

Cyclistic Bike-Share Analysis Notebook

Shrabana Nandi

November 13, 2025

Contents

| | |
|--|-----------|
| BUSINESS TASK | 2 |
| DATA SOURCE | 2 |
| DATA UPLOAD & CLEANING | 2 |
| Importing Datasets | 3 |
| Cleaning & Standardizing Columns | 4 |
| Combining Datasets | 5 |
| Cleaning Nulls, Timestamps, and Duration | 6 |
| Removing Test Rides & Standardizing Labels | 6 |
| Ride Length Summary | 6 |
| DATA ANALYSIS & INSIGHTS | 6 |
| Feature Engineering | 7 |
| Ride Duration Patterns | 7 |
| Round Trips vs One-Way | 9 |
| Temporal Insights | 9 |
| Day + Hour Trends | 9 |
| Seasonal Patterns | 10 |
| Weekly Trends | 11 |
| Spatial Insights | 12 |
| Top 10 Start Stations | 12 |
| STRATEGIC TAKEAWAYS | 13 |
| Recommendations | 14 |
| 1 <input type="checkbox"/> Seasonal Membership Campaign (Late Spring – Early Summer) | 14 |
| 2 <input type="checkbox"/> Public Transport Advertising for Casual Riders | 14 |
| 3 <input type="checkbox"/> Add Extra Ride-Time Benefits to Annual Membership | 14 |
| CONCLUSION | 15 |

BUSINESS TASK

Goal: Convert casual riders into annual members

Core Questions: 1. How do annual members and casual riders use Cyclistic bikes differently?

2. Why might casual riders buy annual memberships?

3. How can Cyclistic use digital media to influence conversion?

□ *The plan:* Explore patterns in usage (time, distance, frequency, and stations) and identify behavioral cues that could drive marketing actions.

DATA SOURCE

Cyclistic bike share trip data from Q1 2019 and Q1 2020.

DATA UPLOAD & CLEANING

First, load the essentials. These will handle wrangling, cleaning, and plotting.

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 4.5.2
```

```
## Warning: package 'readr' was built under R version 4.5.2
```

```
## Warning: package 'lubridate' was built under R version 4.5.2
```

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

```
## v dplyr      1.1.4      v readr      2.1.5
```

```
## v forcats   1.0.1      v stringr   1.5.2
```

```
## v ggplot2    4.0.0      v tibble    3.3.0
```

```
## v lubridate  1.9.4      v tidyr     1.3.1
```

```
## v 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(lubridate)
```

```
library(janitor)
```

```
## Warning: package 'janitor' was built under R version 4.5.2
```

```
##
```

```
## Attaching package: 'janitor'
```

```
##
## The following objects are masked from 'package:stats':
##
##   chisq.test, fisher.test

library(skimr)

## Warning: package 'skimr' was built under R version 4.5.2

library(readr)
library(dplyr)
```

□ *Thoughts:* No need to go overboard — **tidyverse** covers most cleaning and analysis. **janitor** will clean column names, and **skimr** gives a quick glance at structure.

Importing Datasets

```
divvy_2019 <- read_csv("D:/COURSERA/GOOGLE ANALYTICS/COURSE 8/Divvy_Trips_2019_Q1 - Divvy_Trips_2019_Q1.csv")

## Rows: 365069 Columns: 12
## -- Column specification -----
## Delimiter: ","
## chr (6): start_time, end_time, from_station_name, to_station_name, usertype,...
## dbl (5): trip_id, bikeid, from_station_id, to_station_id, birthyear
## num (1): tripduration
##
## 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.

divvy_2020 <- read_csv("D:/COURSERA/GOOGLE ANALYTICS/COURSE 8/Divvy_Trips_2020_Q1 - Divvy_Trips_2020_Q1.csv")

## Rows: 426887 Columns: 13
## -- Column specification -----
## Delimiter: ","
## chr (7): ride_id, rideable_type, started_at, ended_at, start_station_name, e...
## dbl (6): start_station_id, end_station_id, start_lat, start_lng, end_lat, en...
##
## 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.

head(divvy_2019)

## # A tibble: 6 x 12
##   trip_id start_time      end_time bikeid tripduration from_station_id
##   <dbl> <chr>          <chr>    <dbl>      <dbl>          <dbl>
## 1 21742443 2019-01-01 0:04:37 2019-01-01 0:~    2167          390          199
```

```
## 2 21742444 2019-01-01 0:08:13 2019-01-01 0:~ 4386 441 44
## 3 21742445 2019-01-01 0:13:23 2019-01-01 0:~ 1524 829 15
## 4 21742446 2019-01-01 0:13:45 2019-01-01 0:~ 252 1783 123
## 5 21742447 2019-01-01 0:14:52 2019-01-01 0:~ 1170 364 173
## 6 21742448 2019-01-01 0:15:33 2019-01-01 0:~ 2437 216 98
## # i 6 more variables: from_station_name <chr>, to_station_id <dbl>,
## # to_station_name <chr>, usertype <chr>, gender <chr>, birthyear <dbl>
head(divvy_2020)
```

```
## # A tibble: 6 x 13
## ride_id rideable_type started_at ended_at start_station_name start_station_id
## <chr> <chr> <chr> <chr> <chr> <dbl>
## 1 EACB191~ docked_bike 2020-01-2~ 2020-01~ Western Ave & Lel~ 239
## 2 8FED874~ docked_bike 2020-01-3~ 2020-01~ Clark St & Montro~ 234
## 3 789F3C2~ docked_bike 2020-01-0~ 2020-01~ Broadway & Belmon~ 296
## 4 C9A388D~ docked_bike 2020-01-0~ 2020-01~ Clark St & Randol~ 51
## 5 943BC3C~ docked_bike 2020-01-3~ 2020-01~ Clinton St & Lake~ 66
## 6 6D9C8A6~ docked_bike 2020-01-1~ 2020-01~ Wells St & Hubbar~ 212
## # i 7 more variables: end_station_name <chr>, end_station_id <dbl>,
## # start_lat <dbl>, start_lng <dbl>, end_lat <dbl>, end_lng <dbl>,
## # member_casual <chr>
```

□ *Quick check:* Just verifying both files load properly — better to catch encoding or column mismatches early.

Cleaning & Standardizing Columns

```
divvy_2019 <- clean_names(divvy_2019)
divvy_2020 <- clean_names(divvy_2020)

divvy_2019 <- divvy_2019 %>%
  rename(
    ride_id = trip_id,
    started_at = start_time,
    ended_at = end_time,
    start_station_id = from_station_id,
    start_station_name = from_station_name,
    end_station_id = to_station_id,
    end_station_name = to_station_name,
    member_casual = usertype
  ) %>%
  mutate(ride_id = as.character(ride_id))
```

□ *Rationale:* 2019 uses older column names, so aligning them with 2020 ensures consistency before merging. Schema

mismatches are common in multi-year data and can break joins if left uncorrected.

Combining Datasets

```
common_cols <- intersect(colnames(divvy_2019), colnames(divvy_2020))
divvy_data <- bind_rows(
  divvy_2019 %>% select(all_of(common_cols)),
  divvy_2020 %>% select(all_of(common_cols))
)

rm(divvy_2019, divvy_2020)
skim_without_charts(divvy_data)
```

Table 1: Data summary

| | |
|------------------------|------------|
| Name | divvy_data |
| Number of rows | 791956 |
| Number of columns | 8 |
| Column type frequency: | |
| character | 6 |
| numeric | 2 |
| Group variables | None |

Variable type: character

| skim_variable | n_missing | complete_rate | min | max | empty | n_unique | whitespace |
|--------------------|-----------|---------------|-----|-----|-------|----------|------------|
| ride_id | 0 | 1 | 7 | 16 | 0 | 791956 | 0 |
| started_at | 0 | 1 | 18 | 19 | 0 | 742287 | 0 |
| ended_at | 0 | 1 | 18 | 19 | 0 | 737899 | 0 |
| start_station_name | 0 | 1 | 5 | 43 | 0 | 636 | 0 |
| end_station_name | 1 | 1 | 5 | 43 | 0 | 636 | 0 |
| member_casual | 0 | 1 | 6 | 10 | 0 | 4 | 0 |

Variable type: numeric

| skim_variable | n_missing | complete_rate | mean | sd | p0 | p25 | p50 | p75 | p100 |
|------------------|-----------|---------------|--------|--------|----|-----|-----|-----|------|
| start_station_id | 0 | 1 | 204.40 | 158.92 | 2 | 77 | 174 | 291 | 675 |
| end_station_id | 1 | 1 | 204.38 | 159.32 | 2 | 77 | 174 | 291 | 675 |

□ *Note:* Merged successfully. Removing old objects frees up memory. `skim_without_charts()` gives a quick data-health summary — missing values, column types, etc.

Cleaning Nulls, Timestamps, and Duration

```
divvy_data <- divvy_data %>%
  drop_na() %>%
  mutate(
    started_at = as.POSIXct(started_at, format = "%Y-%m-%d %H:%M:%S"),
    ended_at   = as.POSIXct(ended_at,   format = "%Y-%m-%d %H:%M:%S"),
    ride_length = as.numeric(difftime(ended_at, started_at, units = "mins"))
  ) %>%
  filter(ride_length > 0)
```

□ *Thinking:* Dropped nulls because they're minimal — incomplete records add noise. Negative or zero durations represent system test rides or logging errors.

Removing Test Rides & Standardizing Labels

```
divvy_data <- divvy_data %>%
  filter(start_station_name != "HQ QR", ride_length < 1440) %>%
  mutate(member_casual = recode(member_casual,
                                "Subscriber" = "member",
                                "Customer"   = "casual"))
```

□ *Observation:* “HQ QR” denotes internal QA rides. A 1-day (1440 min) cutoff removes anomalies. Renaming “Subscriber/Customer” keeps rider-type labels consistent.

Ride Length Summary

```
summary(divvy_data$ride_length)
```

```
##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## 1.667e-02 5.500e+00 8.983e+00 1.372e+01 1.518e+01 1.436e+03
```

□ *Interpretation:* Most rides last **5–15 minutes** (median ≈ 9 min) — short, practical commuting behavior.

DATA ANALYSIS & INSIGHTS

Feature Engineering

```
divvy_data <- divvy_data %>%
  mutate(
    month = month(started_at, label = TRUE, abbr = TRUE),
    day_of_week = wday(started_at, label = TRUE, abbr = TRUE),
    hour = hour(started_at),
    season = case_when(
      month(started_at) %in% c(12,1,2) ~ "Winter",
      month(started_at) %in% c(3,4,5) ~ "Spring",
      month(started_at) %in% c(6,7,8) ~ "Summer",
      TRUE ~ "Fall"
    ),
    week_type = ifelse(day_of_week %in% c("Sat","Sun"), "Weekend", "Weekday"),
    ride_category = case_when(
      ride_length <= 5 ~ "Short (<5 min)",
      ride_length <= 15 ~ "Medium (5-15 min)",
      ride_length <= 30 ~ "Long (15-30 min)",
      TRUE ~ "Extended (>30 min)"
    )
  )
```

□ *Reasoning:* Creating temporal and categorical features early simplifies later exploration. Season, weekday, and duration categories expose behavioral patterns at a glance.

Ride Duration Patterns

```
divvy_data %>%
  group_by(member_casual) %>%
  summarise(
    avg_ride_length = mean(ride_length),
    median_ride_length = median(ride_length)
  )
```

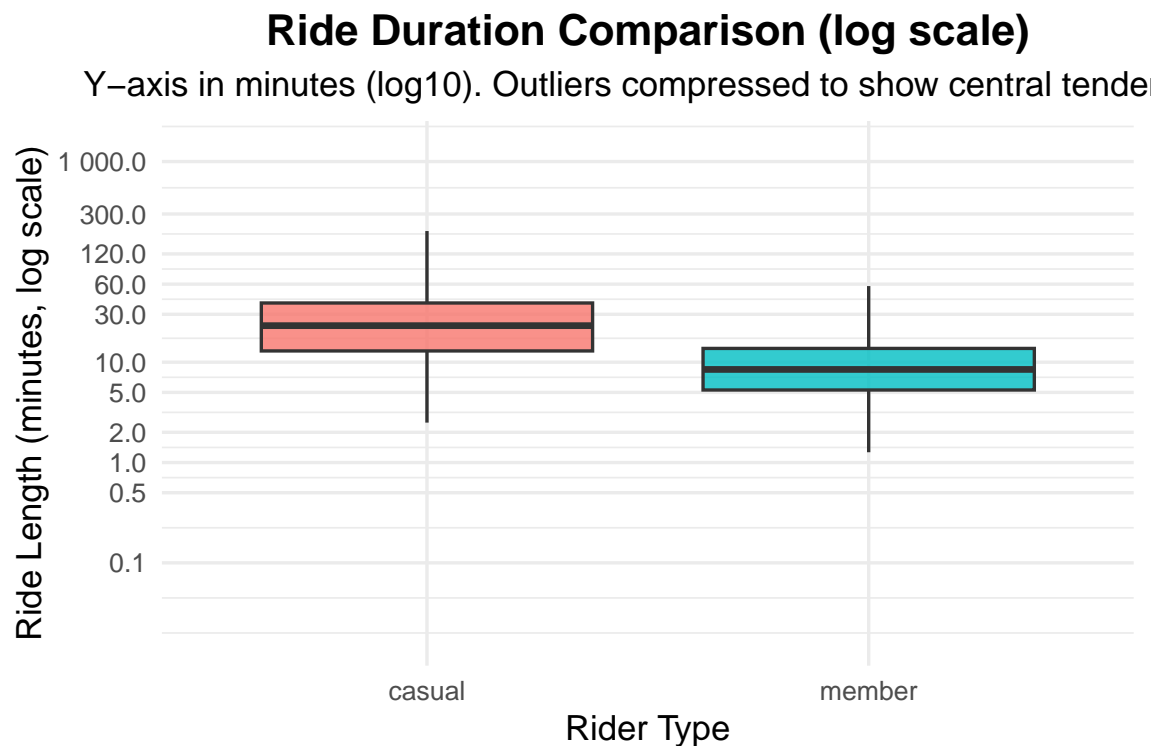
```
## # A tibble: 2 x 3
##   member_casual avg_ride_length median_ride_length
##   <chr>          <dbl>          <dbl>
## 1 casual        38.4            23.1
## 2 member        11.4            8.47
```

□ *Result:* Casual riders' average rides are notably longer — clear leisure orientation.

```
library(ggplot2)
library(dplyr)
```

```
library(scales)

divvy_data %>%
  filter(ride_length > 0) %>%           # remove zeros/negatives
  ggplot(aes(x = member_casual, y = ride_length, fill = member_casual)) +
  geom_boxplot(outlier.shape = NA, alpha = 0.8) + # hide extreme outliers for visibility
  #geom_jitter(width = 0.15, alpha = 0.12, size = 0.6) + # optional: show point cloud
  scale_y_log10(
    labels = label_number(accuracy = 0.1),
    breaks = c(0.1, 0.5, 1, 2, 5, 10, 30, 60, 120, 300, 1000)
  ) +
  labs(
    title = "Ride Duration Comparison (log scale)",
    subtitle = "Y-axis in minutes (log10). Outliers compressed to show central tendency.",
    x = "Rider Type",
    y = "Ride Length (minutes, log scale)"
  ) +
  theme_minimal(base_size = 13) +
  theme(legend.position = "none",
        plot.title = element_text(hjust = 0.5, face = "bold"),
        plot.subtitle = element_text(hjust = 0.5))
```



Notes:

- `scale_y_log10()` compresses the heavy tail so medians and IQRs are more visible.
- I hid extreme outliers in the boxplot (`outlier.shape = NA`) (but added a faint jitter cloud in comment so you can still see distribution density if you remove the #).
- breaks include small and large minute values so axis ticks are meaningful.

□ *Insight:* Casuals display a far wider duration spread — the classic *tourist vs commuter* contrast. Members ride frequently but briefly; casuals ride fewer times but linger longer.

Round Trips vs One-Way

```
divvy_data %>%
  mutate(is_round_trip = start_station_id == end_station_id) %>%
  group_by(member_casual, is_round_trip) %>%
  summarise(total_rides = n())
```

`summarise()` has grouped output by 'member_casual'. You can override using the ## `.groups` argument.

A tibble: 4 x 3

Groups: member_casual [2]

| ## | member_casual | is_round_trip | total_rides |
|------|---------------|---------------|-------------|
| ## | <chr> | <lgl> | <int> |
| ## 1 | casual | FALSE | 57903 |
| ## 2 | casual | TRUE | 9679 |
| ## 3 | member | FALSE | 706273 |
| ## 4 | member | TRUE | 13852 |

□ *Takeaway:* Casuals record more round trips — likely sightseeing or leisure loops. Members' rides are overwhelmingly one-way, consistent with commuting.

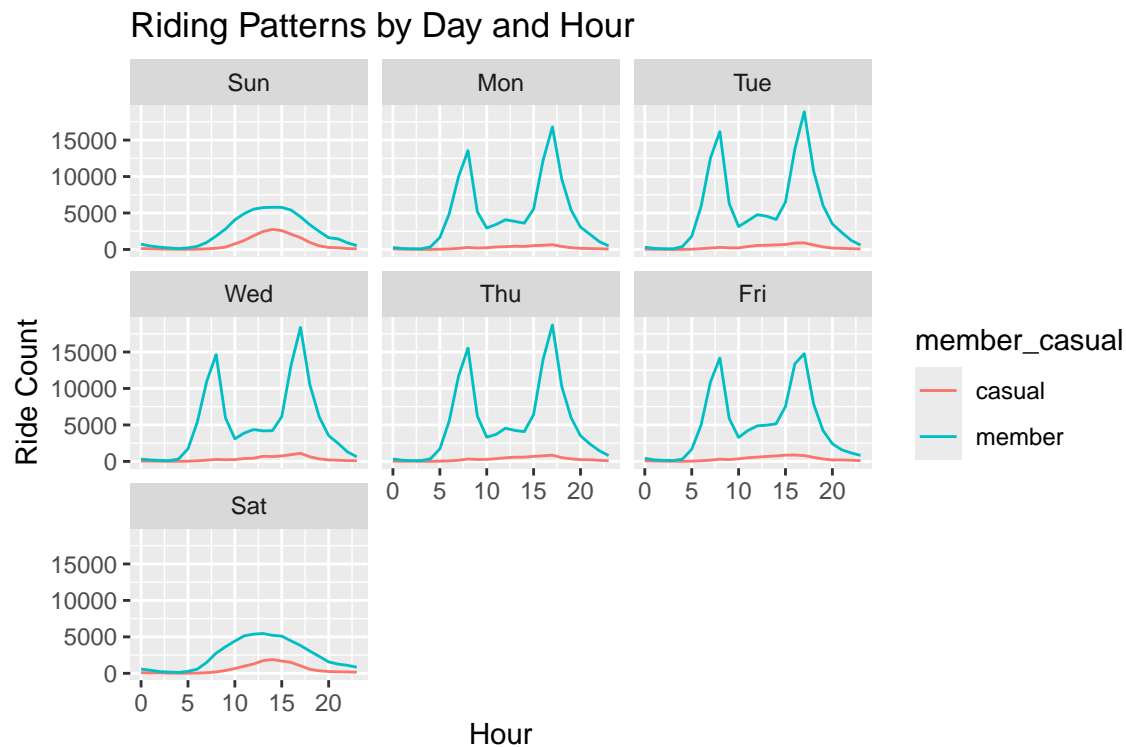
Temporal Insights

Day + Hour Trends

```
divvy_data %>%
  group_by(member_casual, day_of_week, hour) %>%
  summarise(total_rides = n()) %>%
  ggplot(aes(x = hour, y = total_rides, color = member_casual)) +
  geom_line() +
  facet_wrap(~ day_of_week) +
  labs(title = "Riding Patterns by Day and Hour",
       x = "Hour", y = "Ride Count")
```

`summarise()` has grouped output by 'member_casual', 'day_of_week'. You can

override using the ``.groups`` argument.

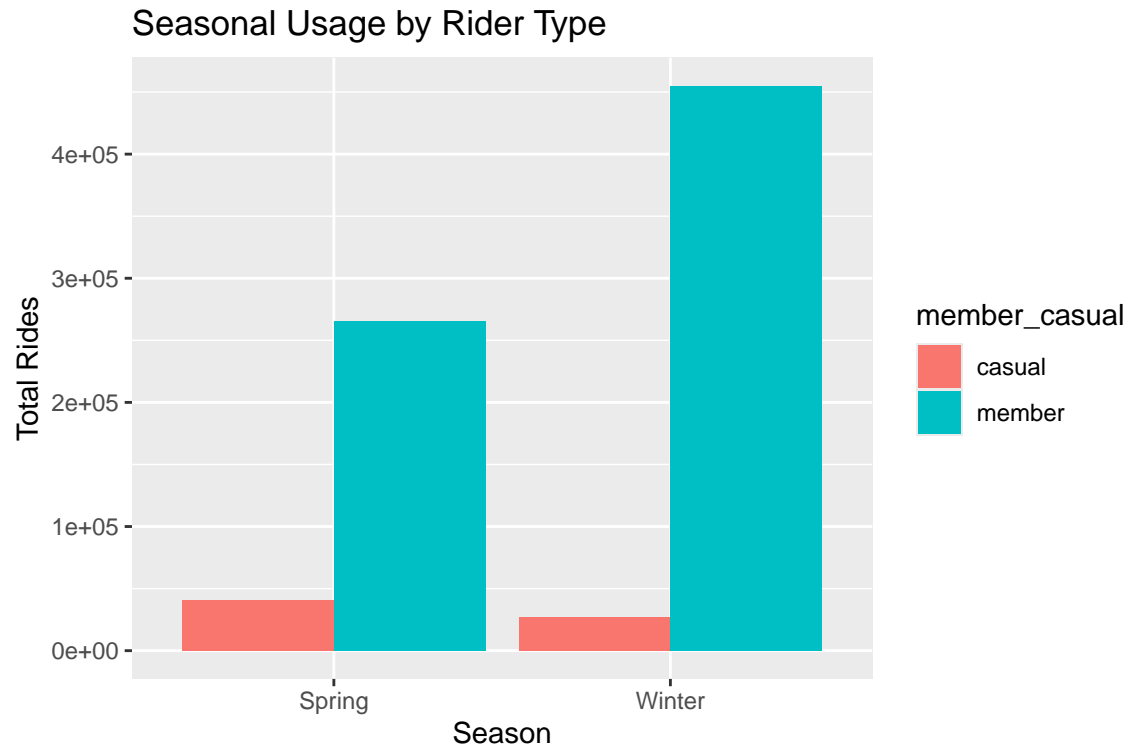


□ *Interpretation:* Members peak at **8–9 AM** and **5–6 PM** — commute hours. Casuals peak mid-day and on weekends — pure leisure timing. Two rhythms emerge: *structured weekdays* vs *spontaneous weekends*.

Seasonal Patterns

```
divvy_data %>%
  group_by(member_casual, season) %>%
  summarise(total_rides = n(),
            avg_ride_length = mean(ride_length)) %>%
  ggplot(aes(x = season, y = total_rides, fill = member_casual)) +
  geom_bar(stat = "identity", position = "dodge") +
  labs(title = "Seasonal Usage by Rider Type",
       x = "Season", y = "Total Rides")
```

``summarise()`` has grouped output by 'member_casual'. You can override using the
``.groups`` argument.

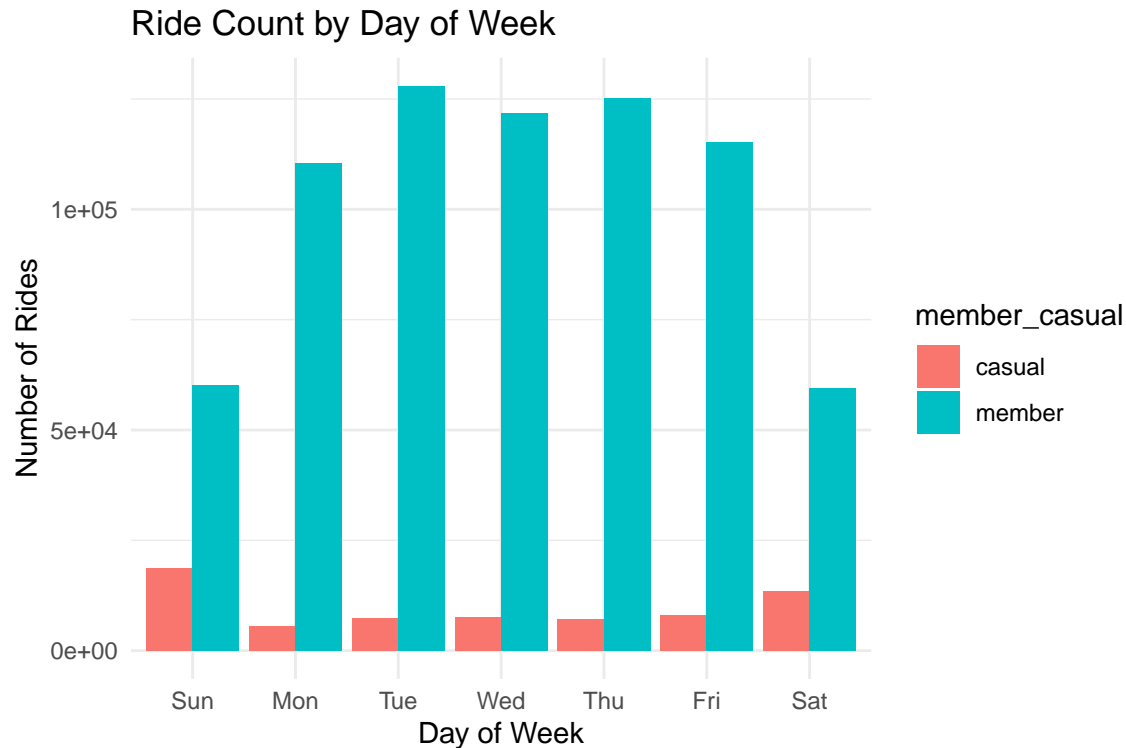


□ *Observation:* Casual ridership soars in **spring & summer**, drops in winter — weather-sensitive. Members stay consistent year-round, reflecting reliability and habit.

Weekly Trends

```
divvy_data %>%
  group_by(member_casual, day_of_week) %>%
  summarise(total_rides = n()) %>%
  ggplot(aes(x = day_of_week, y = total_rides, fill = member_casual)) +
  geom_col(position = "dodge") +
  labs(title = "Ride Count by Day of Week",
       x = "Day of Week", y = "Number of Rides") +
  theme_minimal()
```

`summarise()` has grouped output by 'member_casual'. You can override using the
`groups` argument.



□ *Observation:* Members dominate weekdays; casuals surge on weekends — reinforcing work-vs-leisure behavior.

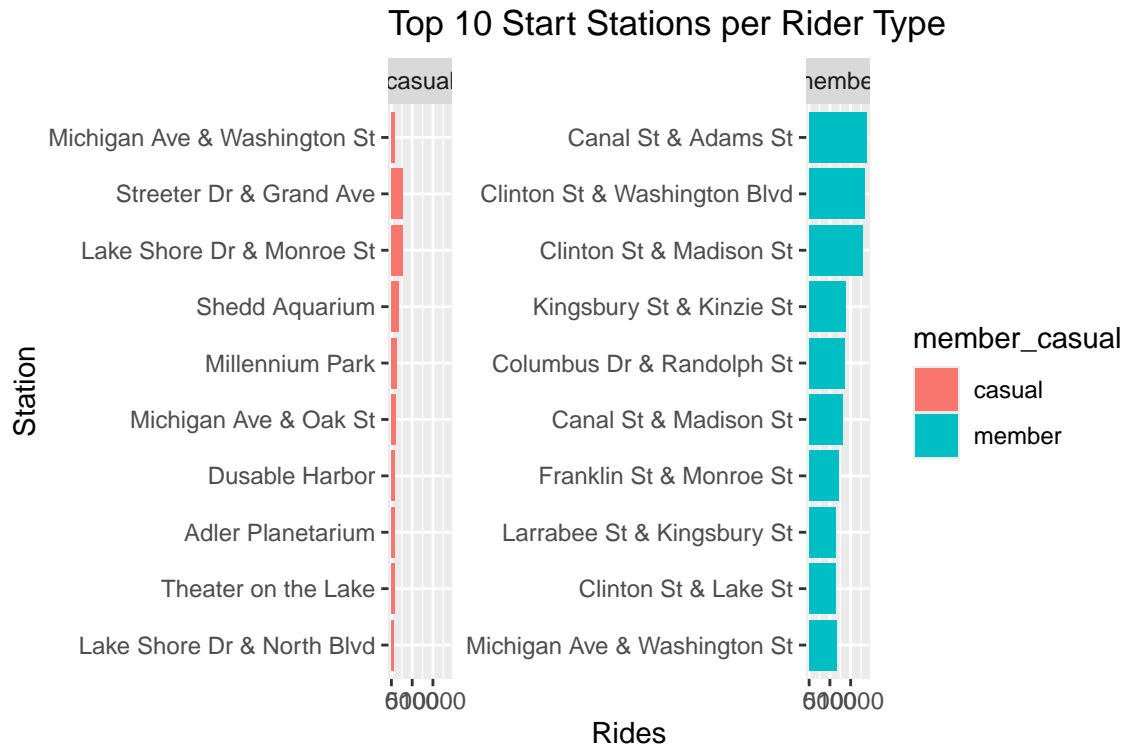
Spatial Insights

Top 10 Start Stations

```
top_stations <- divvy_data %>%
  group_by(member_casual, start_station_name) %>%
  summarise(rides = n()) %>%
  arrange(member_casual, desc(rides)) %>%
  slice_head(n = 10)
```

`summarise()` has grouped output by 'member_casual'. You can override using the
`.groups` argument.

```
ggplot(top_stations,
  aes(x = reorder(start_station_name, rides),
    y = rides, fill = member_casual)) +
  geom_bar(stat = "identity") +
  coord_flip() +
  facet_wrap(~ member_casual, scales = "free_y") +
  labs(title = "Top 10 Start Stations per Rider Type",
    x = "Station", y = "Rides")
```



□ *Insight:* Members cluster around **business districts and transit hubs**. Casuals favor **tourist zones and waterfronts** — ideal spots for day-pass marketing.

STRATEGIC TAKEAWAYS

| Insight Area | Observation | Business Implication |
|----------------------|--|---|
| Ride Duration | Casual rides are $\approx 3\times$ longer | Promote annual memberships for frequent leisure users |
| Peak Times | Members commute mornings/evenings; casuals ride weekends | Offer weekday discounts to casuals |
| Seasonality | Casual rides surge in summer | Launch pre-summer marketing campaigns |
| Top Stations | Casuals use tourist-heavy stations | Place ads / QR banners at these hubs |
| Round Trips | Casuals take more round trips | Introduce “unlimited short-ride” plans |

□ *Final Thoughts:* **Members = commuters; Casuals = tourists and leisure seekers.** Push membership incentives during warmer months when casual ridership peaks — ideally before summer weekends. Highlight **cost savings, convenience, and priority perks** to nudge casual riders into annual membership tiers.

Recommendations

1□ Seasonal Membership Campaign (Late Spring – Early Summer)

What: Launch a targeted membership promotion focused on converting casual riders during peak leisure months.

How:

- Advertise membership perks tied to leisure use (ride discounts, scenic-station offers, partner attraction perks).
- Use QR codes at high-traffic scenic stations for instant sign-up.

When:

- Begin advertising in **late April**
- Run campaign through **July**, when casual ridership is highest.

By Whom:

- **Marketing Team** (campaign design, messaging, creatives)
 - **Partnerships Team** (tourist sites, attractions, scenic zones)
 - **Operations Team** (QR placement at stations)
-

2□ Public Transport Advertising for Casual Riders

What: Promote Cyclistic membership on public transport channels frequently used by casual riders.

How:

- Place ads with QR codes inside buses, trains, and metro stations connected to tourist and weekend-heavy areas.
- Highlight convenience and cost savings for multi-modal travelers.

When:

- Run continuously, with increased frequency on **weekends** and **tourist season** (May–September).

By Whom:

- **Marketing Team** (ad design, CTA optimisation)
 - **Transit Advertising Partners** (bus, metro, rail networks)
 - **Brand Outreach Team** (negotiating placements)
-

3□ Add Extra Ride-Time Benefits to Annual Membership

What: Provide additional free ride-time hours as a built-in benefit for annual members.

How:

- Offer monthly bonus minutes or extended ride duration before overage fees.

- Promote these perks during onboarding and all seasonal campaigns.

When:

- Introduce the benefit at the start of **summer** (June)
- Maintain year-round as a permanent membership feature.

By Whom:

- **Product Team** (membership feature design)
 - **Finance Team** (cost modeling)
 - **Marketing Team** (promotion and messaging)
-

If you'd like, I can also format this as a **one-page executive summary**, a **slide deck**, or insert it cleanly into your final PDF/HTML report.

CONCLUSION

This analysis explored how Cyclistic's members and casual riders differ in their riding behavior across time, season, and location.

Members ride frequently, briefly, and predictably — reflecting daily commuters who depend on Cyclistic for practicality and routine.

Casual riders, on the other hand, ride longer, mostly on weekends and during warmer months, often looping back to where they started — echoing leisure, tourism, and exploration.

These patterns reveal two distinct motivations: **efficiency vs experience**.

Bridging that gap through well-timed seasonal promotions, flexible membership plans, and on-site marketing at tourist stations can transform occasional riders into loyal members.