# Analyzing Momentum and Predicting Match Swings in Tennis: A Comprehensive Study

## Summary

In the realm of tennis, **the term 'momentum' frequently serves as a pivotal gauge of the match flow**, intimately tied to its ultimate outcome. This paper introduces a viable **Dynamic Match Flow Model** aimed at quantifying and forecasting player momentum. It integrates both **static factors**, such as age and the Elo Rating System, and **dynamic factors** like in-match performance, in predicting the probability of winning each point. Moreover, it provides insights into constructing an additional model for predicting key moments of momentum reversal, offering strategic guidance for players facing new opponents.

To quantify momentum, we introduce the concept of 'leverage' to **assess the significance of each scoring point in a match**. Building upon this notion, we factor in a player's in-match performance to gauge momentum. This involves the development of a model for dynamically predicting a player's match-winning rate. Our methodology entails training an **LSTM model** complemented by **counter-factual analysis** on datasets to forecast a player's match performance. Additionally, we devise **an innovative function** to assess the player's efficiency over time and its influence on win rates. Finally, we linearly combine these factors with current momentum, scoring rate, and bias terms to derive updated win probabilities.

We perform **Autocorrelation tests** on momentum across 31 matches in the dataset, revealing a non-random distribution. Through computed **Pearson correlation coefficients** among various variables, we identify pivotal indicators strongly associated with momentum. Moreover, **Students' t-test** is employed to examine player performance and match outcomes. The results highlight significant divergences in indicators such as point_victor, player_points, and player_games in nearly 95% of matches. Additionally, approximately 40% of matches exhibit notable differences between these metrics and in-game player performances, including metrics like ace and unf_err. These findings underscore the robust relationship between momentum and players' performance across diverse matches, underscoring its substantial influence on game results.

By analyzing momentum trends in each match, pivotal moments of momentum reversal can be identified. At these critical junctures, we employ **Lasso regression** to determine highly correlated features, which are displayed in Figure 11. Then, we extract these features and build a linear regression model, allowing us to utilize these key indicators to predict crucial moments of momentum reversal. Results from $R^2$ **tests and mean square error** calculations indicate accurate prediction outcomes.

To assess the model's generalization and accuracy, we conduct cross-validation, calculating True Positive Rate (TPR) and False Positive Rate (FPR) of predicted win rates against actual outcomes. This analysis yields a **Receiver Operating Characteristic** (ROC) curve and corresponding **Area Under the Curve** (AUC) values, demonstrating the model's excellent predictive performance.

**Keywords**: Tennis momentum; Counter-factual analysis; Sports analysis; Assessment system

# Contents

# 1 Introduction

## 1.1 Background

In recent years, there has been a notable surge in academic attention toward sports forecasting models, evidenced by a growing number of studies dedicated to predictive analytics in this domain. These models primarily target exploiting information inefficiencies within betting markets, offering savvy bettors a competitive edge. Moreover, predictive assessment models serve to benefit athletes and coaches directly by refining rating systems and monitoring athletes' performance capabilities.[2]

However, existing academic research primarily concentrates on forecasting football match results, with limited attention given to predicting tennis momentum despite the sport's popularity among bettors. Quantitative definitions and scholarly inquiries into the dynamics of momentum or strength in tennis contests are notably lacking.

Therefore, it is significant to develop an evaluation and forecasting model for momentum. This model will aid players in utilizing their understanding of momentum during competitions to enhance their performance and achieve victories.



Figure 1: In the 2023 Wimbledon Gentlemen's final, 20-year-old Spanish rising star Carlos Alcaraz defeated 36-year-old Novak Djokovic.

## 1.2 Restatement of the Problem

Given the background information and identified constraints outlined in the problem statement, we must address the following issues:

1. Develop a model to assess the performance, or momentum, of each player during a match at any given moment. Apply this model to one or more matches and generate a visualization illustrating the flow of play based on the model's output.

2. Utilize the previously developed model to evaluate whether momentum exhibits randomness during play, thus establishing its significance in influencing match outcomes.

3. Develop a predictive model to anticipate shifts in match dynamics and identify key factors influencing these changes. Provide strategic recommendations to players on how to adapt their tactics when facing opponents experiencing momentum swings.

4. Evaluate the model's performance on additional matches to assess its accuracy and identify areas for improvement. Test the model's generalizability to ensure its effectiveness across various match scenarios.

5. Write a memo to summarize the conclusions of the paper and give advice about momentum.

## 1.3 Literature Review

In the realm of tennis, numerous predictive models are employed to anticipate match outcomes. Regression models stand out as a prevalent approach, encompassing a range of predictive factors, including player rankings (Klaassen and Magnus, 2003), player seedings (Boulier and Stekler, 1999), player demographics (Del Corral and Prieto-Rodriguez, 2010), and prize earnings (Gilsdorf and Sukhatme, 2008). Another popular approach is the Point-based model, distinguished by its focus on predicting the points won in serving and subsequently deriving the conditional probabilities of other events based on this estimated value. In addition, when considered over a representative set of professional matches, Elo scoring predictions come closest to the performance of betting markets. A notable drawback common to most tennis prediction models is their lack of real-time updates as the match progresses. Only point-level models possess the capability to adjust win-loss predictions during play, yet their limitation lies in disregarding in-match information when updating expected scores, and relying on a pre-defined notion of a key-moment only being a break point or game, set, or match point –it could include other points, and these points depend on the relative strengths of the players as well as what has occurred already during the match[1].



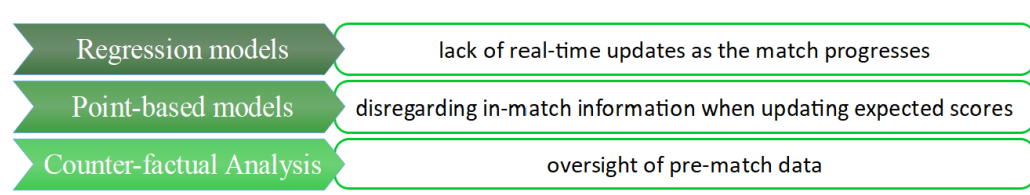| Regression models | lack of real-time updates as the match progresses |
| Point-based models | disregarding in-match information when updating expected scores |
| Counter-factual Analysis | oversight of pre-match data |

Figure 2: Literature review

The recent trend in machine learning emphasizes the adoption of interpretable techniques, such as counter-factual analysis, to explicate predictions for individual events. In the context of sports, this approach can be applied to clarify the impact of a particular match or event on the overall competition outcome. Within this framework, several metrics, including 'leverage,' 'clutch,' and 'momentum,' have been quantified. A chain of counter-factual prediction methods, designed to generate these metrics, effectively captures the significance of specific points, automatically highlighting them before pivotal moments occur, and establishes connections between players' performances on a seasonal level. This approach demonstrates effectiveness comparable to cutting-edge methodologies, addressing the static nature inherent in traditional methods, albeit with a focus on in-match information and an oversight of pre-match data.[3] Figure 2 illustrates these methods and their weaknesses.

## 1.4   Our Work

In sporting events, swings on the field often occur due to the dynamic nature of player performance. Even in scenarios where one player holds a substantial lead, unexpected turnarounds, termed as 'momentum shifts,' can occur, often attributed to psychological or situational factors. Consequently, we define 'momentum' as the current performance level of a player, wherein performance extends beyond mere technical skill to encompass a relative evaluation. A player who initially demonstrates clear dominance may experience a gradual decline in performance as the game progresses, especially if they have built a significant lead. Conversely, their opponent may become more focused and determined, leading to a reversal of fortunes on the field. Identifying pivotal indicators of momentum shifts is imperative for coaches and athletes to devise effective strategies in response. Broadly speaking, heightened momentum signifies consistent and superior performance, thus increasing the likelihood of victory. Although quantifying momentum shifts presents challenges and may involve concealed variables, examining correlations in empirical data can elucidate influential factors at play.[6]

We propose a novel prediction model that integrates elements from previous methodologies, leveraging counterfactual prediction techniques to derive indicators such as 'leverage' and 'momentum'. These indicators, along with static match information, are dynamically updated and synthesized to estimate the probabilities of player victories and losses, and to discern the flow of play. Our model operates at a fine-grained level, focusing on predicting win/loss probabilities for each point, subsequently extrapolating to game, set, and overall match outcomes. Notably, our approach is hierarchical, with each layer's output serving as input for the subsequent layer, culminating in a comprehensive prediction. Additionally, our model is highly adaptable and transferable, with flexibility in data selection and utilization tailored to specific match contexts. This adaptability extends its applicability across diverse match types. Our methodology combines dynamic programming principles at the macro level with a fusion of traditional machine learning and deep learning techniques at the micro level, ensuring precise and robust predictions. Notably, our model accurately anticipates match dynamics by considering a multitude of factors and data inputs.
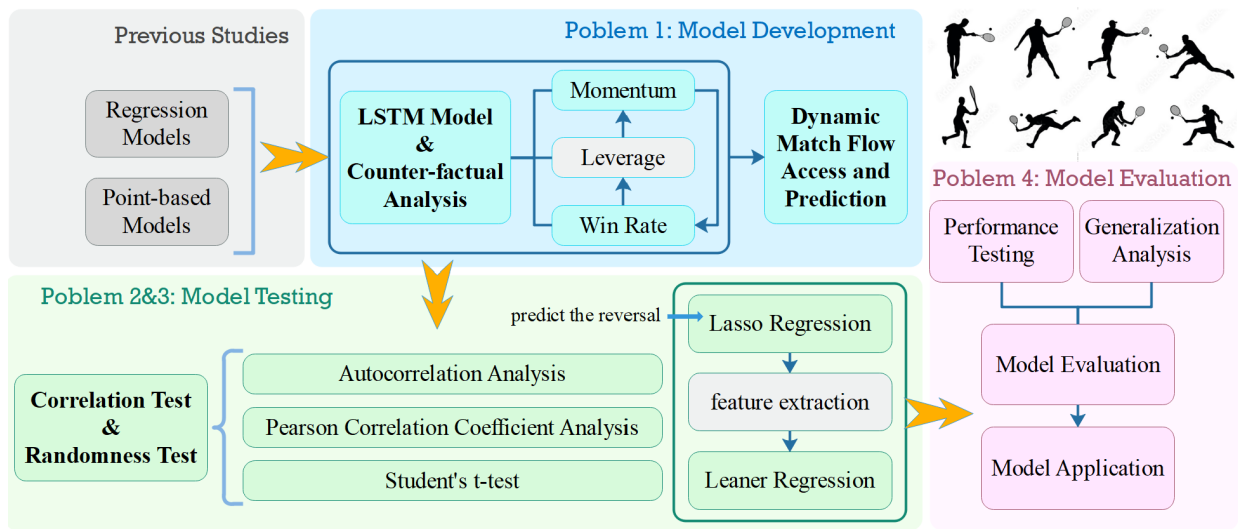


Figure 3: Flow chart of our work

## 2 Assumptions and Justifications

Given the inherent complexity of practical problems, our initial step involves making reasonable assumptions to simplify the model. Each hypothesis is meticulously followed by its corresponding explanation, ensuring clarity and coherence in our approach.

- **Assumption 1**: When computing the winning probabilities for individual games, sets, and overall matches based on each point in a match, we assume that the winning rates for subsequent points, both with and without the service advantage, remain constant. In other words, when calculating the win/loss ratio for a game, set, and match, it is solely contingent on the current model's presumed win rates for points with serve and points without serve.

- **Justification**: This consideration is necessary because, in calculating leverage, we must evaluate the significance of a point, namely its influence on the overall match-winning probability. Consequently, we assess its effect on the fluctuation of the winning probability for each game, set, or match. For this objective, any subsequent deviations stemming from non-independent distributions are considered negligible and thus not accounted for at the current point.

- **Assumption 2**: In a certain range of momentum, whether it is high or low, the performance characteristics of players in a game follow a normal distribution.

- **Justification**: The rationale for this consideration is as follows: Intuitively, when a player exhibits strong momentum, their performance tends to be superior compared to when their momentum is weak. Nevertheless, performance inevitably varies and can be presumed to conform to a normal distribution. In assessing the influence of momentum on the game, we utilize the t-test method grounded on this assumption.

## 3 Notations

Some important mathematical notations used in this paper are listed in table 1.

Table 1: Notations used in this paper

| Symbol | Definition |
| --- | --- |
| $m$ | momentum |
| $l$ | leverage, representing the importance of a point |
| $L$ | variable that represent short-term advantages |
| $P_f$ | variable that represent player's current state |
| $p_1$ | win rate of points when serving |
| $p_2$ | win rate of points when not serving |
| $p_e$ | the inferred probability of winning based on the Elo rating system |
| $p_g$ | the deviation in the probability of winning caused by other factors |
| $h$ | the probability of winning predicted by in-match performance |
| $u$ | efficiency difference |
| $q_1$ | scoring rate by the player when serving |
| $q_2$ | scoring rate by the player when not serving |

Note: Some variables are not listed here and will be discussed in detail in each section.

# 4   Data

Due to non-numeric and missing data in the provided dataset (`Wimbledon_featured_matches.csv`), we conducted data preprocessing before model construction. The steps include:

- **serve_width, serve_depth, return_depth:**   These non-numeric data were encoded with numerical labels (0, 1, ...).

- **speed_mph, return_depth:** Missing values denoted by 'NA' were addressed through statistical analysis on non-missing data.  Imputation was performed by randomly filling missing values based on calculated probabilities.

- **p1_score, p2_score:** In tennis scoring, scores 15, 30, 40, and AD were translated into number of points won (1, 2, 3, 4, ...), facilitating computational processes.

- **elapsed_time:** The elapsed time of the competition was converted from hours, minutes, and seconds into a numerical representation in hours for ease of computation.

Additionally, it should be noted that in this paper, apart from utilizing the official dataset provided by COMAP, supplementary data on players' age and Elo ratings was also obtained. The source of this additional data is from.`https://www.wheeloratings.com/tennis_atp_ratings.html`.

The objective of the Elo rating is to act as an initial value for computing the probability of winning based on pre-match information, whereas the function of age is to ascertain the influence of 'efficiency' decline on win probability.

# 5   Dynamic Match Flow Model

## 5.1   Model Establishment

### 5.1.1   Overview

In the realm of sports competitions, 'momentum' can be defined as the force generated from short-term advantages or favorable conditions, such as a player's current state, which may motivate the player to produce an enhanced performance in the near future. Accurately quantifying this impetus is paramount in our research to adequately capture and analyze the notion of momentum.

Short-term match advantage can be quantified by examining the number of points a player wins in the recent sequence of points. However, it's essential to recognize that tennis differs from sports like soccer or basketball, where the overall score is evaluated after the match. In tennis, assessment occurs at a more detailed level, considering individual games and sets, with the ultimate determination of the winner based on sets and games. Consequently, the significance of each point in a tennis match varies considerably. For example, a set point, which determines the outcome of that set, carries more weight than a point scored at the beginning of the set. To accommodate this variability in importance, we introduce the leverage variable $l$ to gauge the criticality of each point:

$$l = c \cdot (\pi_w - \pi_l)$$

$$c = \begin{cases} 1 & \text{, if win this point,} \\ -1 & \text{, otherwise.} \end{cases} \tag{1}$$

$\pi_w$ indicates how much winning the point would increase the probability of winning, and $\pi_l$ indicates how much losing the store would decrease the probability of winning and c represents whether this point is actually won. The gap between them represents the leverage associated with this point.

As short-term match advantage focuses on the recent points won, it becomes necessary to apply a weighted treatment to the leverage of the recent points:

$$L = \frac{l_t + (1 - \alpha)l_{t-1} + (1 - \alpha)^2 l_{t-2} + \cdots + (1 - \alpha)^t l_0}{1 + (1 - \alpha) + (1 - \alpha)^2 + \cdots + (1 - \alpha)^t} \tag{2}$$

where $[l_0, l_1, \cdots, l_t]$ are gained leverages of the last $t$ points with smoothing factor $\alpha$ and $L$ the weighted leverage at point $t$. In this work, we utilize $\alpha = 0.33$. The closer the leverage is to the current point, the greater the weight it holds, exerting a more significant influence on the weighted leverage.

Next, we need to consider the positive or negative feedback generated by the player's state factors during the match. This includes on-court events that influence the player's state factors, such as negative feedback from a double fault and positive feedback from an untouchable shot. Additionally, we should account for the impact of the current match environment and the macroscopic perspective of the player's season career. This can be expressed as:

$$P_f = f(\S) \tag{3}$$

where $\S$ is a feature vector that dynamically updates throughout the match. Generally, $P_f$ can be considered as a function following a linear relationship:

$$P_f = \omega \cdot \S + b \tag{4}$$

We have researched and evaluated the impact of each feature on the player's momentum through literature review and data collection. Based on this, we classified these features as rewarding or penalizing items and assigned them a score given in $P_f$. The example of player p1 is shown in table 2.

Table 2: Score of features in $P_f$

| Feature | Score | Feature | Score |
|---|---|---|---|
| p1_points | 1 | p1_break_pt | 2.5 |
| p1_games | 5 | p1_break_pt_won | 7.5 |
| p1_sets | 30 | p1_net_pt | 0.3 |
| p1_ace | 0.7 | p1_double_fault | -0.65 |
| p1_winner | 0.5 | p1_unf_err | -0.45 |

In addition, we have calculated the score for winning streaks through a weighted approach:

$$k = \frac{value}{1 + e^{-sx}} \tag{5}$$

where $x$ is a variable representing the winning streak, $value$ is the ideal point value for winning two consecutive sets, set to 90 in this paper, and $s$ is a parameter determining the steepness of the impact function, set to 0.15. This setting allows for a comprehensive consideration of the cumulative effects resulting from various scenarios of consecutive victories, and as $x$ approaches 20, the function asymptotically approaches 1. Equation 6 illustrates the calculation of $x$, where $x_1$ represents the number of consecutive points won, $x_2$ represents the number of consecutive games won, and $x_3$ represents the number of consecutive sets won. The three are weighted to obtain the variable $x$, with weights $\alpha$, $\beta$, and $\gamma$ set to 0.5, 1, and 3, respectively:

$$x = \alpha \cdot x_1 + \beta \cdot x_2 + \gamma \cdot x_3 \tag{6}$$

By taking the difference between the $P_f$ values for the two players, we obtain the overall match's $\Delta P_f$.

Thus, we have the formula for momentum $m_t$, derived from the combination of two factors: weighted leverage and player performance. :

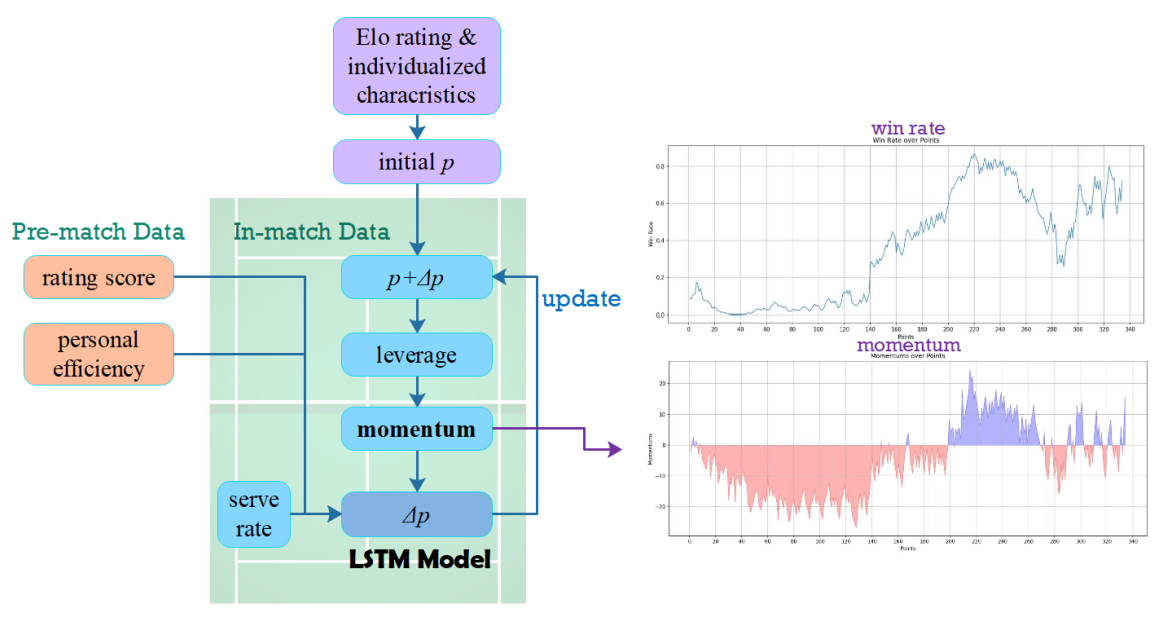$$m_t = \alpha \cdot L + \beta \cdot \Delta P_f \tag{7}$$



Figure 4: Flow chart

In equation (1) for the leverage, the probability is determined by our model based on various factors including momentum, technical prowess, and physical fitness. As per the definition of momentum outlined earlier, it is inherently dependent on probabilities within the leverage. Consequently, this scenario presents itself as a dynamic programming challenge, where we initially establish a probability distribution and continually update the winning probabilities for each scoring point based on this distribution. Concurrently, we also obtain the evolving trend of momentum.

The model's essence is encapsulated in two pivotal aspects: the dynamic evolution and fine-tuning of variable $p$ over time, and the initialization strategy for $p$. For the former, we employ machine learning techniques to craft a model that discerns the requisite functional relationship. Regarding the latter, we draw upon insights from the Elo rating system. By combining this rating with other pertinent factors, we compute the initial probability of winning via a well-defined mathematical function.

### 5.1.2 Initialization of $p$

Based on the weighted Elo rating method for tennis match prediction proposed by Giovanni Angelinia, Vincenzo Candilab, and Luca De Angelisc et al.[3], we can obtain a pre-match rating score for each player, based on which we can roughly infer the probability of winning. The most straightforward approach involves deriving the probability of winning based on the ratio of the Elo ratings. In our model, we also provide an optional deviation value ,denoted as $p_g$. The composition of $p$ is as follows:

$$p = p_0 + p_g \tag{8}$$

where $p_0$ represents the inferred probability of winning based on the Elo rating system, and $p_g$ is the deviation in the probability of winning caused by other factors.

$$p_0 = \alpha \cdot \left(\frac{elo_1}{elo_1 + elo2_2}\right) \tag{9}$$

where $elo_1$ is the rating score of player 1, $elo_2$ is the rating score of player 2, and $p_0$ represents the probability of player 1 winning.

$$p_g = \beta \cdot g \tag{10}$$

where one possible logistic regression form for $g$ in equation (10) is given by:

$$g = \frac{1}{1 + e^{-(\beta_0 + \beta_1 z_1 + \beta_2 z_2 + \beta_3 z_3 + \ldots + z_n)}} \tag{11}$$

The function is used to calculate the probability deviation of winning caused by factors other than the rating score. Specifically, we can consider factors such as the current match environment, player age, event type, etc., to determine the deviation generated by their combined effect. Once the form and parameters of the model are established, statistical methods can be employed to estimate these parameters. We can either pre-set hyperparameters and weights based on empirical knowledge or use machine learning methods, such as maximum likelihood estimation, to obtain $g$. Subsequently, the obtained model is utilized to predict the probability of winning the initial point and analyze the impact of various factors on the probability of winning. Combining the results from the Elo rating system and the initial winning probability calculated by this model provides the initial values for dynamic programming.

### 5.1.3 Updating Methods for p

Updating the probability of winning the next point based on the current momentum and game situation can be challenging. To accurately represent a player's probability of winning, two probabilities will be used for each match: $p_1$ for when the player has the serve and $p_2$ for when they do not. It is

important to note that a player is typically more likely to win when they have the serve. If the player is serving, we update the win rate $p_1$ for that match. If not, we update $p_2$.

$$p_1 = \mu_1 p_1 + \alpha_1 m_t + \beta_1 h + \gamma_1 u + \delta_1 q_{t_1}$$
$$p_2 = \mu_2 p_2 + \alpha_2 m_t + \beta_2 h + \gamma_2 u + \delta_2 q_{t_2}$$

(12)

When updating $p_1$ and $p_2$, we use a linear model to calculate the probability of winning the next point. The calculation of momentum has already been explained above, where $m_t$ represents the difference in momentum between the two players at the current point in time. The values of $h$ and $u$ represent the probability of winning, respectively, as predicted by in-game performance, and the efficiency difference. The win probability predicted by in-game performance is obtained through training an LSTM model. Efficiency is reflected in the detracting effect of the athlete's individualized characteristics as time progresses. For instance, the reduction in physical strength over time may significantly affect older players, but not younger ones. $q_{t_1}$ and $q_{t_2}$ represent the probability of winning of points scored by the player when serving and not serving, respectively. This metric is more instructive, so a separate item is used to update $p_1$ and $p_2$, which can be determined empirically or through machine learning methods.

$$h = LSTM(\varphi_0, \varphi_1, \varphi_2, ...)$$

(13)

As previously stated, the win probability predicted by the player's performance in the match is calculated by inputting their performance features during the inning into the LSTM computation. The LSTM model is chosen due to its ability to handle time series data effectively. To ensure that the predicted winning probability based on in-game performance aligns with the initial value predicted by Elo scoring before the game, a penalty term was added to the loss function which helps to maintain accuracy and consistency in our predictions. Specifically, the loss function is as follow:

$$Loss = \sum [y_i \log \hat{p} + (1 - y_i) \cdot \log(1 - \hat{p})] + \sigma(\hat{p} - p_0)^2$$

(14)

Where $\hat{p}$ is the probability of winning predicted by the LSTM model, $p_0$ is the initial value of the probability of winning based on the Elo score and $y_i$ indicates whether the point is winning or not. $y_i$ then indicates whether the point wins or not, 1 if it wins, 0 otherwise. $\sigma$ is the strength of the self-defined penalty term, the larger it is, the less the predicted probability deviates from the initial value.

To represent the relationship between efficiency and time variation, a function is employed:

$$u = \phi(t, u_0)$$

(15)

**The modeling of the function of efficiency over time is an intriguing and innovative aspect**, where we consider that intermissions will contribute to the recovery of physical strength. Therefore, we analyze the relationship between efficiency over time for the game between breaks. It's important

to note that breaks do not result in a complete recovery of physical strength. The variable $u_0$ represents the initial value after each break, decreasing with the number of breaks $i$. It is worth mentioning that $u_0$ and the function $phi$ may vary among players. In practice, we do not take into account the effect of $u_0$ decreasing with the number of breaks. Instead, the default $u$ consistently decreases over time from the initial $u_0$. This is because the positive feedback from breaks is insignificant in the tennis system, which has short and infrequent breaks. We use a specific function:

$$\phi(t, u_0) = \frac{2 \cdot u_0}{1 + e^{\theta \cdot t}} \tag{16}$$

This function was chosen due to its monotonically decreasing trend, which aligns with the anticipated decrease in efficiency over time. With appropriate adjustments, it can gradually converge to 0 as $t$ approaches 4, which is the average time spent in most races. The idea is intuitive, older players may experience a sharper decrease in stamina over time. The symbol $\theta$ can be used to control the steepness of the function, i.e., the rate of decrease in stamina.In this case, stamina can be used as part of the efficiency of the expression.

We devised a function to associate ages with parameters (denoted as $\theta$). This mapping function exhibits a monotonically increasing convex shape, with its peak values skewed towards regions where age distribution is more concentrated. This design choice aims to better capture the concentrated age distribution of players. The parameters derived from this mapping are treated as the mean of a normal distribution function. Notably, the mapping of regions with concentrated age distributions corresponds to a smaller standard deviation due to the larger population in these areas, which tends to be younger according to research findings.[1]
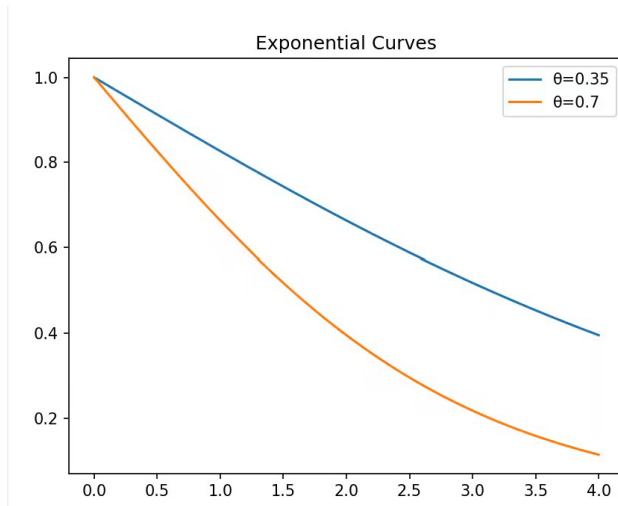


Figure 5: Exponential curves

We utilize the resulting normal distribution for sampling, randomly selecting values within the range of $\alpha$ and $\beta$ and employing the sampled result as the parameter $\theta$. We set $\alpha$ and $\beta$ for this function to 0.35 and 0.7, respectively, as the steepness of this interval better reflects the effect of monotonically decreasing efficiency over time.The primary purpose of introducing randomness is to accommodate the potential distorting effects of factors other than age on efficiency. Even players of the same age may exhibit variations in their $\theta$ values. In summary, we observe a notable trend of decreasing efficiency over time.Using eq. (11) it is possible to find the gain probability, which is then added to the probability value of the previous moment to obtain the updated probability.

---

[1]https://www.itftennis.com/en/players/

### 5.1.4 Hierarchical Progression of $p$

The probability of winning a given scoring point is obtained, so we can use that probability to obtain the probability of winning a game, a set, or a match and update it incrementally. This part of the approach can be considered by a simple permutation using the current model predictions for the probability with and without serve, denoted as $p_1$ and $p_2$.

Consider the following problem: Two players, denoted as A and B, participate in a game where each point is won at a rate of $p$ for A and $q = 1 - p$ for B. Each point is an independent event during the game. The game is configured with an $n$-point format where the first player to score $n$ points and lead by at least 2 points wins. If the score is tied at $n - 1 : n - 1$, the first player to lead by 2 points wins. To determine the probability of both A and B winning the game, we use the notation $w(i, j)$ to represent the probability of A winning a game given that A has scored $i$ points and B has scored $j$ points.According to the rules of the game and the probability assumptions

$$
\begin{aligned}
w(i, j) &= 1 \qquad i \geq n, i - j \geq 2 \\
w(i, j) &= 0 \qquad j \geq n, j - i \geq 2 \\
w(n - 2, n - 2) &= w(n - 1, n - 1) = w(n, n) = \ldots
\end{aligned}
\tag{17}
$$

By simple permutations and combinations, we can determine the probability of $w(m - 2, m - 2)$:

$$
\begin{aligned}
w(m - 2, m - 2) &= p^2 + p^2(2p(1 - p)) + p^2(2p(1 - p))^2 + \ldots \\
&= \frac{p^2}{1 - 2p(1 - p)}
\end{aligned}
\tag{18}
$$

By the law of total probability, the following relationship exists:

$$
w(i, j) = p \cdot w(i + 1, j) + (1 - p) \cdot w(i, j + 1)
\tag{19}
$$

Under the above recursive relationship, employing a dynamic programming approach allows us to calculate the probability of winning this set.

To calculate the winning probability of each subsequent set, we can firstly use either $p_1$ or $p_2$ to get the winning probability of the newly started game. The server alternate in each set, so the winning probability of a set can be found with slight modifications to the above method. Finally, the probability of winning the match is calculated.

## 5.2 Results and Analysis

### 5.2.1 Model training process

The LSTM model underwent training using K-fold cross-validation. In the validation stage, the final match considered was the encounter between Carlos Alcaraz and Novak Djokovic. Through this method, we were able to track the momentum fluctuations as points were scored, thereby quantifying the superiority of each player at distinct moments and determining the extent of their advantage.

### 5.2.2 Match flow results and explanations

Figure 6 shows the momentum trend of both players throughout the match as a function of points. It is evident that Djokovic's momentum was dominant from about the beginning of the match until about 200 points. He maintained a peak during the 80 to 130 point, followed by a sharp decline. In the subsequent period, Djokovic's momentum remained lower in absolute terms, although it was still higher most of the time. At approximately 200, Alcaraz countered Djokovic's momentum. Alcaraz maintained a higher momentum for some time before being countered again by Djokovic at approximately 270. The momentum then shifted back and forth between the two players until the end of the match, when Alcaraz ultimately won with the higher momentum.
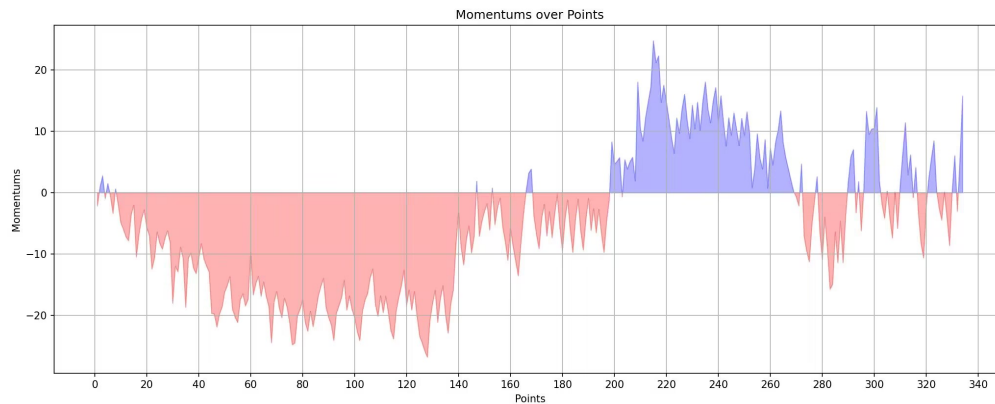


Figure 6: The momentum of both players is **symmetrical** about the point axis. Therefore,the momentum graph for player 1 can represent the momentum of the entire match. A positive value of momentum indicates that Player 1 performed better, while a negative value indicates that Player 2 performed better. **The absolute value of the momentum indicates the quality of the performance.**
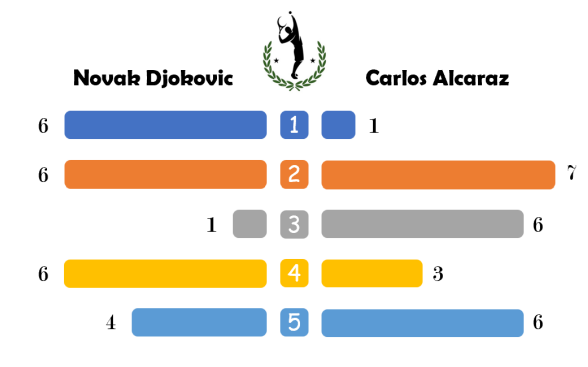


Figure 7: Score

This can be attributed to the fact that Djokovic's momentum peaked just as he won the first set and almost won the second. The momentum shifted after he lost the second set in a tiebreak, but Djokovic still had an advantage due to his previous good performances. At around 200 points, Alcaraz gained momentum and took a 4-1 lead in the third set. During this period with higher momentum, he had a 2-1 advantage.However, at approximately 270 points, Alcaraz became the underdog in the fourth set, which Djokovic eventually won, reversing the momentum.

With the score tied at 2-2 sets, the fifth set tiebreaker was a tense affair with momentum constantly shifting between the two players. In the end, Alcaraz emerged victorious in the fifth set and won the match.

The term momentum accurately describes the performance of both players and the situation on the court. We analyzed the momentum of several other matches and compared them with the actual scores. The results showed that momentum can describe the players' performance well.

Additionally, the model can be expanded to selectively incorporate various information. For instance, if information about the playing environment needs to be considered, it can be added to the momentum's $p_g$. For $u$ in the update method of $p$, we can determine the corresponding $u_0$ value and its decreasing relationship over time based on the personal information of the statistical players, and introduce a more precise function to describe the relationship and changes.

# 6 Correlation Test and Randomness Validation

## 6.1 Overview

We aim to determine whether momentum affects the outcome of the game. In Problem 1, we explicitly defined 'momentum' as a factor that provides feedback on the winning rate of a match. Based on this definition, we proposed a method for updating the winning percentage, and thus, we conclude that momentum plays a role in the match. Therefore, we posit that momentum is a factor in the game, and that fluctuations in the game and the success of players are not random events. To verify this, we conducted an Autocorrelation test to check for the presence of autocorrelation in the data. Autocorrelation reflects the association of data points in time or space. If the data is random, there should be no significant autocorrelation between them.

To evaluate the impact of momentum on the game, we can use $t$-test. The method allows us to determine the correlation between the eigenvectors of each score point and momentum, as well as the degree of correlation between them.

## 6.2 Results And Analysis

### 6.2.1 Autocorrelation Analysis

We calculate the autocorrelation coefficients for the obtained 'momentum data' and then plot the autocorrelation function in Figure 8:
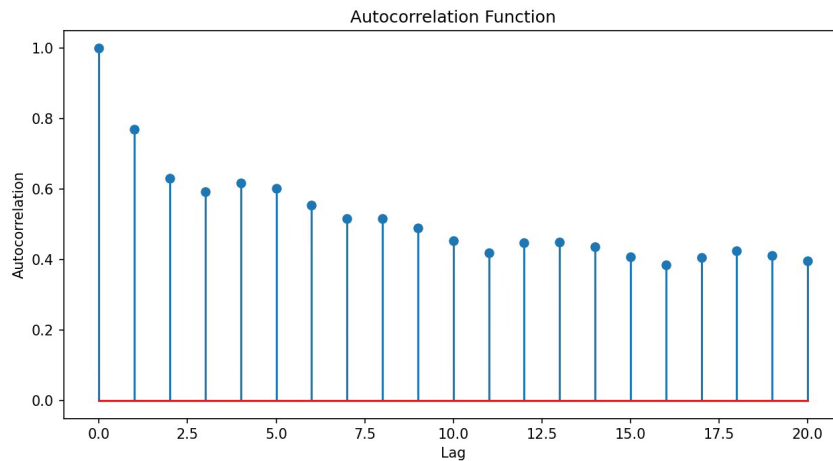


Figure 8: Autocorrelation function

- The horizontal axis (X-axis) represents the lag order, which is the time interval between the current observation and the previous observation. A positive integer indicates the correlation between the current observation and a past point in time, while a negative integer indicates the correlation between the current observation and a future point in time.

- The vertical axis (Y-axis) represents the autocorrelation coefficient, which indicates the degree of correlation at a given lag order. This coefficient typically ranges from -1 to 1, with 0 indicating no correlation, 1 indicating perfect positive correlation, and -1 indicating perfect negative correlation.

The autocorrelation coefficient is significantly non-zero at certain lag orders, indicating significant autocorrelation in the data. This suggests that the observed result is not random.

### 6.2.2 Pearson Correlation Coefficient Analysis

On the other hand, we calculated the Pearson correlation coefficients between the different variables and obtained the heatmap shown in Figure 9. It is evident that the variable 'momentum' has a high correlation with several other variables, indicating that the 'momentum' we obtained is not random and will have a greater impact on the race.
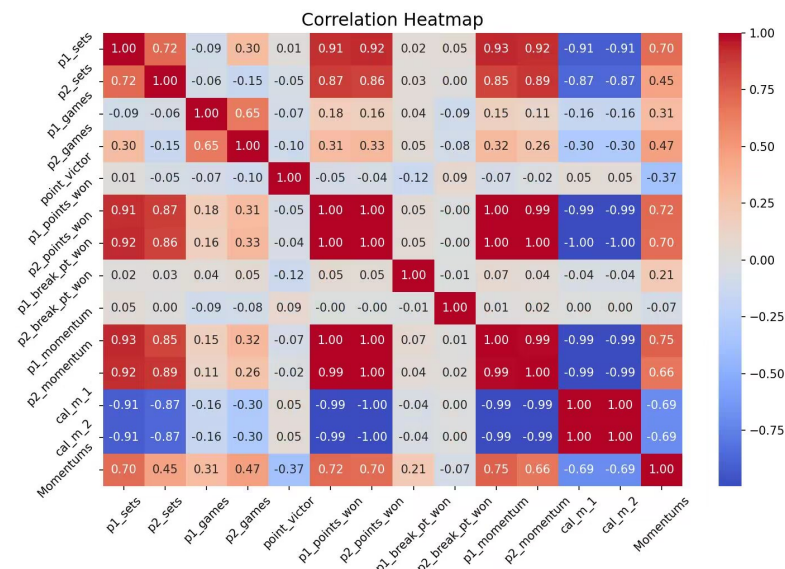


Figure 9: Correlation Heatmap

### 6.2.3 Student's t-test

We analyzed all 31 games, testing performances and results in points with momentum less than half of the minimum and greater than half of the maximum using t-tests. A difference was considered significant if the p-value was less than 0.05. The study results indicated a significant difference in point victor performance in 29 out of 31 matches, suggesting a strong correlation between momentum and point victories. Most of the $t$-test results for the other performances also indicate significant differences of more than 5 matches.

$$t = \frac{\bar{X}_1 - \bar{X}_2}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \tag{20}$$

Figure 10 displays the t-test results for some of the features in four randomly selected games, with the values on the radar plot indicating the magnitude of the $1 - p$ values.
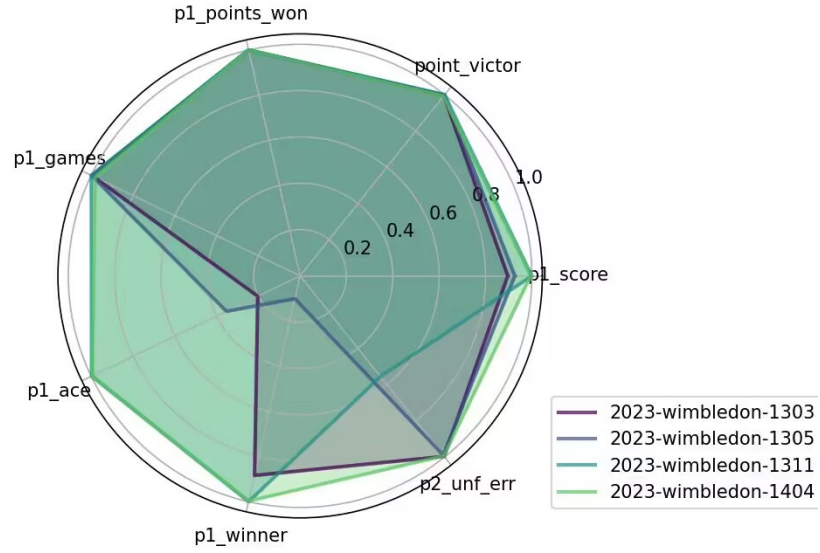


Figure 10: t-test Radarmap

Figure 10 shows that all of these features are correlated with momentum in one or more matches. Specifically, indicators such as point victor, p1 score, and p1 winner are significantly related to momentum in almost every game, indicating that momentum plays a crucial role in the game.

# 7    Match Momentum Swing Prediction Model

## 7.1    Model Establishment

The momentum trend graph from our dynamic match flow model can help determine when one player has an advantage over the other. The zero point where the sign of the momentum changes of the momentum variable indicates when the situation reverses. These key moments are characterized as a function of momentum. Plotting the trend of momentum as a function of points scored, allows for easy identification of fluctuations throughout the game and key points of change.

To identify the factors that influence these changes, we utilize the Lasso regression method. The mathematical expression for Lasso regression is:

$$\min_{\beta} \left( \frac{1}{2N} \|y - X\beta\|_2^2 + \alpha \|\beta\|_1 \right) \tag{21}$$

- y represents the vector of observed target variables.

- X is the matrix that contains the characteristics of each observation.

- β is the vector of coefficients for a linear regression model.

- N represents the number of observations.

- α is the regularization parameter that controls the strength of the regularization

Lasso regression, employing L1 regularization to induce sparsity in the coefficient vector β, achieves feature selection. The L1 regularization term penalizes the sum of absolute values of the coefficient vector, driving many coefficients to zero during optimization. Consequently, models trained using Lasso regression often yield features with zero coefficients, indicating a smaller impact on the target variable, while non-zero coefficients suggest greater impact. Examining the coefficient vector identifies features with significant impact on the target variable. This aids in determining factors associated with changes in the on-field situation[4].

Next, we extract the feature vectors with high correlation and build a new model. We use a simple linear regression fitting strategy to predict the point at which the momentum reverses. This information can be used to provide players with advice on how to play against different opponents.

## 7.2 Results And Analysis

### 7.2.1 Identification of Key Indicators for Momentum Reversal

Based on the results obtained through Lasso regression, it is evident that certain indicators in Figure 11 played a crucial role in the critical point where the momentum shifted.
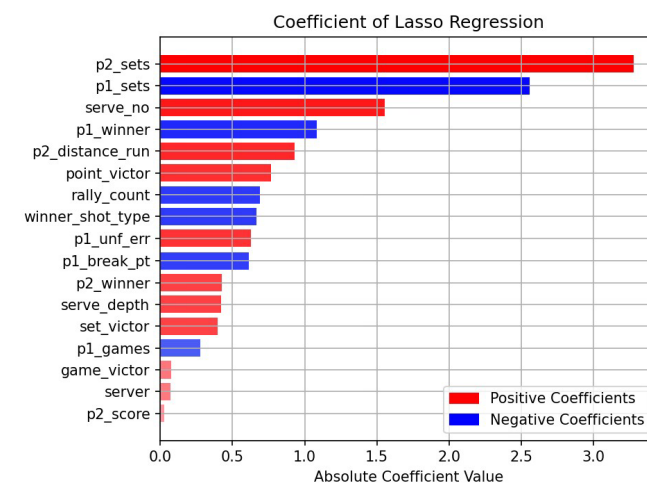


Figure 11: Coefficient of Lasso Regression

Specifically, "sets", "serve_no", and "winner" had a significant impact on the reversal of momentum. Next, we utilize key indicators to map momentum reversal points in a linear regression. We then use the trained model to predict potential momentum reversal points.

### 7.2.2 A Linear Regression Model for Predicting Momentum Reversal Points

The predicted results are highly consistent with the actual values, as demonstrated in Figure 12. The linear regression fit is valid as the $R^2$ value of the model is close to 1 and the mean square error is small.
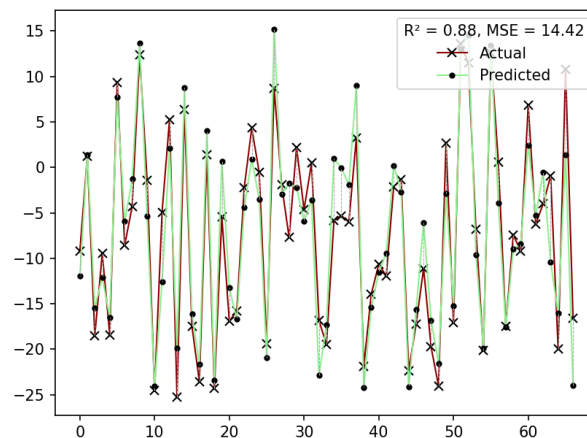
Figure 12: The result of Liner Regression

### 7.2.3 Application of the Model for Individual Players

Based on our findings, we can practically apply the model to benefit individual players by tailoring specific outcomes for each player. To achieve this, we can analyze the "momentum" trajectory of a player's performance against various opponents, pinpointing the pivotal moments where the player transitions from advantage to disadvantage. Subsequently, we can identify the key metrics, as described earlier, and utilize them to train the model in predicting these critical "momentum" turning points. Armed with knowledge about the influence of each indicator on the on-field scenario, we can monitor these indicators to formulate a strategic response when the player encounters a new opponent.

For instance, if we know that the number of strokes, distance traveled, and number of serves affect momentum, we should select the optimal measures to avoid reversing or further increase the current momentum.

# 8 Methodology for model evaluation and application

## 8.1 Calculation of performance indicators

During the model training process, we perform cross-validation on 31 matches. For each validation, we select 5 matches and calculate the True Positive Rate (TPR) and False Positive Rate (FPR). We then use this information to generate an ROC curve and determine the AUC value, which serves as an evaluation metric for our model. It is important to note that our model predicts 'momentum', although there is no actual label for this metric. Our methodology involves converting momentum into the corresponding probability of winning. The probability and actual game results are then used to calculate the AUC value and ROC curve, which are used to evaluate the model's effectiveness.

As mentioned above, the core of our model is based on a point-level prediction method, which means that for such hierarchical tournament events as table tennis, badminton, and tennis, it can be adapted to different scenarios with only minor modifications. On the other hand, since we take into account various factors such as pre-match information and assessment of players' personal situations,

this scoring system can be useful for any tournament. At this point, our model is very general and can be extended according to different tournament types. We have not yet applied it to other types of tournaments because we have not found the right data.

For different levels of tournament types for the same sport, we suspect that the biggest influencing factor is the difference in the technical ability of the players and the athlete rating scores, an aspect that plays directly into the initial probability of winning that we predict. To deal with this, we consider the possibility of adding to the (8) equation for each level of participation an initial event rating score, such as the average of the rating scores of the competitors who participated in the event, in case the competitors have less pre-tournament information. This score is dynamically adjusted during the tournament to deal with the situation of facing a fresh new opponent (without enough information about him or her to evaluate), i.e. using dynamic data from the time of the tournament for evaluation and updating by technical ability, score, etc., which allows a more accurate prediction of win/loss ratio even against a first-time participant, even if the tournament has not yet ended. Notably, In sensitivity tests, we observed that minor inaccuracies in the initial win rates could be promptly rectified by the model over a relatively short period.

## 8.2 Results And Analysis

A total of 10 rounds of validation were conducted, with matches randomly selected as the validation set each time. Finally, we obtained the AUC values ranging from 0.82 to 0.93. Figure 13 shows the ROC curve graph corresponding to an AUC value of 0.93.
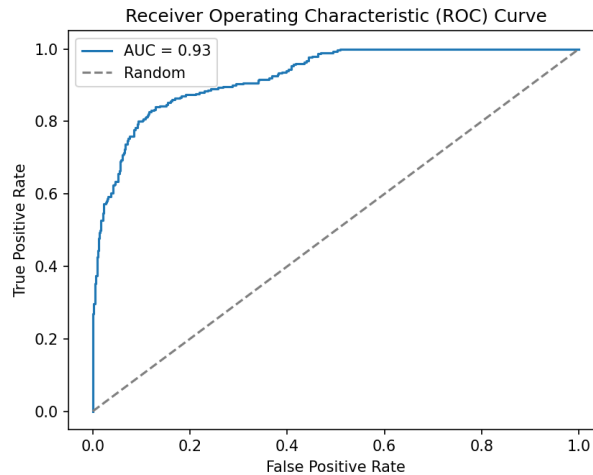


Figure 13: ROC Curve

By analyzing the prediction results, we found that the accuracy of the model is very high. If new factors need to be introduced to improve the model in the future, we have designed various extensible functions and expressions in the model. For example, the (8) equation mentioned earlier initializes p, which is critical to the model's predictions. For $p_g$ in this equation, to maximize its accuracy, we might consider introducing some other factors that take into account information about the game itself (not just the players' ratings), such as the effect of environmental factors on the outcome of the game, since different players may perform more adaptively in different environments. In formula (15), we simulate

the extent to which different players are affected as the game progresses through a simple model. This part of the design can be further refined to ensure that it is more accurate and efficient when updating the p.

# 9    Sensitivity Analysis

In this section, we conducted a sensitivity analysis specifically for our dynamic match flow model. For the subsequent models derived from the dynamic match flow model, if the sensitivity analysis of our model proves effective, the robustness of the new models generated from it will likely be high. Therefore, our analysis primarily focused on two core components of the dynamic match flow model.

## 9.1    Sensitivity to Probability Updating Coefficient

As described in Section 5, we adopt a linear combination approach for updating the probability $p$, where the adjustments for momentum and efficiency are minor, exerting a relatively small impact on the updated results. In contrast, the coefficients $\mu$, $\beta$, and $\gamma$ in the linear combination of the other three probability components have a more substantial influence on the updated probabilities. To assess the robustness of the results, we conducted a sensitivity analysis on $\mu$, $\beta$, and $\gamma$. Figure 14 illustrates the changes in the predicted win probabilities and momentum under different values for these coefficients.
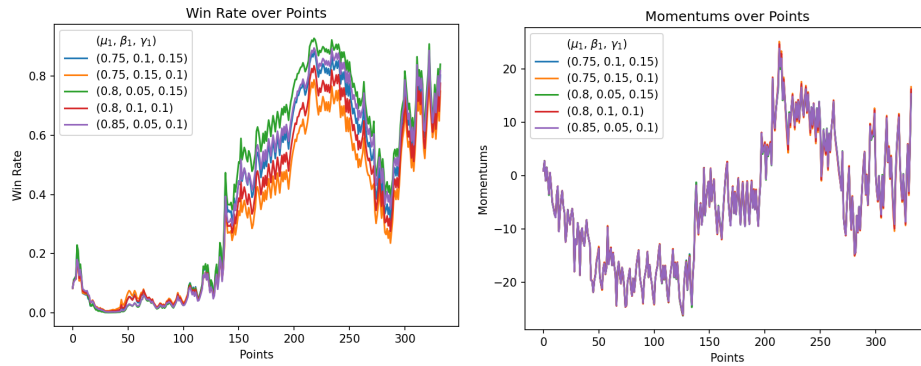


Figure 14: Win rate and momentums over points under different updating coefficient

As evident from the results, the three coefficients in the linear combinations exhibit minimal influence on the momentum. Despite having some impact on the model's predicted probability of winning, the overall trajectory of the predicted win probability remains consistent. This observation highlights the robust performance of the model in sensitivity tests conducted on these parameters.

## 9.2    Sensitivity to Initial Probability

In our dynamic match flow model , the initialization of a player's win probability in a point is determined through elo scores and other fine-tuning factors. We aim to evaluate whether our model can promptly adjust the predicted probability if the initialized probability is not sufficiently accurate. To

assess this, we conducted a sensitivity test on the initialized probability values, specifically by adjusting the weights of fine-tuning factors to modify the model's initialized point winning probability. Figure 15 illustrates the changes in the model's predicted win probability and momentum under different initialized win probability values.
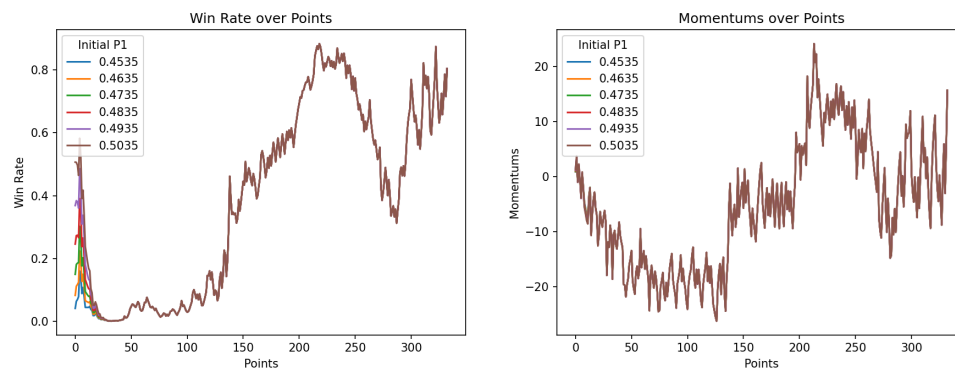


Figure 15: Win rate and momentums over points under different Initial Probability

It can be observed that despite the initial differences in probabilities, after a relatively small number of points, the predicted win probabilities of the model converge closely, and the changes in momentum are nearly identical across various initial probability settings. This indicates that the model is capable of promptly correcting minor deviations in the initialized probabilities, showcasing its robustness. This ability may stem from our comprehensive consideration of player performance during the match when updating probabilities. Consequently, the model can promptly adjust predicted probabilities based on players' on-field performances and scoring outcomes. Such findings underscore the model's resilience and adaptability.

# 10 Strength and Weakness

## 10.1 Strength

1. **Versatility and Scalability:** Our model boasts a high level of generality and scalability, making it compatible with various factors. By integrating static pre-game information with dynamic in-game data during the model construction process, we ensure its versatility. What sets our model apart is its comprehensive consideration of the diverse effects stemming from random factors during matches, rather than relying solely on historical data for assessment. With the insights offered by our model, it can be readily applied to diverse match types. Furthermore, refinement of some initial considerations and prototypes within our model can enhance its adaptability and predictive prowess.

2. **Robustness and Resilience:** Our model exhibits a remarkable robustness to realistic data due to its consideration of multiple factors. Even when individual factors experience perturbations, the overall results remain largely unaffected. Should outliers emerge in the data, our model possesses the capability to mitigate or even eliminate their influence, owing to its complexity and comprehensiveness.

## 10.2   Weakness

1. **Cursory consideration:** While our model adequately considers a plethora of factors, the handling and modeling of individual factors may be somewhat rudimentary. For instance, the prediction of the initial pre-match probability of winning may lack precision, and the formulation of the function describing the change in efficiency over match duration may not accurately reflect practical scenarios. Additionally, our scoring system for in-match predictions may be insufficiently refined, resulting in potentially significant impacts.

2. **Lack of validation in different scenarios** In theory, our model should be applicable to a wide range of match types and scenarios. However, practical verification was hindered by limited datasets and resources. Furthermore, the small dataset utilized may compromise the model's generalization performance across different playing fields.

# References

[1] Kovalchik and Reid, 2019 S.A. Kovalchik, M. Reid "A calibration method with dynamic updates for within-match forecasting of wins in tennis International Journal of Forecasting", 35 (2) (2019), pp. 756-766,

[2] Angelini, Candila & De Angelis, 2022 G. Angelini, V. Candila, L. De Angelis "Weighted Elo rating for tennis match predictions".European Journal of Operational Research, 297 (1) (2022), pp. 120-132,

[3] Robert Seidl & Patrick Lucey ."Live Counter-Factual Analysis in Women ' s Tennis using Automatic Key-Moment Detection".MIT Sloan Sports Analytics Conference, 2022At: Boston, USA

[4] R. Tibshirani (1996) ."Regression Shrinkage and Selection via the Lasso".Journal of the royal statistical society series b-methodological. Mathematics, Computer Science

[5] Jonathan Manuel."Capturing Momentum in Tennis". `https://theanalyst.com/eu/2022/03/capturing-momentum-in-tennis/`

[6] Dietl, Helmut; Nesseler, Cornel (2017). Momentum in tennis: controlling the match. UZH Business Working Paper Series 365, University of Zurich.

[7] Fitzpatrick, A., Stone, J.A., Choppin, S. and Kelley, J., 2019. A simple new method for identifying performance characteristics associated with success in elite tennis. International Journal of Sports Science & Coaching, 14(1), pp.43-50.

# Appendices

In the context of our model establishment, training, and validation, all relevant code has been made publicly accessible on GitHub. Interested parties can access the code repository through the following link: `https://github.com/Jia040223/MCM`.

Dear Sir or Madam:

With the improvement of statistical analysis methods, sports analytics for athletes has become increasingly detailed and scientific. Indeed, there remains a paucity of research on momentum in tennis, compounded by the absence of a definitive and universally accepted definition of momentum. Current interpretations are often grounded in consistent yet abstract intuition, leaving considerable ambiguity regarding the extent to which momentum truly impacts the game.

In the context of sports competitions, 'momentum' can be intuitively understood as the force derived from short-term advantages and favorable factors such as the current state of the player. This force reflects the player's performance and, to some extent, serves as motivation for the player to perform better in the subsequent phases of the game.
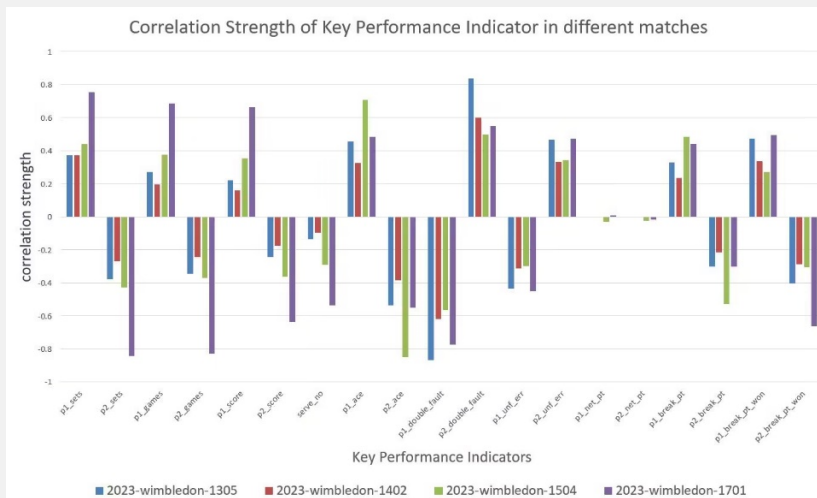
Our team has proposed a model to assess match flow and predict key moments of momentum reversal. We aim for these research findings to assist coaches and players.

Extensive research has underscored the significant correlation between momentum and match outcomes, emphasizing the importance of sustaining high momentum levels on the field. However, momentum reversals are frequent occurrences during play, necessitating a keen awareness of these pivotal moments for players. Our investigation has unveiled compelling correlations between specific key indicators and momentum reversals. Given this insight, coaches can strategically guide athletes in exploiting these indicators when faced with adversity, facilitating a turnaround, or preemptively mitigating momentum shifts when in a favorable position.

The following are specific conclusions and recommendations:

- **Ace**: The research findings reveal a highly significant correlation between aces and momentum. This aligns with our intuition. Therefore, when coaching players, emphasis can be placed on training the speed and placement of serves. This enables players to achieve higher accuracy in delivering ace shots on the court, leading to enhanced momentum.

- **Break/Hold**: The research findings indicate that breaking the opponent's serve is a crucial factor influencing momentum and even the outcome of the match. This highlights the importance of players maintaining composure when facing a break and underscores the significance of successful holds.

- **Serve Count**: Research indicates that the number of serves is a key factor influencing momentum shifts. When a player successfully scores with the first serve, they often experience better momentum since there is no pressure from a potential double-fault. This underscores the importance of training to improve serving skills. Additionally, it suggests that moments of double-faults on the first serve are vulnerable to a momentum reversal, emphasizing the significance of maintaining composure during such situations.

- **Unforced Errors**: Research suggests that unforced errors significantly impact the momentum of a match. This implies that players should strive to minimize unforced errors. At the same time, capitalizing on the opponent's unforced errors is a key factor in influencing momentum.

- **Net Approaches**: Our research findings did not identify a strong correlation between making it to the net and momentum, which is somewhat counterintuitive. This may be due to the infrequent nature of net approaches. Additionally, net approaches, while offering high rewards, often come with substantial risks, making their success less impactful on momentum. In other words, when both players have comparable momentum, it is not advisable to approach the net recklessly. When the opponent approaches the net, there may not be significant pressure to respond.

- Additionally, when the momentum is entirely in favor of the player, the opponent is more likely to change tactics. The player in the lead should also aim to diversify their game.



Correlation Strength of Key Performance Indicator in different matches

This graph illustrates the autocorrelation coefficients of various indicators at the turning points of momentum reversal in four matches. Intuitively, the larger the values, the more significant these indicators are. Positive indicators are advantageous for players, while negative indicators are unfavorable.

We believe that, with your recognition of the significance of momentum in matches and a comprehensive understanding of how various indicators affect game dynamics, you can scientifically plan training content and match strategies for the players. If you have any suggestions or ideas, please feel free to contact our team. Wishing you and your players smooth and successful matches.

Yours Sincerely,
Team# 2429435

# Report on Use of AI

We use AI in three main ways: to provide ideas, to write code and to translate corrections. The specific AI models used are Chatgpt and the Copilot model from Github.

## Part one :The Use to provide ideas

The GPT is used to provide us with ideas, mainly focusing on the use of the GPT to provide us with a choice of solutions and formulas when faced with a particular problem, and it is up to us to decide which method to use and to judge its reasonableness.

ChatGPT 3.5

**Query1**: If I need to verify that a series of data is random, what approach can I take?
**Output**: To verify whether a series of data is random, you can employ statistical tests designed for randomness. One commonly used test is the runs test, which checks whether the sequence of data alternates between two or more states in a random manner. Other tests include:

- Chi-square test: This test compares the observed frequency of events with the expected frequency in different categories. If the observed frequencies match the expected frequencies, the data may be considered random.

- Kolmogorov-Smirnov test: This test compares the cumulative distribution function of the data with a reference distribution (e.g., uniform distribution). If the two distributions are similar, the data may be considered random.

- Autocorrelation test: This test examines whether there is any correlation between the values of a variable and its lagged values. In a random sequence, there should be no significant autocorrelation.

- Serial correlation test: Similar to autocorrelation, this test checks for correlation between consecutive observations in the sequence. A lack of serial correlation suggests randomness.

- Frequency test: This test examines the frequency of occurrence of certain patterns or digits within the data. If the frequencies are consistent with what would be expected in a random sequence, the data may be considered random.

When conducting these tests, it's important to set a significance level (e.g., 0.05) to determine whether the results are statistically significant. Additionally, the choice of test may depend on the nature of your data and the specific characteristics you want to test for randomness.
    I hope this helps!

**Query2**:What charts can I use if I want to visualize the correlation of a series of data?
**Output**:When visualizing the correlation of a series of data, you have several options depending on the nature of your data and the level of detail you want to convey. Here are some common types of charts for visualizing correlations:

1. Scatter plot: This is the most straightforward way to visualize the relationship between two variables. Each data point is plotted on a graph, with one variable on the x-axis and the other on the y-axis. The pattern of the points can give you a sense of the correlation between the variables.

2. Correlation matrix heatmap: If you have multiple variables and want to visualize the correlations between all pairs of variables, a correlation matrix heatmap is useful. In this chart, each cell in the matrix represents the correlation coefficient between two variables, and the colors or intensity of the cells indicate the strength and direction of the correlation.

3. Line chart: If you're interested in how the correlation between two variables changes over time or another continuous variable, you can use a line chart. Plot one variable on the x-axis and the other on the y-axis, and connect the data points with lines.

4. Bubble chart: Similar to a scatter plot, a bubble chart adds a third variable by varying the size or color of the points based on a third variable. This can help visualize correlations between three variables simultaneously.

5. Correlation network graph: For more complex datasets with many variables, a correlation network graph can be helpful. In this type of chart, each variable is represented as a node, and the strength of the correlation between variables is represented by the thickness or color of the edges connecting the nodes.

Choose the chart type that best suits your data and the insights you want to convey. Experimenting with different visualizations can help you better understand the relationships in your data.

**Query3**:Are you able to summarize my recent use of you as the subject
**Output**: Of course. Recently, you've been primarily focused on exploring LaTeX, asking about various aspects such as inserting images, setting fonts, and modifying page layouts. You've also inquired about topics related to artificial intelligence, data analysis, and statistics. Additionally, you've discussed writing methods for reports on AI usage and some other technical subjects. Overall, you've been delving into and learning about LaTeX and related topics while also showing interest in artificial intelligence and data analysis fields.

...

## Part Two :Write Codes

The primary purpose of utilizing GitHub Copilot is to streamline the resolution of straightforward yet laborious tasks, particularly those involving graphing and data processing.

GitHub Copilot (Nov 5, 2024 version)
**Query1**: I want to visualise the results of the lasso regression.
**Output**: To visualise the results of a Lasso regression, you can create a graph containing the original target values and the predicted target values. Here is the Python code for how to do this using matplotlib:

```
import matplotlib.pyplot as plt

plt.figure(figsize=(10, 6))
plt.plot(y_test.values, label='Actual')
plt.plot(y_pred, label='Predicted')
plt.title('Lasso Regression Results')
plt.legend()
plt.show()
```

In this code, we first create a new graph and then plot the original and predicted target values in the graph. We use the label parameter to add a label to each data series so that a legend is displayed in the upper right corner of the plot. Finally, we use plt.show() to display the plot.

ChatGPT 3.5

**Query2**: I have an array that represents the win rate change graph, and I would like to visualize it as a line chart. How can I do that?

**Output**: To visualize the win rate changes as a line chart in Python, you can use the matplotlib library. Here's an example code snippet:

```python
import matplotlib.pyplot as plt

win_rate_changes = [...]

x_values = list(range(len(win_rate_changes)))

plt.plot(x_values, win_rate_changes, marker='o', linestyle='-', color='b')
plt.xlabel('Data Points')
plt.ylabel('Win Rate')
plt.title('Win Rate Changes Over Time')
plt.show()
```

This code uses matplotlib.pyplot.plot to create a line chart. You can customize the appearance of the chart, such as choosing markers, linestyle, and color according to your preferences.

...

## Part Three : Translation and Revision

We use ChatGPT for translation work, mainly to correct errors and to make the terminology more in line with the norms of scientific papers.

ChatGPT3.5

**Query1**: Please help me revise the language in this passage to make it more coherent and adhere to the language conventions of academic writing: As statistical analysis methods have improved, sports analysis of athletes has become increasingly detailed and scientific. However, there remains a paucity of research on momentum in tennis, compounded by the lack of a definitive and universally accepted definition of momentum. Current interpretations are often based on consistent yet abstract intuition, leaving considerable ambiguity about the extent to which momentum actually affects the game.

**Output**: Certainly, here's a revised version: With the improvement of statistical analysis methods, sports analytics for athletes has become increasingly detailed and scientific. Indeed, there remains a paucity of research on momentum in tennis, compounded by the absence of a definitive and universally accepted definition of momentum. Current interpretations are often grounded in consistent yet abstract intuition, leaving considerable ambiguity regarding the extent to which momentum truly impacts the game.

This is just an example; in fact, we have essentially polished the entire text using this approach.

...