

Ride Membership Analytics

Toyin Olapade

2021-11-08

This is a capstone project undertaken as part of the Google Data Analytics course offered on Coursera. It involves working with Cyclistic, a bike-share company in Chicago, as a Junior Data Analyst. The director of marketing at Cyclistic believes the company's future success depends on maximizing the number of annual memberships. Therefore, the marketing team wants to understand how casual riders and annual members use Cyclistic bikes differently. From these insights, the team will design a new marketing strategy to convert casual riders into annual members. In order to proceed with this goal the executives at Cyclistic require compelling data insights and professional data visualizations, and must approve recommendations resulting from this analysis.

Introduction

In 2016, Cyclistic¹ launched a successful bike-share offering. Since then, the program has grown to a fleet of 5,824 bicycles that are geo-tracked and locked into a network of 692 stations across Chicago. The bikes can be unlocked from one station and returned to any other station in the system anytime.

Until now, Cyclistic's marketing strategy relied on building general awareness and appealing to broad consumer segments. One approach that helped make these things possible was the flexibility of its pricing plans: single-ride passes, full-day passes, and annual memberships. Customers who purchase single-ride or full-day passes are referred to as **casual** riders. Customers who purchase annual memberships are Cyclistic **members**².

Cyclistic's finance analysts have concluded that **annual** members are much more profitable than **casual** riders. Although the pricing flexibility helps Cyclistic attract more customers, Moreno³ believes that maximizing the number of annual members will be key to future growth. Rather than creating a marketing campaign that targets all-new customers, she believes there is a very good chance to convert casual riders into members. She notes that casual riders are already aware of the Cyclistic program and have chosen Cyclistic for their mobility needs.

Moreno has set a clear goal: Design marketing strategies aimed at converting casual riders into annual members. In order to do that, however, my team⁴ needs to better understand how annual members and casual riders differ, why casual riders would buy a membership, and how digital media could affect their marketing tactics. Moreno

¹ **Cyclistic** is bike-share company resident in Chicago.

² For the purpose of clear description during data analysis we have relabelled the **member** to **annual**.

³ Lily Moreno is the director of marketing at Cyclistic

⁴ The Marketing Analytics Team

and her team are interested in analyzing the Cyclistic historical bike trip data to identify trends.

Three questions will guide the future marketing program: How do annual members and casual riders use Cyclistic bikes differently? Why would casual riders buy Cyclistic annual memberships? How can Cyclistic use digital media to influence casual riders to become members? I have been assigned the first question, and this note⁵ reports the following deliverable:

1. A clear statement of the business task
2. A description of all data sources used
3. Documentation of any cleaning or manipulation of data
4. A summary of your analysis
5. Supporting visualizations and key findings
6. Your top three recommendations based on your analysis

Business Task

In order to convert casual riders to annual members, Cyclistic wants to understand how annual members and casual riders differs.

Data Collection

This data is collected from Cyclistic's marketing strategy at Google case study, from 05/2020 to 04/2021.

Connecting to Postgres

We used the DBI package to facilitate database connectivity. It has the core functionality of connecting R to database servers. We load, `RPostgres`, a package that implements the core functionality of DBI for `PostgreSQL`, create and test a `dbConnect` object to hold the connection to `Postgres` database.⁶

Data Dictionary

We use Cyclistic's historical trip data⁷ to analyze and identify trends. This is public data that you can use to explore how different customer types are using Cyclistic bikes. But note that data-privacy issues prohibit you from using riders' personally identifiable information. This means that you won't be able to connect pass purchases to credit card numbers to determine if casual riders live in the Cyclistic service area or if they have purchased multiple single passes.

⁵ This note is prepared in RStudio. See Irene Steve's [Using SQL in RStudio](https://irene.rbind.io/post/using-sql-in-rstudio/#passing-variables-tofrom-sql-chunks) at <https://irene.rbind.io/post/using-sql-in-rstudio/#passing-variables-tofrom-sql-chunks> and [Creating Dynamic Documents with RMarkdown and Knitr](https://rstudio-pubs-static.s3.amazonaws.com/180546_e2d5bf84795745ebb5cd3be3dab71fca.html) at https://rstudio-pubs-static.s3.amazonaws.com/180546_e2d5bf84795745ebb5cd3be3dab71fca.html for guidance. For the backend DBMS we used Postgres (pgAdmin 14).

⁶ It should be noted that this code was generated on my local machine connected to a local copy of the database. Your connection details will be different. I also have permissions to modify this database.

⁷ See [data](https://divvy-tripdata.s3.amazonaws.com/index.html) at <https://divvy-tripdata.s3.amazonaws.com/index.html>. The datasets have a different name because Cyclistic is a fictional company. For the purposes of this case study, the datasets are appropriate and will enable you to answer the business questions. The data has been made available by Motivate International Inc. under this [license](#)

Table 1: Cyclistic Data Dictionary.

var_nm	description	data_type
ride_id	Unique identifier for each ride	text
rideable_type	Bike types rideable (classic, docked, electric)	text
started_at	Time at start of trip	datetime
ended_at	Time at end of trip	datetime
start_station_name	Trip-start station name	text
start_station_id	Unique trip-start station identifier	text
end_station_name	Trip-end station name	text
end_station_id	Unique trip-end station identifier	text
start_lat	Latitude of trip-start geolocation	text
start_lng	Longitude of trip-start geolocation	text
end_lat	Latitude of trip-end geolocation	numeric
end_lng	Longitude of trip-end geolocation	numeric
member_tpy	Type of riders (casual, annual)	text

Data Preparation

We want to take a look at 12 months of data ranging from 2020-04-01 to 2021-03-31. We notice that for easier description it may be better to rename the `member_casual` field to `member_tpy` and its values from 'member' to 'annual'⁸.

⁸ See Appendix for SQL queries used to achieve this.

```
select ride_id, rideable_type, started_at, member_tpy
from divvy_tripdata_
where date(started_at) >= cast('2020-04-01' as date)
and date(started_at) <= cast('2021-03-31' as date)
```

```
## # A tibble: 6 x 4
##   ride_id      rideable_type started_at      member_tpy
##   <chr>        <chr>      <dtm>      <chr>
## 1 07F8B580615B90E1 classic_bike 2020-12-17 16:53:24 casual
## 2 6B4F9AEFDD1C7451 classic_bike 2020-12-07 18:53:13 annual
## 3 B0841F69454DB588 electric_bike 2020-12-14 10:28:21 annual
## 4 1B882D036F1934B3 electric_bike 2020-12-01 10:54:42 annual
## 5 90BDD2BD6CD16182 classic_bike 2020-12-21 15:52:04 annual
## 6 A7B686BFC96E3A1D classic_bike 2020-12-06 12:40:21 annual
```

Questions to explore are:

- Which station is patronized the most, and what type of members patronized that station? In other words, does a member type tend to prefer a station over another?

- Do members tend to drop their bike at the same place they picked it from?

Furthermore, one question that is more striking is how many riders patronizes Cyclistics within some period and how is membership type distributed among these riders? However, the data collected by Cyclistic only describe rides and not riders. Each ride has a unique ride id.

```
select member_typ, count(distinct ride_id) as n_ride
from divvy_tripdata_
where date(started_at) >= cast('2020-04-01' as date)
      and date(started_at) <= cast('2021-03-31' as date)
group by member_typ

dt_01 %>%
  group_by(member_typ) %>%
  summarise (n=n()) %>%
  mutate(pct=paste0(round(100 * n/sum(n), 0), "%"))

## # A tibble: 2 x 3
##   member_typ      n pct
##   <chr>         <int> <chr>
## 1 annual       2059372 59%
## 2 casual       1430376 41%

## # A tibble: 6 x 4
##   ride_id      rideable_type started_at      member_typ
##   <chr>         <chr>      <dtm>         <chr>
## 1 07F8B580615B90E1 classic_bike 2020-12-17 16:53:24 casual
## 2 6B4F9AEFDD1C7451 classic_bike 2020-12-07 18:53:13 annual
## 3 B0841F69454DB588 electric_bike 2020-12-14 10:28:21 annual
## 4 1B882D036F1934B3 electric_bike 2020-12-01 10:54:42 annual
## 5 90BDD2BD6CD16182 classic_bike 2020-12-21 15:52:04 annual
## 6 A7B686BFC96E3A1D classic_bike 2020-12-06 12:40:21 annual
```

1226304 casual ride that use classic_bike

How do annual members and casual riders use Cyclistic bikes differently

Assumption: We are assuming that a ride day is identified by the ride start date-time.

We're trying to compute number of rides per day by ride membership type. We wanted to do that for just one day and we realized that we needed to convert text to date.

```
select member_typ,
  to_char(date(started_at), 'yy-mm') as ride_mth,
```

```

    date(started_at) as ride_dt,
    count(distinct ride_id) as ride_cnt
from divvy_tripdata_
where date(started_at) >= cast('2020-04-01' as date)
    and date(started_at) <= cast('2021-03-31' as date)
group by member_typ, to_char(date(started_at), 'yy-mm'), date(started_at)

## # A tibble: 726 x 4
##   member_typ ride_mth ride_dt   ride_cnt
##   <chr>      <chr>    <date>    <int64>
## 1 annual    20-04    2020-04-01    1895
## 2 annual    20-04    2020-04-02    2060
## 3 annual    20-04    2020-04-03    2570
## 4 annual    20-04    2020-04-04    1758
## 5 annual    20-04    2020-04-05    1982
## 6 annual    20-04    2020-04-06    1935
## 7 annual    20-04    2020-04-07    3374
## 8 annual    20-04    2020-04-08    1765
## 9 annual    20-04    2020-04-09    1563
## 10 annual   20-04    2020-04-10    2056
## # ... with 716 more rows

```

We now have ride counts per day but we want to visualize its trend by ride membership type so we can compare patterns between them.

```

dt_01 %>%
  mutate(ride_cnt_=as.numeric(ride_cnt)) %>%
  ggplot(mapping=aes(x=ride_mth, y=ride_cnt_, color=member_typ)) +
  geom_boxplot() +
  theme(legend.position = 'bottom')

select member_typ, to_char(ride_dt, 'yy-mm') as ride_mth,
       avg(ride_cnt) as ride_avg
from (
  select member_typ, date(started_at) as ride_dt,
         count(distinct ride_id) as ride_cnt
  from divvy_tripdata_
  where date(started_at) >= cast('2020-04-01' as date)
        and date(started_at) <= cast('2021-03-31' as date)
  group by member_typ, date(started_at)
) as a
group by member_typ, to_char(ride_dt, 'yy-mm')
order by member_typ, ride_mth

dt_02 %>%
  ggplot(mapping=aes(x=ride_mth, y=ride_avg, color=member_typ)) +
  geom_point()

```

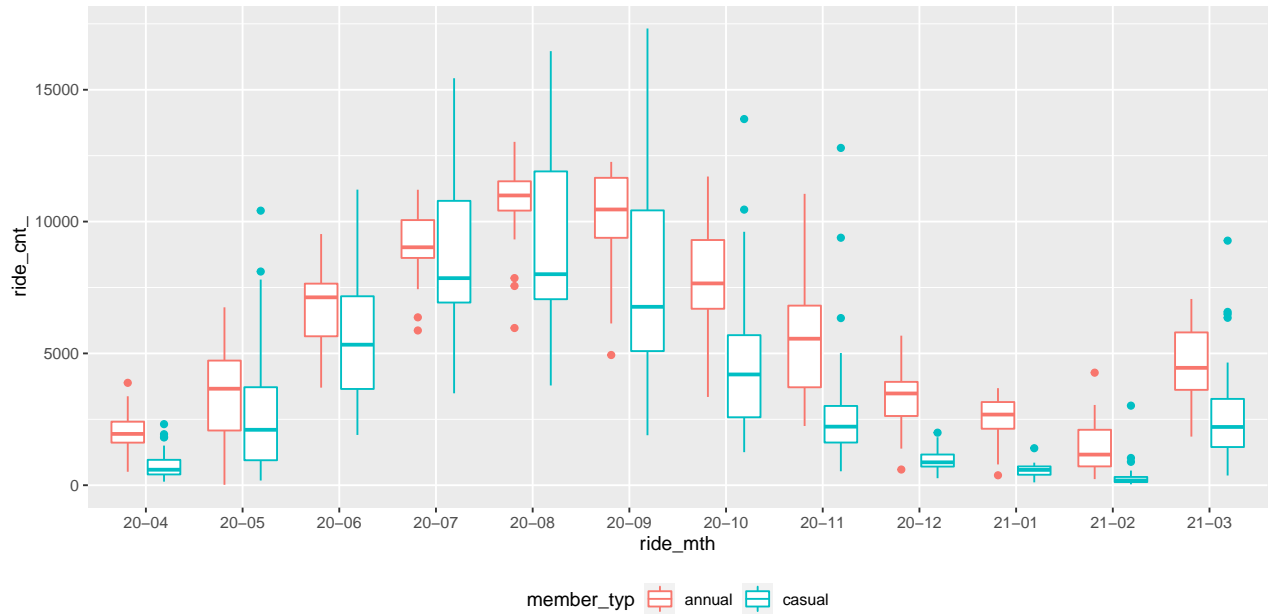
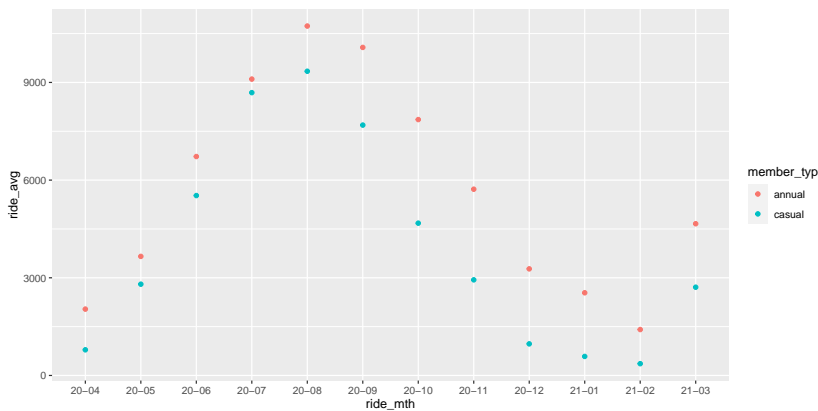


Figure 1: MPG vs horsepower, colored by transmission.

Figure 2: MPG vs horsepower, colored by transmission.



```
dt_01 %>%
  mutate(ride_cnt_=as.numeric(ride_cnt)) %>%
  ggplot(mapping=aes(x=ride_mth, y=ride_cnt_, color=member_typ)) +
  geom_boxplot(mapping=aes(color=member_typ)) +
  stat_summary(fun.y=mean, geom="point", shape=1)
```

Warning: `fun.y` is deprecated. Use `fun` instead.

```
# geom_point(data=dt_02, mapping=aes(x=ride_mth, y=ride_avg, color=member_typ))
```

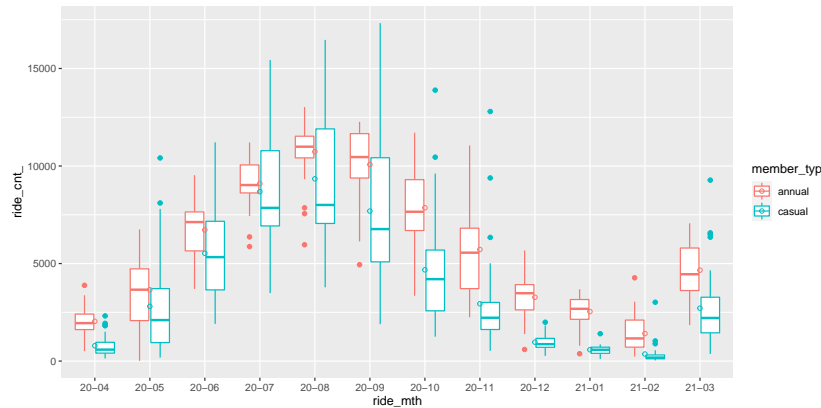


Figure 3: MPG vs horsepower, colored by transmission.

```
dt_01 %>%
  mutate(ride_cnt_=as.numeric(ride_cnt)) %>%
  ggplot(mapping=aes(x=ride_mth, y=ride_cnt_, fill=member_typ)) +
  stat_summary(fun.y=median, geom="point", shape=21, size=4)

## Warning: `fun.y` is deprecated. Use `fun` instead.
```

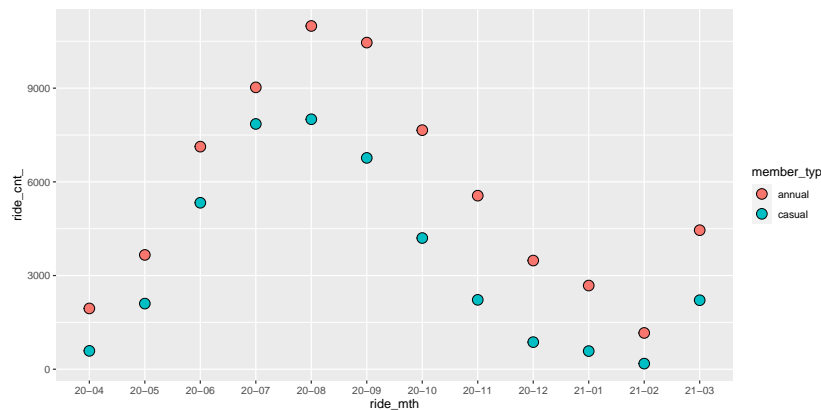


Figure 4: MPG vs horsepower, colored by transmission.

```
dt_01 %>%
  mutate(ride_cnt_=as.numeric(ride_cnt)) %>%
  ggplot(mapping=aes(x=ride_mth, y=ride_cnt_, fill=member_typ)) +
  stat_summary(fun.y=median, geom="point", shape=21, size=4)+
  labs(title="divvy_cyclist:ride_cnt_ vs.ride_mth", subtitle = "sample of the three rideable_type",
       caption = "Data collected by divvy cyclyst")+
  annotate("text", x=5,y=40,label="The annual member are the highest user")

## Warning: `fun.y` is deprecated. Use `fun` instead.
```

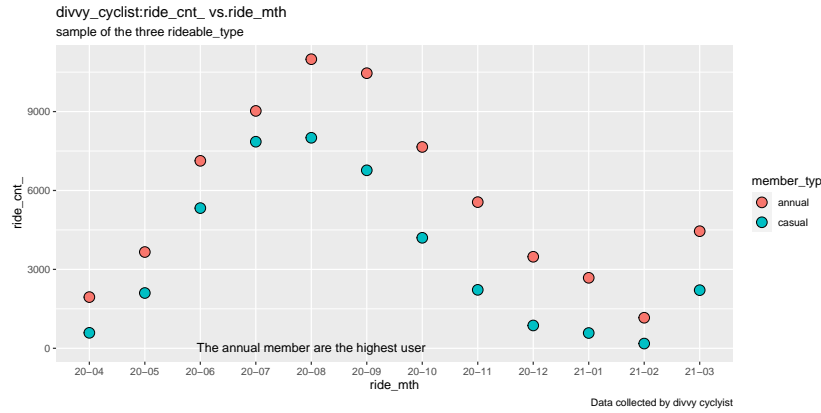


Figure 5: MPG vs horsepower, colored by transmission.

```
select member_type, rideable_type,
  to_char(date(started_at), 'yy-mm') as ride_mth,
  date(started_at) as ride_dt,
  count(distinct ride_id) as ride_cnt
from divvy_tripdata_
where date(started_at) >= cast('2020-04-01' as date)
  and date(started_at) <= cast('2021-03-31' as date)
group by member_type, rideable_type, to_char(date(started_at), 'yy-mm'), date(started_at)

dt_03 %>%
  mutate(ride_cnt_=as.numeric(ride_cnt)) %>%
  ggplot(., mapping=aes(x=ride_mth, y=ride_cnt_, color=member_type)) +
  geom_boxplot() +
  facet_wrap(~rideable_type)
```

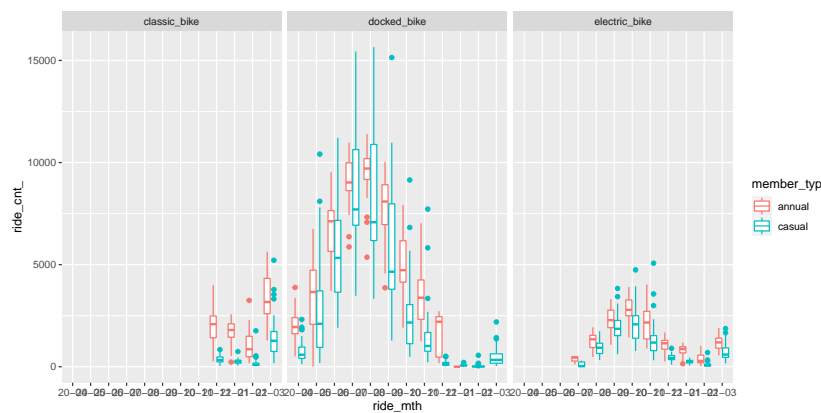


Figure 6: MPG vs horsepower, colored by transmission.

From the plot, the annual member has the high ride per day.


```
dt_03 %>%
  mutate(ride_cnt_=as.numeric(ride_cnt)) %>%
  ggplot(., mapping=aes(x=ride_mth, y=ride_cnt_, color=member_typ, fill=member_typ)) +
  stat_summary(fun.y=median, geom="point", shape=21, size=4) +
  facet_wrap(~rideable_type)+
  labs(title="divvy_cyclist:ride_cnt_ vs.ride_mth", subtitle = "sample of the three rideable_type",
       caption = "Data collected by divvy cyclyist")+
  annotate("text", x=5,y=40,label="The annual member are the highest user",angle=15)

## Warning: `fun.y` is deprecated. Use `fun` instead.
```

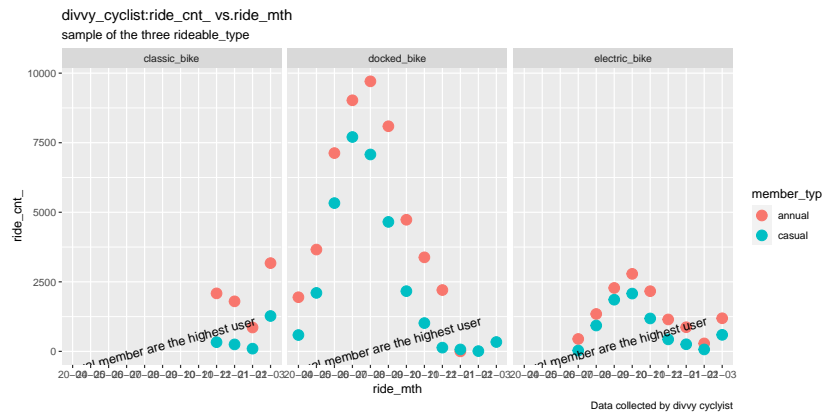


Figure 7: MPG vs horsepower, colored by transmission.