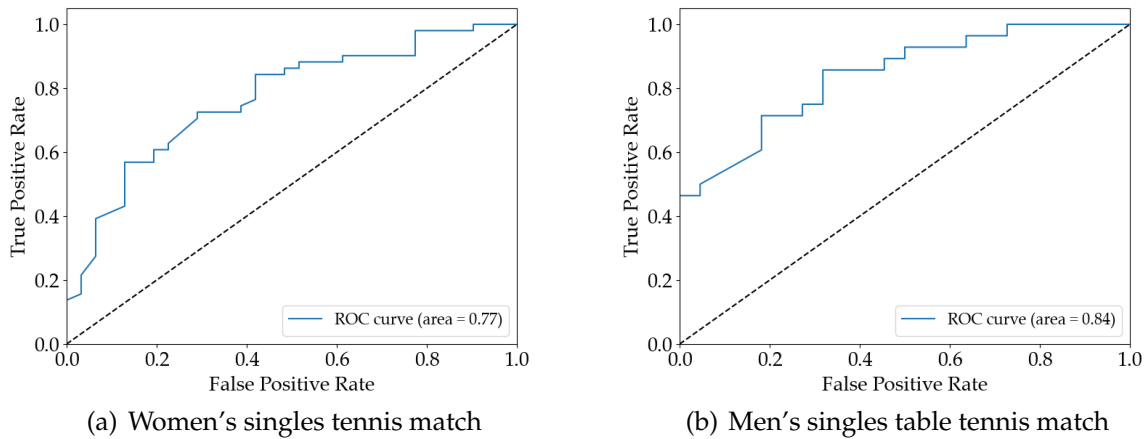


Figure 17: AUC figures

Figure 17(a) on the left shows the AUC curve for the women's singles tennis match, which is calculated as $AUC = 0.77$. The confusion matrix is calculated as $[[24, 7], [12, 49]]$, and $F1 = 0.716$. This shows that the model still has high generalisability when the match categories are the same but the genders are different.

Figure 17(b) on the right shows the AUC curve for the table tennis men's singles tournament, which is calculated with $AUC = 0.84$. The confusion matrix is calculated as $[[15, 7], [5, 23]]$, and $F1 = 0.714$. This shows that when the categories of the tournaments are different but similar, the model still has a high degree of generalisability.

9 Sensitivity Analysis

In tennis, the player who serves is much more likely to win the points. In Problem 1, the winning percentage of the server is calculated to be 67.31%, which is a constant. Therefore, the winning percentage of the server is chosen to perform a sensitivity analysis of the player performance evaluation model.

By constantly changing the winning percentage of the server, the performance of the player is drawn as shown below:

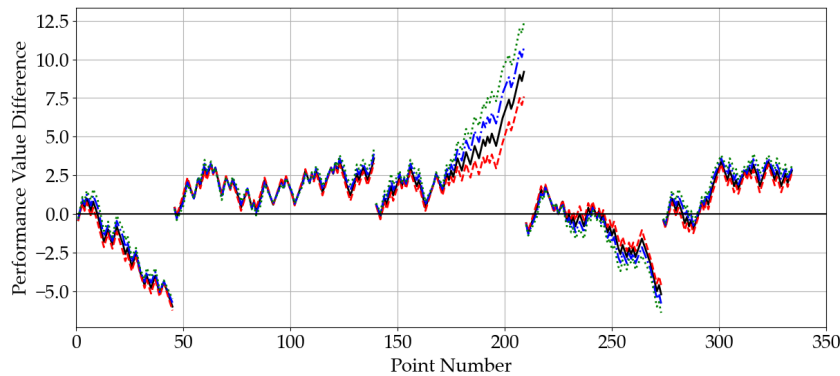


Figure 18: Sensitivity analysis of the winning percentage of the server

Where the four lines represent 64%, 65%, 67%, and 69% of the win rate of the server, which shows that when the win rate of the server is changed, the results of the model do not fluctuate massively. This indicates that our model has good stability.

10 Evaluation of Strengths and Weaknesses

10.1 Strengths

- The establishment of each model is rooted in a deep understanding of the problem, combining objectivity and rationality.
- The selection of indicators affecting momentum is scientifically rigorous and supported by the literature. This allows for a more accurate analysis of momentum swings and match flows, and makes the paper's findings more convincing.
- The use of the CRITIC method to objectively assign weights to the indicators affecting momentum is accurate and robust.
- Random Forest is capable of operating on large data sets. Using it to make predictions about momentum swings is extremely accurate,

10.2 Weaknesses

- Simply using points to define a player's performance may be unreasonable.
- Although six indicators are selected as impact factors for momentum, there may still be important impact indicators that have not been considered.

MEMORANDUM

To: tennis coach

From: Team #2420494 of 2024 MCM

Subject: Findings on Player Momentum and its Effects

Date: February 6, 2024

In tennis, it often happens that players who seem to have an advantage tend to suddenly lose points in quick succession, leading to losses in games and sets. This incredible swing is usually attributed to momentum. Here are the results of our team's research on momentum.

We first quantify the performance of the player by constructing a performance evaluation model based on an exponential function, and visualize the performance of the player in terms of sets and the match flow. Next, we define momentum and explore the influencing factors of it. By reviewing the literature, we find that ace, volley points, break points, double faults, unforced errors, and break points missed affect the player's momentum, and we determine the weights of each influencing index through the CRITIC method, completing the quantitative treatment of the player's momentum. Then, we conduct a correlation analysis between the player's momentum and performance. Finally, we use the random forest algorithm based on grid search to construct the match swings prediction model to predict the match swings.

Based on our research, in response to the role of momentum in the match, we offer you the following advice:

- By analyzing the correlation between a player's momentum and performance, we determined that a strong positive correlation exists between a player's momentum and performance. The higher the player's momentum, the better the performance.
- Since momentum has a big impact on a player's performance, it is helpful to hold steady and try to get ahead of the momentum during the course of a match, no matter what the situation is.
- Since improving momentum is beneficial to the game, the six factors affecting momentum need to be taken into account. In daily training, the practice of serving should be strengthened while unforced errors should be avoided.

In addition, based on the Match Swings Prediction Model we built, we categorize match situations into five categories based on differences in the momentum swings of previous matches, and propose strategies for each situation:

- When the opponent's momentum is overwhelming, the disadvantaged player is advised to adopt a prudent strategy, which requires them to calm down and draw on experience in turning around unfavorable situations.
- When the momentum of the opponent is gradually moving in the direction of being stronger than oneself, the disadvantaged player is advised to adopt a competitive strategy, concentrate on increasing energy, and compete aggressively to reverse the disadvantage in a short time.

- When both players have roughly the same momentum, it is recommended that players use a winning with stability strategy, using their preparation patterns on serve and competing for control of the match.
- When a player has better momentum in comparison to the opponent, the player should use a feedback strategy. Think about the reasons why the advantage is established and apply the experience to the rest of the match.
- When a player is in complete control of the match and the momentum is overwhelming, the advantage player is advised to adopt a pursuit strategy. They should aggressively adjust his tactics and keep compressing the opponent's room to win the match.

The above are the results of our team's research on momentum in the match. If you still need more detailed information about momentum, please feel free to contact us.

References

- [1] https://www.espn.com/tennis/story/_/id/38020644/carlos-alcaraz-stops-novak-djokovic-win-1st-wimbledon-title
- [2] <https://weibo.com/1663312090/NacnvhkFb>
- [3] Vallerand, R. J., Colavecchio, P. G., and Pelletier, L. G. (1988). Psychological momentum and performance inferences: a preliminary test of the antecedent consequences psychological momentum model. *Journal of Sport and Exercise Psychology*, 10, 92-108.
- [4] Taylor, J. and Demick, A. (1994). A multidimensional model of momentum in sports. *Journal of Applied Sport Psychology*, 6, 51-70.
- [5] Cornelius, A. E., Silva, J. M., Conroy, D. E. and Petersen, G. (1997). The projected performance model: relating cognitive and performance antecedents of psychological momentum. *Perceptual and Motor Skills*, 84, 475-485.
- [6] Cooper, M. (2003). *Existential therapies*. London: Sage.
- [7] Jin, Jiasheng. "A Comparative Study of Three Common Methods for Comprehensive Sports Evaluation—Ranking of the Athletic Performance of Outstanding Centers in the 2003-2004 NBA Season." *Journal of Beijing Sport University*, 28(2), 213-216.
- [8] Zhu, Wenfu. "An Exploratory Study on Predictive Models for Sports Competition Outcomes." *Journal of Chongqing Technology and Business University (Natural Science Edition)*, 28(3), 318-321.
- [9] Huang, Yi. "Predicting Soccer Match Results Using Neural Networks." *Journal of Microcomputer Applications*, 37(11), 137-140.
- [10] Ban, Yang, et al. "Predicting Soccer Match Results Based on Machine Learning Algorithms." *12th National Sports Science Conference*, Pages: 104-105.
- [11] BREIMAN L. Random forests—random features. *Machine Learning*, 45(1) 5-32.