



# ANALYSIS OF WOMEN'S TENNIS ASSOCIATION MATCH DATA FOR THE YEAR 2023 AND 2024

AMANI ALBARAZI, MARAM ALHUSAMI, BUSHRA  
ALSHEHRI, AND SARA IMRAN  
SUPERVISED BY: DR.ZAIN Balfagih  
FALL-2024  
CS 3072

# CONTENTS

INTRODUCTION

PROBLEM STATEMENT

APPLYING BASIC PANDAS FUNCTIONS

DATA VISULATION

CLEANING DATA

MERGING THE TWO DATASETS

APPLYING MODELS

CONCLUSION

# INTRODUCTION

This project focuses on analyzing match data from the Women's Tennis Association (WTA) to develop predictive models for match outcomes and player rankings.

# PROBLEM STATEMENT

## Research Question:

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

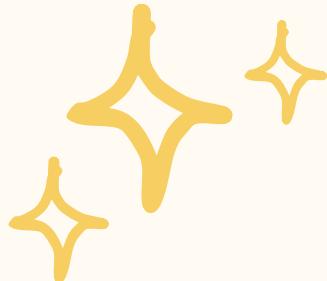
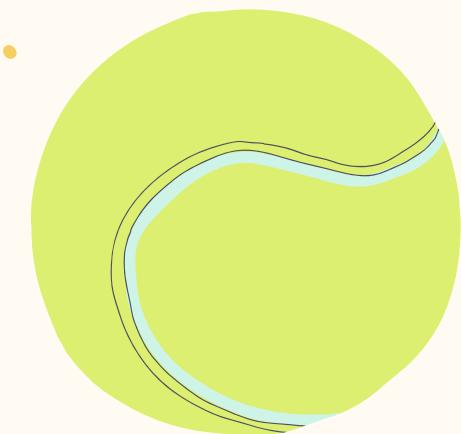
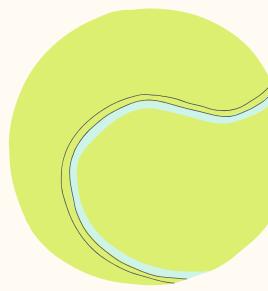
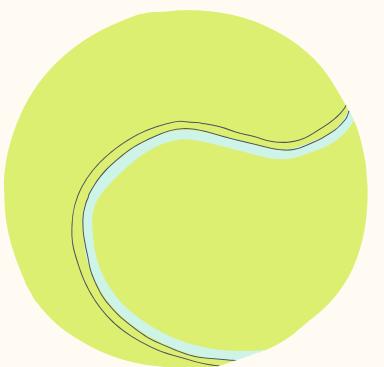
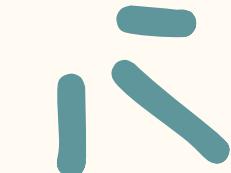
## Aim:

This project aims to predict two variables: winner rank and match outcome.

## Objectives:

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.

# APPLYING BASIC PANDAS FUNCTIONS



# WTA\_MATCHES\_QUAL\_ITF\_2023

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	...
	tourney_id	tourney_name	surface	draw_size	tourney_level	tourney_date	match_num	winner_id	winner_seed	winner_entry	winner_name	...
1	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	102	220578	NaN	Q	Julia Konishi	Camargo Silva
2	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	103	223367	NaN	NaN	Maria Fernanda Navarro	
3	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	106	206037	NaN	NaN	Daniela Seguel	
4	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	107	216075	NaN	Q	Jessica Bertoldo	
5	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	110	214605	NaN	Q	Agustina Chipac	
6	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	111	221657	NaN	Q	Camilla Zanolini	
7	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	114	223336	NaN	NaN	Weronika Bazzak	
8	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	115	214713	NaN	Q	Camila Romero	
9	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	118	211646	NaN	NaN	Oana Georgeata Simion	
10	2024-W-ITF-ARG-01A-2024	W35 Buenos Aires	Clay	48	35	20240115	119	214238	NaN	NaN	Verena Meliss	

10 rows × 48 columns

# WTA\_MATCHES\_QUAL\_ITF\_2023

```
wta_matches_2023.shape
```

```
(14319, 48)
```

```
wta_matches_2023.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
```

# WTA\_MATCHES\_QUAL\_ITF\_2024

```
wta_matches_2024.head()
```

	tourney_name	surface	draw_size	tourney_level	tourney_date	match_num	winner_id	winner_seed	winner_entry	winner_name	...
tourney_id											
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

# WTA\_MATCHES\_QUAL\_ITF\_2024

wta\_matches\_2024.shape

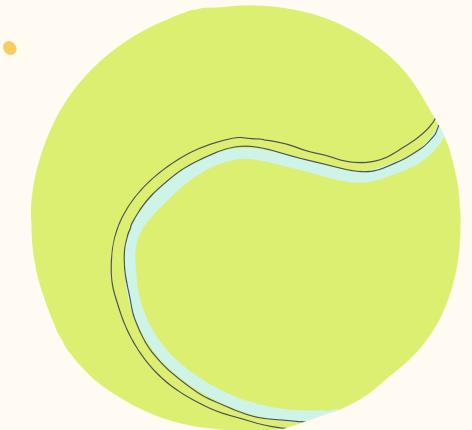
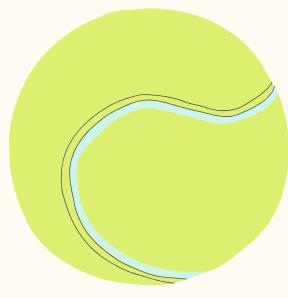
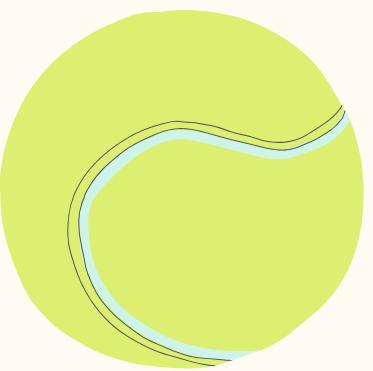
(14319, 48)

wta\_matches\_2024.describe()

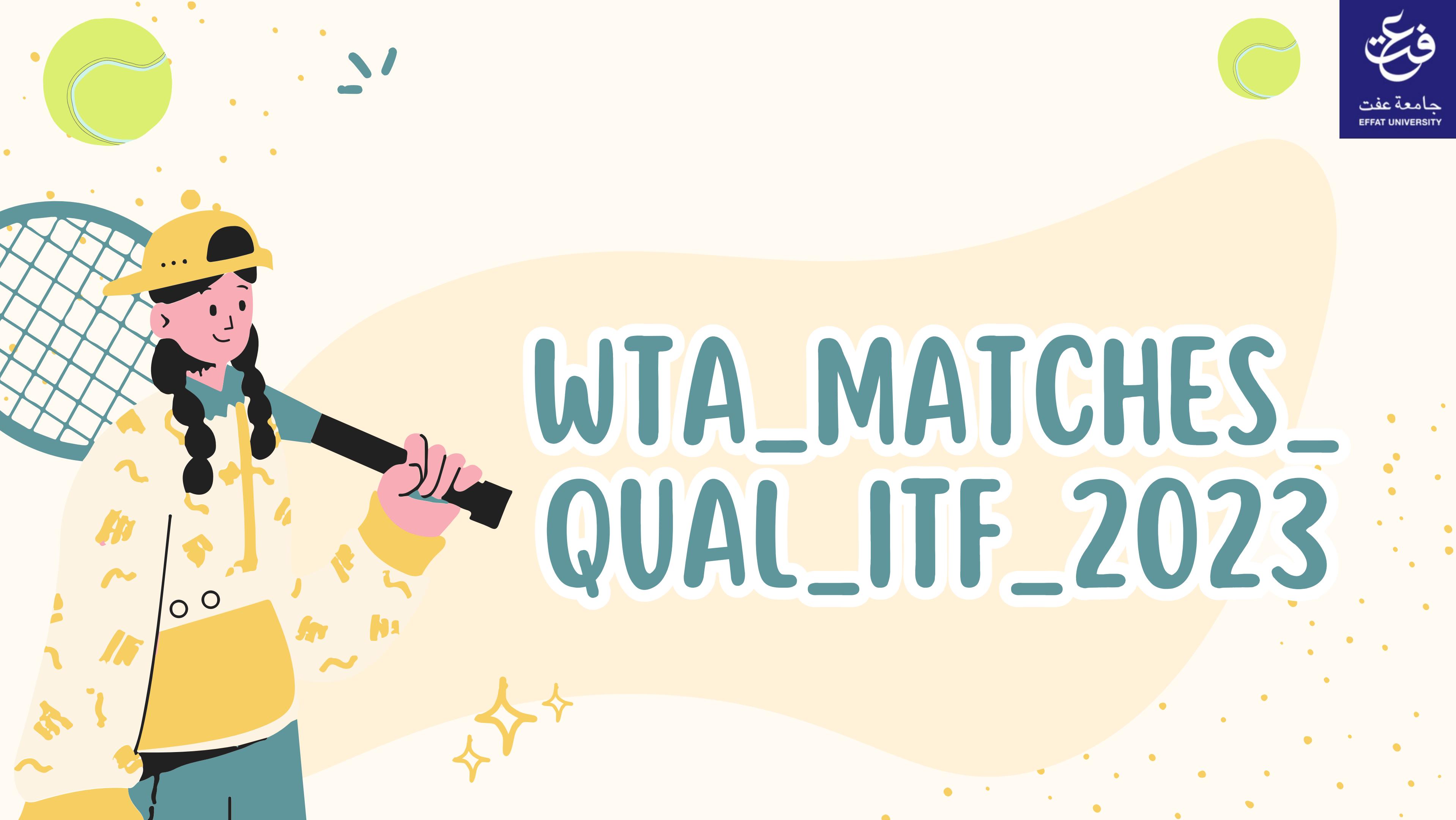
	draw_size	tourney_date	match_num	winner_id	winner_ht	winner_age	loser_id	loser_ht	loser_age	best_of
count	14319.000000	1.431900e+04	14319.000000	14319.000000	1744.000000	12136.000000	14319.000000	1367.000000	10782.000000	14319.0
mean	34.505203	2.024033e+07	185.131643	228164.922550	171.943234	23.615985	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

# DATA VISUALIZATION



# WTA\_MATCHES QUAL\_ITF\_2023

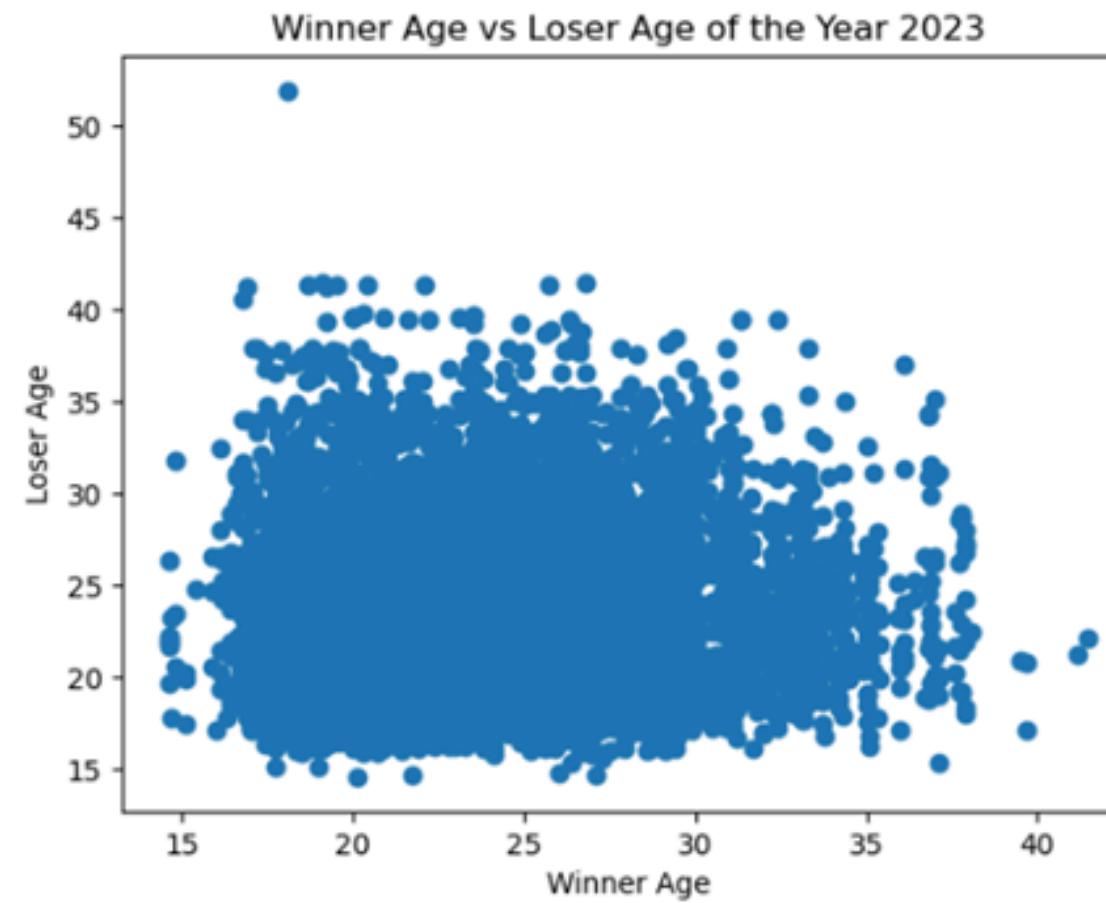


# A SCATTER PLOT OF WINNER AGE VS LOSER AGE

```
# 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')
```



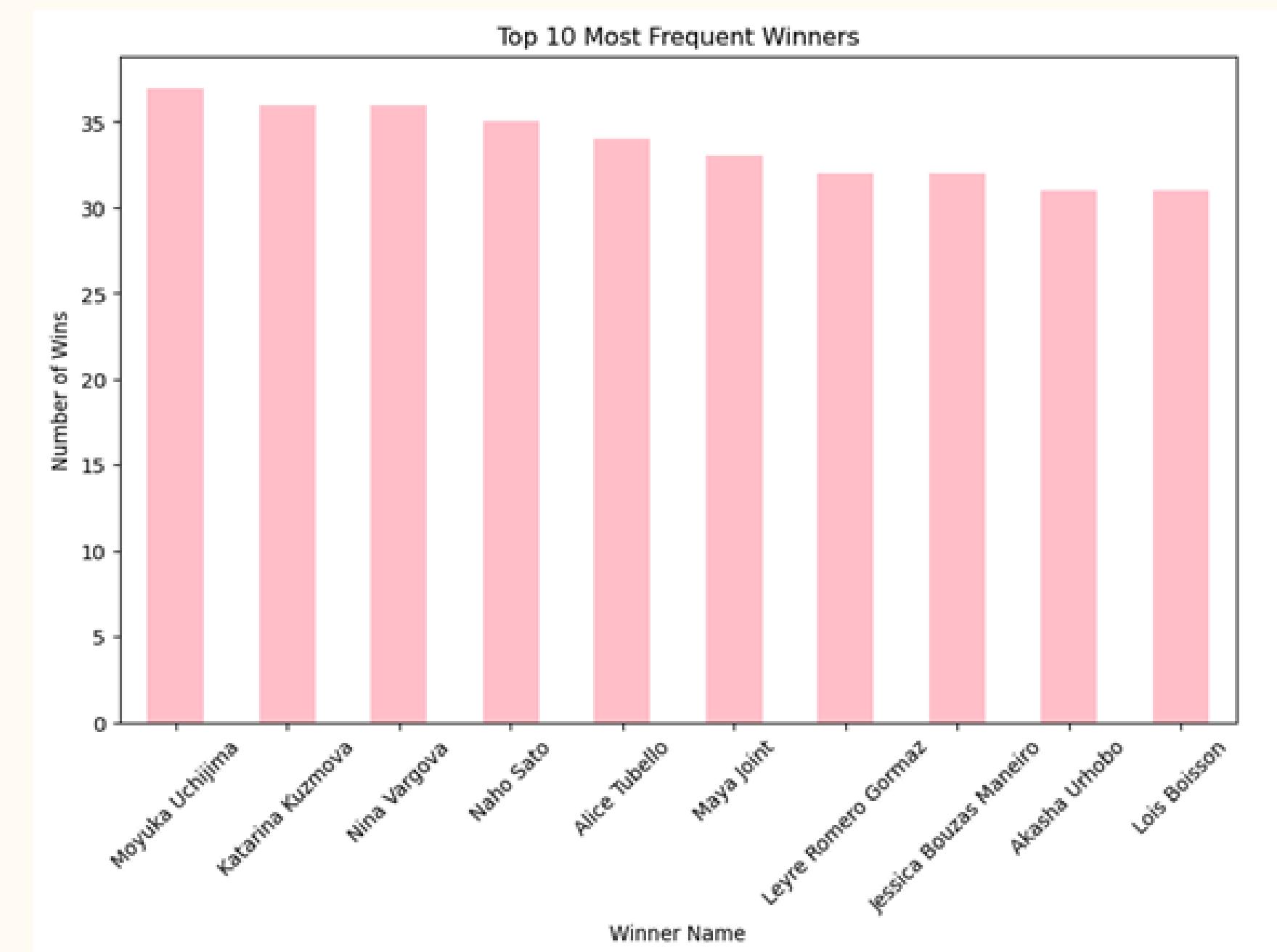
# TOP 10 MOST FREQUENT WINNERS

```
# 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()
```

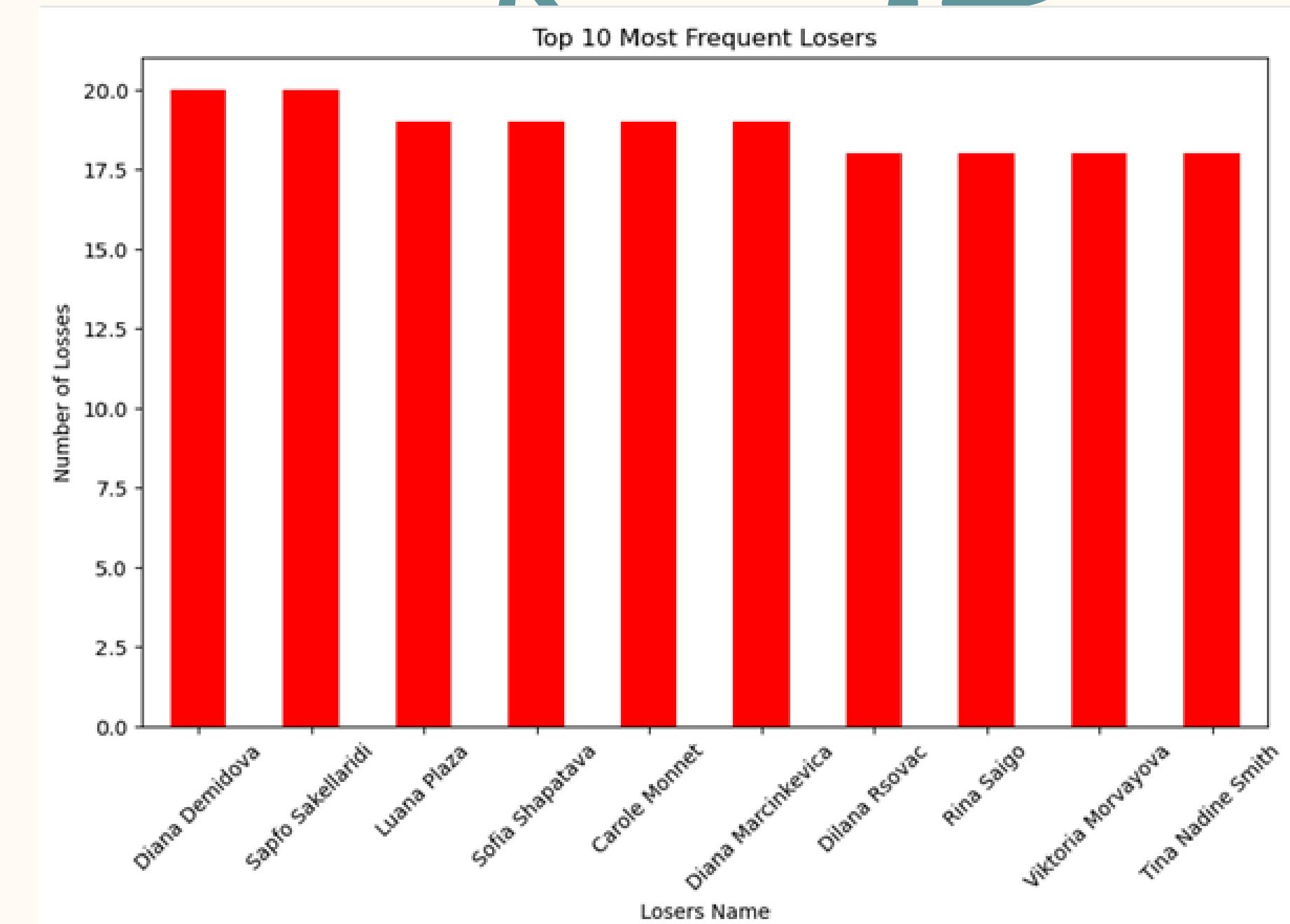


# MATCH DURATION TRENDS

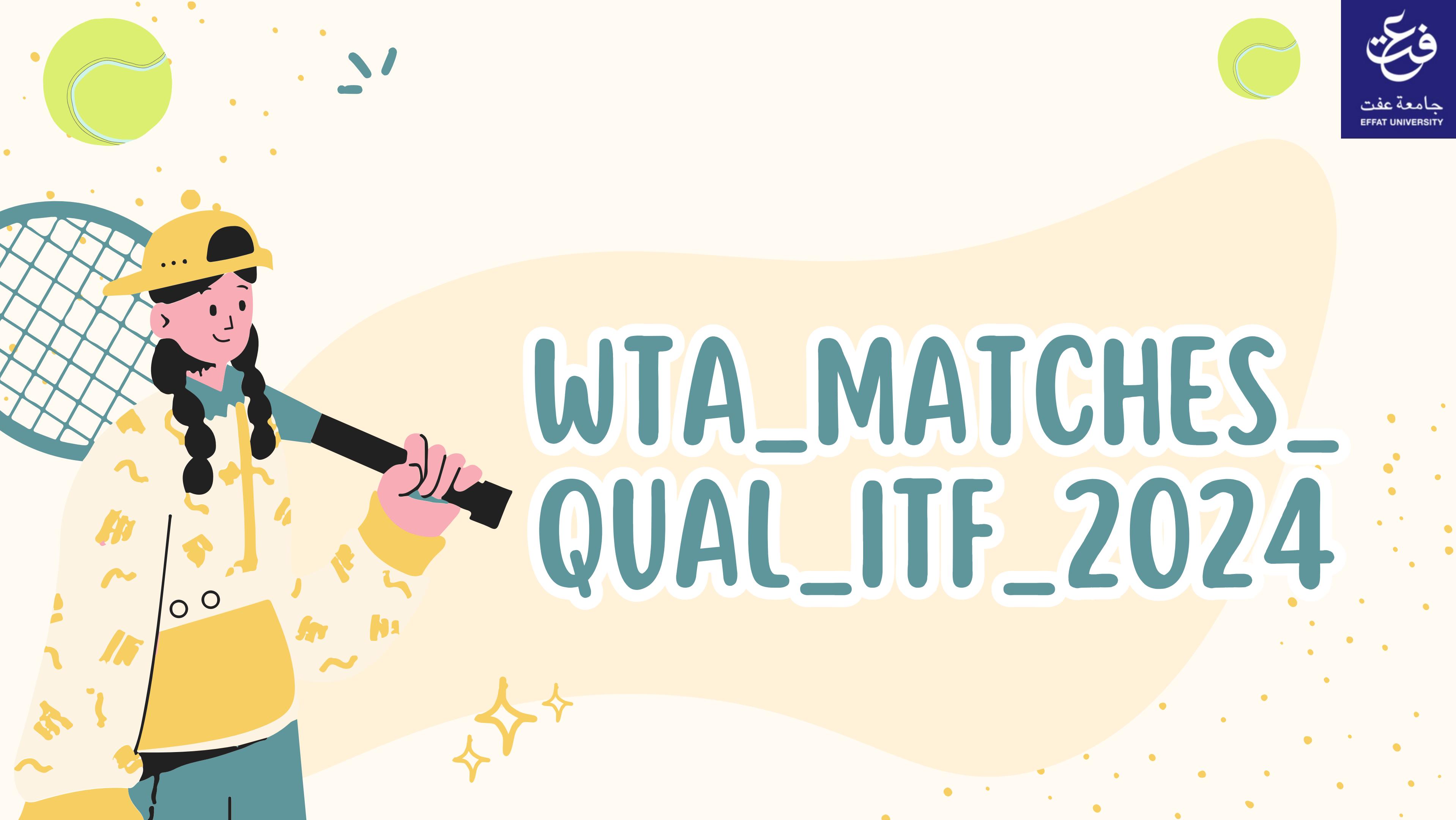
```
# Top 10 Most Frequent Losers
# Count top 10 losers by frequency
top_10_losers = wta_matches_2023['loser_name'].value_counts().head(10)

# 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()
```

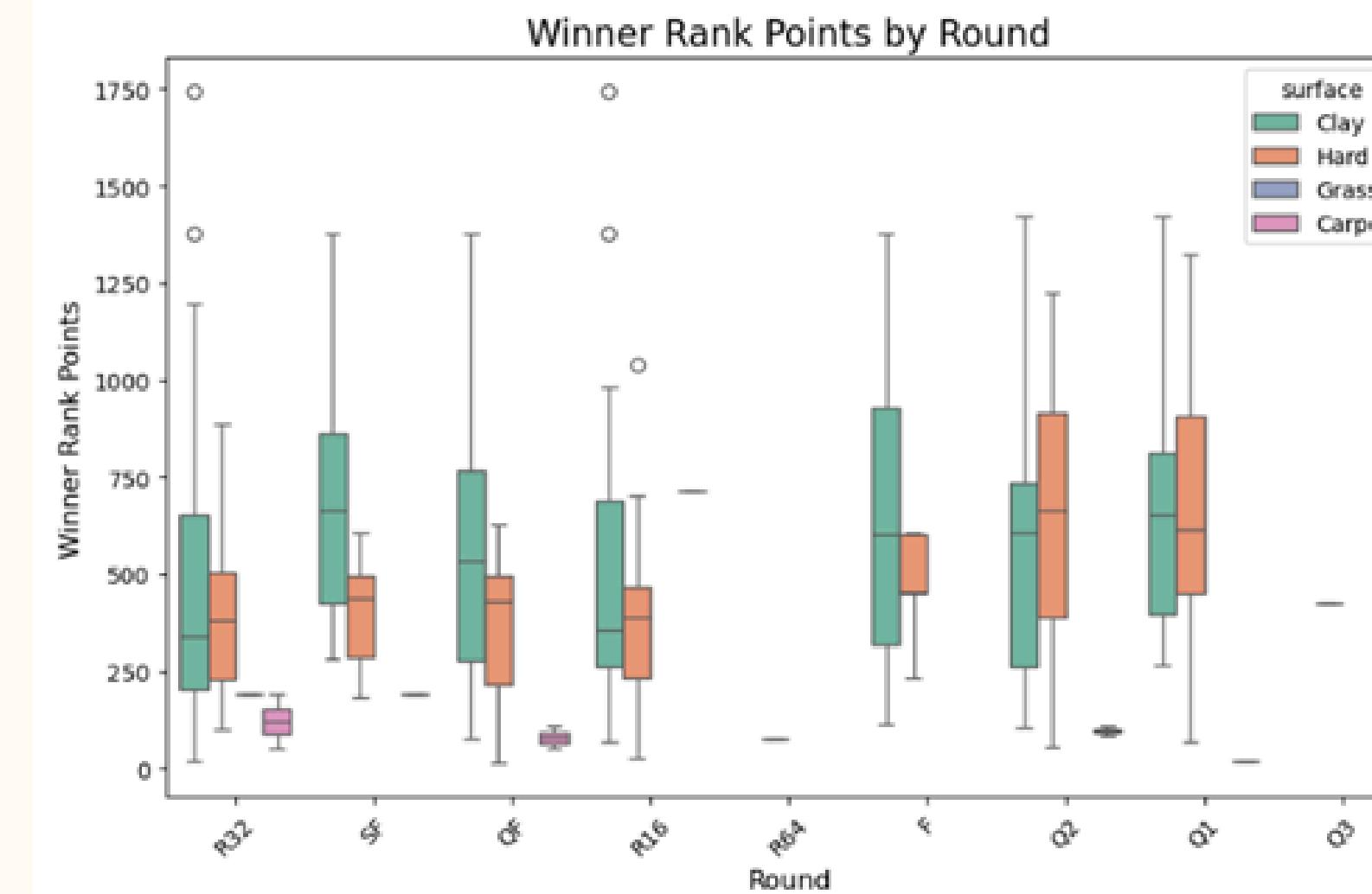


# WTA\_MATCHES QUAL\_ITF\_2024



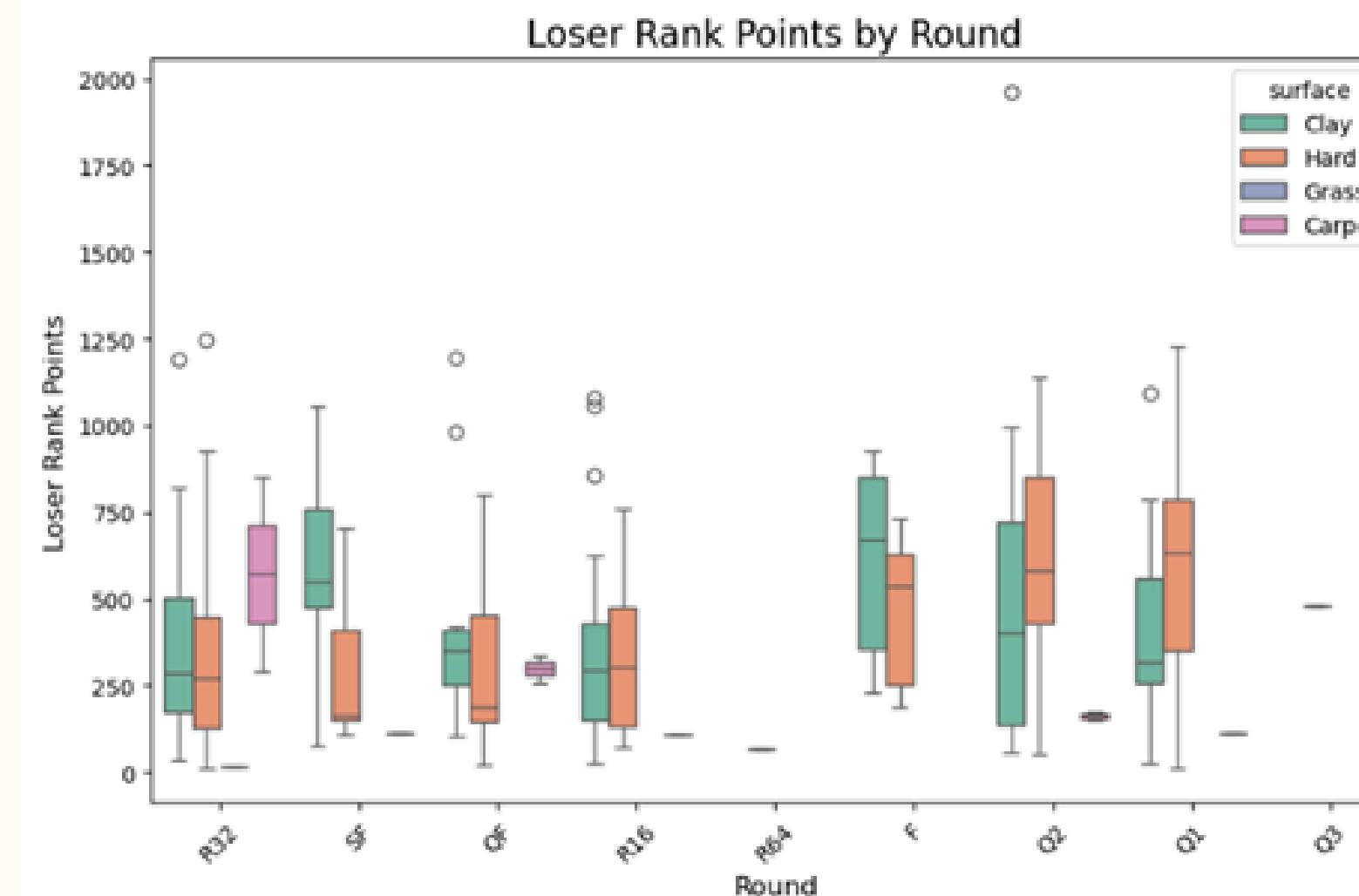
# BOXPLOT OF WINNER RANK POINTS BY ROUND

```
# 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()
```

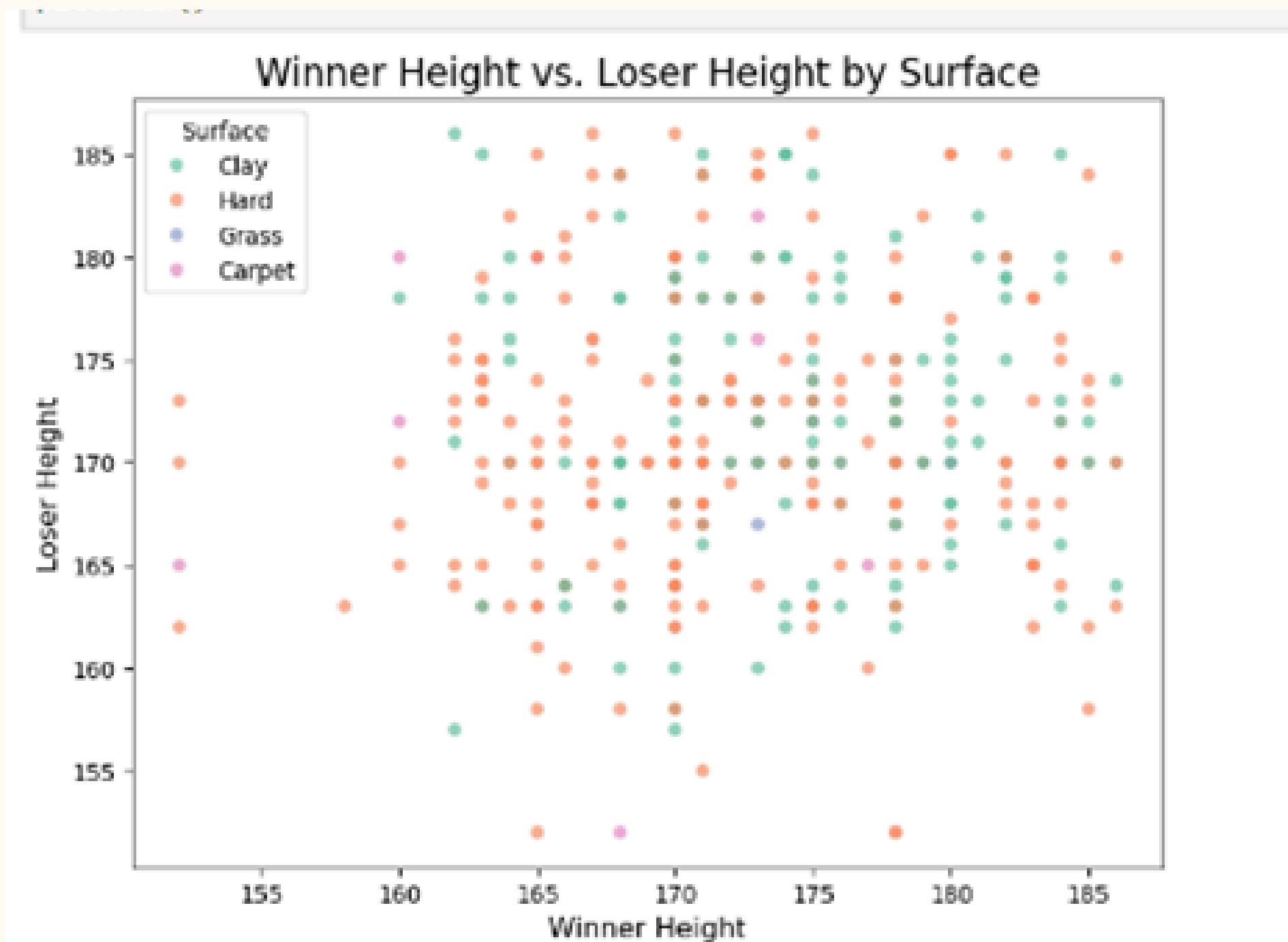


# BOXPLOT OF LOSER RANK POINTS BY ROUND

```
# 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()
```



# WTA\_MATCHES\_QUAL\_ITF\_2024



# CLEANING DATA

# WTA\_MATCHES\_QUAL\_ITF\_2023&24

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

wta_matches_2024.isnull().sum()	
tourney_name	0
surface	0
draw_size	0
tourney_level	0
tourney_date	0
match_num	0
winner_id	0
winner_seed	157
winner_entry	346
winner_name	0
winner_hand	0
winner_ht	0
winner_ioc	0
winner_age	0
loser_id	0
loser_seed	259
loser_entry	322
loser_name	0
loser_hand	0
loser_ht	0
loser_ioc	0
loser_age	0
score	0
best_of	0
round	0
minutes	210
w_ace	211
w_df	211
w_svpt	211
w_1stIn	211
w_1stWon	211
w_2ndWon	211
w_SvGms	211
w_bpSaved	211
w_bpFaced	211
l_ace	211
l_df	211
l_svpt	211
l_1stIn	211

# WTA\_MATCHES\_QUAL\_ITF\_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 unneccery 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')
```

# WTA\_MATCHES\_QUAL\_ITF\_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 unnecessery 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 Joao 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_ioc             object
winner_age              float64
loser_id              int64
loser_seed           object
loser_entry          object
loser_name            object
loser_hand            float64
loser_ht              float64
loser_ioc             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' 'WC' '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' 'WC' '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' 'WC' '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' 'WC' '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
```

# MERGING THE TWO DATASETS

٦٦

```
# 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

#Amani
# 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)
```



# APPLYING MODELS

# RANDOM FOREST REGRESSOR MODEL

## PROCESSING, SPLITTING, TRAINING, TESTING, AND SCALING

```

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_wta_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

```

# COMPARING THE RESULT

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

# RANDOM FOREST CLASSIFIER MODEL

## PROCESSING, SPLITTING, TRAINING, TESTING, AND SCALING

```

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: 0.9953758292005649

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

# LOGISTIC REGRESSION MODEL

```
[205]: # Apply Logistic Regression Algorithm

[209]: 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

[211]: # 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
```

DATA  
PREPROCESSING,  
DATA SPLITTING,  
MODEL TRAINING,  
OUTCOME  
PREDICTION

# LINEAR REGRESSION MODEL FOR WINNER RANK

```
[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_s
# 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

R<sup>2</sup> 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]

## DATA NORMALIZATION, MODEL TRAINING AND EVALUATION, FEATURE IMPORTANCE

# LOGISTIC REGRESSION MODEL FOR MATCH OUTCOME

```
[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, random_state=42)

# 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]]
```

DATA  
SPLITTING,  
MODEL  
TRAINING AND  
EVALUATION,  
FEATURE  
IMPORTANCE

# COMPARING THE RESULT

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

# DECISION TREE CLASSIFIER MODEL FOR 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)

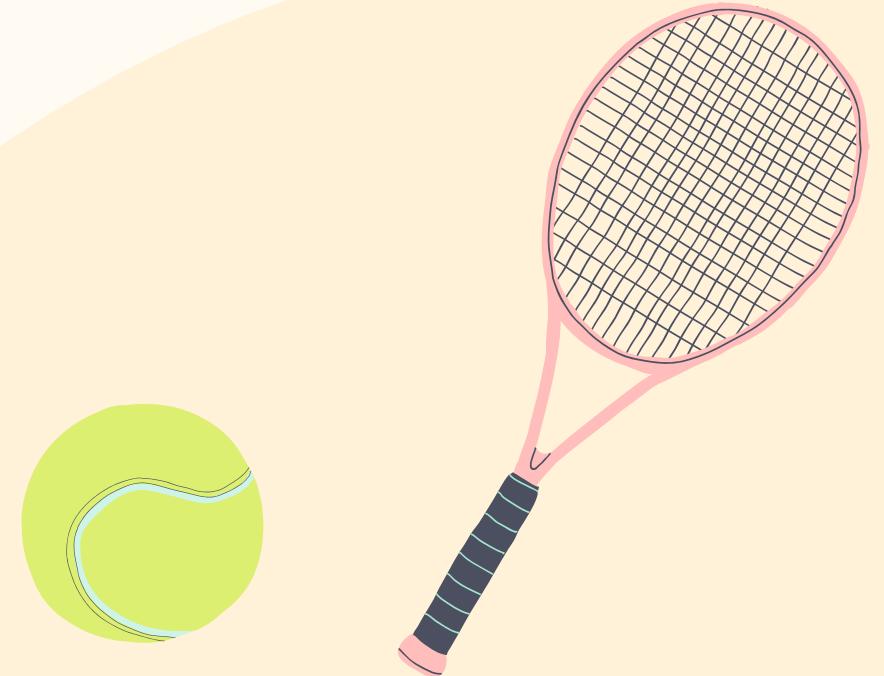
# Train the model
classifier = DecisionTreeClassifier(max_depth=5, random_state=42)
classifier.fit(X_train, y_train)

# 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
```

## MODEL TRAINING AND EVALUATION



# DECISION TREE REGRESSOR MODEL FOR 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)

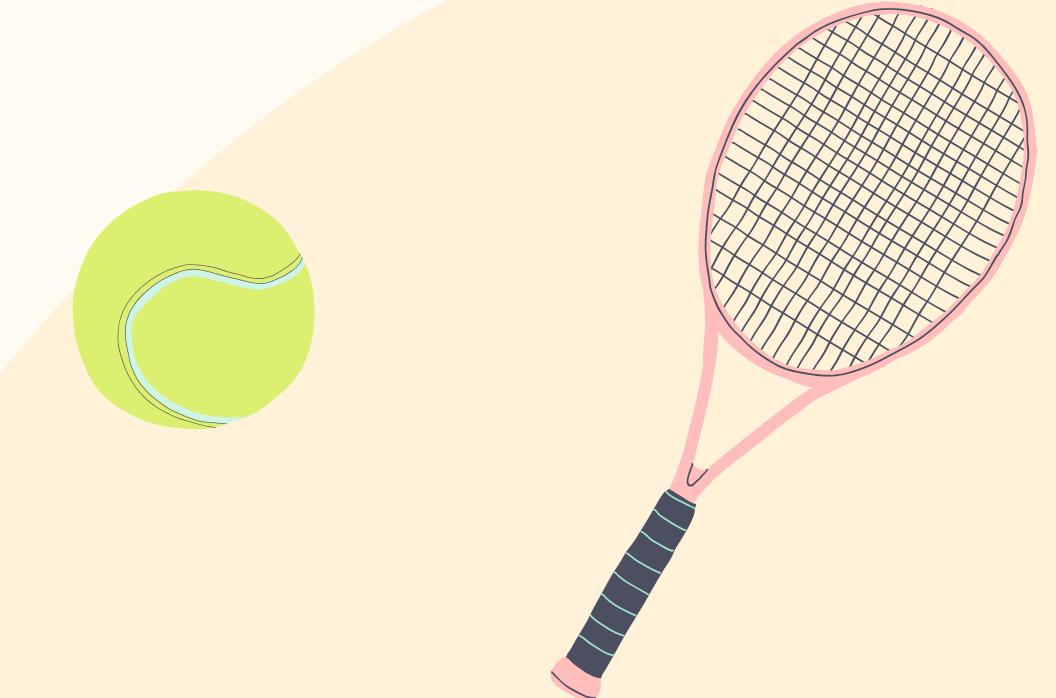
# Evaluate the model's performance using Mean Absolute Error (MAE)
mae = mean_absolute_error(y_test, y_pred)

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
```

## MODEL TRAINING AND EVALUATION



# COMPARING THE RESULTS

## FOR MATCH OUTCOME

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	Achieved high Recall for the minority 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.

## FOR WINNER RANK

Comparison Factor	Our Model	Paper 2's Model
Winner Rank 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 SUCCESSFULLY PREDICTED WTA MATCH OUTCOMES AND WINNER RANKS USING PLAYER PERFORMANCE DATA (E.G., RANK, HEIGHT, AGE, AND SURFACE).
- MODELS LIKE RANDOM FOREST ACHIEVED EXCELLENT ACCURACY.
- MISSING VALUES AND CLASS IMBALANCES
- OMISSION OF CONTEXTUAL FACTORS (E.G., PSYCHOLOGICAL STATES, CROWD IMPACT)
- EXPAND THE DATASET ACROSS MORE YEARS
- INCORPORATE DEEP LEARNING (E.G., LSTM) AND EXTERNAL VARIABLES (E.G., WEATHER)
- WE SURPASSED PREVIOUS RESEARCH, ACHIEVING HIGH ACCURACY, DEMONSTRATING THE POWER OF MACHINE LEARNING IN SPORTS ANALYTIC





THANK YOU

