

EDA Proposal Statistical

Anant Patel - 0866771

```
#libraries  
library(lubridate)
```

Warning: package 'lubridate' was built under R version 4.4.2

Attaching package: 'lubridate'

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

date, intersect, setdiff, union

```
library(dplyr)
```

Warning: package 'dplyr' was built under R version 4.4.2

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(ggplot2)
```

Warning: package 'ggplot2' was built under R version 4.4.2

```
library(GGally)
```

Warning: package 'GGally' was built under R version 4.4.2

```
Registered S3 method overwritten by 'GGally':  
  method from  
+ .gg      ggplot2
```

```
tuesdata <- tidyTuesdayR::tt_load('2023-04-25')
```

```
---- Compiling #TidyTuesday Information for 2023-04-25 ----  
--- There are 2 files available ---
```

```
-- Downloading files -----  
  
1 of 2: "winners.csv"  
2 of 2: "london_marathon.csv"
```

```
tuesdata <- tidyTuesdayR::tt_load(2023, week = 17)
```

```
---- Compiling #TidyTuesday Information for 2023-04-25 ----  
--- There are 2 files available ---
```

```
-- Downloading files -----  
  
1 of 2: "winners.csv"  
2 of 2: "london_marathon.csv"
```

```
winners <- tuesdata$winners  
london_marathon <- tuesdata$london_marathon
```

```
View(winners)
View(london_marathon)
```

```
winners$Time.Seconds <- period_to_seconds(hms(winners$Time))

str(winners)
```

```
spc_tbl_ [163 x 6] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
 $ Category      : chr [1:163] "Men" "Men" "Men" "Men" ...
 $ Year          : num [1:163] 1981 1981 1982 1983 1984 ...
 $ Athlete       : chr [1:163] "Dick Beardsley (Tie)" "Inge Simonsen (Tie)" "Hugh Jones" "Mike
 $ Nationality   : chr [1:163] "United States" "Norway" "United Kingdom" "United Kingdom" ...
 $ Time          : 'hms' num [1:163] 02:11:48 02:11:48 02:09:24 02:09:43 ...
 ..- attr(*, "units")= chr "secs"
 $ Time.Seconds: num [1:163] 7908 7908 7764 7783 7797 ...
 - attr(*, "spec")=
 .. cols(
 ..   Category = col_character(),
 ..   Year = col_double(),
 ..   Athlete = col_character(),
 ..   Nationality = col_character(),
 ..   Time = col_time(format = "")
 .. )
 - attr(*, "problems")=<externalptr>
```

```
# Factoring the variables
winners$Category <- factor(winners$Category)
winners$Athlete <- factor(winners$Athlete)
winners$Nationality <- factor(winners$Nationality)

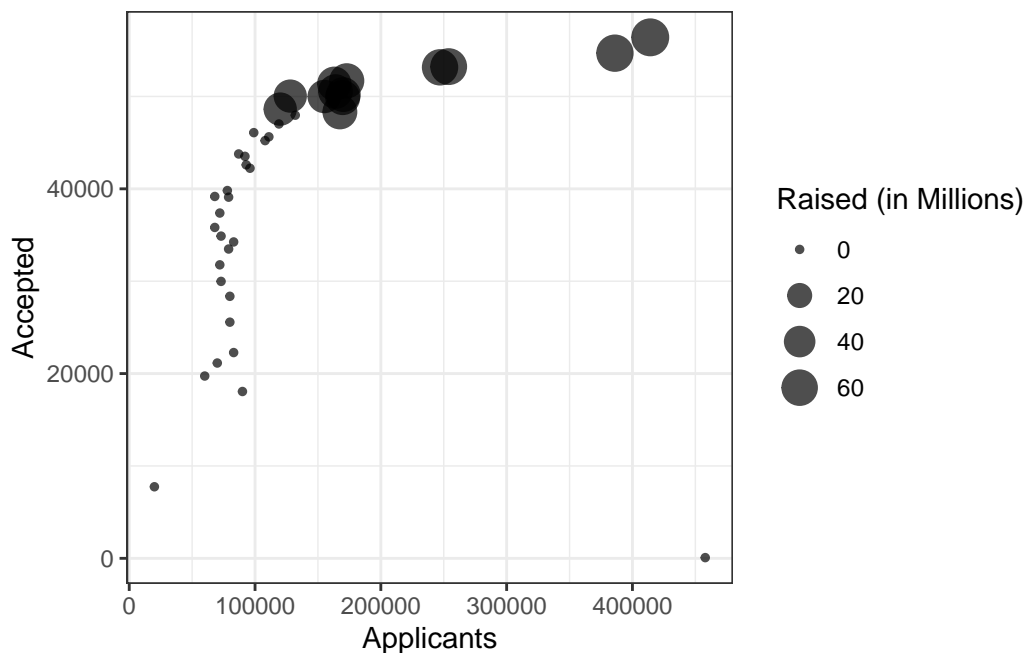
str(london_marathon)
```

```
spc_tbl_ [42 x 8] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
 $ Date          : Date[1:42], format: "1981-03-29" "1982-05-09" ...
 $ Year          : num [1:42] 1981 1982 1983 1984 1985 ...
 $ Applicants    : num [1:42] 20000 90000 60000 70000 83000 80000 80000 73000 72000 73000
 $ Accepted      : num [1:42] 7747 18059 19735 21142 22274 ...
 $ Starters      : num [1:42] 7055 16350 16500 16992 17500 ...
 $ Finishers     : num [1:42] 6255 15116 15793 15675 15873 ...
 $ Raised        : num [1:42] NA NA NA NA NA NA NA NA NA ...
 $ Official charity: chr [1:42] NA NA NA NA ...
```

```
- attr(*, "spec")=
.. cols(
..   Date = col_date(format = ""),
..   Year = col_double(),
..   Applicants = col_double(),
..   Accepted = col_double(),
..   Starters = col_double(),
..   Finishers = col_double(),
..   Raised = col_double(),
..   `Official charity` = col_character()
.. )
- attr(*, "problems")=<externalptr>
```

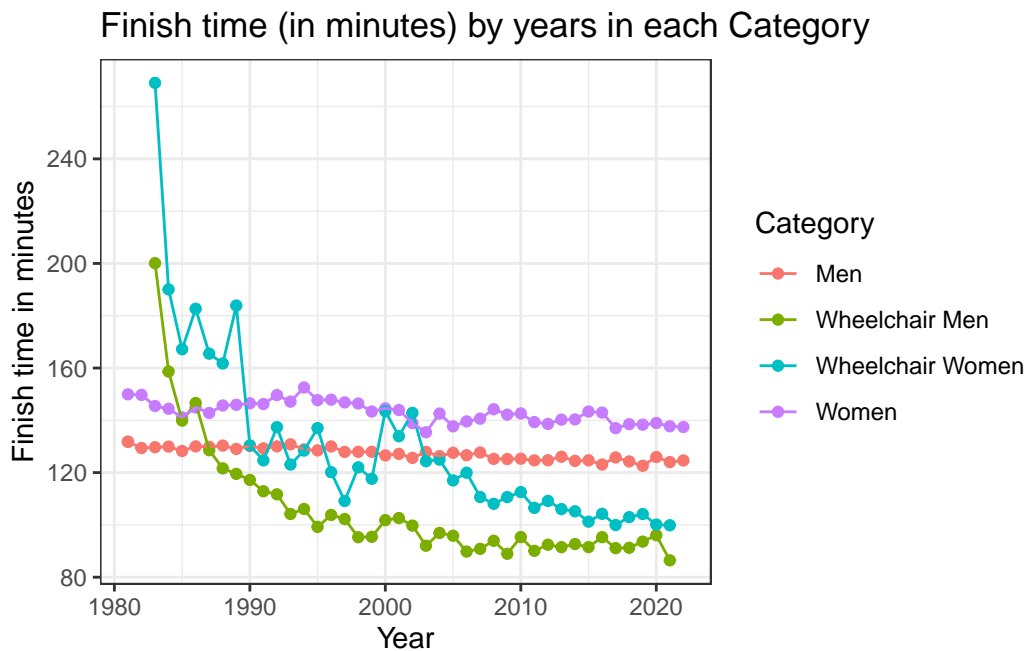
```
# Handling the NA in Raised
london_marathon$Raised[is.na(london_marathon$Raised)] = 0
london_marathon <- london_marathon[rowSums(is.na(london_marathon)) <= 2,]
```

```
options(scipen = 999)
# Accepted participants vs finishers by amount raised
london_marathon %>% ggplot(aes(x=Applicants, y = Accepted, size = Raised)) +
  geom_point(alpha = 0.7) +
  scale_size_continuous(name = "Raised (in Millions)") +
  theme_bw()
```



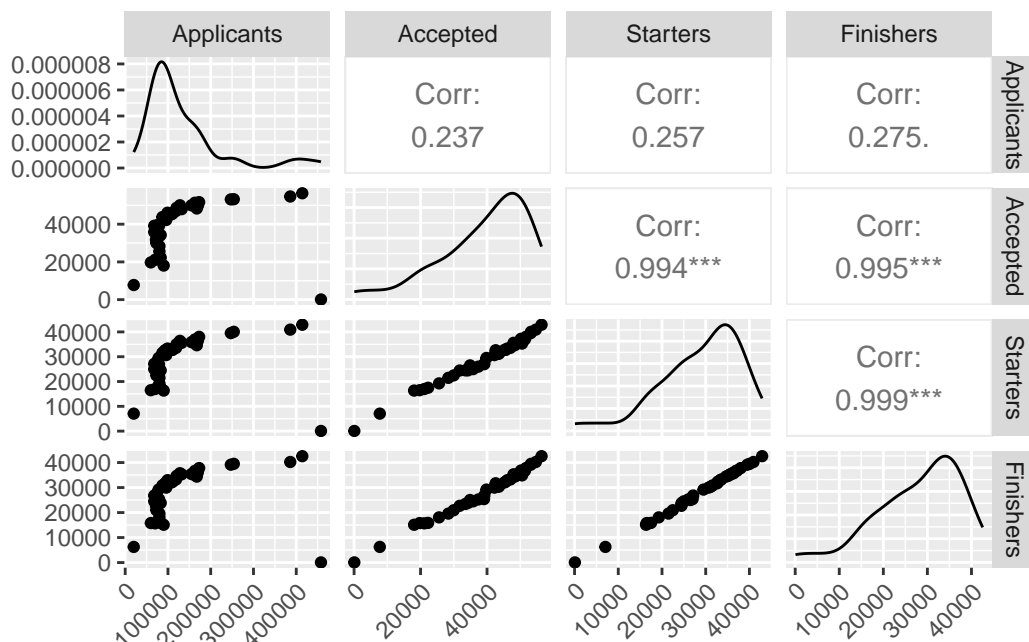
Question: Does the amount gets raised when the applicants are accepted more? **OR** Claim: when the applicants are accepted more the amount is raised.

```
# Year vs Time by Category
winners %>%
  ggplot(aes(x = Year, y = Time.Seconds / 60, color = Category)) +
  geom_point() +
  geom_line() +
  labs(
    title = "Finish time (in minutes) by years in each Category",
    x = "Year",
    y = "Finish time in minutes"
  ) +
  theme_bw()
```



Question: Wheelchair individuals have some correlation with time to finish the race

```
options(scipen=10)
ggpairs(london_marathon[,c("Applicants", "Accepted", "Starters", "Finishers")]) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1)) # chat gpt helped me
```



OpenAI. (2024). ChatGPT [Large language model]. <https://chatgpt.com>

(OpenAI, 2024)

Question: Can we predict the Finishers based on the Starters and Accepted. (Linear Model)

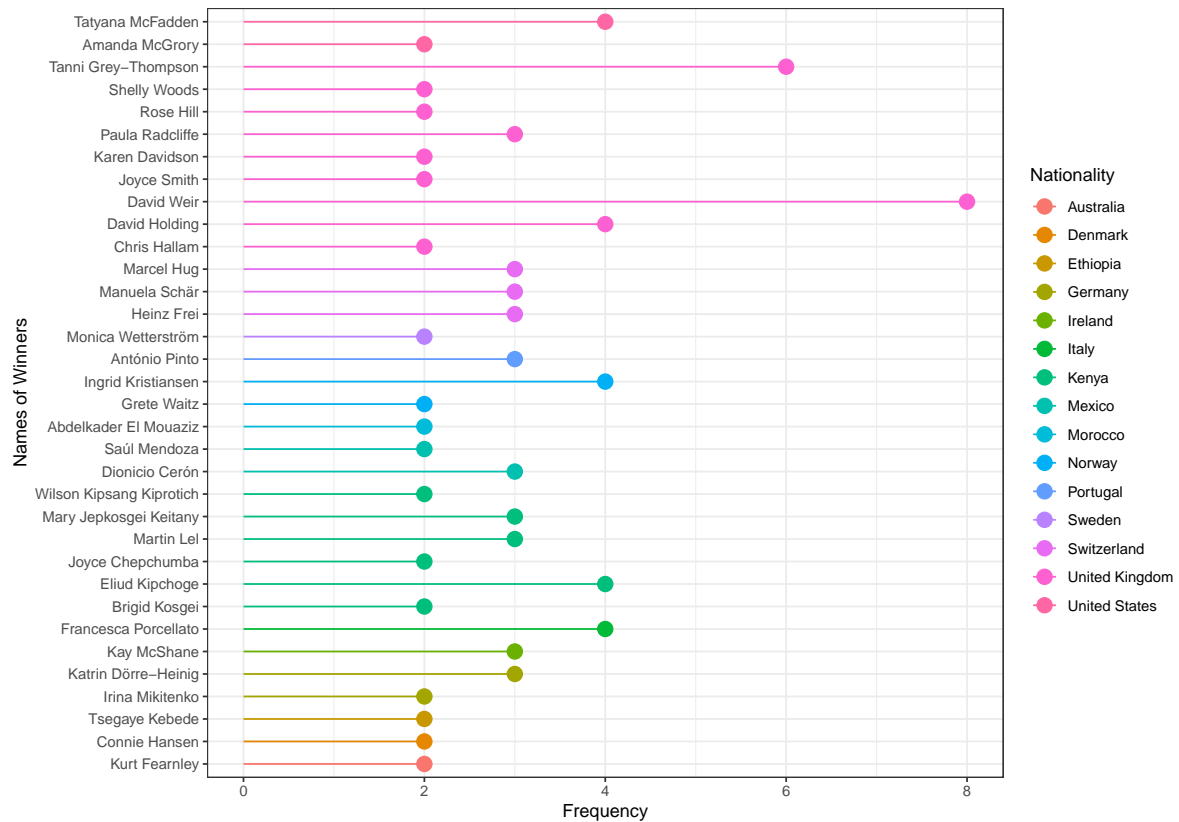
```
winners_count <- data.frame(table(winners$Athlete))
names(winners_count) <- c("Athlete", "Frequency")

winners_nationality <- unique(left_join(winners_count, winners[,c("Athlete", "Nationality")],

# ggplot(winners_nationality[winners_count$Frequency > 1,], aes(x=Athlete, y=Frequency)) +
#   geom_segment( aes(x=Athlete, xend=Athlete, y=0, yend=Frequency), color="skyblue") +
#   geom_point( color="blue", size=4, alpha=0.6) +
#   theme_light() +
#   coord_flip() +
#   theme(
#     panel.grid.major.y = element_blank(),
#     panel.border = element_blank(),
#     axis.ticks.y = element_blank()
#   )

winners_nationality[winners_nationality$Frequency > 1,] %>%
  arrange(Nationality) %>%      # First sort by val. This sort the dataframe but NOT the f
```

```
mutate(name=factor(Athlete, levels=Athlete)) %>% # This trick update the factor levels
ggplot( aes(x=name, y=Frequency, color=Nationality)) +
geom_segment( aes(x=name, xend=name, y=Frequency,yend=0)) +
geom_point( size=4) +
coord_flip() +
theme_bw() +
xlab("Names of Winners")
```



Question: Is united kingdom at advantage of winning the marathon? What does the proportions say?