

## 1. Opis zbioru

Dane dotyczą wynajmu rowerów przez fikcyjną firmę Cyclistic w okresie od kwietnia 2021 do marca 2022.

Zbiory zawierają informacje o numerach identyfikacyjnych przejazdów, typach wypożyczonych rowerów, datach wypożyczenia i oddania rowerów oraz nazwy, numery identyfikacyjne i współrzędne geograficzne stacji początkowych oraz końcowych.

Łącznie w 12 tabelach znajdują się 5723532 wpisy.

Dane pochodzą z obserwacji.

Źródło: <https://www.kaggle.com/datasets/evangower/cyclistic-bike-share> (<https://www.kaggle.com/datasets/evangower/cyclistic-bike-share>)

## 2. Analiza eksploracyjna

Dwoma głównymi typami wypożyczanych rowerów są **rower klasyczny** (56.7 % wszystkich wpisów) i **rower elektryczny** (37.9 %).

Przeprowadziłam analizę średnich czasów wypożyczenia obydwu rodzajów roweru, a następnie przetestowałam, czy jeden z nich ma średnio dłuższy czas.

Najpierw wczytuję dane oraz przygotowuję funkcję zwracającą dla danego miesiąca wektor czasów wypożyczeń.

```
In [1]: 1 april_21 <- read.csv("202104-divvy-tripdata.csv")
2 may_21 <- read.csv("202105-divvy-tripdata.csv")
3 june_21 <- read.csv("202106-divvy-tripdata.csv")
4 july_21 <- read.csv("202107-divvy-tripdata.csv")
5 august_21 <- read.csv("202108-divvy-tripdata.csv")
6 september_21 <- read.csv("202109-divvy-tripdata.csv")
7 october_21 <- read.csv("202110-divvy-tripdata.csv")
8 november_21 <- read.csv("202111-divvy-tripdata.csv")
9 december_21 <- read.csv("202112-divvy-tripdata.csv")
10 january_22 <- read.csv("202201-divvy-tripdata.csv")
11 february_22 <- read.csv("202202-divvy-tripdata.csv")
12 march_22 <- read.csv("202203-divvy-tripdata.csv")
```

```
In [2]: 1 calculate_time <- function(month) {
2   end <- as.POSIXct(month$ended_at, format="%Y-%m-%d %H:%M:%S", tz="UTC")
3   start <- as.POSIXct(month$started_at, format="%Y-%m-%d %H:%M:%S", tz="UTC")
4   time <- as.numeric(difftime(end, start), units="secs")
5 }
```

Następnie tworzę dwie listy zawierające czasy w konkretnych miesiącach dla roweru klasycznego i elektrycznego.

Na ich podstawie liczę średnie z każdego miesiąca. Tworzę też wektory liczby wypożyczeń każdego z dwóch rodzajów roweru.

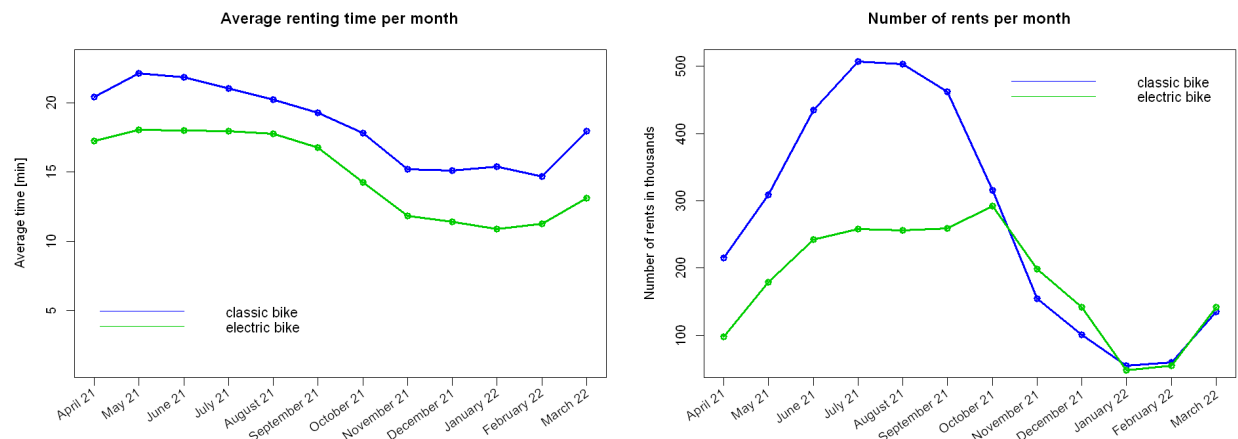
```
In [3]: 1 total_list <- list(april_21, may_21, june_21, july_21, august_21, september_21, october_21, november_21,
2   december_21, january_22, february_22, march_22)
3
4 classic_times <- list()
5 for (i in 1:12) {
6   classic_times[[i]] <- calculate_time(subset(total_list[[i]],
7   rideable_type == "classic_bike"))
8 }
9
10 electric_times <- list()
11 for (i in 1:12) {
12   electric_times[[i]] <- calculate_time(subset(total_list[[i]],
13   rideable_type == "electric_bike"))
14 }
```

```
In [4]: 1 classic_monthly_avg <- vector()
2 for (i in 1:12) {
3   classic_monthly_avg <- c(classic_monthly_avg, mean(unlist(classic_times[i])))
4 }
5
6 electric_monthly_avg <- vector()
7 for (i in 1:12) {
8   electric_monthly_avg <- c(electric_monthly_avg, mean(unlist(electric_times[i])))
9 }
```

```
In [5]: 1 classic_monthly_rents <- vector()
2 for (i in 1:12) {
3   classic_monthly_rents <- c(classic_monthly_rents, length(unlist(classic_times[i])))
4 }
5
6 electric_monthly_rents <- vector()
7 for (i in 1:12) {
8   electric_monthly_rents <- c(electric_monthly_rents, length(unlist(electric_times[i])))
9 }
```

Korzystając z wyliczonych wektorów tworzę wykresy porównujące średnie czasy wypożyczenia oraz liczbę wypożyczeń roweru klasycznego i elektrycznego.

```
In [6]: 1 months = c("April 21", "May 21", "June 21", "July 21", "August 21", "September 21",
2           "October 21", "November 21", "December 21", "January 22", "February 22", "March 22")
3
4 par(mfrow=c(1,2))
5 options(repr.plot.width=16, repr.plot.height=6)
6
7 plot(classic_monthly_avg/60, type = "o", col = "blue", main = "Average renting time per month",
8      xaxt = "n", xlab = "", ylab = "Average time [min]", lwd = 3, ylim = c(1,23))
9 lines(electric_monthly_avg/60, type = "o", col = "green3", lwd = 3)
10 axis(1, 1:12, labels = FALSE)
11 text(x = 1:12,
12      y = par("usr")[3] - 1,
13      labels = months,
14      xpd = NA,
15      srt = 35,
16      adj = 1)
17 legend(0.2,6,legend=c("classic bike","electric bike"),col=c("blue", "green3"),lty=1:1, cex=1.1,
18       bty='n', box.lwd=0)
19
20 plot(classic_monthly_rents/1000, type = "o", col = "blue", main = "Number of rents per month",
21      xaxt = "n", xlab = "", ylab = "Number of rents in thousands", lwd = 3)
22 lines(electric_monthly_rents/1000, type = "o", col = "green3", lwd = 3)
23 axis(1, 1:12, labels = FALSE)
24 text(x = 1:12,
25      y = par("usr")[3] - 22,
26      labels = months,
27      xpd = NA,
28      srt = 35,
29      adj = 1)
30 legend(6.5,500,legend=c("classic bike","electric bike"),col=c("blue", "green3"),lty=1:1,
31       cex=1.1, bty='n', box.lwd=0)
```



Jak widać, rower klasyczny jest średnio wypożyczany na dłużej.  
Z wyjątkiem miesięcy zimowych jest również wybierany znacznie częściej od elektrycznego.

Do sprawdzenia poprawności wniosków z testu statystycznego potrzebować będę średniego czasu przejazdu rowerem klasycznym i elektrycznym.

```
In [7]: 1 classic_renting_times <- unlist(classic_times)
2 electric_renting_times <- unlist(electric_times)
```

Funkcja `sd()` wylicza odchylenie standardowe w próbce, dlatego w przypadku populacji mnożę wynik przez  $\sqrt{\frac{n-1}{n}}$ .

```
In [8]: 1 classic_avg = mean(classic_renting_times)
2 classic_sd <- sd(classic_renting_times)
3          *(sqrt((length(classic_renting_times)-1)/length(classic_renting_times)))
4 cat("Średni czas wypożyczania roweru klasycznego:", classic_avg,
5     "min z odchyleniem standardowym", classic_sd, "min.\n")
6
7 electric_avg = mean(electric_renting_times)
8 electric_sd <- sd(electric_renting_times)
9          *(sqrt((length(electric_renting_times)-1)/length(electric_renting_times)))
10 cat("Średni czas wypożyczania roweru elektrycznego:", electric_avg,
11     "min z odchyleniem standardowym", electric_sd, "min.")
```

Średni czas wypożyczania roweru klasycznego: 1184.347 min z odchyleniem standardowym 3297.516 min.  
 Średni czas wypożyczania roweru elektrycznego: 939.2228 min z odchyleniem standardowym 1202.942 min.

### 3. Test statystyczny

Zamierzam sprawdzić, czy średni czas jazdy jest dłuższy dla roweru klasycznego względem roweru elektrycznego.

W tym celu przeprowadzę test Z dla dwóch populacji niezależnych.

Jako, że posiadam już dane o całej populacji, wylosuję z niej po 1000 wpisów dla każdego rodzaju roweru i na jej podstawie przeprowadzę test. Rozmiar próby jest na tyle duży, że na podstawie Centralnego Twierdzenia Klasycznego można uznać, że liczona średnia pochodzi z rozkładu normalnego.

```
In [9]: 1 n <- 1000
2 classic_sample <- sample(classic_renting_times, n)
3 electric_sample <- sample(electric_renting_times, n)
```

Do przeprowadzenia testu konieczne będzie wyznaczenie średnich i odchyłeń standardowych obydwu próbek.

```
In [10]: 1 classic_sample_avg <- mean(classic_sample)
2 classic_sample_sd <- sd(classic_sample)
3
4 electric_sample_avg <- mean(electric_sample)
5 electric_sample_sd <- sd(electric_sample)
6
7 cat("Średni czas wypożyczenia roweru klasycznego:", classic_sample_avg,
8     "s z odchyleniem standardowym", classic_sample_sd, "s.\n")
9 cat("Średni czas wypożyczenia roweru elektrycznego:", electric_sample_avg,
10     "s z odchyleniem standardowym", electric_sample_sd, "s.")
```

Średni czas wypożyczenia roweru klasycznego: 1273.37 s z odchyleniem standardowym 4123.939 s.  
 Średni czas wypożyczenia roweru elektrycznego: 854.32 s z odchyleniem standardowym 868.4347 s.

Zacznę od sformułowania hipotezy zerowej i alternatywnej.

$H_0 : \mu_c = \mu_e$  - średni czas wypożyczenia roweru klasycznego i elektrycznego jest równy

$H_A : \mu_c > \mu_e$  - średni czas wypożyczenia roweru klasycznego jest większy od elektrycznego

Następnie, wyznaczę zbiór krytyczny dla poziomu istotności  $\alpha = 0.01$

Zbiór krytyczny: (2.326;  $+\infty$ )

```
In [11]: 1 qnorm(0.01)
```

-2.32634787404084

Pozostaje obliczenie wartości statystyki Z, co zrobię korzystając ze wzoru:

$$Z_{\bar{X}_1 - \bar{X}_2} = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

```
In [12]: 1 Z <- (classic_sample_avg - electric_sample_avg)/sqrt((classic_sample_sd^2 + electric_sample_sd^2)/n)
2 cat("Wartość statystyki Z:", round(Z,3))
```

Wartość statystyki Z: 3.144

#### Wnioski

Otrzymana wartość Z należy do zbioru krytycznego, w związku z czym **możemy odrzucić hipotezę zerową i zaakceptować hipotezę alternatywną**.

Jako, że wartości średnich czasów dla obydwu typów roweru są nam znane, możemy zauważyć, że przeprowadzony test doprowadził do poprawnych wniosków.