

# Sprint 6 Task 1 (S06\_T01)

Author: Alberto Achaual

## Level 1

### Practice 1

Grab a sports theme dataset you like and select an attribute from the dataset. Calculate the mode, median, standard deviation, and arithmetic mean.

To solve this practice we start by **importing** all the required libraries:

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats
import statistics
```

For this practice, I'm going to use data from the ATP (Association of Professional Tennis) matches from the year 2018, the same one I used for Sprint 5.

In this data set, the winner and loser are coded with an ID. We have another dataset that references this ID with player information.

For the matches dataset we have the following column reference:

- ace = absolute number of aces
- df = number of double faults
- svpt = total serve points
- 1stIn = 1st serve in
- 1stWon = points won on 1st serve
- 2ndWon = points won on 2nd serve
- SvGms = serve games
- bpSaved = break point saved
- bpFaced = break point faced

```
source: https://www.kaggle.com/datasets/pablodroca/atp-matches-2000-2019?select=atp_matches_2019.csv
```

```
In [2]: ranking = pd.read_csv('atp_matches_2018.csv') # matches 2018
         pd.read_csv('atp_players.csv') # players information
```

```
In [3]: pd.set_option('display.max_columns', 50)
         ranking.sample(5)
```

	tourney_id	tourney_name	tourney_date	surface	winner_id	loser_id	score	best_of	round	minutes	w_ace	w_df	w_svpt	w_1stIn	w
183	2018-S80	Australian Open	20180115	Hard	111202	104999	6-2 4-1 7-6(4)	5	R128	48.0	3.0	2.0	35.0	21.0	
2768	2018-0328	Basel	20181022	Hard	104668	105208	7-6(4) 6-3 6(0)	3	R16	136.0	4.0	4.0	76.0	38.0	
2734	2018-0429	Stockholm	20181015	Hard	104259	104269	4-6 7-5 6(5)	3	R16	165.0	10.0	0.0	112.0	67.0	
2045	2018-0319	Kitzbuehl	20180730	Clay	104797	104259	5-7 6-3 6-7(5)	3	R16	129.0	6.0	6.0	93.0	64.0	
92	2018-0301	Auckland	20180108	Hard	105449	132482	7-5 6-3 6(4)	3	R32	104.0	12.0	1.0	78.0	43.0	

```
In [4]: ranking.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2889 entries, 0 to 2888
Data columns (total 32 columns):
 #   Column                Non-Null Count  Dtype  ---
 0   tourney_id            2889 non-null   object
 1   tourney_name          2889 non-null   object
 2   tourney_date          2889 non-null   object
 3   surface               2889 non-null   object
 4   winner_id             2889 non-null   int64
 5   loser_id              2889 non-null   object
 6   score                 2889 non-null   object
 7   best_of               2889 non-null   int64
 8   round                2889 non-null   object
 9   minutes               2889 non-null   float64
10  w_ace                 2863 non-null   float64
11  w_df                  2863 non-null   float64
12  w_svpt                2863 non-null   float64
13  w_1stIn               2863 non-null   float64
14  w_1stWon              2863 non-null   float64
15  w_2ndWon              2863 non-null   float64
16  w_SvGms               2863 non-null   float64
17  w_bpSaved              2863 non-null   float64
18  w_bpFaced              2863 non-null   float64
19  w_ace                 2863 non-null   float64
20  w_df                  2863 non-null   float64
21  w_svpt                2863 non-null   float64
22  w_1stIn               2863 non-null   float64
23  w_1stWon              2863 non-null   float64
24  w_2ndWon              2863 non-null   float64
25  w_SvGms               2863 non-null   float64
26  w_bpSaved              2863 non-null   float64
27  w_bpFaced              2863 non-null   float64
28  winner_rank           2882 non-null   float64
29  winner_rank_points     2882 non-null   float64
30  loser_rank             2858 non-null   float64
31  loser_rank_points      2858 non-null   float64
dtypes: float64(23), int64(4), object(5)
memory usage: 722.4+ KB
```

```
In [5]: players.head()
```

	player_id	name	first_name	last_name	hand	birthdate	country
0	100001	Gardnar	Mulloy	R	19131122.0	USA	
1	100002	Pancho	Segura	R	19210620.0	ECU	
2	100003	Frank	Sedgman	R	19271020.0	AUS	
3	100004	Giuseppe	Merlo	R	19271011.0	ITA	
4	100005	Richard Pancho	Gonzales	R	19280509.0	ITA	

```
In [6]: players.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 54067 entries, 0 to 54066
Data columns (total 6 columns):
 #   Column                Non-Null Count  Dtype  ---
 0   player_id             54067 non-null  int64
 1   name_first            53890 non-null  object
 2   name_last             54025 non-null  object
 3   hand                  48289 non-null  object
 4   birthdate             43099 non-null  float64
 5   country               54067 non-null  object
dtypes: float64(1), int64(1), object(4)
memory usage: 2.5+ MB
```

I want to merge both dataframes in order to replace in **matches** the **winner\_id** and **loser\_id** with actual name and contry from **players**.

```
In [7]: # new column for players with joined name
players['player'] = players['name_last'] + ', ' + players['name_first']
```

```
In [8]: # merge dataframes, first winner and then loser
winners = pd.merge(ranking, players.loc[:, ['player_id', 'country', 'player']],
                  left_on='winner_id', right_on='player_id').drop(columns = ['winner_id', 'player_id'])
atp_2018 = pd.merge(winners, players.loc[:, ['player_id', 'country', 'player']], how='inner',
                  left_on='loser_id', right_on='player_id').drop(columns = ['loser_id', 'player_id'])
```

```
In [9]: # reorder and rename columns
atp_2018 = atp_2018[['tourney_id', 'tourney_name', 'tourney_date', 'surface', 'player_x', 'country_x', 'player_y', 'score', 'best_of', 'round', 'minutes', 'w_ace', 'w_df', 'w_svpt', 'w_1stIn', 'w_1stWon', 'w_2ndWon', 'w_SvGms', 'w_bpSaved', 'w_bpFaced', 'w_ace', 'w_df', 'w_svpt', 'w_1stIn', 'w_1stWon', 'w_2ndWon', 'w_SvGms', 'w_bpSaved', 'w_bpFaced', 'winner_rank', 'winner_rank_point', 'loser_rank', 'loser_rank_points']]
atp_2018.rename(columns={'player_x': 'winner', 'country_x': 'country_winner', 'player_y': 'loser', 'country_y': 'country_loser'}, inplace=True)
```

```
In [10]: atp_2018.head()
```

	tourney_id	tourney_name	tourney_date	surface	winner	country_winner	loser	country_loser	score	best_of	round	minutes	w_ace
0	2018-M020	Brisbane	20180101	Hard	Harrison, Ryan	USA	Mayer, Leonardo	ARG	6-4 3-6 6-2	3	R32	123.0	
1	2018-6242	Winston Salem	20180820	Hard	Edmund, Kyle	GBR	Mayer, Leonardo	ARG	6-3 6-3 7(7)	3	R32	72.0	
2	2018-M016	Eastbourne	20180625	Grass	Millman, John	AUS	Mayer, Leonardo	ARG	6(7) 6-4	3	R16	121.0	
3	2018-M007	Miami Masters	20180319	Hard	Coric, Borna	CRO	Mayer, Leonardo	ARG	7(5) 6-3 6-4	3	R64	176.0	
4	2018-0414	Hamburg	20180723	Clay	Basilasch, Nikoloz	RUS	Mayer, Leonardo	ARG	6-4 0-6 7-5	3	F	132.0	

```
In [11]: atp_2018.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 2889 entries, 0 to 2888
Data columns (total 34 columns):
 #   Column                Non-Null Count  Dtype  ---
 0   tourney_id            2889 non-null   object
 1   tourney_name          2889 non-null   object
 2   tourney_date          2889 non-null   object
 3   surface               2889 non-null   object
 4   winner                2889 non-null   object
 5   country_winner         2889 non-null   object
 6   loser                 2889 non-null   object
 7   country_loser          2889 non-null   object
 8   score                 2889 non-null   object
 9   best_of               2889 non-null   int64
10   round                 2889 non-null   object
11   minutes               2854 non-null   float64
12   w_ace                 2863 non-null   float64
13   w_df                  2863 non-null   float64
14   w_svpt                2863 non-null   float64
15   w_1stIn               2863 non-null   float64
16   w_1stWon              2863 non-null   float64
17   w_2ndWon              2863 non-null   float64
18   w_SvGms               2863 non-null   float64
19   w_bpSaved              2863 non-null   float64
20   w_bpFaced              2863 non-null   float64
21   w_ace                 2863 non-null   float64
22   w_df                  2863 non-null   float64
23   w_svpt                2863 non-null   float64
24   w_1stIn               2863 non-null   float64
25   w_1stWon              2863 non-null   float64
26   w_2ndWon              2863 non-null   float64
27   w_SvGms               2863 non-null   float64
28   w_bpSaved              2863 non-null   float64
29   w_bpFaced              2863 non-null   float64
30   winner_rank           2882 non-null   float64
31   winner_rank_points     2882 non-null   float64
32   loser_rank            2858 non-null   float64
33   loser_rank_points      2858 non-null   float64
dtypes: float64(23), int64(2), object(9)
memory usage: 790.0+ KB
```

Now we have the dataframe as we wanted! (tt! Let's do the **mode**, **median**, **standard deviation**, and **arithmetic mean** for the number of aces per minute from the winner of the match. It will be new column **w\_ace/hour** (since we have matches "to the best of 3 games" and "to the best of 5 games").

Let's see about NaN values in **w\_ace** and **minutes** columns:

```
In [12]: atp_2018[atp_2018['w_ace'].isnull() | atp_2018['minutes'].isnull()]
```

	tourney_id	tourney_name	tourney_date	surface	winner	country_winner	loser	country_loser	score	best_of	round	minutes	w_ace
61	2018-M001	Sydney	20180108	Hard	Medvedev, Daniil	RUS	De Minaur, Alex	AUS	1-6 6-4 7-5	3	F	133.0	
276	2018-5014	Shanghai Masters	20181008	Hard	Cecchinato, Marco	ITA	Chung, Hyeon	KOR	4-6 6-7 6(5)	3	R32	NaN	
301	2018-5014	Shanghai Masters	20181008	Hard	Bautista Agut, Roberto	ESP	Mcdonald, Mackenzie	USA	3-6 3-6 6-4	3	R32	NaN	
309	2018-6242	Winston Salem	20180820	Hard	Harrison, Ryan	USA	Simon, Gilles	FRA	W/O	3	R32	NaN	
392	2018-5014	Shanghai Masters	20181008	Hard	Zverev, Alexander	GER	Basilasch, Nikoloz	RUS	7-5 6-4	3	R32	NaN	
395	2018-M-DC-2018-G2-IPA-M-101	Davis Cup G2 R1: GEO vs MAR	20180203	Clay	Ouahab, Lamine	ALG	Basilasch, Nikoloz	RUS	6-1 6-3	3	R16	NaN	
456	2018-M006	Indian Wells	20180305	Hard	Racic, Milos	CAN	Baghdatis, Marcos	CYP	W/O	3	R16	NaN	
498	2018-5014	Shanghai Masters	20181008	Hard	Querrey, Sam	USA	Fritz, Taylor Harry	USA	6-3 6-7 7-5	3	R32	NaN	
640	2018-5014	Shanghai Masters	20181008	Hard	Edmund, Kyle	GBR	Seppi, Andreas	ITA	6-3 6-4	3	R32	NaN	
727	2018-0311	London	20180618	Grass	Lopez, Feliciano	ESP	Racic, Milos	CAN	W/O	3	R16	NaN	
728	2018-0410	Monte Carlo Masters	20180416	Clay	Cilic, Marin	CRO	Racic, Milos	CAN	W/O	3	R16	NaN	
738	2018-0352	Paris Masters	20181029	Hard	Federer, Roger	SUI	Racic, Milos	CAN	W/O	3	R32	NaN	
779	2018-5014	Shanghai Masters	20181008	Hard	De Minaur, Alex	AUS	Paire, Benoit	FRA	6-4 6-3	3	R32	NaN	
888	2018-M035	Washington	20180730	Hard	De Minaur, Alex	AUS	Murray, Andy	GBR	W/O	3	QF	NaN	
900	2018-0451	Doha	20180101	Hard	Monfils, Gael	FRA	Rublev, Andrey	RUS	6-2 6-3	3	F	60.0	
1057	2018-5014	Shanghai Masters	20181008	Hard	Anderson, Kevin	RSA	Kukushkin, Mikhail	KAZ	6-2 6-3	3	R32	NaN	
1394	2018-0451	Doha	20180101	Hard	Monfils, Gael	FRA	Thiem, Dominic	AUT	W/O	3	SF	NaN	
1476	2018-5014	Shanghai Masters	20181008	Hard	Federer, Roger	SUI	Medvedev, Daniil	RUS	4-6 6-2 6-4	3	R32	NaN	
1545	2018-0407	Rotterdam	20180212	Hard	Goffin, David	BEL	Berdych, Tomas	CZE	W/O	3	QF	NaN	
1725	2018-0414	Hamburg	20180723	Clay	Jarry, Nicolas	CHI	Gasquet, Richard	FRA	W/O	3	R16	NaN	
1729	2018-5014	Shanghai Masters	20181008	Hard	Del Potro, Juan Martin	ARG	Gasquet, Richard	FRA	7-5 7-6(7)	3	R32	NaN	
1752	2018-0352	Paris Masters	20181029	Hard	Fognini, Fabio	ITA	Fucsovics, Marton	HUN	W/O	3	R32	NaN	
1833	2018-7485	Antwerp	20181015	Hard	Edmund, Kyle	GBR	Ivashka, Ilya	BLR	W/O	3	QF	NaN	
1916	2018-M015	Beijing	20181001	Hard	Del Potro, Juan Martin	ARG	Fognini, Fabio	ITA	W/O	3	SF	NaN	
2040	2018-M014	Moscow	20181015	Hard	Basir, Mirza	BIH	Kyrgios, Nick	AUS	W/O	3	R16	NaN	
2472	2018-5014	Shanghai Masters	20181008	Hard	Tsitipas, Stefanos	GRE	Khachanov, Karen	RUS	7-6 6-2 6(8)	3	R32	NaN	
2494	2018-M-DC-2018-G2-IPA-M-SLO-POL-01	Davis Cup G2 R1: SLO vs POL	20180203	Hard	Bedene, Aljaz	SLO	Majchrzak, Kamil	POL	6-3 6-4	3	R	NaN	
2637	2018-0337	Vienna	20181022	Hard	Anderson, Kevin	RSA	Melzer, Jurgen	AUT	W/O	3	R16	NaN	
2678	2018-5014	Shanghai Masters	20181008	Hard	Nishikori, Kei	JPN	Wu, Yibing	CHN	3-6 6-0 6-3	3	R32	NaN	
2757	2018-M-DC-2018-G1-AM-M-COL-BRA-01	Davis Cup G1 R2: COL vs BRA	20180406	Hard	Gonzalez, Alejandro	COL	Sorri, Jose Pedro	BRA	6-3 6-7 6(0)	3	RR	NaN	
2776	2018-M-DC-2018-G2-AM-M-BOL-PER-01	Davis Cup G2 R1: BOL vs PER	20180203	Clay	Panta Herreros, Jorge Brian	PER	Mendoza, Alejandro	BOL	6-2 6-4	3	RR	NaN	
2777	2018-M-DC-2018-G2-IPA-M-BOL-PER-01	Davis Cup G2 R1: BOL vs PER	20180203	Clay	Chavez, Villalpando, Luis Diego	ESP	Echazuz, Mauricio	PER	6-3 6-6 [10-8]	3	RR	NaN	
2792	2018-M-DC-2018-IRL-DEN-01	Davis Cup G2 R1: IRL vs DEN	20180203	Hard	Hess Olesen, Soren	DEN	Carr, Simon	IRL	2-6 6-6 [10-8]	3	RR	NaN	
2804	2018-M-DC-2018-G2-AM-M-LBN-TPE-01	Davis Cup G2 R1: LBN vs TPE	20180203	Hard	Chen, Ti	TPE	Ballout, Jad	TPE	6-1 6-1	3	RR	NaN	
2826	2018-M-DC-2018-G2-IPA-M-SLO-POL-01	Davis Cup G2 R1: SLO vs POL	20180203	Hard	Hurkacz, Hubert	POL	Kocovar Desman, Tom	SLO	6-3 6-2	3	RR	NaN	
2840	2018-M-DC-2018-G2-IPA-M-TUN-FIN-01	Davis Cup G2 R1: TUN vs FIN	20180203	Hard	Hellouvaara, Harri	FIN	Chargui, Moez	TUN	3-6 6-6 6-4	3	RR	NaN	
2843	2018-M-DC-2018-G2-IPA-M-TUR-ZIM-01	Davis Cup G2 R1: TUR vs ZIM	20180203	Hard	Agabigun, Sarp	TUR	Sibahidi, Mehdi Don Ayanda	ZIM	7-5 6-4	3	RR	NaN	

Half of these rows are W/O (walk over) matches in **score** column. This means matches that have not been played or finished early due to some problem. Some other rows are from matches in early stages of the tournament where no data has been collected. We could fill these NaN values with the mean from the rest, but early stages of a tournament indicates ""not so good players"" (among a lot of quotes) and it might not be a good idea to replace them with the average aces of winners in advanced stages of the tournament. Both these rows are missing **w\_ace** and **minutes** data.

Conclusion: I'll drop these rows and the remaining ones were only the **minutes** are missing I will fill them with mean. In this case only in matches "to the best of 3" are missing the minutes.

```
In [13]: atp_2018.dropna(subset = ['w_ace'], how = 'any', inplace = True) # rows with NaN values in w_ace= column
```

```
In [14]: atp_2018[atp_2018['w_ace'].isnull() | atp_2018['minutes'].isnull()] # rows with NaN in -minutes- column
```

	tourney_id	tourney_name	tourney_date	surface	winner	country_winner	loser	country_loser	score	best_of	round	minutes	w_ace
276	2018-5014	Shanghai Masters	20181008	Hard	Cecchinato, Marco	ITA	Chung, Hyeon	KOR	4-6 7-6(5) 6(5)	3	R32	NaN	
301	2018-5014	Shanghai Masters	20181008	Hard	Bautista Agut, Roberto	ESP	Mcdonald, Mackenzie	USA	3-6 6-4 6-1	3	R32	NaN	
392	2018-5014	Shanghai Masters	20181008	Hard	Zverev, Alexander	GER	Basilasch, Nikoloz	RUS	7-5 6-4	3	R32	NaN	
498	2018-5014	Shanghai Masters	20181008	Hard	Querrey, Sam	USA	Fritz, Taylor Harry	USA	6-3 6-7 7(4)	3	R32	NaN	
640	2018-5014	Shanghai Masters	20181008	Hard	Edmund, Kyle	GBR	Seppi, Andreas	ITA	6-3 6-4	3	R32	NaN	
779	2018-5014	Shanghai Masters	20181008	Hard	De Minaur, Alex	AUS	Paire, Benoit	FRA	6-4 6-3	3	R32	NaN	
1057	2018-5014	Shanghai Masters	20181008	Hard	Anderson, Kevin	RSA	Kukushkin, Mikhail	KAZ	6-3 6-2	3	R32	NaN	
1476	2018-5014	Shanghai Masters	20181008	Hard	Federer, Roger	SUI	Medvedev, Daniil	RUS	4-6 6-4 6-4	3	R32	NaN	
1729	2018-5014	Shanghai Masters	20181008	Hard	Del Potro, Juan Martin	ARG	Gasquet, Richard	FRA	7-5 7-6(7)	3	R32	NaN	
2472	2018-5014	Shanghai Masters	20181008	Hard	Tsitipas, Stefanos	GRE	Khachanov, Karen	RUS	6-4 7-6 6(8)	3	R32	NaN	
2678	2018-5014	Shanghai Masters	20181008	Hard	Nishikori, Kei	JPN	Wu, Yibing	CHN	3-6 6-0 6-3	3	R32	NaN	

```
In [15]: grouped_mins = atp_2018.groupby('best_of')['minutes'].mean() # excluding missing values by default
grouped_mins
```

```
best_of
3      99.371666
5     156.341593
Name: minutes, dtype: float64
```

```
In [16]: mean_bo3 =
```