

Final Project

Alexia Crisologo and Olivia Encarnacion

```
library(dplyr)
```

```
Attaching package: 'dplyr'
```

```
The following objects are masked from 'package:stats':
```

```
filter, lag
```

```
The following objects are masked from 'package:base':
```

```
intersect, setdiff, setequal, union
```

```
library(tidyverse)
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
```

```
v forcats 1.0.0    v readr      2.1.4
```

```
v ggplot2  3.4.3    v stringr    1.5.0
```

```
v lubridate 1.9.2    v tibble     3.2.1
```

```
v purrr     1.0.2    v tidyr      1.3.0
```

```
-- Conflicts ----- tidyverse_conflicts() --
```

```
x dplyr::filter() masks stats::filter()
```

```
x dplyr::lag()     masks stats::lag()
```

```
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
library(ggplot2)
tennis <- read.csv("data/wta_matches_qual_itf_2023.csv")
```

Introduction and Data

Research Question

Does the court type impact the duration of a women's single tennis match, considering the winner's height, age, and playing hand?

The data we chose was the Women's Tennis Association data set from Awesomedata's GitHub repository (https://github.com/JeffSackmann/tennis_wta/blob/master/wta_matches_qual_itf_2023.csv). The data was created in 2023 and was collected from the International Tennis Federation. The data contains 34,323 observations and 49 variables. The variables of interest in our research include surface, the winner and loser id, winner and loser IOC (country), winner and loser height, winner and loser age, and winner and loser hand. These variables will help us answer the question of the court types impact on the duration of tennis matches. There are some NA variables corresponding to the players height. Since this is our variable of interest, we will be dropping all NA values corresponding to height. This will leave us with 1,256 observations. Additionally, since grand slams are played on either clay, grass, or hard courts, we will be dropping the matches that were played on carpet. This will then leave us with 1,239 observations. The motivation behind this project is to help Women's Tennis players better prepare for the length of their match based on the surface they will be playing on. With the Olympics coming up, this data will help the tennis player better prepare for a match.

```
tennis <- tennis %>%
  filter(!is.na(winner_ht) & !is.na(loser_ht) & surface != "Carpet")
```

Variables of Interest

Surface: Surface the match was played on (clay, grass, or hard)

Winner_id: Identification of the winner

Winner_hand: Dominant playing hand of the winner (right, left, undecided)

Winner_ht: Height of the winner in centimeters (cm)

Winner_ioc: Country of the winner as assigned by the International Olympic Committee

Winner_age: Age of the winner

Loser_id: Identification of the loser

Loser_hand: Dominant playing hand of the loser (right, left, or undecided)

Loser_ht: Height of the loser in centimeters (cm)

Loser_ioc: Country of the loser as assigned by the International Olympic Committee

Loser_age: Age of the loser

Minutes: Duration of the tennis match in minutes

```
winner_hand_counts <- tennis |>
  count(winner_hand)

loser_hand_counts <- tennis |>
  count(loser_hand)

print(winner_hand_counts)
```

	winner_hand	n
1	L	116
2	R	1110
3	U	13

```
print(loser_hand_counts)
```

	loser_hand	n
1	L	107
2	R	1107
3	U	25

Methodology

```
ggplot(data = tennis, aes(x = surface, y = minutes)) +
  geom_boxplot() +
  labs(title = "Match Duration by Court Surface",
       x = "Surface",
       y = "Match Duration (minutes)")
```

Warning: Removed 620 rows containing non-finite values (``stat_boxplot()``).

