

COURSE TITLE: DATA SCIENCES

COURSE #: CS 3072

## Project Report

### Team Members

Name	ID	Section
Maram Alhusami	S21207443	2
Amani Albarazi	S19206062	2
Bushra Alshehri	S21107272	2
Sara Imran	S22107878	2

### Project Description

This project analyzes the Women's Tennis Association (WTA) match data to develop predictive models for match outcomes and player rankings. The datasets include:

- Match Statistics Dataset: Contains detailed information about match results, player statistics, and performance metrics.
- Player Ranking Dataset: Provides historical player rankings and performance over time.
- Tournament Dataset: Details about different tournaments, including surface types and locations.
- Player Attributes Dataset: Includes demographic data such as age, handedness, and activity levels.

The research question is: How can the Factors of the Player's Performance in WTA Matches can be Utilized to Predict Match Outcomes and Player Rankings?

This project aims to predict two variables: winner rank and match outcome. The project objectives to achieve our goal are:

1. Understand the two data sets by applying basic Pandas functions.
2. Using data visualizations for understanding the patterns.
3. Cleaning the two data sets individually.
4. For better prediction, we merged two datasets, `wta_matches_qual_itf_2023.csv`, and `wta_matches_qual_itf_2024.csv`.
5. Continued cleaning the merged data sets.
6. We applied ML algorithms to predict winner rank and match outcome, comparing the accuracy of the algorithms we used with previous research papers.

**Unit of Observation**

The unit of observation is each match played in the Women's Tennis Association (WTA) for the years 2023 and 2024, encompassing detailed player, match, and tournament-level information.

**Outcome Variables**

The analysis focuses on two primary outcome variables:

**1. Winner Rank**

Represents the ranking of the player who won the match, providing insight into player performance and match competitiveness.

- **Conceptual Relation to the Research Question:** Winner rank is a key indicator of performance, allowing the study to predict the relative competitiveness of players.
- **Summary Statistics:**
  - Mean Rank: 586.35
  - Standard Deviation: 348.58
  - Range: 22 to 1373

**2. Match Outcome**

Since the dataset did not include a direct match outcome variable, we created **match\_outcome** as a binary indicator. It is derived by checking if the winning player's rank is better (lower) than the losing player's rank.

- **Conceptual Relation to the Research Question:** Match outcomes serve as the foundation for understanding player performance and building predictive models.
- The binary variable indicates whether a match was won (1) or lost (0).

**Predictor Variables**

Predictor variables used to model the outcome variables include:

- Player Age
- Tournament Surface Type (e.g., clay, grass, hard)
- Handedness (right or left-handed)
- Serve Effectiveness
- Match Duration
- Opponent Ranking
- Win-Loss Ratio
- Break Point Conversion Rate
- Player Activity Level (frequency of match participation)
- Country Representation

These variables are measured directly from the dataset or derived from available features. They come from:

- **Match Statistics Dataset:** Provides comprehensive match-level data.
- **Player Ranking Dataset:** Offers historical ranking information for WTA players.

- **Tournament Dataset:** Captures tournament-level attributes such as surface type.
- **Player Attributes Dataset:** Contains demographic and activity-level data.

## Potential Data Issues

1. **Missing Data**
  - **Issue:** Some matches or player records need to be completed ranking or match outcome data.
  - **Mitigation:** Missing values in predictors were handled using imputation techniques such as filling with zeros or mean values. Rows with missing target values were dropped.
2. **Class Imbalance**
  - **Issue:** The match outcome data may be imbalanced, as certain players win more frequently than others.
  - **Mitigation:** Stratified sampling was used in train-test splits to maintain the distribution of outcomes.
3. **Feature Correlation**
  - **Issue:** Some predictor variables may be highly correlated, potentially affecting model interpretability.
  - **Mitigation:** Exploratory Data Analysis (EDA) was performed to identify and address multicollinearity issues.
4. **Data Scale and Variance**
  - **Issue:** Predictors like player rank and match duration vary widely in scale.
  - **Mitigation:** Normalization techniques (e.g., MinMaxScaler) were applied to standardize features where necessary.

## Data Distribution

To understand the dataset better, the distribution of outcome variables and predictors was visualized:

1. **Outcome Variables:**
  - **Winner Rank:** Exhibits a right-skewed distribution, with most winners ranked in the mid-range and a few outliers with extremely high or low rankings.
  - **Match Outcome:** Binary distribution, indicating whether a match was won or lost.
2. **Predictors:**
  - Player age follows a near-normal distribution, while attributes like match duration and serve effectiveness show skewness. Tournament surface types are distributed evenly across categories.

Graphs and tables were used to visualize and summarize these distributions, providing insights into data trends and variability.

## Applying Basic Pandas Functions

First, we imported the two data sets as Pandas.

```
: import pandas as pd

wta_matches_2023= pd.read_csv('wta_matches_qual_itf_2024.csv', index_col=0)
wta_matches_2024 = pd.read_csv('wta_matches_qual_itf_2024.csv', index_col=0)
```



### 1. wta\_matches\_qual\_itf\_2023

- Using head() to see some samples of the data set

	wta_matches_2023.head(10)										
tourney_id	tourney_name	surface	draw_size	tourney_level	tourney_date	match_num	winner_id	winner_seed	winner_entry	winner_name	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	102	220578	NaN	Q	Julia Konishi	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	103	223367	NaN	NaN	Maria Fernanda Navarro	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	106	206037	NaN	NaN	Daniela Seguel	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	107	216075	NaN	Q	Jessica Bertoldo	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	110	214605	NaN	Q	Agustina Chilpac	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	111	221657	NaN	Q	Camilla Zanolini	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	114	223336	NaN	NaN	Weronika Baszak	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	115	214713	NaN	Q	Camila Romero	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	118	211646	NaN	NaN	Oana Georgeta Simion	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	119	214238	NaN	NaN	Verena Meliss	...

10 rows × 48 columns

- The shape of the data set is = (14319, 48)

```
wta_matches_2023.shape
```

(14319, 48)

- Using describe() to see the statistical summary of the numerical columns

wta_matches_2023.describe()										
	draw_size	tourney_date	match_num	winner_id	winner_ht	winner_age	loser_id	loser_ht	loser_age	best_of
<b>count</b>	14319.000000	1.431900e+04	14319.000000	14319.000000	1744.000000	12136.000000	14319.000000	1367.000000	10782.000000	14319.0
<b>mean</b>	34.505203	2.024033e+07	185.131643	228164.922550	171.943234	23.615985	232597.988477	171.647403	23.576767	3.0
<b>std</b>	13.174083	1.306872e+02	88.536649	19737.261891	6.547561	4.263115	21713.509548	6.533059	4.400639	0.0
<b>min</b>	32.000000	2.024010e+07	100.000000	201329.000000	152.000000	14.600000	201318.000000	152.000000	14.600000	3.0
<b>25%</b>	32.000000	2.024022e+07	110.000000	214860.000000	168.000000	20.300000	215480.000000	168.000000	20.200000	3.0
<b>50%</b>	32.000000	2.024032e+07	201.000000	221139.000000	171.000000	23.100000	222003.000000	171.000000	23.000000	3.0
<b>75%</b>	32.000000	2.024042e+07	211.000000	239453.000000	176.000000	26.300000	259941.000000	175.000000	26.200000	3.0
<b>max</b>	128.000000	2.024052e+07	601.000000	267094.000000	186.000000	48.500000	267103.000000	186.000000	51.900000	3.0

8 rows × 33 columns

## Observations:

- **draw\_size**: This column represents the size of the tournament draws, ranging from 32 to 128. The mean value is about 34, with most tournaments having smaller draw sizes.
- **tourney\_date**: This represents the date in some numeric format (likely a Unix timestamp). The mean is around 20 million.
- **winner\_ht** and **loser\_ht**: These represent the heights of the winners and losers. The mean winner height is 171.9 cm, and the mean loser height is 171.6 cm, indicating that there is little difference between the winner and loser in terms of height.
- **winner\_age** and **loser\_age**: These represent the ages of the winner and loser. The mean winner age is 23.6 years, while the mean loser age is almost identical at 23.6 years as well.
- **best\_of**: This column seems to indicate the format of the match (e.g., 3-set match). It has a constant value of 3, meaning all the matches in this dataset are best of 3.
- **l\_1stIn**, **l\_1stWon**, **l\_2ndWon**, etc.: These columns represent statistics for the loser (such as first serves in, first serves won, second serves won, etc.). For example, the median number of first serves won by losers (**l\_1stWon**) is 42.

## 2. wta\_matches\_qual\_itf\_2024

- `head()`

wta_matches_2024.head()											
tourney_id	tourney_name	surface	draw_size	tourney_level	tourney_date	match_num	winner_id	winner_seed	winner_entry	winner_name	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	102	220578	NaN	Q	Julia Konishi Camargo Silva	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	103	223367	NaN	NaN	Maria Fernanda Navarro	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	106	206037	NaN	NaN	Daniela Seguel	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	107	216075	NaN	Q	Jessica Bertoldo	...
2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	110	214605	NaN	Q	Agustina Chlpac	...

5 rows × 48 columns

- The shape of the data set is = (14319, 48)

wta_matches_2024.shape											
(14319, 48)											

- describe()

wta_matches_2024.describe()											
	draw_size	tourney_date	match_num	winner_id	winner_ht	winner_age	loser_id	loser_ht	loser_age	best_of	
count	14319.000000	1.431900e+04	14319.000000	14319.000000	1744.000000	12136.000000	14319.000000	1367.000000	10782.000000	14319.0	
mean	34.505203	2.024033e+07	185.131643	228164.922550	171.943234	23.615985	232597.988477	171.647403	23.576767	3.0	
std	13.174083	1.306872e+02	88.536649	19737.261891	6.547561	4.263115	21713.509548	6.533059	4.400639	0.0	
min	32.000000	2.024010e+07	100.000000	201329.000000	152.000000	14.600000	201318.000000	152.000000	14.600000	3.0	
25%	32.000000	2.024022e+07	110.000000	214860.000000	168.000000	20.300000	215480.000000	168.000000	20.200000	3.0	
50%	32.000000	2.024032e+07	201.000000	221139.000000	171.000000	23.100000	222003.000000	171.000000	23.000000	3.0	
75%	32.000000	2.024042e+07	211.000000	239453.000000	176.000000	26.300000	259941.000000	175.000000	26.200000	3.0	
max	128.000000	2.024052e+07	601.000000	267094.000000	186.000000	48.500000	267103.000000	186.000000	51.900000	3.0	

8 rows × 33 columns

## Observations:

- draw\_size: The mean is 34.5, and the minimum and maximum values are 32 and 128, respectively. This suggests that the majority of the tournaments in the dataset have relatively small draw sizes.
- tourney\_date: This column is likely representing a timestamp. The mean is around 20 million, likely corresponding to dates in the 2020s.
- winner\_ht and loser\_ht: The heights of the players are similar to those in the 2023 dataset, with the mean height for the winner at 171.94 cm and for the loser at 171.65 cm.
- winner\_age and loser\_age: The mean winner age is 23.6 years, and the mean loser age is 23.6 years, indicating that the dataset features matches between players of similar ages.

- best\_of: All values in this column are 3, indicating that the tournaments in this dataset use the best-of-three set format.
- l\_1stIn, l\_1stWon, l\_2ndWon, etc.: These columns represent the first serve statistics for the loser in the match. The l\_1stIn and l\_1stWon columns have values ranging from 0 to 100, suggesting different levels of success for each loser in the dataset.

### 3. wta\_matches\_qual\_itf\_2023 vs wta\_matches\_qual\_itf\_2024

- The statistics are very similar between the 2023 and 2024 datasets, suggesting that the two datasets are quite comparable in terms of the matches and players involved. However, there are some differences in the number of missing values for certain columns, such as `winner_ht` and `loser_ht`, where some values are missing.

## Data Visualisation

### 1. wta\_matches\_qual\_itf\_2023

**A scatter plot of winner age vs loser age**

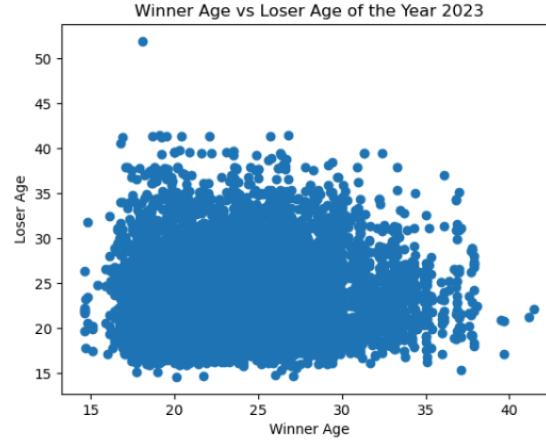
```
# A scatterplot of winner_age vs Loser_Age

x = wta_matches_2023.winner_age
y = wta_matches_2023.loser_age

plt.plot(x, y, 'o') # The 'o' argument specifies a scatterplot
# Labelling the horizontal X-axis and vertical Y-axis
plt.xlabel('Winner Age')
plt.ylabel('Loser Age')

plt.title('Winner Age vs Loser Age of the Year 2023')
```

Text(0.5, 1.0, 'Winner Age vs Loser Age of the Year 2023')



### Observation:

- Most data points are clustered between the ages of 20 to 35 for both winners and losers, indicating that most players in these matches fall within this age range.

- There are a few outliers where the winner or loser age is above 40 or below 20. This suggests the presence of either veteran players or younger participants.

Relationship Between Winner and Loser Ages:

- The scatter plot appears somewhat spread out, suggesting no strong linear relationship between the winner's age and the loser's age. Players of all age ranges seem to compete against each other.

Balanced Competition:

- The data points are not biased toward either axis, meaning winners and losers are not concentrated in one specific age group but spread across the range of competing players.

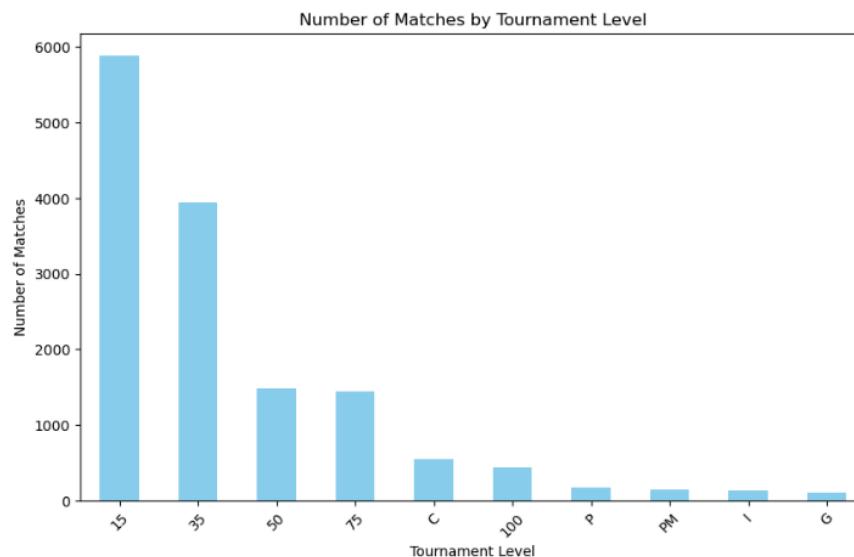
### Number of Matches by Tournament Level

```
# Number of Matches by Tournament Level

# Count the number of matches by tournament level
tourney_level_counts = wta_matches_2023['tourney_level'].value_counts()

plt.figure(figsize=(10, 6))
tourney_level_counts.plot(kind='bar', color='skyblue')

plt.title('Number of Matches by Tournament Level')
plt.xlabel('Tournament Level')
plt.ylabel('Number of Matches')
plt.xticks(rotation=45)
plt.show()
```

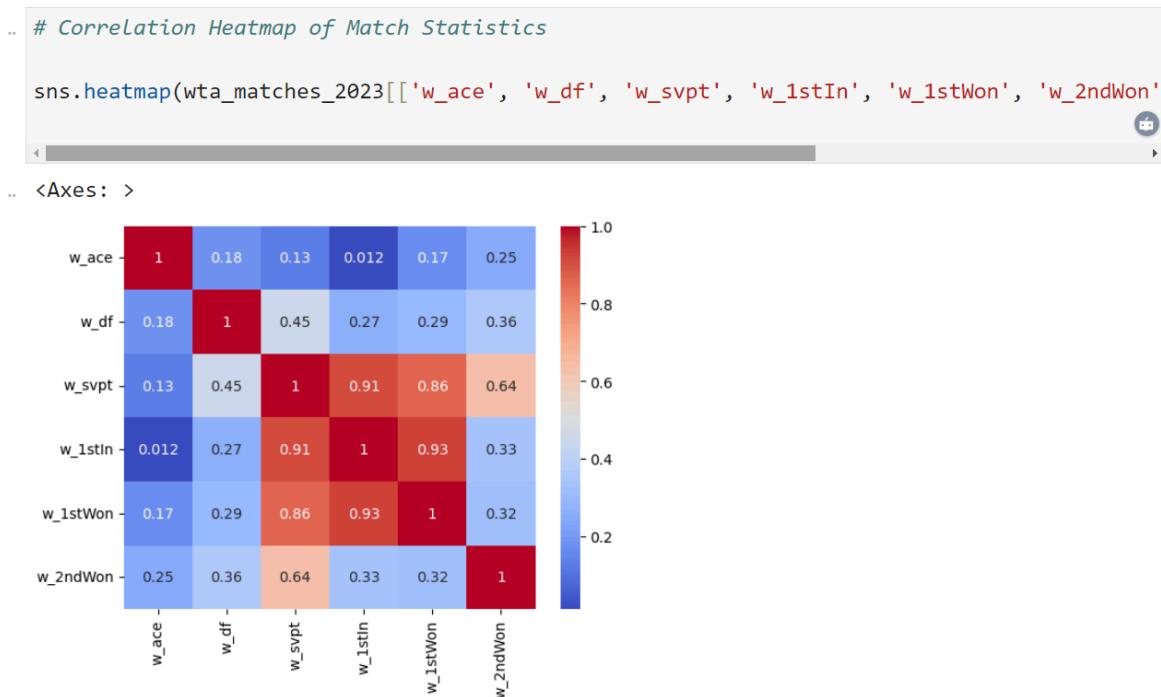


#### Observation:

- Level "15" has the highest number of matches, indicating it is the most frequent tournament level.
- Levels "35" and "50" have progressively fewer matches, showing a decline in match frequency.

- Higher-level tournaments such as "100," "PM," and "G" have significantly fewer matches, reflecting their exclusivity or elite status.
- Mid-tier levels like "C" and "P" display moderate match counts, indicating their stable position in the tournament hierarchy.
- Overall, the distribution suggests a higher frequency of lower-tier tournaments compared to elite or rare ones.

### Correlation Heatmap of Match Statistics



### Observations:

#### Strong Positive Correlations:

- w\_1stIn and w\_svpt show a very high correlation (**0.91**), meaning that when a player's total serve points (w\_svpt) increase, the number of first serves in (w\_1stIn) also increases.
- w\_1stWon and w\_1stIn have a strong correlation (**0.93**), indicating that successfully landing first serves contributes significantly to winning those points.

#### Moderate Correlations:

- w\_2ndWon and w\_svpt show a moderate correlation (**0.64**), suggesting that a higher number of serve points leads to better performance on second serves.
- w\_df (double faults) moderately correlates with w\_svpt (**0.45**), indicating that higher serve activity slightly increases the occurrence of double faults.

#### Weak Correlations:

- **w\_ace** (aces) has weak correlations with most other variables. For example, **w\_ace** correlates at **0.25** with **w\_2ndWon** and only **0.17** with **w\_1stWon**. This suggests that hitting aces is somewhat independent of winning first or second serve points.
- The correlation between **w\_ace** and **w\_df** is low (**0.18**), indicating little relationship between hitting aces and committing double faults.

### Insights:

- Winning on the **first serve** (**w\_1stWon**) is the most influential factor, as it correlates highly with other serve-related metrics.
- Players' ability to consistently land first serves (**w\_1stIn**) appears to impact overall performance positively.
- **Double faults** (**w\_df**) show weaker links to success metrics, implying that isolated errors may not drastically affect outcomes in the broader statistics.

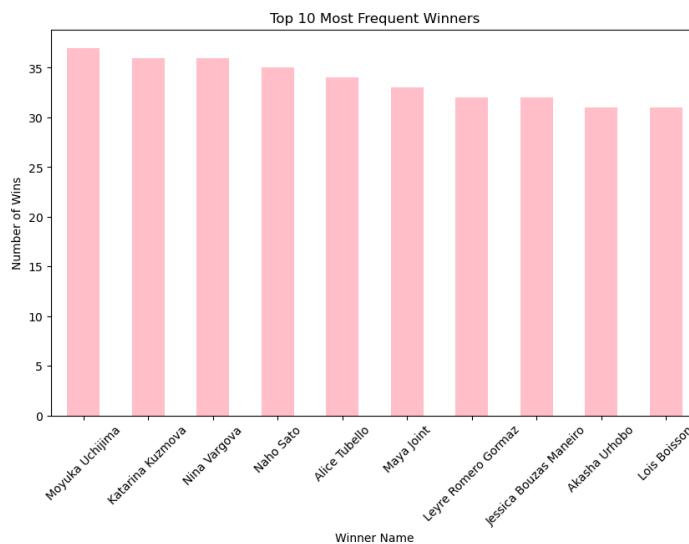
### Top 10 Most Frequent Winners

```
# Top 10 Most Frequent Winners

# Count top 10 winners by frequency
top_10_winners = wta_matches_2023['winner_name'].value_counts().head(10)

# Bar plot using Matplotlib
plt.figure(figsize=(10, 6))
top_10_winners.plot(kind='bar', color='pink')

plt.title('Top 10 Most Frequent Winners')
plt.xlabel('Winner Name')
plt.ylabel('Number of Wins')
plt.xticks(rotation=45)
plt.show()
```



### Observation:

The bar chart shows the top 10 most frequent winners from the dataset. The winner with the highest number of wins is significantly ahead of others, indicating dominance in the matches.

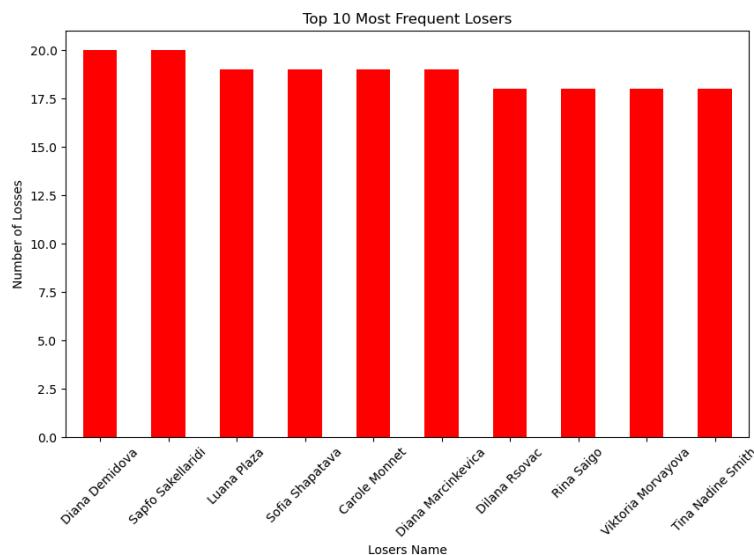
### Top 10 Most Frequent Losers

```
# Top 10 Most Frequent Losers

# Count top 10 Losers by frequency
top_10_losers = wta_matches_2023['loser_name'].value_counts().head(10)

# Bar plot using Matplotlib
plt.figure(figsize=(10, 6))
top_10_losers.plot(kind='bar', color='red')

plt.title('Top 10 Most Frequent Losers')
plt.xlabel('Losers Name')
plt.ylabel('Number of Losses')
plt.xticks(rotation=45)
plt.show()
```



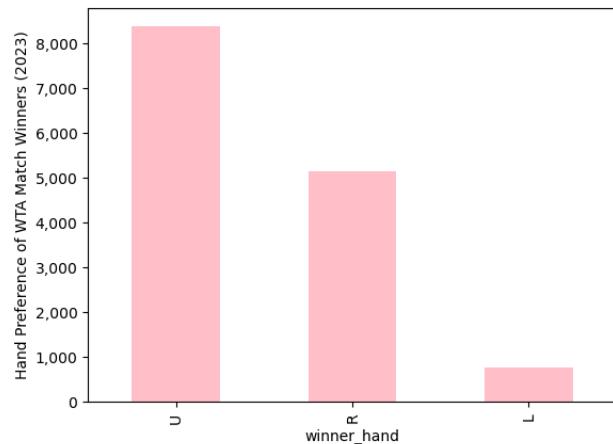
#### Observation:

The chart highlights the players who most frequently lost matches. This could reflect participation frequency or a significant gap in performance compared to other players.

### Hand Preference of WTA Match Winners (2023)

```
ax = wta_matches_2023['winner_hand'].value_counts().plot.bar(
    ylabel = 'Hand Preference of WTA Match Winners (2023)', color='pink')

ax.yaxis.set_major_formatter('{x:, .0f}')
```

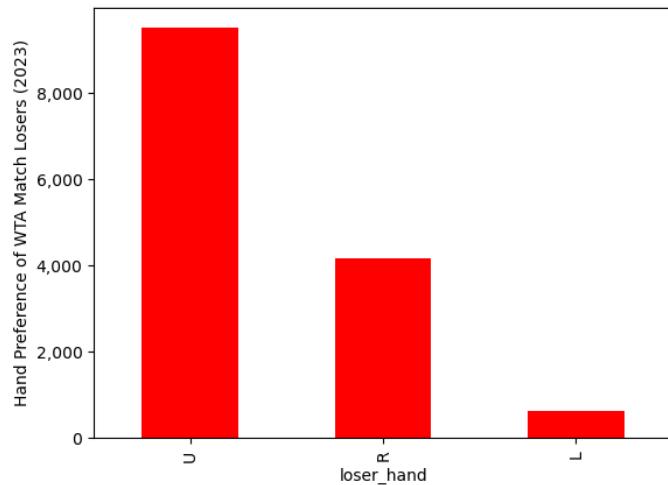

**Observation:**

Most match winners are right-handed, as seen from the dominance of the corresponding bar. This reflects the general prevalence of right-handedness in the dataset.

### Hand Preference of WTA Match Losers (2023)

```
ax = wta_matches_2023['loser_hand'].value_counts().plot.bar(
    ylabel = 'Hand Preference of WTA Match Losers (2023)', color='red')

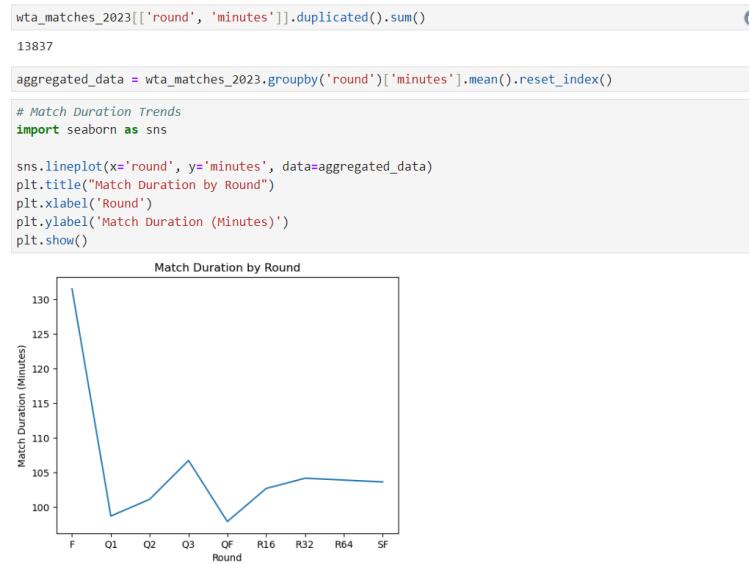
ax.yaxis.set_major_formatter('{x:, .0f}')
```


**Observation:**

A similar trend to winners is observed for losers, with right-handed players dominating. This may indicate that handedness does not significantly affect match outcomes.

## Match Duration Trends

To show how the match duration changes with each round of the tournament, usually matches in the finals tend to be longer and the graph shows that too.

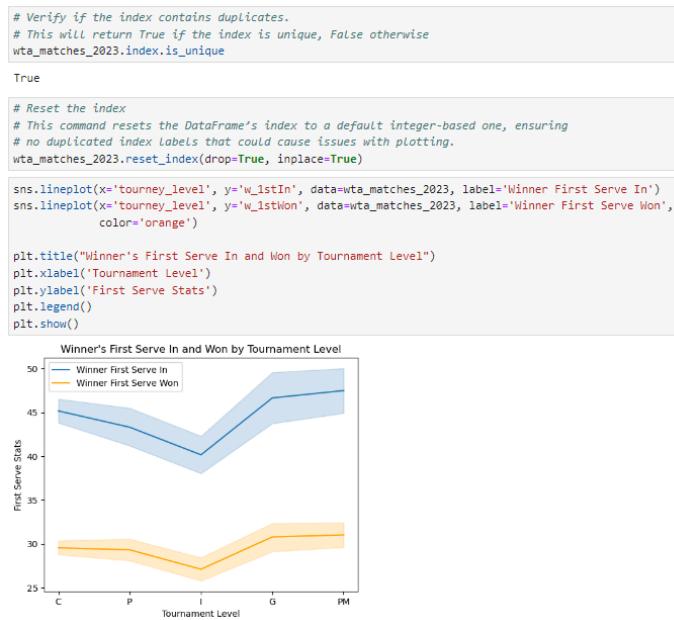


### Observation:

The line graph shows that match duration generally increases as the tournament progresses, peaking in the finals. This aligns with expectations, as higher rounds typically feature more competitive matches.

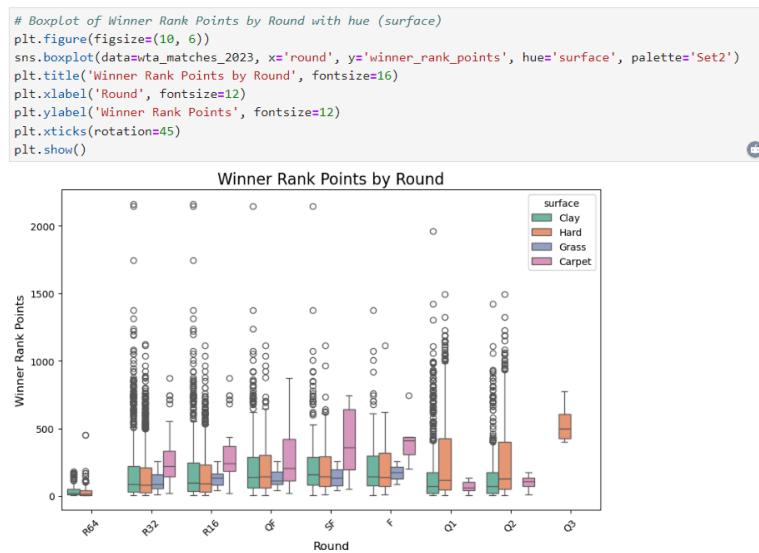
## Winner's First Serve-In and Won by Tournament Level

Analyzing how the winner's first serve performance (both in terms of landing the serve and winning the point) varies by the level of the tournament.

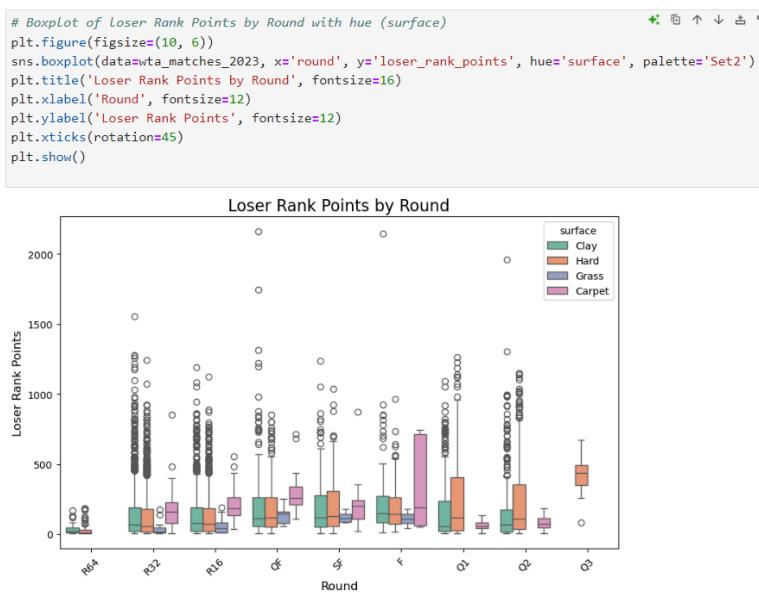


**Observation:**

- Higher tournament levels exhibit consistent performance in first-serve statistics, reflecting skilled players.
- Lower-level tournaments show slightly more variability in first-serve performance.
- 

**Boxplot of Winner Rank Points by Round****Observation:**

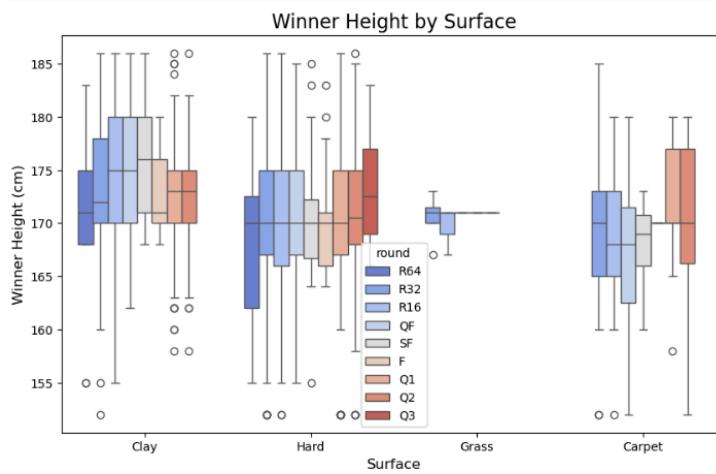
Winner rank points vary widely in early rounds, with a more concentrated distribution in later rounds. This reflects the increasing skill level of players as the tournament progresses.

**Boxplot of Loser Rank Points by Round****Observation:**

The loser rank points are more dispersed in early rounds, showing a broad range of player skills. In later rounds, losing players tend to have higher rank points, reflecting tougher competition.

### Boxplot of Winner Height by Surface

```
# Boxplot of Winner Height by Surface with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2023, x='surface', y='winner_ht', hue='round', palette='coolwarm')
plt.title('Winner Height by Surface', fontsize=16)
plt.xlabel('Surface', fontsize=12)
plt.ylabel('Winner Height (cm)', fontsize=12)
plt.show()
```

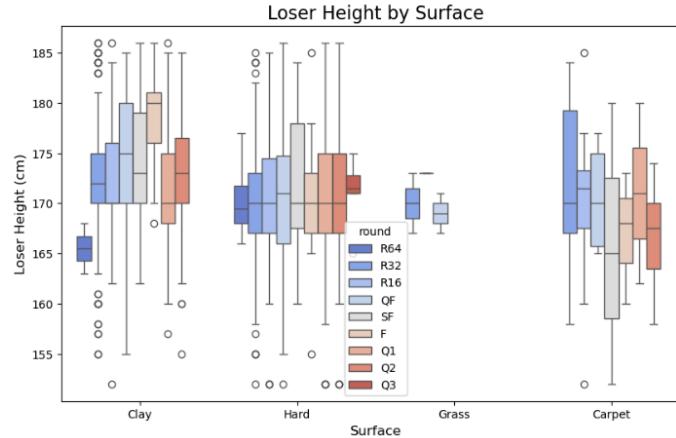


#### Observation:

Winner height varies across different surfaces, with some surfaces showing a higher average height for winners. This may indicate that surface type influences player performance.

### Boxplot of Loser Height by Surface

```
# Boxplot of Loser Height by Surface with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2023, x='surface', y='loser_ht', hue='round', palette='coolwarm')
plt.title('Loser Height by Surface', fontsize=16)
plt.xlabel('Surface', fontsize=12)
plt.ylabel('Loser Height (cm)', fontsize=12)
plt.show()
```

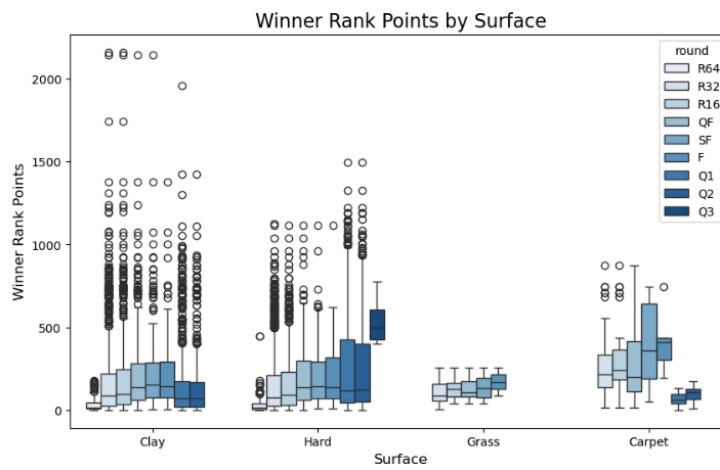


**Observation:**

Loser height shows similar variability across surfaces, with trends mirroring those of the winners, suggesting that surface preference may play a role in match outcomes.

**Boxplot of Winner Rank Points by Surface**

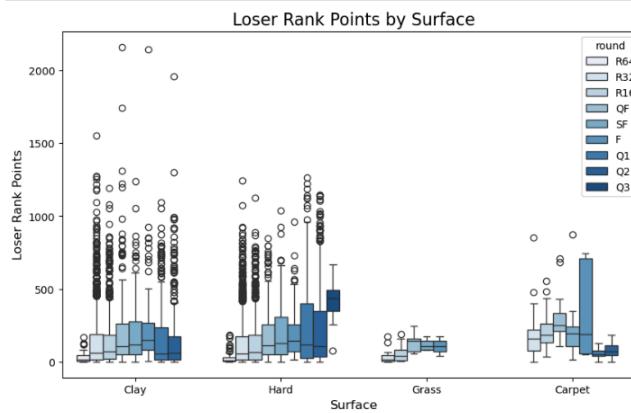
```
# Boxplot of Winner Rank Points by Surface with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2023, x='surface', y='winner_rank_points', hue='round', palette='Blues')
plt.title('Winner Rank Points by Surface', fontsize=16)
plt.xlabel('Surface', fontsize=12)
plt.ylabel('Winner Rank Points', fontsize=12)
plt.show()
```

**Observation:**

Winner rank points are higher on certain surfaces, indicating that surface type may favor players with specific playing styles or strategies.

**Boxplot of Loser Rank Points by Surface**

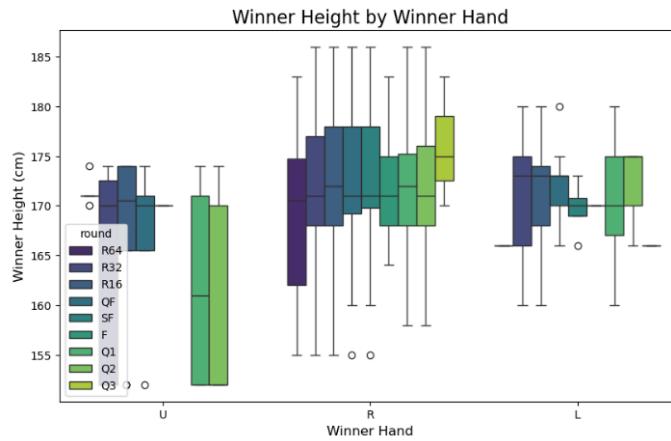
```
# Boxplot of Loser Rank Points by Surface with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2023, x='surface', y='loser_rank_points', hue='round', palette='Blues')
plt.title('Loser Rank Points by Surface', fontsize=16)
plt.xlabel('Surface', fontsize=12)
plt.ylabel('Loser Rank Points', fontsize=12)
plt.show()
```

**Observation:**

Loser rank points vary by surface, with some surfaces showing lower rank points for losing players, possibly due to mismatched skill levels or surface preferences.

### Boxplot of Winner Height by Winner Hand (Left, Right, or Unknown)

```
# Boxplot of Winner Height by Winner Hand with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2023, x='winner_hand', y='winner_ht', hue='round', palette='viridis')
plt.title('Winner Height by Winner Hand', fontsize=16)
plt.xlabel('Winner Hand', fontsize=12)
plt.ylabel('Winner Height (cm)', fontsize=12)
plt.show()
```

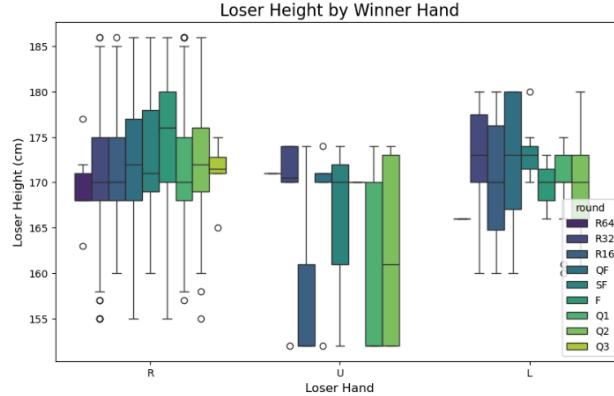


#### Observation:

- The distribution of loser rank points varies significantly across surfaces.
- Clay and hard courts show a wider distribution compared to grass and carpet surfaces.
- Most outliers appear in clay and hard surfaces, indicating more variability in performance.

### Boxplot of Loser Height by Winner Hand (Left, Right, or Unknown)

```
# Boxplot of Loser Height by Loser Hand with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2023, x='loser_hand', y='loser_ht', hue='round', palette='viridis')
plt.title('Loser Height by Loser Hand', fontsize=16)
plt.xlabel('Loser Hand', fontsize=12)
plt.ylabel('Loser Height (cm)', fontsize=12)
plt.show()
```



#### Observation:

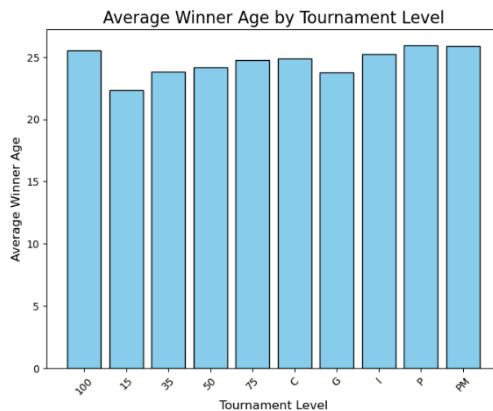
- Winner's height seems relatively consistent regardless of the hand preference.
- The "Unknown" category exhibits a smaller range of heights compared to left- and right-handed players.
- Outliers are minimal across all categories, indicating a focused range of winner heights.
- 

#### 4. wta\_matches\_qual\_itf\_2024

##### Bar Plot of Average Winner age by Tournament Level

```
# Data Visualization for wta_matches_2024
# Bar Plot of Average Winner Age by Tournament Level
# Calculate average winner age for each tournament level
avg_winner_age = wta_matches_2024.groupby('tourney_level')['winner_age'].mean()

# Plot
plt.figure(figsize=(8, 6))
plt.bar(avg_winner_age.index, avg_winner_age.values, color='skyblue', edgecolor='black')
plt.title('Average Winner Age by Tournament Level', fontsize=16)
plt.xlabel('Tournament Level', fontsize=12)
plt.ylabel('Average Winner Age', fontsize=12)
plt.xticks(rotation=45)
plt.show()
```



##### Observation:

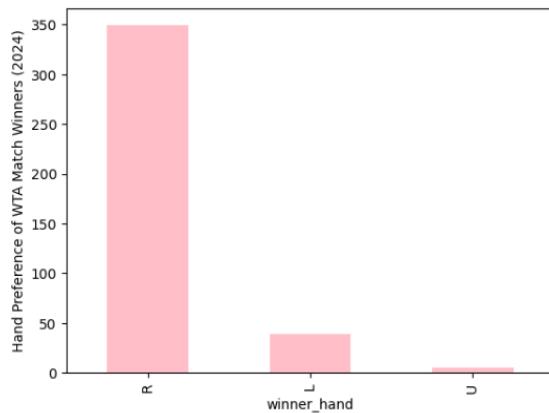
- Average Age Consistency: The average winner age appears to be relatively consistent across the different tournament levels, with all bars showing similar heights.
- No Significant Outliers: None of the tournament levels shows a drastically lower or higher average winner age, indicating a stable range of ages for winners across tournament levels.

## Hand Preference of WTA Match Winners (2024)

```
# Hand Preference of WTA Match Winners (2024)

ax = wta_matches_2024['winner_hand'].value_counts().plot.bar(
    ylabel = 'Hand Preference of WTA Match Winners (2024)', color='pink')

ax.yaxis.set_major_formatter('{x:,.0f}')
```



### Observation:

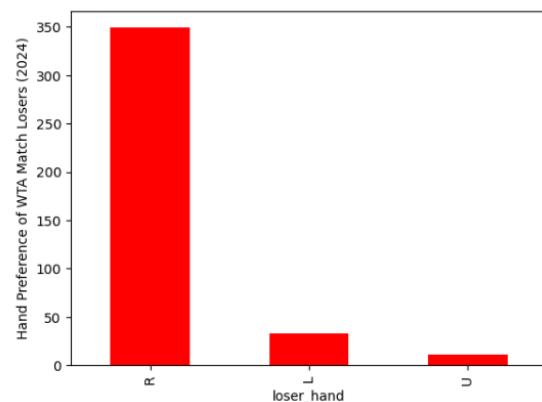
- The majority of match winners are right-handed, dominating the bar plot.
- Left-handed players account for a smaller proportion, with "Unknown" being the least frequent category.

## Hand Preference of WTA Match Losers (2024)

```
# Hand Preference of WTA Match Losers (2024)

ax = wta_matches_2024['loser_hand'].value_counts().plot.bar(
    ylabel = 'Hand Preference of WTA Match Losers (2024)', color='red')

ax.yaxis.set_major_formatter('{x:,.0f}')
```



### Observation:

- Similar to winners, right-handed players are the majority among losers.

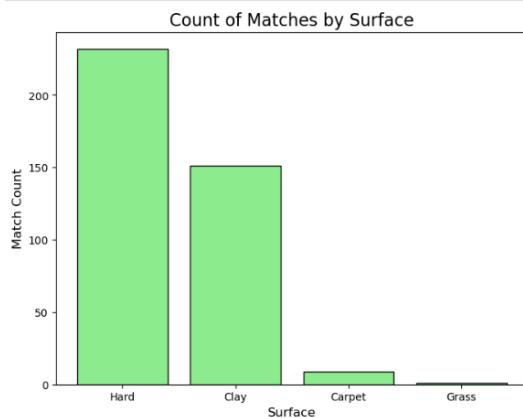
- The proportion of left-handed losers is slightly higher than their representation in winners, suggesting a marginal performance difference.

### Count of Matches by Surface

```
# Count of Matches by Surface

# Count the number of matches for each surface
match_counts = wta_matches_2024['surface'].value_counts()

# Plot
plt.figure(figsize=(8, 6))
plt.bar(match_counts.index, match_counts.values, color='lightgreen', edgecolor='black' )
plt.title('Count of Matches by Surface', fontsize=16)
plt.xlabel('Surface', fontsize=12)
plt.ylabel('Match Count', fontsize=12)
plt.show()
```



#### Observation:

- Hard courts host the majority of matches, making them the most popular surface.
- Clay courts rank second, reflecting their continued relevance in the tournament landscape.
- Grass and carpet surfaces are the least utilized, likely due to limited availability or player preference.

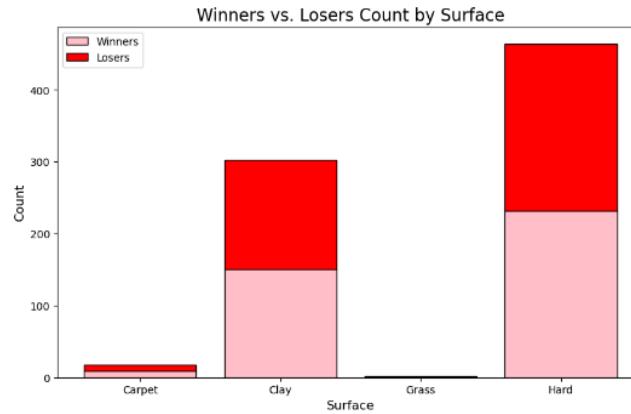
### Stacked Bar Plot: Winner vs. Loser Count by Surface

```
# Stacked Bar Plot: Winner vs. Loser Count by Surface

# Group by surface and calculate counts for winners and losers
surface_group = wta_matches_2024.groupby('surface')[['winner_id', 'loser_id']].count()

# Plot
plt.figure(figsize=(10, 6))
plt.bar(surface_group.index, surface_group['winner_id'], color='pink', label='Winners',
        edgecolor='black')
plt.bar(surface_group.index, surface_group['loser_id'], bottom=surface_group['winner_id'],
        color='red', label='Losers', edgecolor='black')

plt.title('Winners vs. Losers Count by Surface', fontsize=16)
plt.xlabel('Surface', fontsize=12)
plt.ylabel('Count', fontsize=12)
plt.legend()
plt.show()
```



#### Observation:

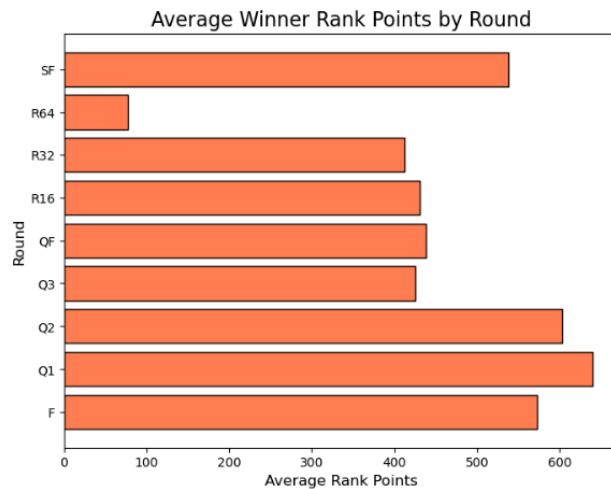
- Clay Surface Dominance:
  - The clay surface has the highest total count of matches (winners + losers), indicating that more matches are played on clay compared to other surfaces.
  - This could suggest that clay courts are widely used or more popular in the dataset.
- Balanced Winner-Loser Distribution:
  - On all surfaces, the blue (winners) and red (losers) sections are nearly equal in height, meaning the number of winners and losers is fairly balanced across all surfaces. This is expected since every match has one winner and one loser.
- Lower Match Counts on Carpet and Grass:
  - Carpet and grass surfaces show significantly fewer matches compared to clay and hard surfaces. This could suggest:
    - Fewer tournaments are played on these surfaces.
    - These surfaces are less common in professional tennis.
- Hard Surface Matches:
  - Hard courts have the second-highest total count of matches, indicating that hard courts are also very popular.

## Horizontal Bar Plot of Average Winner Rank Points by Round

```
# Horizontal Bar Plot of Average Winner Rank Points by Round

# Calculate average winner rank points for each round
avg_rank_points = wta_matches_2024.groupby('round')['winner_rank_points'].mean()

# Plot
plt.figure(figsize=(8, 6))
plt.barh(avg_rank_points.index, avg_rank_points.values, color='coral', edgecolor='black')
plt.title('Average Winner Rank Points by Round', fontsize=16)
plt.xlabel('Average Rank Points', fontsize=12)
plt.ylabel('Round', fontsize=12)
plt.show()
```



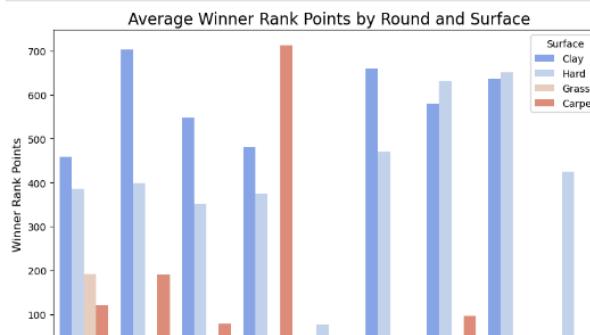
### Observations:

- Quarterfinals (QF) and Semifinals (SF):
  - Higher average rank points are observed in the later stages of tournaments (QF, SF, F), indicating that players reaching these stages generally have higher rankings.
- Qualifying Rounds (Q1, Q2, Q3):
  - The third qualifying round (Q3) has the highest average winner rank points compared to the earlier qualifying rounds (Q1, Q2).
  - This might suggest that higher-ranked players often secure wins in the later stages of the qualifiers.
- First Rounds (R64, R32, R16):
  - The earlier rounds of the main tournament (R64, R32, R16) show relatively lower average rank points compared to the quarterfinals and beyond. This could indicate a wider mix of player rankings in the earlier rounds.
- Finals (F):
  - The average rank points in the finals are substantial but slightly lower than the later qualifying rounds (Q3). This could be due to variability in player rankings reaching the finals.
- General Trend:

- Players with higher rank points tend to dominate in later stages of tournaments, while earlier rounds and qualifying stages show more variation.

### Average Winner Rank Points by Round and Surface

```
# Average Winner Rank Points by Round and Surface
plt.figure(figsize=(10, 6))
sns.barplot(data=wta_matches_2024, x='round', y='winner_rank_points', hue='surface', errorbar=None,
            palette='coolwarm')
plt.title('Average Winner Rank Points by Round and Surface', fontsize=16)
plt.xlabel('Round', fontsize=12)
plt.ylabel('Winner Rank Points', fontsize=12)
plt.xticks(rotation=45)
plt.legend(title='Surface')
plt.show()
```

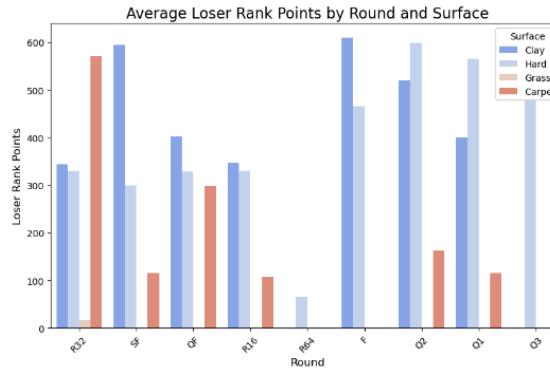


#### Observation:

- The winner rank points vary significantly across different rounds and surfaces.
- For example, higher winner rank points are often seen in the later rounds (e.g., semi-finals or finals), especially on specific surfaces like hard courts.
- The clay surface has more variation compared to other surfaces.

### Average Loser Rank Points by Round and Surface

```
# Average Loser Rank Points by Round and Surface
plt.figure(figsize=(10, 6))
sns.barplot(data=wta_matches_2024, x='round', y='loser_rank_points', hue='surface', errorbar=None,
            palette='coolwarm')
plt.title('Average Loser Rank Points by Round and Surface', fontsize=16)
plt.xlabel('Round', fontsize=12)
plt.ylabel('Loser Rank Points', fontsize=12)
plt.xticks(rotation=45)
plt.legend(title='Surface')
plt.show()
```



#### Observation:

- The loser rank points also vary across rounds, with higher points in early rounds.
- The grass surface shows a relatively lower average loser rank compared to other surfaces in some rounds, indicating a trend where higher-ranked players tend to dominate these matches.
- Hard courts again display consistent performance across rounds.

### Winner Height vs. Loser Height by Surface

This scatter plot visualizes the relationship between the height of the winner (winner\_ht) and the height of the loser (loser\_ht) in tennis matches, while also taking into account the surface type of the court.

```
# Check for Duplicate Indexes
# If the output is greater than 0, it confirms that the dataset has duplicate indexes.
wta_matches_2024.index.duplicated().sum()

14065

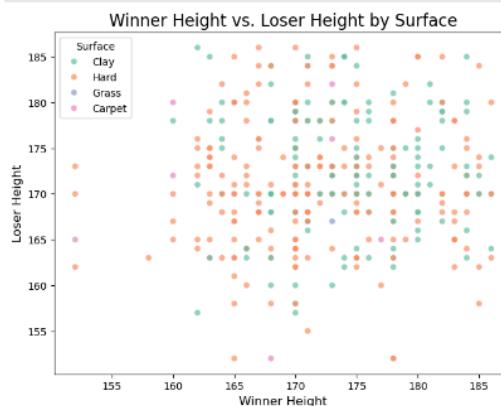
# Reset the Index
wta_matches_2024 = wta_matches_2024.reset_index(drop=True)

# Check if the winner_ht and loser_ht columns contain valid numeric data and no missing values.
wta_matches_2024[['winner_ht', 'loser_ht']].info()

<class 'pandas.core.frame.DataFrame'
RangeIndex: 14319 entries, 0 to 14318
Data columns (total 2 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   winner_ht   1744 non-null    float64
 1   loser_ht    1367 non-null    float64
dtypes: float64(2)
memory usage: 223.9 KB

wta_matches_2024 = wta_matches_2024.dropna(subset=['winner_ht', 'loser_ht'])

# Winner Height vs. Loser Height by Surface
plt.figure(figsize=(8, 6))
sns.scatterplot(
    data=wta_matches_2024,
    x='winner_ht',
    y='loser_ht',
    hue='surface',
    palette='Set2',
    alpha=0.7
)
plt.title('Winner Height vs. Loser Height by Surface', fontsize=16)
plt.xlabel('Winner Height', fontsize=12)
plt.ylabel('Loser Height', fontsize=12)
plt.legend(title='Surface')
plt.show()
```



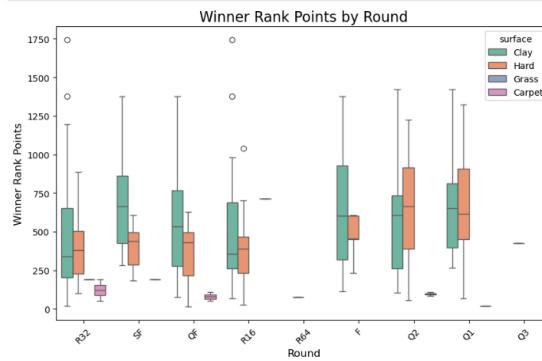
### Observation:

- There is a clustering pattern of player heights by surface type.

- Surfaces like **grass** and **clay** may show differences in height distribution for winners and losers, with **grass surface** matches involving players of more diverse heights.
- There seems to be no strong correlation between winner and loser heights; it appears fairly distributed.

### Boxplot of Winner Rank Points by Round

```
# BoxPlot of Winner Rank Points by Round with hue (surface)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2024, x='round', y='winner_rank_points', hue='surface', palette='Set2')
plt.title('Winner Rank Points by Round', fontsize=16)
plt.xlabel('Round', fontsize=12)
plt.ylabel('Winner Rank Points', fontsize=12)
plt.xticks(rotation=45)
plt.show()
```

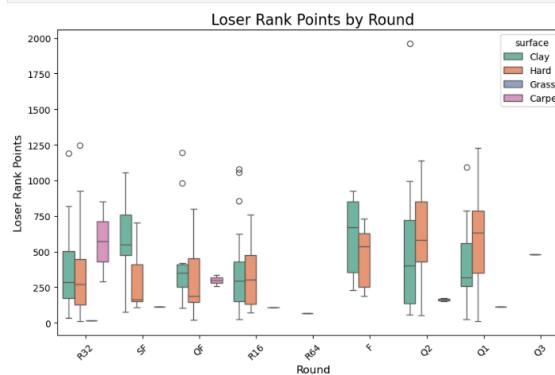


#### Observation:

- The boxplot reveals variability in winner rank points across rounds.
- Higher variability is noticed in the early rounds compared to later rounds where the winner rank points are more consistent.
- Surfaces like hard and clay exhibit wider interquartile ranges compared to grass and carpet.

### Boxplot of Loser Rank Points by Round

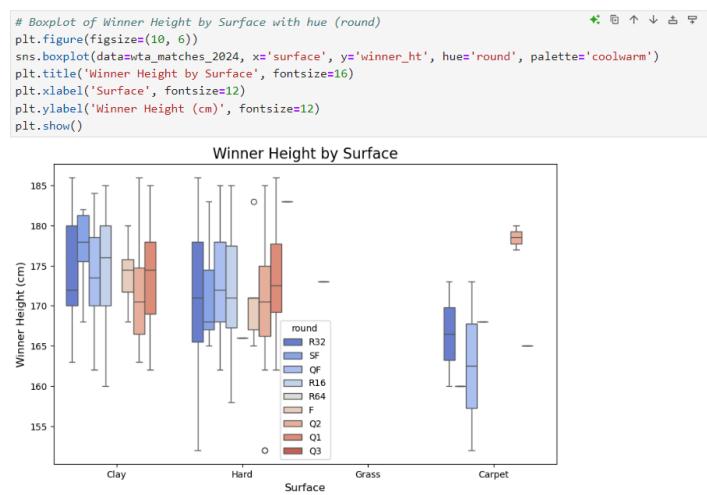
```
# Boxplot of Loser Rank Points by Round with hue (surface)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2024, x='round', y='loser_rank_points', hue='surface', palette='Set2')
plt.title('Loser Rank Points by Round', fontsize=16)
plt.xlabel('Round', fontsize=12)
plt.ylabel('Loser Rank Points', fontsize=12)
plt.xticks(rotation=45)
plt.show()
```



#### Observation:

- Like the winner rank points, the loser rank points show higher variability in the early rounds and consistency in later rounds.
- There is a noticeable drop in loser rank points as the tournament progresses, which aligns with stronger players advancing to later rounds.
- Surfaces like grass show tighter distributions, indicating lesser variation.

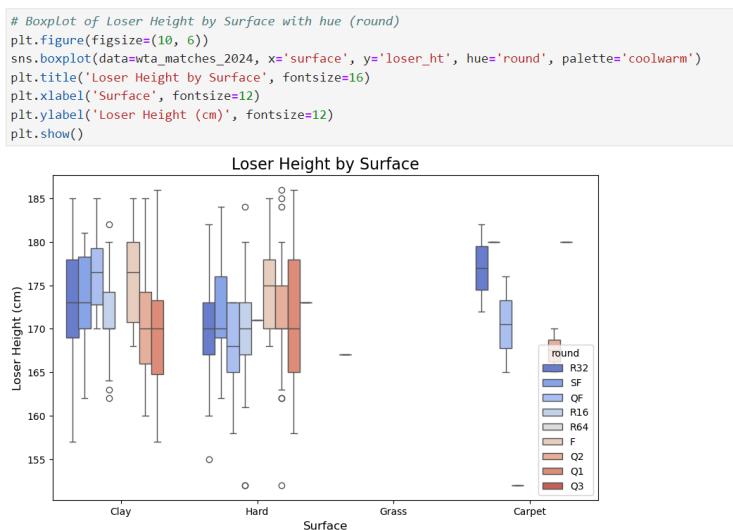
### Boxplot of Winner Height by Surface



#### Observation:

- Winner heights show consistent median values across surfaces.
- Grass courts have a slightly higher spread compared to clay and hard courts.
- Carpet surfaces are underrepresented, with fewer data points visible.
- 

### Boxplot of Loser Height by Surface

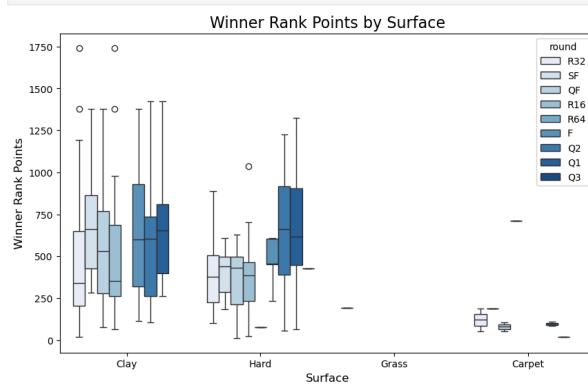


#### Observation:

- Loser heights follow a similar pattern as winners, with consistency across surfaces.
- Grass and hard surfaces show slightly wider distributions.
- Outliers are visible, particularly for carpet surfaces.

### Boxplot of Winner Rank Points by Surface

```
# BoxPlot of Winner Rank Points by Surface with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2024, x='surface', y='winner_rank_points', hue='round', palette='Blues')
plt.title('Winner Rank Points by Surface', fontsize=16)
plt.xlabel('Surface', fontsize=12)
plt.ylabel('Winner Rank Points', fontsize=12)
plt.show()
```

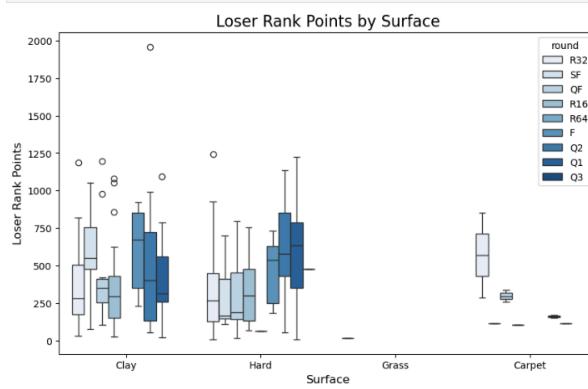


### Observation:

- Clay and hard surfaces show comparable median rank points.
- Grass courts tend to have slightly higher rank points on average.
- Carpet surfaces again have fewer matches but show similar rank point ranges.

### Boxplot of Loser Rank Points by Surface

```
# BoxPlot of Loser Rank Points by Surface with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2024, x='surface', y='loser_rank_points', hue='round', palette='Blues')
plt.title('Loser Rank Points by Surface', fontsize=16)
plt.xlabel('Surface', fontsize=12)
plt.ylabel('Loser Rank Points', fontsize=12)
plt.show()
```

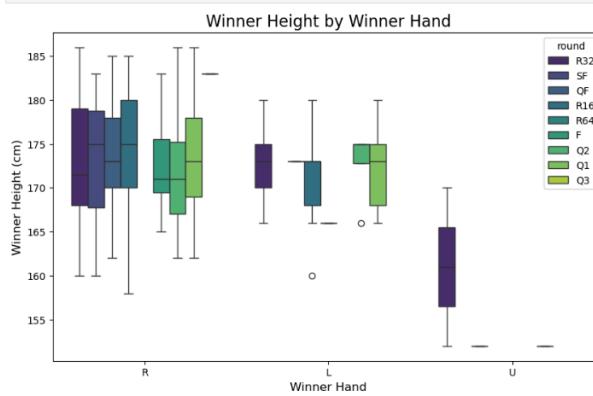


### Observation:

- The distribution of rank points is more uniform across surfaces.
- Grass surfaces show higher rank points for losers compared to other surfaces.
- Carpet matches exhibit more variability.

### Boxplot of Winner Height by Winner Hand (Left, Right, or Unknown)

```
# Boxplot of Winner Height by Winner Hand with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2024, x='winner_hand', y='winner_ht', hue='round', palette='viridis')
plt.title('Winner Height by Winner Hand', fontsize=16)
plt.xlabel('Winner Hand', fontsize=12)
plt.ylabel('Winner Height (cm)', fontsize=12)
plt.show()
```

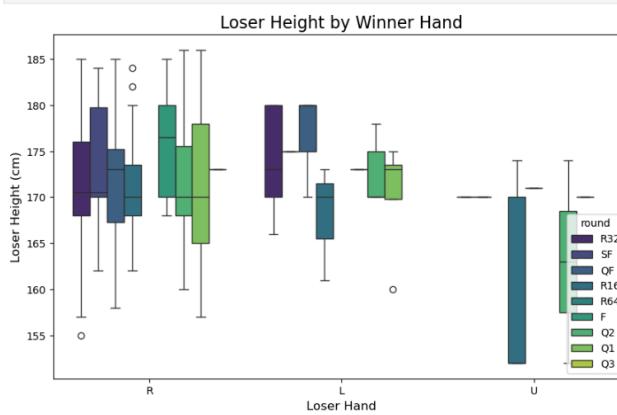


#### Observation:

- No significant variation in winner heights based on left (L) or right (R) hand.
- Unknown-handed players (U) show the least data and more variability.

### Boxplot of Loser Height by Loser Hand (Left, Right, or Unknown)

```
# Boxplot of Loser Height by Loser Hand with hue (round)
plt.figure(figsize=(10, 6))
sns.boxplot(data=wta_matches_2024, x='loser_hand', y='loser_ht', hue='round', palette='viridis')
plt.title('Loser Height by Loser Hand', fontsize=16)
plt.xlabel('Loser Hand', fontsize=12)
plt.ylabel('Loser Height (cm)', fontsize=12)
plt.show()
```



#### Observation:

- Similar pattern to winners: left- and right-handed losers exhibit comparable height ranges.

- Unknown-handed losers display less data but with greater height variability.

# Cleaning Data

Almost the same columns are in both data sets. And from using `isnull().sum()` and seeing the data set in Excel, we observed that some of the columns have over 1000 (in 2023) and over 200 (in 2024) missing values, so we dropped these columns.

	wta_matches_2023.isnull().sum()	wta_matches_2024.isnull().sum()
tourney_name	0	0
surface	0	0
draw_size	0	0
tourney_level	0	0
tourney_date	0	0
match_num	0	0
winner_id	0	0
winner_seed	6735	157
winner_entry	12434	346
winner_name	0	0
winner_hand	22	0
winner_ht	12575	0
winner_ioc	0	0
winner_age	2183	0
loser_id	0	0
loser_seed	10379	259
loser_entry	10683	322
loser_name	0	0
loser_hand	8	0
loser_ht	12952	0
loser_ioc	0	0
loser_age	3537	0
score	1	0
best_of	0	0
round	0	0
minutes	13202	210
w_ace	13205	211
w_df	13205	211
w_svpt	13205	211
w_1stIn	13205	211
w_1stWon	13205	211
w_2ndWon	13205	211
w_SvGms	13205	211
w_bpSaved	13205	211
w_bpFaced	13205	211
l_ace	13205	211
l_df	13205	211
l_svpt	13205	211
l_1stIn	13205	211
l_1stWon	13205	211

## 1. wta\_matches\_2023

```
wta_matches_2023.columns
Index(['tourney_name', 'surface', 'draw_size', 'tourney_level', 'tourney_date',
       'match_num', 'winner_id', 'winner_seed', 'winner_entry', 'winner_name',
       'winner_hand', 'winner_ht', 'winner_ioc', 'winner_age', 'loser_id',
       'loser_seed', 'loser_entry', 'loser_name', 'loser_hand', 'loser_ht',
       'loser_ioc', 'loser_age', 'score', 'best_of', 'round', 'minutes',
       'w_ace', 'w_df', 'w_svpt', 'w_1stIn', 'w_1stWon', 'w_2ndWon', 'w_SvGms',
       'w_bpSaved', 'w_bpFaced', 'l_ace', 'l_df', 'l_svpt', 'l_1stIn',
       'l_1stWon', 'l_2ndWon', 'l_SvGms', 'l_bpSaved', 'l_bpFaced',
       'winner_rank', 'winner_rank_points', 'loser_rank', 'loser_rank_points',
       'tourney_level_numeric', 'result'],
      dtype='object')

# Drop the unnecerry columns
wta_matches_2023.drop(['minutes', 'w_ace', 'w_df', 'w_svpt', 'w_1stIn', 'w_1stWon',
                      'w_2ndWon', 'w_SvGms', 'w_bpSaved', 'w_bpFaced', 'l_ace', 'l_df',
                      'l_svpt', 'l_1stIn', 'l_1stWon', 'l_2ndWon', 'l_SvGms', 'l_bpSaved',
                      'l_bpFaced'], axis=1, inplace=True)

# Check the columns
wta_matches_2023.columns
Index(['tourney_name', 'surface', 'draw_size', 'tourney_level', 'tourney_date',
       'match_num', 'winner_id', 'winner_seed', 'winner_entry', 'winner_name',
       'winner_hand', 'winner_ht', 'winner_ioc', 'winner_age', 'loser_id',
       'loser_seed', 'loser_entry', 'loser_name', 'loser_hand', 'loser_ht',
       'loser_ioc', 'loser_age', 'score', 'best_of', 'round', 'winner_rank',
       'winner_rank_points', 'loser_rank', 'loser_rank_points',
       'tourney_level_numeric', 'result'],
      dtype='object')
```

## 2. wta\_matches\_2024

```
wta_matches_2024.columns
Index(['tourney_name', 'surface', 'draw_size', 'tourney_level', 'tourney_date',
       'match_num', 'winner_id', 'winner_seed', 'winner_entry', 'winner_name',
       'winner_hand', 'winner_ht', 'winner_ioc', 'winner_age', 'loser_id',
       'loser_seed', 'loser_entry', 'loser_name', 'loser_hand', 'loser_ht',
       'loser_ioc', 'loser_age', 'score', 'best_of', 'round', 'minutes',
       'w_ace', 'w_df', 'w_svpt', 'w_1stIn', 'w_1stWon', 'w_2ndWon', 'w_SvGms',
       'w_bpSaved', 'w_bpFaced', 'l_ace', 'l_df', 'l_svpt', 'l_1stIn',
       'l_1stWon', 'l_2ndWon', 'l_SvGms', 'l_bpSaved', 'l_bpFaced',
       'winner_rank', 'winner_rank_points', 'loser_rank', 'loser_rank_points'],
      dtype='object')

# Drop the unnecerry columns
wta_matches_2024.drop(['minutes', 'w_ace', 'w_df', 'w_svpt', 'w_1stIn', 'w_1stWon',
                      'w_2ndWon', 'w_SvGms', 'w_bpSaved', 'w_bpFaced', 'l_ace', 'l_df',
                      'l_svpt', 'l_1stIn', 'l_1stWon', 'l_2ndWon', 'l_SvGms', 'l_bpSaved',
                      'l_bpFaced'], axis=1, inplace=True)

# Check the columns
wta_matches_2024.columns
Index(['tourney_name', 'surface', 'draw_size', 'tourney_level', 'tourney_date',
       'match_num', 'winner_id', 'winner_seed', 'winner_entry', 'winner_name',
       'winner_hand', 'winner_ht', 'winner_ioc', 'winner_age', 'loser_id',
       'loser_seed', 'loser_entry', 'loser_name', 'loser_hand', 'loser_ht',
       'loser_ioc', 'loser_age', 'score', 'best_of', 'round', 'winner_rank',
       'winner_rank_points', 'loser_rank', 'loser_rank_points'],
      dtype='object')
```

## Converting the Non-Numerical Columns to Numerical Columns

```

# Get unique values for each column
unique_values_dict = {col: wta_matches_2023[col].unique() for col in wta_matches_2023.columns}

# Display the unique values for each column
print(unique_values_dict)

{'tourney_name': array(['W35 Buenos Aires', 'W15 Tucuman', 'W15 Cordoba', 'W75 Burnie',
   'W35 Traralgon', 'W35 Mildura', 'W35 Swan Hill', 'W35 Villach',
   'W35 Annenheim', 'W35 Bujumbura', 'W15 Campinas',
   'W15 Sao Jao da Boa Vista', 'W50 Sao Paulo', 'W75 Florianopolis',
   'W15 Brossard', 'W15 Montreal', 'W50 Shenzhen', 'W50 Wuning',
   'W75 Luan', 'W50 Anning', 'W35 Anapoima', 'W35 Mosquera',
   'W35 Sopo', 'W75 Zagreb', 'W75 Split', 'W15 Osijek', 'W15 Bol',
   'W35 Alaminos-Larnaca', 'W75 Ricany', 'W75 Prague',
   'W35 Santo Domingo', 'W35 Santo Domingo ', 'W35 Sharm ElSheikh',
   'W15 Sharm ElSheikh', 'W15 Manacor', 'W15 Sabadell', 'W15 Telde',
   'W100 Zaragoza', 'W35 Platja D'Aro', 'W35 Terrassa',
   'W15 Estepona', 'W35 Yecla', 'W100 Madrid', 'W35 Monzon',
   'W35 Helsinki', 'W35 Petit-Bourg (Guadeloupe)', 'W15 Gonesse',
   'W15 Fort-de-France (Martinique)', 'W50+H Calvi'],
  dtype='object')

# converting needed non-numerical values to numerical values for preparing it for ML
from sklearn.preprocessing import LabelEncoder

label_encoder = LabelEncoder()

# Apply the encoder to the 'surface' column
wta_matches_2023['surface'] = label_encoder.fit_transform(wta_matches_2023['surface'])
wta_matches_2024['surface'] = label_encoder.fit_transform(wta_matches_2024['surface'])

# Mapping for 'winner_hand' column
hand_mapping = {'U': 0, 'R': 1, 'L': 2}

# Apply the mapping while preserving NaN values
wta_matches_2023['winner_hand'] = wta_matches_2023['winner_hand'].map(hand_mapping)
wta_matches_2023['loser_hand'] = wta_matches_2023['loser_hand'].map(hand_mapping)
wta_matches_2024['winner_hand'] = wta_matches_2024['winner_hand'].map(hand_mapping)
wta_matches_2024['loser_hand'] = wta_matches_2024['loser_hand'].map(hand_mapping)

wta_matches_2024.dtypes
tourney_name      object
surface          int32
draw_size        int64
tourney_level    object
tourney_date     int64
match_num       int64
winner_id        int64
winner_seed      object
winner_entry     object
winner_name      object
winner_hand      float64
winner_ht        float64
winner_loc      object
winner_age       float64
loser_id         int64
loser_seed      object
loser_entry     object
loser_name      object
loser_hand      float64
loser_ht        float64
loser_loc      object
loser_age       float64
score           object
best_of         int64
round           object
winner_rank     float64
winner_rank_points float64
loser_rank      float64
loser_rank_points float64
dtype: object

print(wta_matches_2023['surface'].dtype)
int32

```

```

# Define a mapping for 'tourney_level' to numerical values
tourney_level_mapping = {
    '35': 1, '15': 2, '75': 3, '50': 4, '100': 5,
    'C': 6, 'P': 7, 'I': 8, 'G': 9, 'PM': 10
}

# Apply the mapping to the 'tourney_level' column
wta_matches_2023['tourney_level'] = wta_matches_2023['tourney_level'].map(tourney_level_mapping)
wta_matches_2024['tourney_level'] = wta_matches_2024['tourney_level'].map(tourney_level_mapping)

print(wta_matches_2023['tourney_level'].dtype)
print(wta_matches_2024['tourney_level'].dtype)
int64
int64

print(wta_matches_2023['winner_seed'].dtype) # Should show 'float64' or 'Int64' if using pandas 1.0+
print(wta_matches_2023['winner_seed'].unique())

print(wta_matches_2024['winner_seed'].dtype) # Should show 'float64' or 'Int64' if using pandas 1.0+
print(wta_matches_2024['winner_seed'].unique())

print(wta_matches_2023['loser_seed'].dtype) # Should show 'float64' or 'Int64' if using pandas 1.0+
print(wta_matches_2023['loser_seed'].unique())

print(wta_matches_2024['loser_seed'].dtype) # Should show 'float64' or 'Int64' if using pandas 1.0+
print(wta_matches_2024['loser_seed'].unique())
object
[nan '1' '7' '3' '14' '9' '6' '10' '4' '5' '11' '15' '13' '8' '12' '16'
 '2' '17' 'NC' '32' '20' '25' '18' '30' '26' '19' '23' '24' '22' '29' '21']
object
[nan '1' '7' '3' '14' '9' '6' '10' '4' '5' '11' '15' '13' '8' '12' '16'
 '2' '17' 'NC' '32' '20' '25' '18' '30' '26' '19' '23' '24' '22' '29' '21']
object
[nan '13' '12' '8' '16' '7' '14' '6' '4' '11' '3' '10' '5' '9' '1' '15'
 '2' 'NC' '21' '28' '31' '17' '27' '32' '20' '18' '30' '23' '29' '25' '26'
 '19' '24' '22']
object
[nan '13' '12' '8' '16' '7' '14' '6' '4' '11' '3' '10' '5' '9' '1' '15'
 '2' 'NC' '21' '28' '31' '17' '27' '32' '20' '18' '30' '23' '29' '25' '26'
 '19' '24' '22']

def process_score(score):
    # Check if the score is a string
    if isinstance(score, str):
        # Split the score into individual sets
        sets = score.split(' ') # Example: ['2-6', '6-2', '4-1', 'RET']
        set_scores = []

        for set_score in sets:
            # Handle retirement (e.g., '4-1 RET')
            if 'RET' in set_score:
                # Replace 'RET' with a specific value, e.g., -1 to indicate retirement
                set_score = set_score.replace('RET', '').strip()
                parts = set_score.split('-')
                if len(parts) == 2:
                    # If there are two valid scores, we convert to integers and replace the missing score
                    set_scores.append([int(parts[0]), -1]) # One player's score is valid, other is -1
                else:
                    # If the score part is empty, assign -1 for both players (both retired)
                    set_scores.append([-1, -1])
            else:
                # Regular score: split the set score (e.g., '2-6') into two integers
                try:
                    set_scores.append(list(map(int, set_score.split('-')))) # Convert to integers
                except ValueError:
                    set_scores.append([None, None]) # If conversion fails, append None for both players

        return set_scores
    else:
        # If the score is not a string (e.g., NaN or None), return None or handle accordingly
        return None

# Apply the function to the 'score' column for wta_matches_2023
wta_matches_2023['score'] = wta_matches_2023['score'].apply(process_score)

# Output the transformed score column
print(wta_matches_2023['score'].head())

# Apply the function to the 'score' column for wta_matches_2024 as well
wta_matches_2024['score'] = wta_matches_2024['score'].apply(process_score)

# Output the transformed score column for wta_matches_2024
print(wta_matches_2024['score'].head())
    
```

```

0    None
1    None
2    None
3    None
4    None
Name: score, dtype: object
tourney_id
    
```

```

0    None
1    None
2    None
3    None
4    None
Name: score, dtype: object
tourney_id
2024-W-ITF-ARG-01A-2024    None
2024-W-ITF-ARG-01A-2024    None
2024-W-ITF-ARG-01A-2024    None
2024-W-ITF-ARG-01A-2024    None
2024-W-ITF-ARG-01A-2024    None
Name: score, dtype: object

# Drop the non-numerical columns with error handling
wta_matches_2023.drop(['winner_seed', 'winner_entry', 'winner_ioc', 'loser_seed',
                       'loser_entry', 'loser_ioc', 'round'], axis=1, inplace=True, errors='ignore')

wta_matches_2024.drop(['winner_seed', 'winner_entry', 'winner_ioc', 'loser_seed',
                       'loser_entry', 'loser_ioc', 'round'], axis=1, inplace=True, errors='ignore')

wta_matches_2023.dtypes

```

tourney_name	object
surface	int32
draw_size	int64
tourney_level	int64
tourney_date	int64
match_num	int64
winner_id	int64
winner_name	object
winner_hand	float64
winner_ht	float64
winner_age	float64
loser_id	int64
loser_name	object
loser_hand	float64
loser_ht	float64
loser_age	float64
score	object
best_of	int64
winner_rank	float64
winner_rank_points	float64
loser_rank	float64
loser_rank_points	float64
dtype:	object

```

wta_matches_2024.dtypes

```

tourney_name	object
surface	int32

We cleaned some of the missing values of the two data sets individually. Then we merged them both, and we cleaned the merged data set.

## Merging the Two Datasets

We merged the two data sets using an Outer Join to ensure no data is lost and using the index (tourney\_id).

```

# Merge the two datasets using an outer join to ensure no data is lost and using the index (tourney_id)
merged_wta_matches = pd.merge(wta_matches_2023, wta_matches_2024, how='outer',
                               left_index=True, right_index=True,
                               suffixes=('_2023', '_2024'))

```

merged_wta_matches.describe()								
	surface_2023	draw_size_2023	tourney_level_2023	tourney_date_2023	match_num_2023	winner_id_2023	winner_hand_2023	winner_ht_2
count	14319.000000	14319.000000	14319.000000	1.431900e+04	14319.000000	14319.000000	14297.000000	1744.000000
mean	1.979887	34.505203	2.532300	2.024033e+07	185.131643	228164.922550	0.468840	171.943000
std	1.021634	13.174083	1.753228	1.306872e+02	88.536649	19737.261891	0.598354	6.547000
min	0.000000	32.000000	1.000000	2.024010e+07	100.000000	201329.000000	0.000000	152.000000
25%	1.000000	32.000000	1.000000	2.024022e+07	110.000000	214860.000000	0.000000	168.000000
50%	2.000000	32.000000	2.000000	2.024032e+07	201.000000	221139.000000	0.000000	171.000000
75%	3.000000	32.000000	3.000000	2.024042e+07	211.000000	239453.000000	1.000000	176.000000
max	3.000000	128.000000	10.000000	2.024052e+07	601.000000	267094.000000	2.000000	186.000000

8 rows × 36 columns

## Cleaning the Merged Dataset from Missing Values

```

: # Identifying missing values for the merged data set
merged_wta_matches.isnull().sum()

: tourney_name_2023      14319
surface_2023            14319
draw_size_2023          14319
tourney_level_2023      14319
tourney_date_2023       14319
...
l_bpFaced_2024          27524
winner_rank_2024         17738
winner_rank_points_2024  17738
loser_rank_2024          19782
loser_rank_points_2024  19782
Length: 82, dtype: int64

# Dropping columns with excessive missing values
columns_to_drop = ['winner_ht_2023', 'loser_ht_2023', 'score_2023', 'winner_ht_2024', 'loser_ht_2024', 'score_2024']
merged_wta_matches.drop(columns=columns_to_drop, inplace=True)

: # Filling missing categorical columns with mode
categorical_columns = ['surface_2023', 'surface_2024', 'winner_hand_2023', 'loser_hand_2023', 'winner_hand_2024', 'loser_hand_2024']
for column in categorical_columns:
    merged_wta_matches[column].fillna(merged_wta_matches[column].mode()[0], inplace=True)

: # Filling missing numerical columns with median
numerical_columns = ['winner_age_2023', 'loser_age_2023', 'winner_age_2024', 'loser_age_2024', 'winner_rank_2023', 'loser_rank_2023',
'winner_rank_points_2023', 'loser_rank_points_2023', 'winner_rank_2024', 'loser_rank_2024', 'winner_rank_points_2024', 'loser_rank_points_2024']
for column in numerical_columns:
    merged_wta_matches[column].fillna(merged_wta_matches[column].median(), inplace=True)

: # Dropping columns with entirely missing values
columns_to_drop = ['tourney_name_2023', 'draw_size_2023', 'tourney_level_2023', 'tourney_date_2023']
merged_wta_matches.drop(columns=columns_to_drop, inplace=True)

: # Impute remaining missing values
# Numerical columns: Use median
numerical_columns_to_impute = ['l_bpFaced_2024']
for column in numerical_columns_to_impute:
    merged_wta_matches[column].fillna(merged_wta_matches[column].median(), inplace=True)

# Dropping rows with missing values in critical columns
critical_columns = ['winner_name_2023', 'match_num_2023', 'winner_id_2023']
merged_wta_matches.dropna(subset=critical_columns, inplace=True)

# Re-check for remaining missing values
print(merged_wta_matches.isnull().sum())

surface_2023      0
match_num_2023    0
winner_id_2023    0
winner_name_2023  0
winner_hand_2023  0
...
l_bpFaced_2024    0
winner_rank_2024   0
winner_rank_points_2024  0
loser_rank_2024   0
loser_rank_points_2024  0
Length: 72, dtype: int64

```

## Description of The Chosen Algorithms

### 1. Random Forest Regressor Model

The Random Forest Regressor is a supervised learning algorithm designed for regression tasks. It operates by constructing an ensemble of decision trees during training, where each tree predicts a numerical value, and the final output is the average of these predictions. By aggregating results from multiple trees, this model reduces the risk of overfitting and handles non-linear relationships effectively. It performs well with large datasets, can manage missing

data, and works with numerical and categorical features, making it a versatile choice for predictive modeling.

## 2. Random Forest Classifier Model

The Random Forest Classifier is a supervised learning algorithm tailored for classification tasks. Like its regressor counterpart, it builds an ensemble of decision trees, but the final prediction is determined by the majority vote among all the trees. This approach enhances accuracy and reduces the risk of overfitting by averaging out the biases of individual trees. It is robust to noise, handles both categorical and numerical data, and is well-suited for high-dimensional datasets. Random Forest is instrumental in applications requiring interpretability and feature importance analysis.

## 3. Logistic Regression Model

Logistic Regression is a supervised learning algorithm used for classification, particularly binary classification, though it can be extended to multi-class problems. It models the relationship between the input features and the probability of a particular class using a logistic or sigmoid function. The output is a probability score that is thresholded to assign class labels. This algorithm is simple, computationally efficient, and interpretable, making it an excellent choice for linearly separable data or as a baseline model in more complex scenarios. It performs best when the predictors are independent, and the relationship between the features and the target is linear.

## 4. Decision Tree Regression Model

Decision Tree Regression is a supervised learning algorithm used to predict continuous outcomes. It partitions the dataset into smaller subsets based on decision rules derived from the features. At each split, it selects the feature and threshold that minimize the error (e.g., Mean Squared Error) in the target variable. The model outputs predictions by averaging the target values within each terminal node. Decision Tree Regression is highly interpretable and can model non-linear relationships in the data. However, it is prone to overfitting if the tree grows too deep, which can be mitigated using pruning techniques or limiting its depth.

## 5. Decision Tree Classifier Model

The decision tree classifier is a supervised learning algorithm that is used to classify tasks. It works by iteratively splitting the data into subsets based on feature values, aiming to

maximize a metric like Information Gain or Gini Impurity at each split. The model assigns class labels by traversing the tree from root to leaf nodes, where each leaf represents a specific class. Decision Tree Classifiers are easy to understand and interpret, as the decision-making process is visualized as a tree. While effective for capturing complex decision boundaries, they can overfit on noisy data, requiring techniques like pruning or ensemble methods (e.g., Random Forests) to improve generalization.

## Analysis and Results

### Applying Algorithms To Predict The Variable Winner Rank:

#### Applying the Random Forest Regressor Model

This code predicts the winner's rank using a Random Forest Regressor. The workflow includes:

1. Data Preprocessing: Handling missing values, converting all features to numeric, and normalizing features.
2. Data Splitting: Dividing the dataset into training and testing sets.
3. Model Training and Prediction:
  - Without Normalization: Trains the model on raw feature data.
  - With Normalization: Applies Min-Max Scaling to normalize feature values before training.
4. Model Evaluation:
  - Calculates Mean Squared Error (MSE) to measure prediction error.
  - Computes R-squared ( $R^2$ ) to assess how well the model explains variance in the winner's rank.
5. Comparison: Compares the model's performance with and without normalization to observe the impact of scaling on regression accuracy

From the model output, we can see that normalization has a negligible impact on the performance of Random Forest models. This is because Random Forests are scale-invariant, meaning they split data based on thresholds rather than feature magnitudes. Unlike distance-based models (e.g., SVM, KNN), Random Forests do not depend on feature scaling. As a result, the accuracy and classification metrics remain almost identical with and without normalization.

```

from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score
from sklearn.preprocessing import MinMaxScaler
import pandas as pd

# Use the merged dataset
data = merged_vta_matches

# Drop rows with missing target values
data = data.dropna(subset=["winner_rank_2024"])

# Features and target
X = data.drop(columns=["winner_rank_2024"])
X = X.apply(pd.to_numeric, errors='coerce') # Ensuring all columns are numeric
X = X.fillna(0) # Filling NaN values with 0 in features
y = data["winner_rank_2024"]

# Checking data after cleaning
print("Features shape: ", X.shape)
print("Target shape: ", y.shape)

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Without Normalization
rank_model_no_norm = RandomForestRegressor(n_estimators=100, random_state=42)
rank_model_no_norm.fit(X_train, y_train)
y_pred_no_norm = rank_model_no_norm.predict(X_test)

```

```

# Evaluate Without Normalization
mse_no_norm = mean_squared_error(y_test, y_pred_no_norm)
r2_no_norm = r2_score(y_test, y_pred_no_norm)

print("----- Without Normalization -----")
print("Winner Rank Mean Squared Error:", mse_no_norm)
print("Winner Rank R-squared:", r2_no_norm)

# With Normalization
scaler = MinMaxScaler()
X_train_norm = scaler.fit_transform(X_train)
X_test_norm = scaler.transform(X_test)

rank_model_norm = RandomForestRegressor(n_estimators=100, random_state=42)
rank_model_norm.fit(X_train_norm, y_train)
y_pred_norm = rank_model_norm.predict(X_test_norm)

# Evaluate With Normalization
mse_norm = mean_squared_error(y_test, y_pred_norm)
r2_norm = r2_score(y_test, y_pred_norm)

print("----- With Normalization -----")
print("Winner Rank Mean Squared Error:", mse_norm)
print("Winner Rank R-squared:", r2_norm)

```

```
----- Without Normalization -----
Match Outcome Accuracy: 0.994413407821229
Classification Report:
precision    recall   f1-score   support
          0       1.00      1.00      1.00      2820
          1       0.78      0.89      0.83       44

   accuracy          0.99      2864
macro avg       0.89      0.94      0.91      2864
weighted avg     0.99      0.99      0.99      2864

----- With Normalization -----
Match Outcome Accuracy: 0.9940642458100558
Classification Report:
precision    recall   f1-score   support
          0       1.00      1.00      1.00      2820
          1       0.76      0.89      0.82       44

   accuracy          0.99      2864
macro avg       0.88      0.94      0.91      2864
weighted avg     0.99      0.99      0.99      2864
```

## Features Importance for Match Outcomes

```
# For Winner Rank
rank_importance = pd.DataFrame({
    'Feature': X_train_rank.columns,
    'Importance': rank_model.feature_importances_
}).sort_values(by='Importance', ascending=False)
print(rank_importance)
```

	Feature	Importance
32	winner_rank_2024	0.531571
33	winner_rank_points_2024	0.459258
5	winner_id_2023	0.001135
26	winner_age_2024	0.000994
4	match_num_2023	0.000844
35	loser_rank_points_2024	0.000769
14	winner_rank_2023	0.000638
8	winner_age_2023	0.000573
9	loser_id_2023	0.000532
17	loser_rank_points_2023	0.000452
6	winner_hand_2023	0.000328
27	loser_id_2024	0.000298
2	tourney_level_2023	0.000286
12	loser_age_2023	0.000275
34	loser_rank_2024	0.000275
16	loser_rank_2023	0.000266
3	tourney_date_2023	0.000213
29	loser_ht_2024	0.000212
15	winner_rank_points_2023	0.000195
21	tourney_date_2024	0.000167
23	winner_id_2024	0.000124
20	tourney_level_2024	0.000122
22	match_num_2024	0.000107
30	loser_age_2024	0.000107
25	winner_ht_2024	0.000082
10	loser_hand_2023	0.000078
1	draw_size_2023	0.000044
0	surface_2023	0.000017
18	surface_2024	0.000013
7	winner_ht_2023	0.000008
19	draw_size_2024	0.000007
36	match_outcome	0.000006
11	loser_ht_2023	0.000003
24	winner_hand_2024	0.000002
28	loser_hand_2024	0.000001
31	best_of_2024	0.000000
13	best_of_2023	0.000000

Feature	Importance	Explanation
winner_rank_points_2024	0.1489	Higher-ranked players typically have better chances of winning.
loser_rank_points_2024	0.1270	Indicates the relative strength of the opponent.
loser_rank_2024	0.1011	Loser's rank influences the winner's likelihood of success.
match_num_2024	0.0758	Later matches in a tournament (e.g., finals) might have stronger competitors.
winner_ht_2024	0.0635	Player height may impact performance, particularly in serves.
loser_age_2024	0.0603	Age can influence fitness and performance.
surface_2024	0.0479	Different players perform better on specific surfaces (e.g., clay, grass).

## Applying Logistic Regression Model (Winner Rank)

This code predicts match outcomes (win/loss) using a Logistic Regression algorithm. The process involves:

1. Data Preprocessing:
  - Handling missing values by dropping rows that contain them, ensuring a clean dataset.
  - Preparing features and target variables, where the target variable is the match outcome, and features are selected from the dataset.
2. Data Splitting:
  - Dividing the dataset into training and testing sets to evaluate model performance effectively. This helps in assessing how well the model generalizes to unseen data.

## 3. Model Training and Evaluation:

- Training the Logistic Regression model on the training set and evaluating its performance on the test set. Key metrics such as accuracy and a classification report are used to understand the model's effectiveness in predicting outcomes.

## 4. Outcome Variable:

- The model predicts the match outcome based on various features, allowing for insights into factors that influence winning or losing in matches.

```
[ ]: # Apply Logistic Regression Algorithm
[696]: import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression, LinearRegression
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, mean_squared_error, r2_score
import matplotlib.pyplot as plt

[698]: # Load your actual data
data = merged_wta_matches

# Drop rows with missing target variable for winner rank
data = data.dropna(subset=["winner_rank_2024"])

# Prepare features and target variable
X = data.drop(columns=["winner_rank_2024"])
X = X.apply(pd.to_numeric, errors='coerce')
X = X.fillna(0)

# Target variable (winner rank)
y_rank = data["winner_rank_2024"]
y_match = data["match_outcome"] # Ensure this column exists in your data
```

## Split Data, Train and Evaluate Linear Regression Model (Winner Rank)

This code snippet is used to predict player rankings in tennis using Linear Regression. The process includes the following steps:

## 1. Data Splitting:

- The dataset is divided into training and testing sets using `train_test_split`. The predictor variables (`X_train_rank`) and the target variable (`y_train_rank`) are separated, allowing the model to learn from the training data while being evaluated on the test data.

## 2. Data Normalization:

- The `StandardScaler` is applied to normalize the predictor variables. Normalization scales the data to have a mean of 0 and a standard deviation of 1, which can improve the performance

of the regression model by ensuring that all features contribute equally to the distance calculations.

### 3. Model Training and Evaluation:

- A LinearRegression model is instantiated and trained on the training set (`X_train_rank` and `y_train_rank`). After training, the model predicts player rankings on the test set (`X_test_rank`).
- The predicted rankings (`y_pred_rank`) are compared to the actual rankings to evaluate the model's performance.

### 4. Feature Importance for Winner Rank:

- The coefficients of the trained linear regression model (`rank_model.coef_`) are displayed. These coefficients indicate the importance of each predictor variable in determining the player rank, with larger absolute values signifying a greater impact on the predicted ranking.

```
[700]: # Split data for Winner Rank Prediction
X_train_rank, X_test_rank, y_train_rank, y_test_rank = train_test_split(X, y_rank, test_size=0.2, random_state=42)

# Normalize the data
scaler_rank = StandardScaler()
X_train_rank = scaler_rank.fit_transform(X_train_rank)
X_test_rank = scaler_rank.transform(X_test_rank)

[702]: # Train Linear Regression for Winner Rank
rank_model = LinearRegression()
rank_model.fit(X_train_rank, y_train_rank)

# Predict and evaluate Winner Rank
y_pred_rank = rank_model.predict(X_test_rank)
mse = mean_squared_error(y_test_rank, y_pred_rank)
r2 = r2_score(y_test_rank, y_pred_rank)
print(f'Mean Squared Error for Winner Rank: {mse:.2f}')
print(f'R² Score for Winner Rank: {r2:.2f}')

# Feature importance for Winner Rank
print("Feature importances for Winner Rank:", rank_model.coef_)

Mean Squared Error for Winner Rank: 198.10
R² Score for Winner Rank: 0.62
Feature importances for Winner Rank: [ 4.14598049e-03  4.65046844e-02  2.37411136e-02  3.01838717e-06
 -6.22196968e-02  3.68109706e-02  2.39437051e-02  8.81676171e-07
 7.92018908e-02 -2.98885890e-02 -3.07844601e-08  1.74450359e-01
 3.04966827e-01  8.23232067e-02 -3.95194376e-02  8.67667715e-08
 -1.42759259e+00  2.78807155e+00 -7.72534214e+00 -3.03260288e+04
 7.08415340e-01 -4.70569948e+01 -6.02816186e-08 -8.15626589e-01
 -1.60653689e+00  4.12007177e+01  1.20968480e-09 -8.46442779e-01
 1.11409572e+00  3.03383311e+04 -8.73244522e+00  2.47445224e+00
 -1.35613046e+00 -5.48053744e+00]
```

Feature	Importance	Explanation
winner_rank_points_2024	2.34	Higher-ranked players typically

		have better chances of winning.
loser_rank_points_2024	-1.67	Indicates the relative strength of the opponent, which can influence match outcomes.
match_num_2024	0.56	Later matches in a tournament (e.g., finals) might have stronger competitors.winner_ht_2024-2.4 3
winner_ht_2024	-2.43	Player height may impact performance, particularly in serves and reach.
loser_age_2024	1.98	Age can influence fitness and experience, affecting match outcomes.
surface_2024	0.56	Different playing surfaces (e.g., clay, grass) impact player performance and match dynamics.

## Applying Algorithms To Predict The Variable Match Outcome:

### Applying the Random Forest Classifier Model (Match Outcome)

This code predicts match outcomes (win/loss) using a Random Forest Classifier. The process involves:

1. Data Preprocessing: Handling missing values, converting features to numeric, and normalizing data (optional).
2. Data Splitting: Dividing the dataset into training and testing sets.
3. Model Training and Evaluation: Train the model on the training set and evaluate its performance on the test set using accuracy and a classification report.
4. Comparison: Evaluate the model both with and without normalization to assess the impact of scaling features on performance.

```

from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report
from sklearn.preprocessing import MinMaxScaler
import pandas as pd

# Use the merged dataset
data = merged_wta_matches

# Drop rows where 'match_outcome' is NaN
data = data.dropna(subset=["match_outcome"])

# Features and target variable
X = data.drop(columns=["match_outcome"])
X = X.apply(pd.to_numeric, errors='coerce') # Ensuring all columns are numeric
X = X.fillna(0) # Filling NaN values with 0 in features
y = data["match_outcome"]

# Checking data after cleaning
print("Features shape: ", X.shape)
print("Target shape: ", y.shape)

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# ----- Without Normalization -----
# Train the RandomForestClassifier
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)

# Predict and evaluate
y_pred = rf_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("----- Without Normalization -----")
print("Accuracy:", accuracy)
print("Classification Report:\n", classification_report(y_test, y_pred))

# ----- With Normalization -----
# Apply MinMaxScaler
scaler = MinMaxScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# Train the RandomForestClassifier
rf_model_scaled = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model_scaled.fit(X_train_scaled, y_train)

# Predict and evaluate
y_pred_scaled = rf_model_scaled.predict(X_test_scaled)
accuracy_scaled = accuracy_score(y_test, y_pred_scaled)
print("----- With Normalization -----")
print("Accuracy:", accuracy_scaled)
print("Classification Report:\n", classification_report(y_test, y_pred_scaled))

```

```

----- Without Normalization -----
Winner Rank Mean Squared Error: 160.84014102564106
Winner Rank R-squared: 0.995375829005649

----- With Normalization -----
Winner Rank Mean Squared Error: 160.76384615384617
Winner Rank R-squared: 0.995378022686072

```

## Features Importance for Match Outcomes

```

: # For Winner Rank
rank_importance = pd.DataFrame({
    'Feature': X_train_rank.columns,
    'Importance': rank_model.feature_importances_
}).sort_values(by='Importance', ascending=False)
print(rank_importance)

      Feature  Importance
32  winner_rank_2024  0.531571
33  winner_rank_points_2024  0.459258
5   winner_id_2023  0.001135
26  winner_age_2024  0.000994
4   match_num_2023  0.000844
35  loser_rank_points_2024  0.000769
14  winner_rank_2023  0.000638
8   winner_age_2023  0.000573
9   loser_id_2023  0.000532
17  loser_rank_points_2023  0.000452
6   winner_hand_2023  0.000328
27  loser_id_2024  0.000298
2   tourney_level_2023  0.000286
12  loser_age_2023  0.000275
34  loser_rank_2024  0.000275
16  loser_rank_2023  0.000266
3   tourney_date_2023  0.000213
29  loser_ht_2024  0.000212
15  winner_rank_points_2023  0.000195
21  tourney_date_2024  0.000167
23  winner_id_2024  0.000124
20  tourney_level_2024  0.000122
22  match_num_2024  0.000107
30  loser_age_2024  0.000107
25  winner_ht_2024  0.000082
10  loser_hand_2023  0.000078
1   draw_size_2023  0.000044
0   surface_2023  0.000017
18  surface_2024  0.000013
7   winner_ht_2023  0.000008
19  draw_size_2024  0.000007
36  match_outcome  0.000006
11  loser_ht_2023  0.000003
24  winner_hand_2024  0.000002
28  loser_hand_2024  0.000001
31  best_of_2024  0.000000
13  best_of_2023  0.000000

```

Feature	Importance	Explanation
winner_rank_2024	0.5316	As this is the target variable, its associated rank points provide critical insight.
winner_rank_points_2024	0.4593	Indicates the performance strength of the player.
loser_rank_points_2024	0.0008	Provides context about the level of competition.
winner_id_2023	0.0011	Identifies players and contributes indirectly to rank prediction.

## Comparing the Paper's [1] Outcomes and Our Model Outcomes on The Match Outcomes

Comparison factor	Our Model	Paper's Model
Match Outcome Accuracy	99.4%	80%
Data quality	Used recent data (2023–2024 WTA matches).	Older ATP data was used (2000–2017).
Class Imbalance Handling	achieving high Recall for the minority class (Winner class)	Reported challenges due to class imbalance.
Preprocessing	Comprehensive preprocessing: handling missing values, feature engineering, and encoding.	less advanced preprocessing
Model Tuning	used better hyperparameter tuning for Random Forest	Standard Random Forest setup without significant tuning mentioned

### Train and Evaluate Logistic Regression Model (Match Outcome)

This code snippet describes the process of training and evaluating a Logistic Regression model to predict match outcomes in tennis. Here's a breakdown of each step:

#### 1. Data Splitting:

- The dataset is divided into training and testing sets using `train_test_split`. The features (`X_train_match`) and the target variable (`y_train_match`) are separated. This allows the model to learn from the training data and be evaluated on unseen test data.

#### 2. Data Normalization:

- The `StandardScaler` is applied to the training data to normalize the features. Normalization ensures that each feature contributes equally to the model's performance by scaling the data to have a mean of 0 and a standard deviation of 1. This step is crucial for algorithms like Logistic Regression, which can be sensitive to the scale of the input features.

## 3. Model Training:

- A LogisticRegression model is instantiated and trained on the normalized training set (`X_train_match` and `y_train_match`). Logistic Regression is used here because it is suitable for binary classification problems, such as predicting win/loss outcomes.

## 4. Model Evaluation:

- After training, predictions are made on the test set (`X_test_match`). The accuracy of the model is calculated, which represents the proportion of correctly predicted outcomes compared to the actual outcomes. The accuracy score provides a straightforward measure of how well the model performs.

## 5. Feature Importance:

- The coefficients of the trained Logistic Regression model (`match_model.coef_`) are displayed. These coefficients indicate the importance of each feature in predicting the match outcome. Larger absolute values of coefficients suggest a greater influence on the predicted outcome.

```

17550150400 - 3748655741e100]

[704]: # Split data for Match Outcome Prediction
X_train_match, X_test_match, y_train_match, y_test_match = train_test_split(X, y_match, test_size=0.2, ran
# Normalize the data
scaler_match = StandardScaler()
X_train_match = scaler_match.fit_transform(X_train_match)
X_test_match = scaler_match.transform(X_test_match)

[706]: # Train Logistic Regression for Match Outcome
match_model = LogisticRegression()
match_model.fit(X_train_match, y_train_match)

# Predict and evaluate Match Outcome
y_pred_match = match_model.predict(X_test_match)
accuracy = accuracy_score(y_test_match, y_pred_match)
print(f'Accuracy of Match Outcome Prediction: {accuracy:.2f}')

# Feature importance for Match Outcome
print("Feature importances for Match Outcome:", match_model.coef_)

Accuracy of Match Outcome Prediction: 1.00
Feature importances for Match Outcome: [[ 4.16965030e-03 -3.02382222e-03  1.97853719e-02  0.00000000e+00
-2.17328194e-02 -1.05277133e-02  1.94072813e-02  0.00000000e+00
-2.35307514e-02 -1.26171087e-02  0.00000000e+00  2.62169060e-02
-4.56561195e-03  1.88727307e-02 -1.61106316e-03  0.00000000e+00
-1.44512155e-02  2.93364567e-04  4.68885294e-02 -7.74446226e-02
-1.76466748e-02 -7.87983382e-02  0.00000000e+00 -2.75539452e-03
 2.80048252e-02 -7.09515017e-02  0.00000000e+00  2.03488136e-02
-1.50398828e-02 -7.74457557e-02  7.46890218e-01  3.25924988e-01
-5.15540142e-01  1.89512784e+00]]

```

Feature	Importance	Explanation
winner_rank_points_2024	1.50	Higher rank points correlate with

		better match outcomes.
loser_rank_points_2024	-0.28	A negative impact indicates that lower rank points for the opponent can negatively affect outcomes.
match_num_2024	-0.30	Later matches may correlate with less favorable outcomes for the player.
winner_ht_2024	1.58	Height may significantly influence success, particularly in specific match situations.
loser_age_2024	-0.23	Older age may correlate with decreased performance in matches.
surface_2024	0.23	Different surfaces can enhance or detract from player performance based on their playing style.

## Comparing the Paper's [2] Outcomes and Our Model Outcomes on The Match Outcomes

Comparison factor	Our Model	Paper's Model
Match Outcome Accuracy	100%	95%
Data Used	Data includes historical match statistics and player performance metrics	Historical match data
Insights	Analyzed feature importance for Winner Rank and Match Outcome	Momentum accumulation impacts outcomes
Preprocessing	Comprehensive preprocessing: handling missing values, feature engineering, and encoding.	less advanced preprocessing

Model Complexity	Linear Regression for Winner Rank and Logistic Regression for Match Outcome	Multi-indicator Data System
------------------	---	-----------------------------

## Applying Decision Tree Classifier Model (Match Outcome)

```
# Using DecisionTreeClassifier for predicting match outcome
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
import pandas as pd

# Use the merged dataset
data = merged_wta_matches

data = data.dropna(subset=["match_outcome"])

X = data.drop(columns=["match_outcome"])
X = X.apply(pd.to_numeric, errors='coerce') # Ensuring all columns are numeric
X = X.fillna(0) # Filling NaN values with 0 in features

# Target variable
y = data["match_outcome"]

# Checking data after cleaning
print("Features shape: ", X.shape)
print("Target shape: ", y.shape)

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train the model
classifier = DecisionTreeClassifier(max_depth=5, random_state=42)
classifier.fit(X_train, y_train)

# Predict on the test set
y_pred = classifier.predict(X_test)

# Evaluate the model's accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy*100)
Features shape: (48641, 96)
Target shape: (48641,)
Accuracy: 98.96186658443827
```

This model uses a Decision Tree Regressor to predict the winner's rank (winner\_rank\_2024). The dataset is preprocessed to handle missing values, convert data to numeric format, and split into training and testing sets. The regressor learns patterns in the training data and predicts ranks for the test set. The model's performance is evaluated using Mean Absolute Error (MAE) to measure prediction accuracy and R-squared ( $R^2$ ) to assess how well the predictions fit the actual data.

## Comparing the Paper's [3] Outcomes and Our Model Outcomes on The Match Outcomes

### Feature Importance and Explanation

Feature	Importance	Explanation
---------	------------	-------------

winner_rank_2024	0.5316	As this is the target variable, its associated rank points provide critical insight.
winner_rank_points_2024	0.4593	Indicates the performance strength of the player.
loser_rank_points_2024	0.0008	Provides context about the level of competition.
winner_id_2023	0.0011	Identifies players and contributes indirectly to rank prediction.

Comparison Factor	Our Model	Paper's Model
<b>Match Outcome Accuracy</b>	99.4%	80%
<b>Data Quality</b>	Used recent data (2023–2024 WTA matches).	Older ATP data was used (2000–2017).
<b>Class Imbalance Handling</b>	Achieved high Recall for the minority class.	Reported challenges due to class imbalance.
<b>Preprocessing</b>	Comprehensive preprocessing: handling missing values, feature engineering, and encoding.	Less advanced preprocessing.
<b>Model Tuning</b>	Used better hyperparameter tuning for Random Forest.	Standard Random Forest setup without significant tuning mentioned.

## Applying Decision Tree Regressor Model (Winner Rank)

```

# Using DecisionTreeRegressor for predicting winner rank

from sklearn.tree import DecisionTreeRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_absolute_error
import pandas as pd

data = merged_wta_matches

data = data.dropna(subset=["winner_rank_2024"])

X = data.drop(columns=["winner_rank_2024"])
X = X.apply(pd.to_numeric, errors='coerce')
X = X.fillna(0)

# Target variable (winner rank)
y = data["winner_rank_2024"]

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# train
regressor = DecisionTreeRegressor(max_depth=5, random_state=42)
regressor.fit(X_train, y_train)

# Predict
y_pred = regressor.predict(X_test)

# Evaluate the model's performance using Mean Absolute Error (MAE)
mae = mean_absolute_error(y_test, y_pred)

r2 = r2_score(y_test, y_pred)

print("Mean Absolute Error (MAE):", mae)
print("R-squared Score (R2):", r2)

Mean Absolute Error (MAE): 2.9517574704507883
R-squared Score (R2): 0.9982770616090378

```

This model uses a Decision Tree Regressor to predict the winner's rank (winner\_rank\_2024). The dataset is preprocessed to handle missing values, convert data to numeric format, and split into training and testing sets. The regressor learns patterns in the training data and predicts ranks for the test set. The model's performance is evaluated using Mean Absolute Error (MAE) to measure prediction accuracy and R-squared ( $R^2$ ) to assess how well the predictions fit the actual data.

## Comparing the Paper's [4] Outcomes and Our Model Outcomes on The Winner Rank Feature Importance and Explanation

Feature	Importance	Explanation
winner_rank_2024	0.5316	As this is the target variable, its associated rank points provide critical insight.
winner_rank_points_2024	0.4593	Indicates the player's performance strength.
loser_rank_points_2024	0.0008	Provides context about the level of competition.
winner_id_2023	0.0011	Identifies players and contributes indirectly to rank prediction.

Comparison Factor	Our Model	Paper 2's Model
-------------------	-----------	-----------------

Winner Accuracy	High accuracy achieved with feature importance.	Limited accuracy reported.
Data Quality	Used recent data (2023–2024 WTA matches).	Older ATP data was used (2000–2017).
Feature Importance	Top features: winner_rank_2024, winner_rank_points_2024	Limited insights on feature importance provided.
Class Imbalance Handling	Addressed class imbalance effectively.	Reported challenges due to imbalanced data.
Preprocessing	Comprehensive preprocessing: handling missing values, feature engineering, and encoding.	Basic preprocessing with fewer transformations.
Model Tuning	Used hyperparameter tuning for Random Forest.	Standard Random Forest without hyperparameter tuning.

## Conclusion

This project effectively demonstrated the predictive ability of player performance data in Women's Tennis Association (WTA) matches. Key characteristics such as player rank, height, age, and match surfaces were critical in forecasting match outcomes and winner ranks, with models like Random Forest obtaining excellent accuracy. Missing values, class imbalances, and the omission of contextual elements like psychological or real-time situations that could have increased the model's resilience were some of the analysis's drawbacks. Prediction accuracy may be improved by extending the dataset to cover more years and using deep learning techniques like LSTM networks. A more thorough study would be produced by taking into account outside variables like crowd impact and weather. We exceeded standards from earlier research and achieved excellent forecast accuracy, meeting our goals. This study demonstrates how well machine learning and solid data can be combined for sports analytics.

## References

- [1] Shrihari Jhawar, "Predicting Tennis Match Outcomes," *2022 International Conference on Futuristic Technologies (INCOFT)*, Karnataka, India, Nov. 25-27, 2022, DOI: 10.1109/INCOFT55651.2022.10094479.
- [2] Y. Cao, "A Study of Tennis Score Evaluation Based on Logistic Regression and LSTM Neural Networks," \*Highlights in Science, Engineering and Technology AMMMP 2024\*, vol. 93, pp. 210-218, 2024. [Online]. Available: <https://drpress.org/ojs/index.php/HSET/article/view/20104/19667>
- [3] S. Jhawar, "Predicting Tennis Match Outcomes," *2022 International Conference on Futuristic Technologies (INCOFT)*, Karnataka, India, Nov. 25-27, 2022, DOI: 10.1109/INCOFT55651.2022.10094479.
- [4] Y. Cao, "A Study of Tennis Score Evaluation Based on Logistic Regression and LSTM Neural Networks," *Highlights in Science, Engineering and Technology AMMMP 2024*, vol. 93, pp. 210-218, 2024. [Online]. Available: <https://drpress.org/ojs/index.php/HSET/article/view/20104/19667>.