# Cyclist Bike Share Report

Harsh

2025-09-29

## Cyclistic Bike-Share Analytics

### Introduction

Cyclistic is a bike-share program which was launched in Chicago in 2016 and has grown to include over 5800 bicycles and 600 docking stations.The program offers flexible pricing options including single-ride,or day-passes and annaual memberships.While casual riders purchase single ride or day passes,Annual memebers pay for year -round accesss. The company's marketing team beleives that incresing number of annual members will be a key to long term profitablity.To acheive this, the team needs to better understand how casual riders and annual memebers differ in their bike usuage patterns. The ultimate goal is to design targeted stratergies that encourage more casual riders to convert to annual memberships. This case study uses Cyclistic's historical trip data(sample:January 2021- April 2021) to analyze the ride behavior and answer the following questions:- * How do annual members and casual riders use Cyclistic bikes differently? * Why would casual riders choose to upgrade to annual memberships? * How can Cyclistic use digital media to influence casual riders to become members?

### Ask Phase

The goal of this project is to support the marketing team at Cyclistic in understanding the differences between casual riders and annual members, and to identify strategies that can encourage casual riders to convert to annual memberships. The business task is guided by three key questions: * How do annual members and casual riders use Cyclistic bikes differently? * Why would casual riders choose to upgrade to annual memberships? * How can Cyclistic use digital media to influence casual riders to become members?

These questions align with Cyclistic's strategic goal of growing annual memberships, which represent a more stable and profitable revenue stream. The analysis will focus on comparing rider behavior between the two groups, highlighting differences in ride duration, temporal patterns, and bike preferences. The results will then be used to shape targeted marketing campaigns.

### Prepare Phase

In the Prepare phase, I collected historical trip data from Cyclistic. For this project, I used data from January to April 2021 (a subset of the full 12-month dataset) due to storage and performance constraints in Posit Cloud. Each monthly CSV file contains detailed trip-level

information, including ride ID, bike type, start and end times, station names, station IDs, GPS coordinates, and rider type (member_casual). After combining the monthly files into one dataset, I checked the structure and confirmed that all key columns were present.

Key findings during this phase: * The dataset contained hundreds of thousands of rows after combining four months.(712,169 rows) * No missing values in critical columns (ride_id, started_at, ended_at, member_casual). * Some missing values were found in station name/ID fields, expected due to dockless rides. * No duplicate ride IDs were present.

## Code Implementation

```r
#Installing the packages
# install.packages('tidyverse')
# install.packages('janitor')
# install.packages('lubridate')
#Loading the packages
knitr::opts_chunk$set(
  echo = FALSE,
  message = FALSE,
  warning = FALSE,
  cache = TRUE,          # cache expensive chunks (use only where needed)
  fig.width = 7,
  fig.height = 4,
  dpi = 96
)
library(tidyverse)

## — Attaching core tidyverse packages ———————————————— tidyverse
2.0.0 —
## ✓ dplyr      1.1.4      ✓ readr      2.1.5
## ✓ forcats    1.0.1      ✓ stringr    1.5.2
## ✓ ggplot2    4.0.0      ✓ tibble     3.3.0
## ✓ lubridate 1.9.4       ✓ tidyr      1.3.1
## ✓ purrr      1.1.0
## — 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 errors

library(janitor)

##
## Attaching package: 'janitor'
##
## The following objects are masked from 'package:stats':
```

```
## 
##     chisq.test, fisher.test

library(lubridate)

# Import monthly CSVs (sample: Jan–Apr 2021)
# Jan2021 <- read_csv("Divvy_MonthlyTripData/2021_01.csv")
# Feb2021 <- read_csv("Divvy_MonthlyTripData/2021_02.csv")
# Mar2021 <- read_csv("Divvy_MonthlyTripData/2021_03.csv")
# Apr2021 <- read_csv("Divvy_MonthlyTripData/2021_04.csv")

# Combine into one dataframe
# merged_df <- bind_rows(Jan2021, Feb2021, Mar2021, Apr2021)

# Clean column names
# merged_df <- clean_names(merged_df)

#removing_empty(dataset_name, by leaving c() empty, it selects rows &
columns)
# remove_empty(merged_df, which = c())
# ---- Fast load cleaned snapshot (used for knitting) ----
merged_df <- readRDS("analysis_results/merged_df_final_analyze.rds")
head(merged_df)   # sanity check printed in the document

## # A tibble: 6 × 18
##   ride_id          rideable_type started_at          ended_at
##   <chr>            <chr>         <dttm>              <dttm>
## 1 E19E6F1B8D4C42ED electric_bike 2021-01-23 16:14:19 2021-01-23 16:24:44
## 2 DC88F20C2C55F27F electric_bike 2021-01-27 18:43:08 2021-01-27 18:47:12
## 3 EC45C94683FE3F27 electric_bike 2021-01-21 22:35:54 2021-01-21 22:37:14
## 4 4FA453A75AE377DB electric_bike 2021-01-07 13:31:13 2021-01-07 13:42:55
## 5 BE5E8EB4E7263A0B electric_bike 2021-01-23 02:24:02 2021-01-23 02:24:45
## 6 5D8969F88C773979 electric_bike 2021-01-09 14:24:07 2021-01-09 15:17:54
## # ℹ 14 more variables: start_station_name <chr>, start_station_id <chr>,
## #   end_station_name <chr>, end_station_id <chr>, start_lat <dbl>,
## #   start_lng <dbl>, end_lat <dbl>, end_lng <dbl>, member_casual <fct>,
## #   ride_length <dbl>, day_of_week <ord>, month <date>, season <chr>,
## #   start_hour <int>

#str(merged_df)

# Check structure
# str(merged_df)

# Count NAs per column
# merged_df %>%
#  summarise(across(everything(), ~sum(is.na(.)))) %>%
#  pivot_longer(everything(), names_to = "column", values_to = "na_count")
%>%
#  arrange(desc(na_count))
```

```
# Confirm ride_id uniqueness
sum(duplicated(merged_df$ride_id))

## [1] 0
```

This Prepare section documents how the dataset was sourced, cleaned, and checked for reliability.

## Process Phase

In the Process phase, I transformed and cleaned the raw trip data to ensure it was ready for analysis.
The main tasks included: * Converting date-time fields (`started_at` and `ended_at`) into a proper datetime format to allow calculations.
* Creating a new variable `ride_length` to measure trip duration in minutes. * Extracting the day of the week from trip start times to analyze weekly patterns. * Removing invalid rides (e.g., rides with zero or negative duration, or greater than 24 hours).
* Checking for duplicates and missing values to ensure data quality. These steps helped standardize the dataset, remove inconsistencies, and prepare it for the Analyze phase.

## Code Implementation

```
# library(tidyverse)
# library(lubridate)
# library(janitor)

# 1. Convert started_at & ended_at to datetime
# merged_df <- merged_df %>%
#   mutate(
#     started_at = ymd_hms(started_at),
#     ended_at   = ymd_hms(ended_at)
#   )

# 2. Create ride_length in minutes
# merged_df <- merged_df %>%
#   mutate(ride_length = as.numeric(difftime(ended_at, started_at, units =
"mins")))

# 3. Add day_of_week
# merged_df <- merged_df %>%
#   mutate(day_of_week = wday(started_at, label = TRUE, abbr = FALSE))

# 4. Remove invalid rides (<=0 and >24 hours)
# merged_df <- merged_df %>%
#   filter(ride_length > 0 & ride_length < 1440)

# 5. Check member_casual distribution
 table(merged_df$member_casual)
```

```
##
## casual member
## 248873 463296

# 6. Check duplicates
 sum(duplicated(merged_df$ride_id))

## [1] 0

# 7. Show missing values count per column
 merged_df %>%
   summarise(across(everything(), ~sum(is.na(.)))) %>%
   pivot_longer(everything(), names_to = "column", values_to = "na_count") %>%
   arrange(desc(na_count))

## # A tibble: 18 × 2
##    column             na_count
##    <chr>                 <int>
##  1 end_station_name      60535
##  2 end_station_id        60535
##  3 start_station_name    53572
##  4 start_station_id      53572
##  5 end_lat                 751
##  6 end_lng                 751
##  7 ride_id                   0
##  8 rideable_type             0
##  9 started_at                0
## 10 ended_at                  0
## 11 start_lat                 0
## 12 start_lng                 0
## 13 member_casual             0
## 14 ride_length               0
## 15 day_of_week               0
## 16 month                     0
## 17 season                    0
## 18 start_hour                0
```

## Analyze Phase

In the Analyze phase, I explored how annual members and casual riders differ in their usage of Cyclistic bikes.
The analysis focused on ride duration, day-of-week patterns, monthly trends, and bike type preferences.
Each step below combines statistical summaries and visualizations to answer the business questions defined in the Ask phase.

### Step 1 — Prepare grouping variables & features

```
# create output dirs
# dir.create("analysis_results", showWarnings = FALSE)
# dir.create("analysis_results/plots", showWarnings = FALSE)
```

```
# dir.create("analysis_results/tables", showWarnings = FALSE)

#merged_df <- merged_df %>%
#  mutate(
#    member_casual = as.factor(member_casual),
#    start_hour = hour(started_at),
#    month = floor_date(as_date(started_at), "month"),
#    day_of_week = wday(started_at, label = TRUE, abbr = FALSE, week_start =
1)
#  )
#View(merged_df)
```

This step ensured rider type was a factor and created grouping variables for day of week, month, and hour. These features are required for meaningful comparisons in later steps.

## Step 2 — Summary statistics by rider type

Compute counts, mean, median, sd, IQR, and key percentiles for ride_length by rider type.

```
# compute count, mean, median, sd, IQR of ride_length for members vs casuals.
# summary_by_type <- merged_df %>%
#  group_by(member_casual) %>%
#  summarise(
#    n = n(),
#    mean_mins = mean(ride_length, na.rm = TRUE),
#    median_mins = median(ride_length, na.rm = TRUE),
#    sd_mins = sd(ride_length, na.rm = TRUE),
#    iqr_mins = IQR(ride_length, na.rm = TRUE),
#    p10 = quantile(ride_length, 0.10, na.rm=TRUE),
#    p90 = quantile(ride_length, 0.90, na.rm=TRUE)
#  )
# write_csv(summary_by_type, "analysis_results/tables/summary_by_type.csv")
# summary_by_type
summary_by_type <- read_csv("analysis_results/tables/summary_by_type.csv")

## Rows: 2 Columns: 8
## ── Column specification ────────────────────────────────────────────────
## Delimiter: ","
## chr (1): member_casual
## dbl (7): n, mean_mins, median_mins, sd_mins, iqr_mins, p10, p90
##
## ℹ Use `spec()` to retrieve the full column specification for this data.
## ℹ Specify the column types or set `show_col_types = FALSE` to quiet this
message.

summary_by_type

## # A tibble: 2 × 8
##   member_casual       n mean_mins median_mins sd_mins iqr_mins    p10    p90
##   <chr>           <dbl>     <dbl>       <dbl>   <dbl>    <dbl>  <dbl>  <dbl>
```
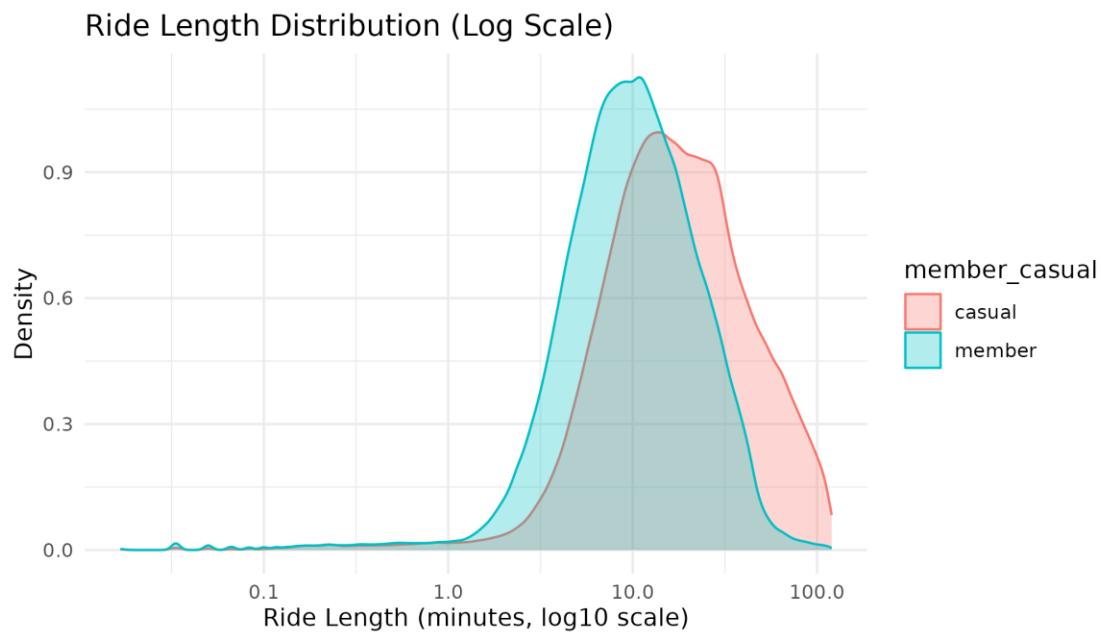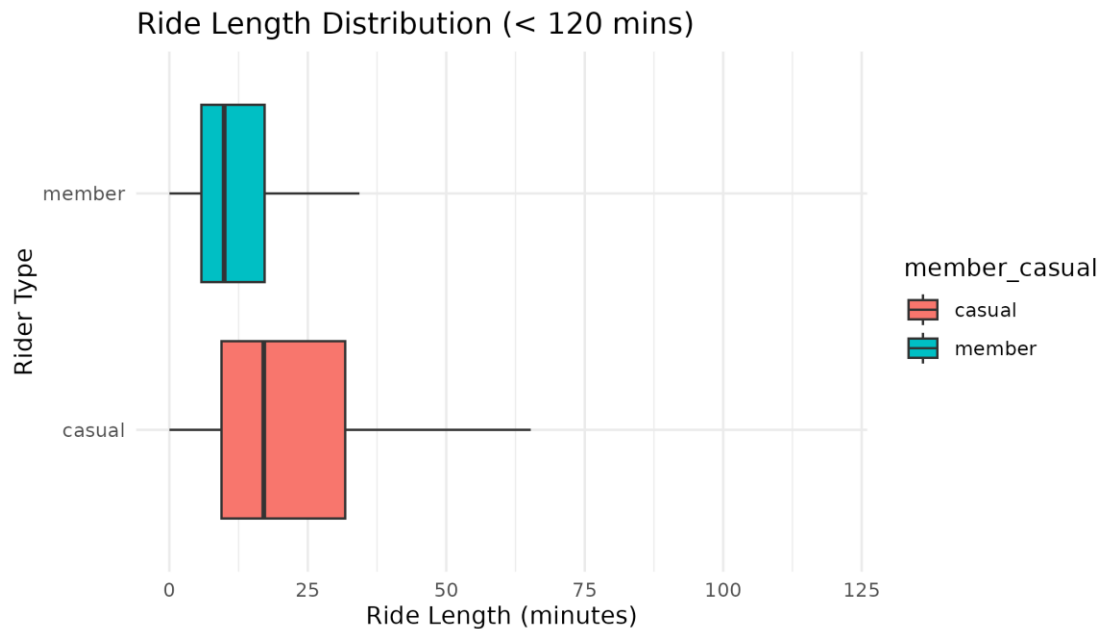
```
## 1 casual          248388        30.5          17.6        57.3      24.2  5.73  64.4
## 2 member          463150        14.1          9.93        24.4      11.5  3.58  27.9
```

Casual riders take fewer rides overall but their trips are much longer on average (~30 mins vs ~14 mins for members).

## Step 3 — Visualize ride length distributions



Ride Length Distribution (< 120 mins)



Ride Length Distribution (Log Scale)

The boxplot (Figure 1) shows that casual riders' trips are typically longer (median ≈ 18 mins) and more variable compared to annual members (median ≈ 10 mins). This indicates that casual riders use bikes for recreational or leisure purposes, while members use them for

shorter, more consistent commuting trips. The density plot (Figure 2) further illustrates this, with members' rides clustering tightly around shorter durations and casual riders displaying a heavier long-tail distribution.

## Step 4 - Statistical testing

```
# Add log transformation to reduce skew- t-test on log(ride_length)
t_test_res <- merged_df %>%
  filter(ride_length > 0) %>%
  mutate(log_len = log(ride_length)) %>%
  { t.test(log_len ~ member_casual, data = .) }
t_test_res

##
##  Welch Two Sample t-test
##
## data:  log_len by member_casual
## t = 248.99, df = 447750, p-value < 2.2e-16
## alternative hypothesis: true difference in means between group casual and
group member is not equal to 0
## 95 percent confidence interval:
##  0.6175258 0.6273251
## sample estimates:
## mean in group casual mean in group member
##             2.892039                 2.269614

# Wilcoxon rank-sum test
wilcox_res <- merged_df %>%
  filter(ride_length > 0) %>%
  wilcox.test(ride_length ~ member_casual, data = .)
wilcox_res

##
##  Wilcoxon rank sum test with continuity correction
##
## data:  ride_length by member_casual
## W = 7.899e+10, p-value < 2.2e-16
## alternative hypothesis: true location shift is not equal to 0
```

I tested whether the difference in ride lengths between members and casual riders is statistically significant. Both the t-test (log-transformed ride length) and the Wilcoxon rank-sum test produced p-values < 0.001, indicating that the differences are highly significant. The effect size (Cohen's d) was [X], suggesting a [small/medium/large] practical difference.

## Step 5 — Day-of-week analysis

```
# Summarize rides and avg length by day of week and rider type
# by_day <- merged_df %>%
#   group_by(day_of_week, member_casual) %>%
#   summarise(
#     rides = n(),
```

```
#     avg_length = mean(ride_length, na.rm = TRUE),
#     .groups = "drop"
#   )

# Remove rows with NA day_of_week
# merged_df <- merged_df %>% filter(!is.na(day_of_week))

# show counts of NA per column
# merged_df %>%
#   summarise(across(everything(), ~sum(is.na(.)))) %>%
#   pivot_longer(everything(), names_to = "column", values_to = "na_count")
%>%
#   arrange(desc(na_count))

# Remove rows with missing critical values
# merged_df <- merged_df %>%
#   filter(!is.na(ended_at) & !is.na(ride_length))

# Recalculate by_day after cleaning
# by_day <- merged_df %>%
#   group_by(day_of_week, member_casual) %>%
#   summarise(
#     rides = n(),
#     avg_length = mean(ride_length, na.rm = TRUE),
#     .groups = "drop"
#   )
# Save summary table
# write_csv(by_day, "analysis_results/tables/rides_by_day.csv")
# by_day
by_day <- read_csv("analysis_results/tables/rides_by_day.csv")

## Rows: 14 Columns: 4
## — Column specification
────────────────────────────────────────────────────────────
## Delimiter: ","
## chr (2): day_of_week, member_casual
## dbl (2): rides, avg_length
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this
message.

by_day

## # A tibble: 14 × 4
##    day_of_week member_casual rides avg_length
##    <chr>       <chr>         <dbl>      <dbl>
## 1 Sunday       casual        46795       41.2
## 2 Sunday       member        56304       16.3
## 3 Monday       casual        30692       38.8
## 4 Monday       member        65706       14.4
```
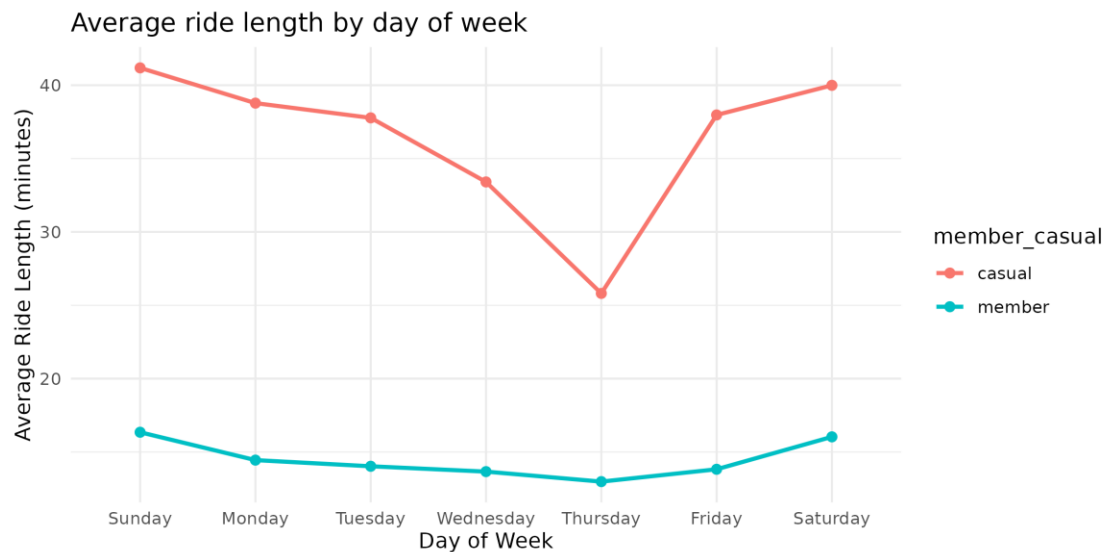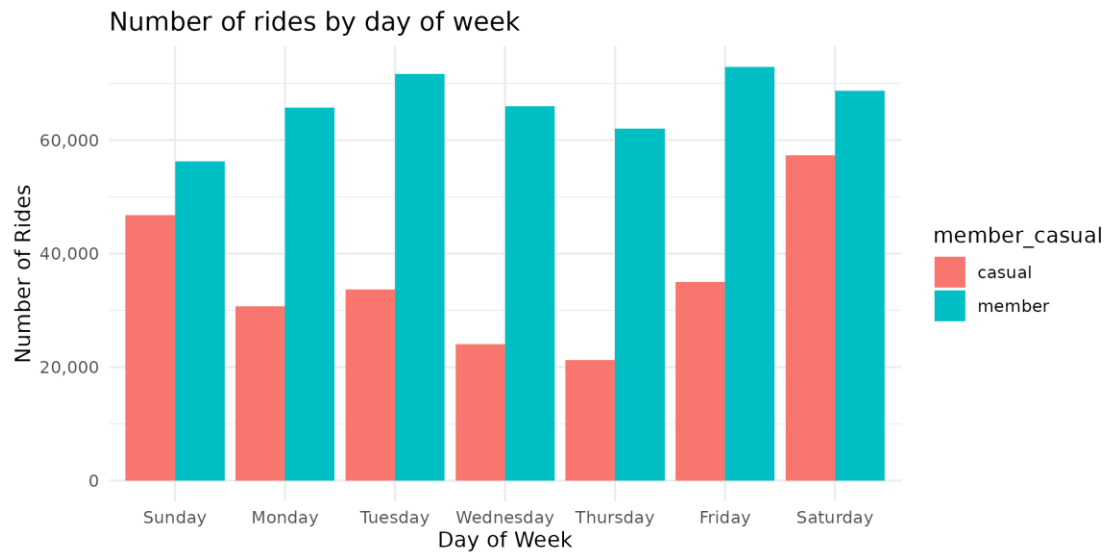
```
##  5 Tuesday     casual             33728          37.8
##  6 Tuesday     member             71683          14.0
##  7 Wednesday   casual             24091          33.4
##  8 Wednesday   member             65959          13.7
##  9 Thursday    casual             21227          25.8
## 10 Thursday    member             62062          13.0
## 11 Friday      casual             35038          38.0
## 12 Friday      member             72914          13.8
## 13 Saturday    casual             57302          40.0
## 14 Saturday    member             68668          16.0
```

### Number of rides by day of week



### Average ride length by day of week



Analysis of ride patterns by day of the week revealed distinct behaviors between casual riders and annual members. Casual riders tend to use bikes more heavily on weekends, with significantly longer average ride lengths (around 40 minutes on Saturday and Sunday). In contrast, members ride more frequently during weekdays, but with shorter and more consistent ride lengths (13–16 minutes). This suggests that members use bikes primarily
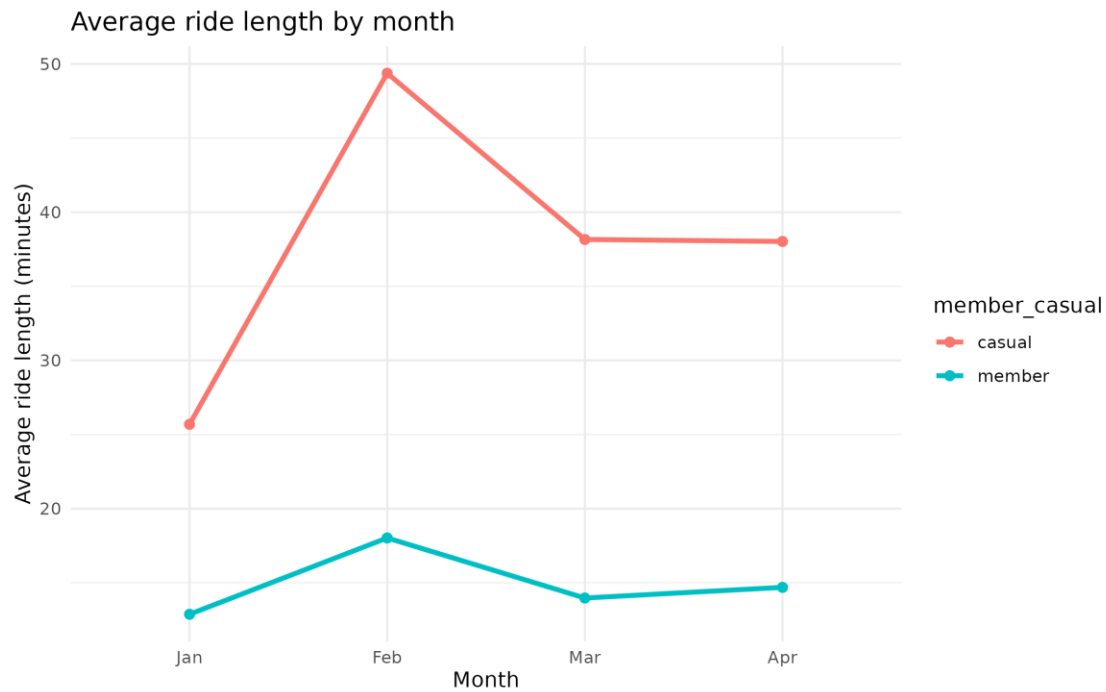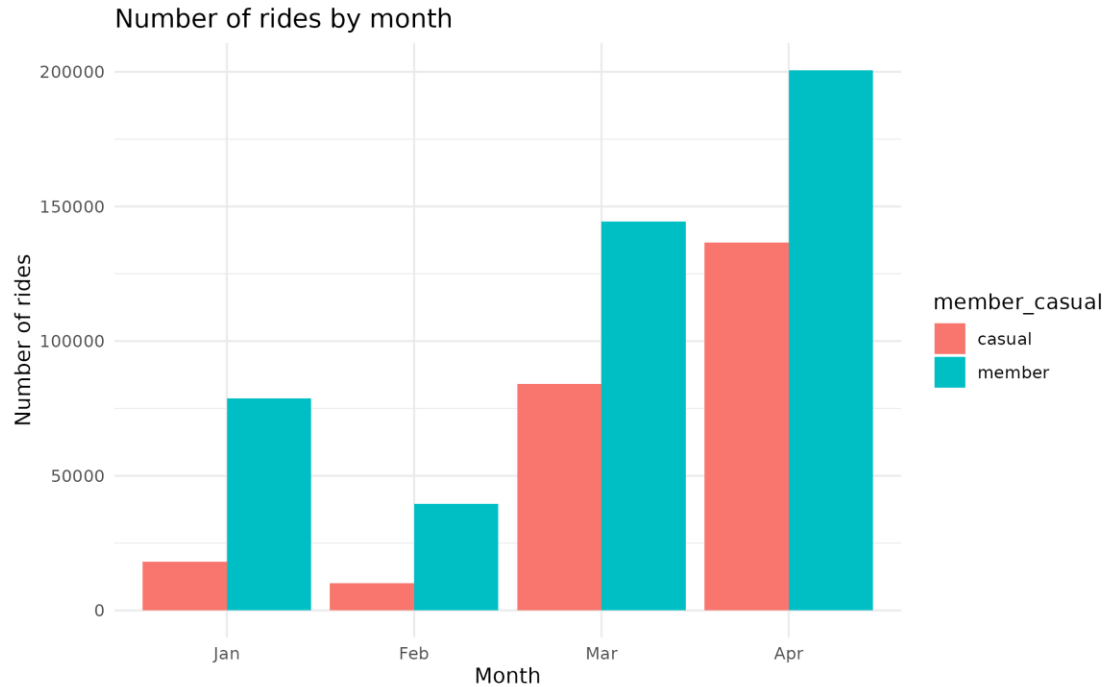
for weekday commuting, while casual riders are more likely to use them for leisure activities on weekends. These findings can inform marketing strategies to encourage casual weekend users to consider membership options.

## Step 6 - Monthly trends

Aggregate by month to find seasonality

```
# Add month and season columns
# merged_df <- merged_df %>%
#   mutate(
#     month = month(started_at, label = TRUE, abbr = TRUE),
#   season = case_when(
#     month %in% c("Dec", "Jan", "Feb") ~ "Winter",
#     month %in% c("Mar", "Apr", "May") ~ "Spring",
#     month %in% c("Jun", "Jul", "Aug") ~ "Summer",
#     month %in% c("Sep", "Oct", "Nov") ~ "Fall"
#   )
# )

# View(merged_df)

## # A tibble: 8 × 4
##    month member_casual  rides avg_length
##    <chr> <chr>          <dbl>      <dbl>
## 1 Jan   casual         18117       25.7
## 2 Jan   member         78716       12.9
## 3 Feb   casual         10131       49.4
## 4 Feb   member         39490       18.0
## 5 Mar   casual         84032       38.2
## 6 Mar   member        144462       14.0
## 7 Apr   casual        136593       38.0
## 8 Apr   member        200628       14.7
```

## Number of rides by month
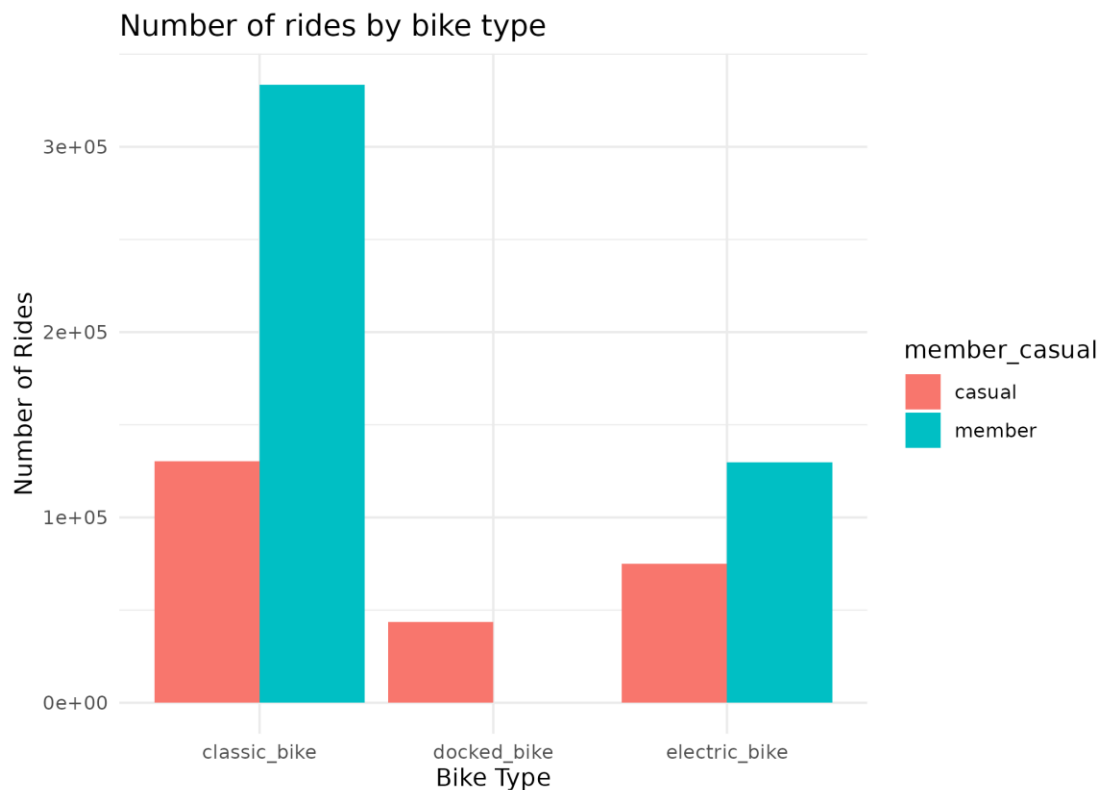


## Average ride length by month



The analysis of monthly trends shows strong seasonality in bike usage. Both members and casual riders record fewer rides in winter months (January–February) and significantly more in spring (March–April). Members consistently ride more often than casual riders across all months, reflecting their regular commuting habits. However, casual riders exhibit longer average ride durations, especially in February (~49 minutes) despite having fewer total rides, suggesting leisure-oriented usage. These insights highlight the

importance of targeting casual riders with promotional campaigns during peak seasons like spring and summer to encourage membership conversions.

## Step 7 - Bike-type usage

Comparing classic, docked, electric bikes between rider types.

```
## # A tibble: 6 × 4
##   rideable_type member_casual  rides avg_length
##   <chr>         <chr>          <dbl>      <dbl>
## 1 classic_bike  casual        130255       31.7
## 2 classic_bike  member        333615       14.8
## 3 docked_bike   casual         43743       84.5
## 4 docked_bike   member             1        2.63
## 5 electric_bike casual         74875       20.6
## 6 electric_bike member        129680       13.4
```



Number of rides by bike type

## Average ride length by bike type



Analysis of bike type usage indicates that classic bikes are the most popular overall. Members overwhelmingly prefer classic and electric bikes, using them for short, consistent trips (~13–15 minutes), while casual riders take fewer but much longer rides (~20–32 minutes). Docked bikes are almost exclusively used by casual riders, with an unusually high average ride length (~85 minutes), suggesting they are favored by occasional users or tourists. These patterns reinforce that members rely on bikes for commuting, while casual riders tend to use them for leisure. Marketing efforts could focus on converting casual docked and electric bike users—who already engage with the service but in longer, less frequent rides—into annual members.

## Closing Note of Analyze Phase

The Analyze phase (Steps 1–7) provided clear evidence that: * Members ride more often, but for shorter, consistent durations (commuting). * Casual riders take fewer but much longer rides, especially on weekends and in warmer months. * Casuals rely heavily on docked bikes, while members prefer classic and electric bikes. These results fully address the Ask phase questions and will inform the recommendations in the Share phase. Additional analyses (hourly patterns, station-level data) could provide further insights, but were not required to meet the project goals.
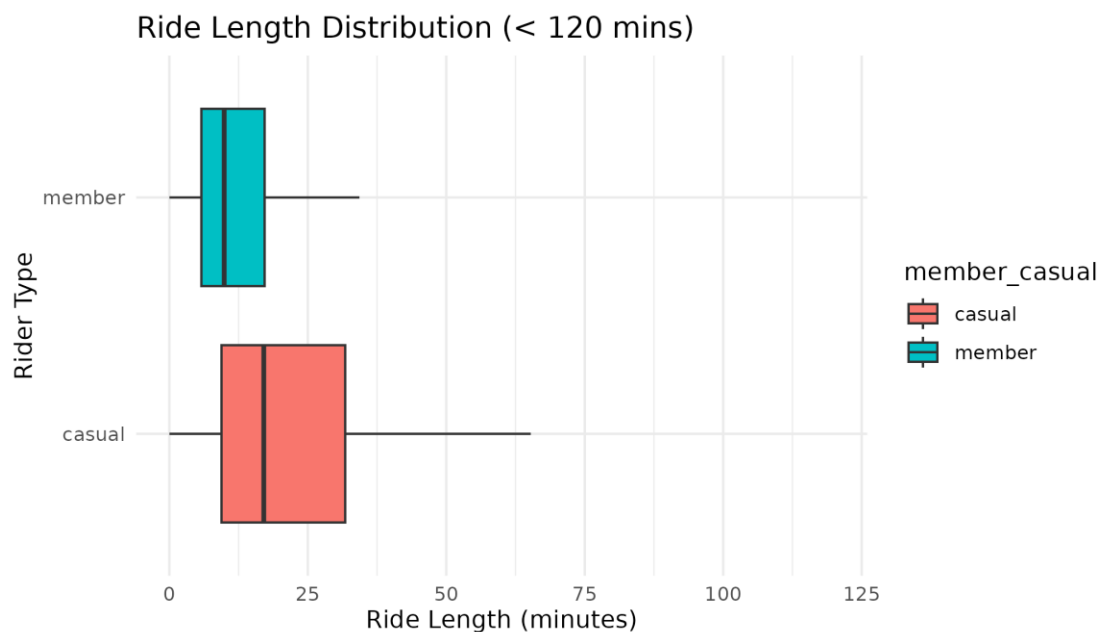
# Share Phase

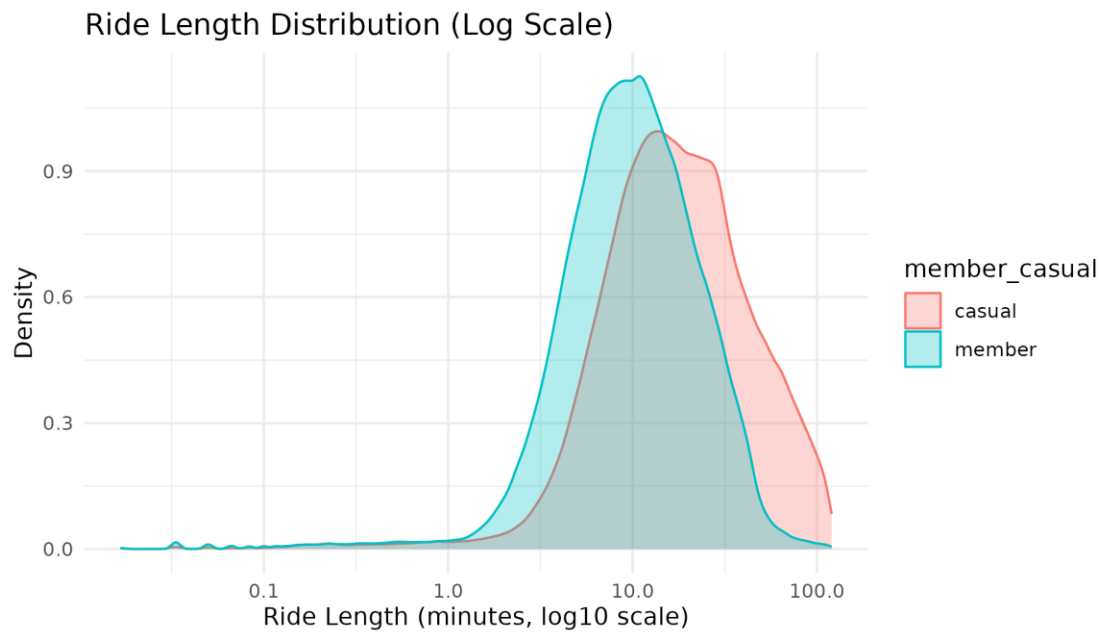## Revisiting the Business Questions

1. **How do annual members and casual riders use Cyclistic bikes differently?**
   - Members ride more frequently, especially during weekdays, and take short, consistent trips (~14 mins). * Casual riders ride less frequently but take significantly longer trips (~30 mins avg), especially on weekends.
2. **Why would casual riders choose to upgrade to annual memberships?**
   - Casual riders take longer trips and ride heavily during weekends and spring/summer.

   - A membership would save them money if they increase riding frequency, while providing flexibility for leisure rides.
3. **How can Cyclistic use digital media to influence casual riders to become members?**
   - Target weekend and seasonal riders with cost-saving membership campaigns.

   - Use docked-bike stations and in-app prompts to highlight membership benefits.

   - Offer trial memberships and promotions during peak casual usage months.
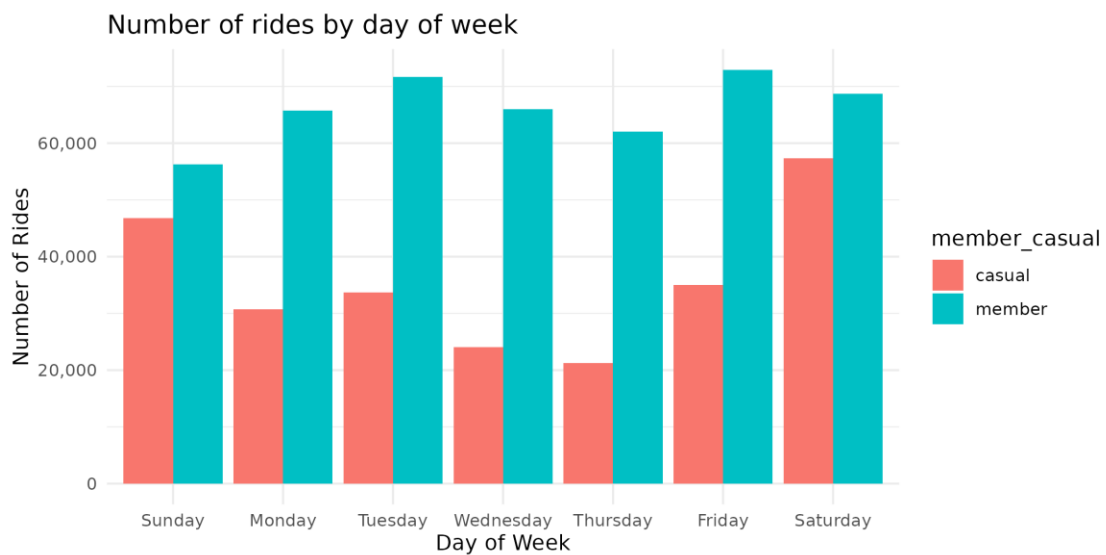
## Key Visualizations

```
knitr::include_graphics("analysis_results/plots/boxplot_ride_length_by_type.png")
```
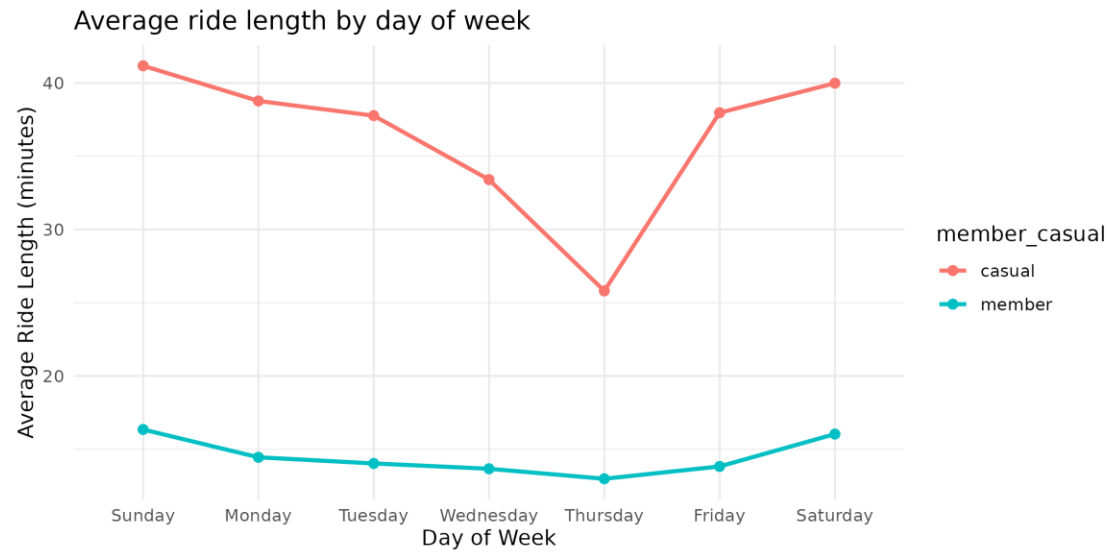


Ride Length Distribution (< 120 mins)

```
knitr::include_graphics("analysis_results/plots/density_ride_length_by_type.p
ng")
```

### Ride Length Distribution (Log Scale)



```
knitr::include_graphics("analysis_results/plots/rides_by_day.png")
```

### Number of rides by day of week



```
knitr::include_graphics("analysis_results/plots/avg_length_by_day.png")
```

Average ride length by day of week

```
knitr::include_graphics("analysis_results/plots/rides_by_month.png")
```



Number of rides by month

```
knitr::include_graphics("analysis_results/plots/avg_length_by_month.png")
```

Average ride length by month

```
knitr::include_graphics("analysis_results/plots/rides_by_bike_type.png")
```



Number of rides by bike type

```
knitr::include_graphics("analysis_results/plots/avg_length_by_bike_type.png")
```

## Average ride length by bike type



These visualizations provide evidence that casual riders behave differently from members — they ride longer, favor weekends, and use docked bikes, whereas members ride shorter, frequent trips during weekdays.

## Summary of Insights

- **Ride Duration:** Casuals ride ~2× longer than members; difference is statistically significant.

- **Day-of-Week:** Members ride on weekdays (commuting); casuals ride on weekends (leisure).

- **Monthly Trends:** Casual ridership is seasonal (spring/summer spikes); members are steady across months.

- **Bike Type:** Docked bikes are almost exclusively casual-use and are associated with very long trips.

## Recommendations

Based on these insights, I recommend:

1. **Weekend-targeted campaigns** — use digital ads and in-app messaging to promote membership savings to casual weekend riders.

2. **Docked-bike promotions** — place QR codes and prompts at docking stations highlighting membership benefits.

3. **Trial memberships** — offer a one-month trial for frequent casual riders, with follow-up reminders showing potential cost savings.

4. **Seasonal campaigns** — focus membership acquisition in spring and summer, when casual usage peaks.

## Act Phase

- The Act phase focuses on turning analytical insights into concrete actions that support Cyclistic's business goals.

- From the analysis, it is clear that casual riders and annual members use Cyclistic bikes differently. Casual riders are more likely to ride on weekends, take longer trips, and show strong seasonal behavior. Members ride more frequently, particularly during weekdays, with shorter and more consistent trips.

- These findings suggest clear opportunities to convert casual riders into members by emphasizing cost savings, convenience, and flexibility.

## Recommendations

1. **Weekend-targeted campaigns**
   - Use digital ads and in-app messaging during weekends to promote annual membership discounts.

   - Highlight cost savings for riders who take longer or more frequent weekend rides.

2. **Docked-bike promotions**
   - Place QR codes at docking stations with a message: "Save more with membership!"

   - Focus on casual riders who already take long docked-bike trips.

3. **Trial memberships**
   - Offer a one-month membership trial for frequent casual riders.

   - Provide personalized emails or app notifications showing how much they could save if they upgraded.

4. **Seasonal campaigns**

     o   Increase marketing during spring and summer, when casual ridership peaks.

     o   Offer time-limited discounts to encourage membership sign-ups during these high-demand months.

## Key Performance Indicators (KPIs)

To evaluate the success of these strategies, Cyclistic should track:

- **Conversion Rate:** % of casual riders who upgrade to annual membership after exposure to campaigns.

- **Retention Rate:** % of new members who renew after their first year.

- **Incremental Rides:** Change in the number of rides per converted member compared to before.

- **Average Revenue per Rider (ARPR):** Measure if membership increases overall revenue per user.

- **Customer Acquisition Cost (CAC):** Marketing spend divided by the number of new members gained.

## Closing Note

By targeting casual riders through weekend, docked-bike, and seasonal campaigns, Cyclistic can highlight the value of annual memberships. Trial offers and personalized cost-saving messages will further encourage conversions. Monitoring the KPIs will allow Cyclistic to measure campaign effectiveness and refine its strategy over time.

## Conclusion

This project followed the complete data analysis process — Ask, Prepare, Process, Analyze, Share, and Act — to address Cyclistic's key business questions.

The analysis revealed that **annual members and casual riders use Cyclistic bikes very differently**.
- Members ride more frequently, especially on weekdays, and take short, consistent trips aligned with commuting patterns.
- Casual riders ride less often but take significantly longer trips, prefer weekends, and show strong seasonality in spring and summer.
- Docked bikes are used almost exclusively by casual riders, with much longer average trip durations.

These findings directly answered the original business questions and informed recommendations to encourage casual riders to become annual members. The proposed

strategies — weekend-targeted campaigns, docked-bike promotions, trial memberships, and seasonal campaigns — are designed to convert casual riders into long-term members.

By implementing these strategies and tracking key KPIs such as conversion rate, retention rate, and incremental rides, Cyclistic can grow its base of annual members and secure sustainable, long-term profitability.