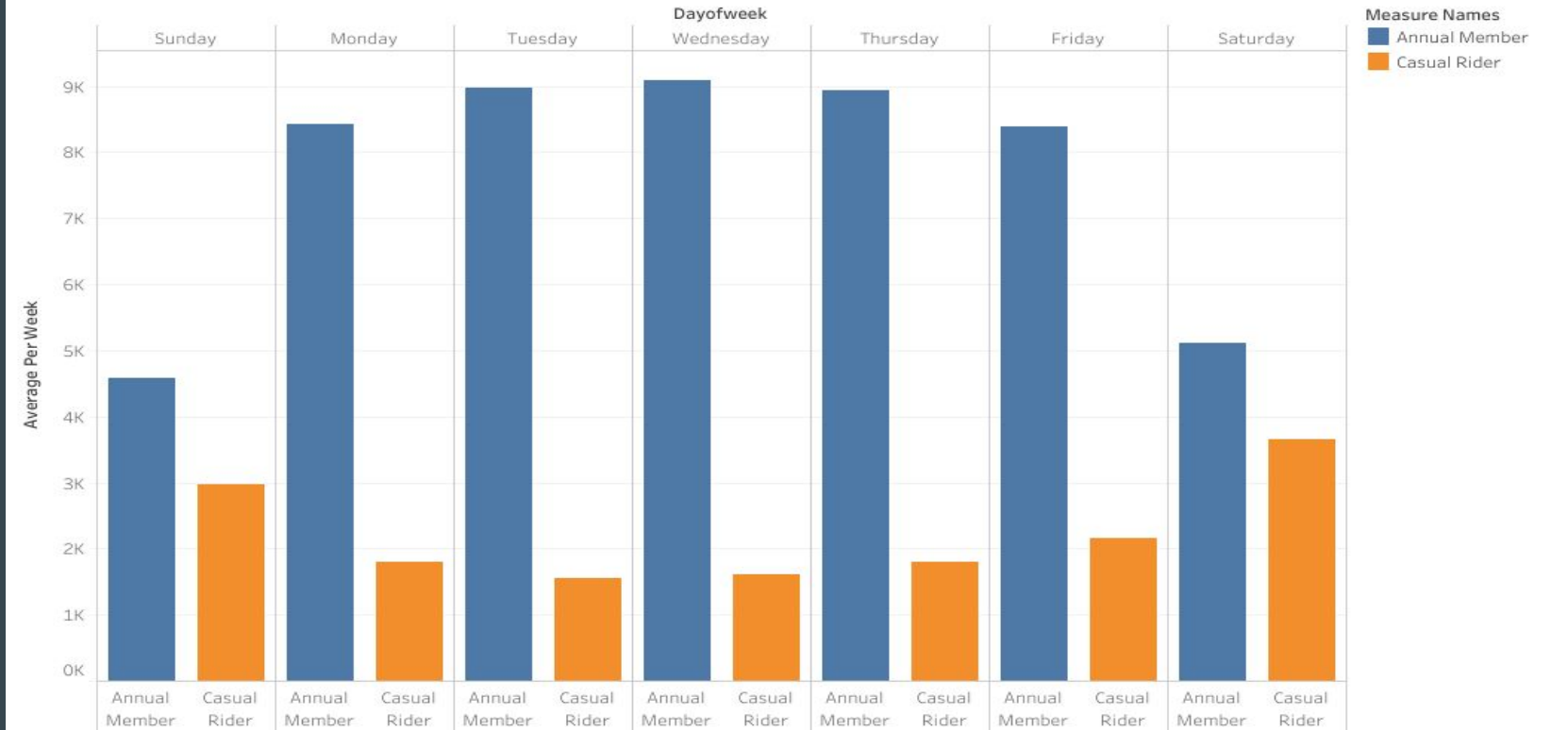# Cyclistic

• • •

# Overview

Created in 2016 and located in Chicago, Illinois, Cyclistic is a successful bike-share company ready to reach new heights. Cyclistic's business model involves casual riders (customers who purchase single-ride or full-day passes) and annual members, and its financial analysts have concluded that annual members drive better profitability. Therefore, the marketing team is now seeking ways to convert casual riders into annual members.

# Key Question

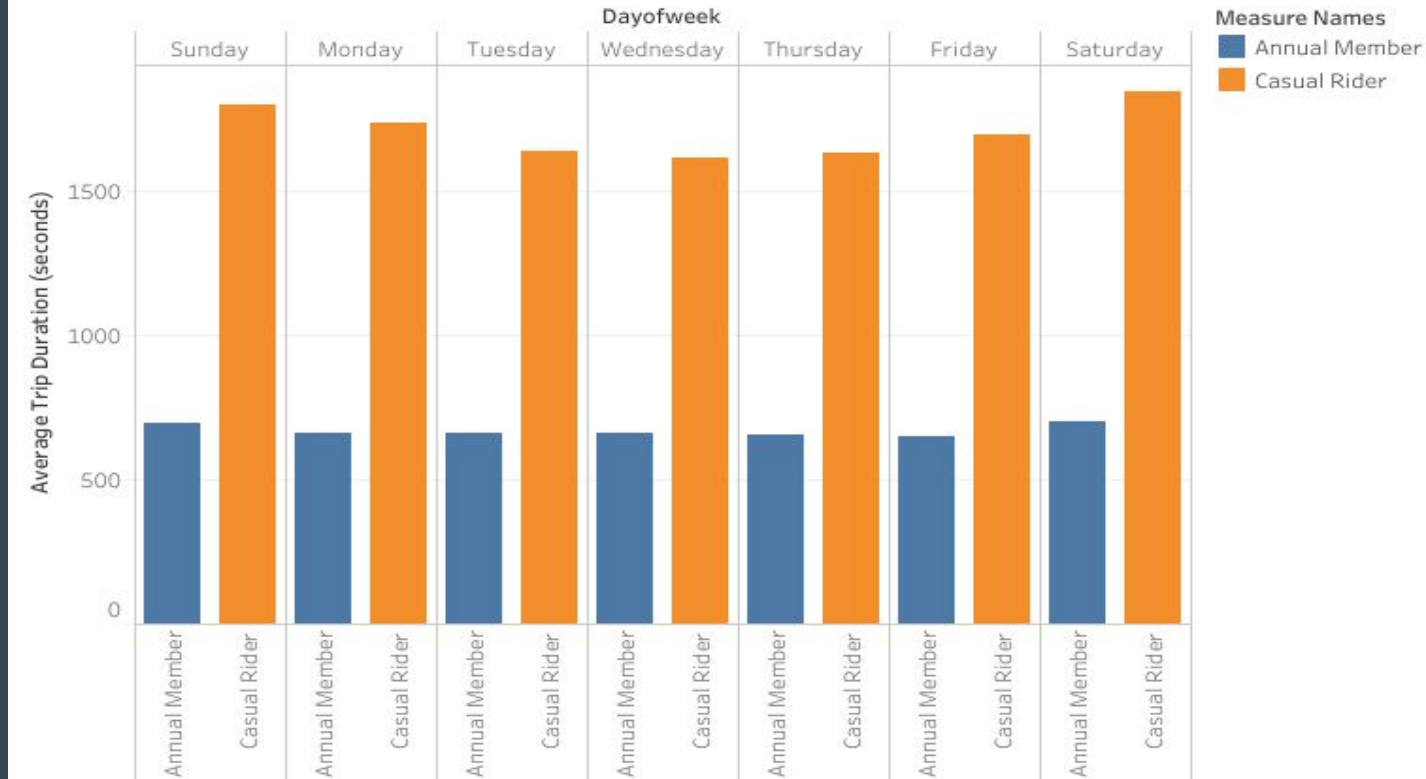How do annual members and casual riders use Cyclistic bikes differently?

We begin our analysis by looking at which days of the week experienced the most number of trips on average. It appears that annual members prefer weekdays whereas casual riders prefer weekends.
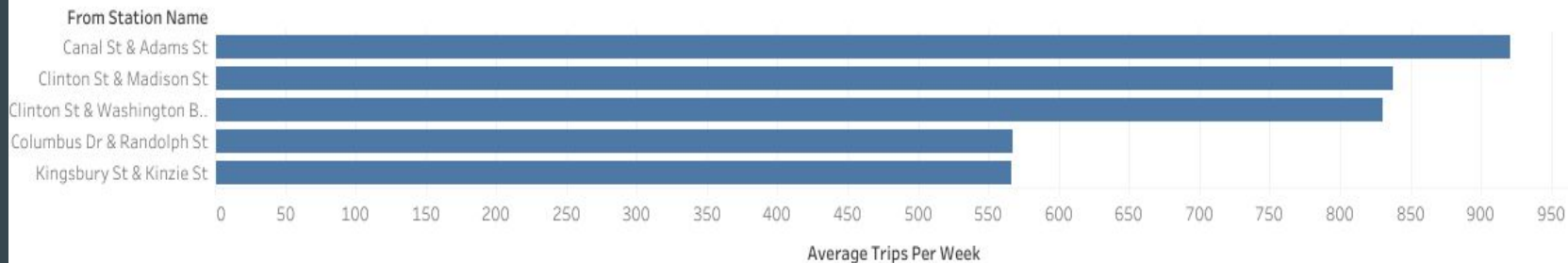
2019 Total Trips Every Week (Chronological)

Through this time series chart, we see that annual members took a larger quantity of trips every week of the year compared to casual riders. The trends of both lines mirror closely which indicates both annual members and casual riders prefer to ride in the summer (week 25 - week 35).
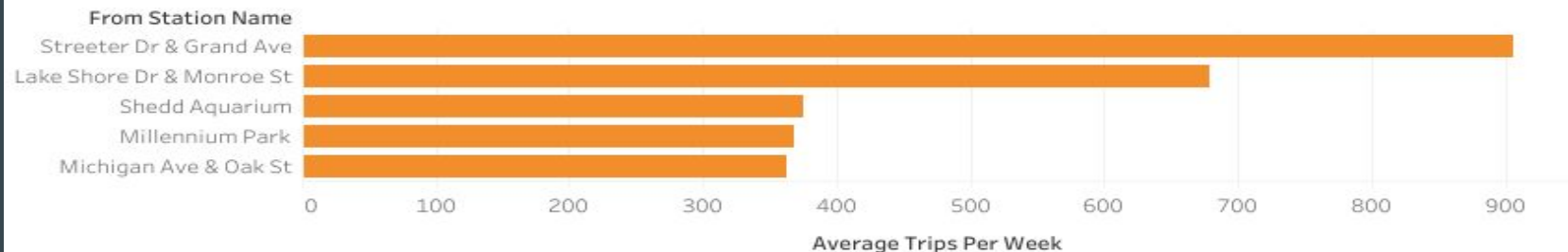
## 2019 Trip Duration Comparison

Above we show the average time length of trips that occurred on each day of the week. Annual members tend to keep a relatively consistent, flat average whereas casual riders tend to ride for longer time lengths on the weekend and overall ride longer than annual members.

## Annual Member Top 5 Stations Weekly

From Station Name

| Station | Average Trips Per Week |
|---|---|
| Canal St & Adams St | ~920 |
| Clinton St & Madison St | ~840 |
| Clinton St & Washington B.. | ~835 |
| Columbus Dr & Randolph St | ~565 |
| Kingsbury St & Kinzie St | ~570 |

Average Trips Per Week

## Casual Rider Top 5 Stations Weekly

From Station Name

| Station | Average Trips Per Week |
|---|---|
| Streeter Dr & Grand Ave | ~910 |
| Lake Shore Dr & Monroe St | ~680 |
| Shedd Aquarium | ~375 |
| Millennium Park | ~370 |
| Michigan Ave & Oak St | ~360 |

Average Trips Per Week

From the two charts, we see the stations with the highest average weekly trips. None of the top five stations are shared by casual riders and annual members.

**Recap of Key Question:** How do annual members and casual riders use Cyclistic bikes differently?

- Annual members ride more frequently per week.

- Annual members prefer weekdays whereas casual riders prefer weekends.

- Casual riders ride for longer durations and show more fluctuation in duration data.

- Casual riders and annual members most popular stations are all in different locations.

## My three recommendations to drive conversion of casual riders to annual members:

1. *Advertise the benefits of annual membership at the most popular casual rider stations, with an emphasis on summertime and weekends.*

2. *Offer an annual membership free trial to casual riders after his or her account shows a high number of trips.*

3. *Set a ride time limit for casual riders after which casual riders are required to watch an advertisement on his or her phone before closing out that trip session.*

# Documentation

# Data Sources Used:

- https://divvy-tripdata.s3.amazonaws.com/Divvy_Trips_2019_Q1.zip
- https://divvy-tripdata.s3.amazonaws.com/Divvy_Trips_2019_Q2.zip
- https://divvy-tripdata.s3.amazonaws.com/Divvy_Trips_2019_Q3.zip
- https://divvy-tripdata.s3.amazonaws.com/Divvy_Trips_2019_Q4.zip

```
 1    create table bicycle_data.2019_full as
 2    select *
 3    from bicycle_data.2019_Q1
 4    union all
 5    select *
 6    from bicycle_data.2019_Q2
 7    union all
 8    select *
 9    from bicycle_data.2019_Q3
10    union all
11    select *
12    from bicycle_data.2019_Q4
```

- I uploaded the prior listed data sources to Google BigQuery because that amount of data overwhelmed my excel, to hone my SQL skills, and to hopefully standout as I believe most completing this case study will follow the R notes given from the Google certificate program.

- I combined the data in BigQuery by creating a new table through the code shown above.

```
1    select birthyear
2    from bicycle_data.2019_full
3    where birthyear is null
4
5    --------
6
7    select gender
8    from bicycle_data.2019_full
9    where gender is null
```

To start the cleaning process, I searched columns for null values in which the gender and birthyear columns each returned close to 500k null values out of 3M rows. I deemed both columns contained too much unreliable data to be used.

```
1   select usertype
2   from bicycle_data.2019_full
3   where usertype <> "Customer" AND usertype <> "Subscriber"
```

Next, I checked for any different values in the Customer and Subscriber columns.

```
1    select trip_id, count(*)
2    from bicycle_data.2019_full
3    group by trip_id
4    having count (*) > 1
```

Here I checked for duplicate trip ids.

Next, I will check for outliers in the data.

```
1   select
2       max(tripduration)
3   from(
4   select
5   tripduration,
6   NTILE(4) OVER (ORDER BY tripduration) AS quartiles
7   from bicycle_data.2019_full
8   where usertype = "Subscriber")
9   where quartiles = 1
10  -- equals 362
11  select
12      max(tripduration)
13  from(
14  select
15  tripduration,
16  NTILE(4) OVER (ORDER BY tripduration) AS quartiles
17  from bicycle_data.2019_full
18  where usertype = "Subscriber")
19  where quartiles = 3
20  -- equals 967
21
22  /* 967-362 = IQR = 605
23     605 * 1.5 = 907.5
24     907.5 + 967 = 1874.5 */
25
26
27
```

```
1   select
2       max(tripduration)
3   from(
4   select
5   tripduration,
6   NTILE(4) OVER (ORDER BY tripduration) AS quartiles
7   from bicycle_data.2019_full
8   where usertype = "Customer")
9   where quartiles = 1
10  -- equals 915
11  select
12      max(tripduration)
13  from(
14  select
15  tripduration,
16  NTILE(4) OVER (ORDER BY tripduration) AS quartiles
17  from bicycle_data.2019_full
18  where usertype = "Customer")
19  where quartiles = 3
20  -- equals 2718
21
22  /* 2718 - 915 = 1803 = IQR
23     1.5 * 1803 = 2704.5
24     2704.5 + 2718 = 5422.5 */
25
26
27
28
29
```

Using the 1.5 * IQR rule I established upper bounds at 1874.5 seconds for "Subscribers" aka Annual Members and 5422.5 seconds for "Customers" aka Casual Riders.

Above these bounds are suspected outliers which I removed from the data. The lower bounds for both "Subscribers" and "Customers" would have been well below zero and no data was found below a value of zero.

```
1
2   create table bicycle_data.2019_full_clean as
3   select
4   * except (birthyear,gender)
5   from bicycle_data.2019_full
6   where usertype = "Customer" AND tripduration <= 5422.5
7   union all
8   select
9   * except (birthyear,gender)
10  from bicycle_data.2019_full
11  where usertype = "Subscriber" AND tripduration <= 1874.5
12
```

To complete the cleaning process, I combined my cleaning techniques into a new table to be used for analysis.
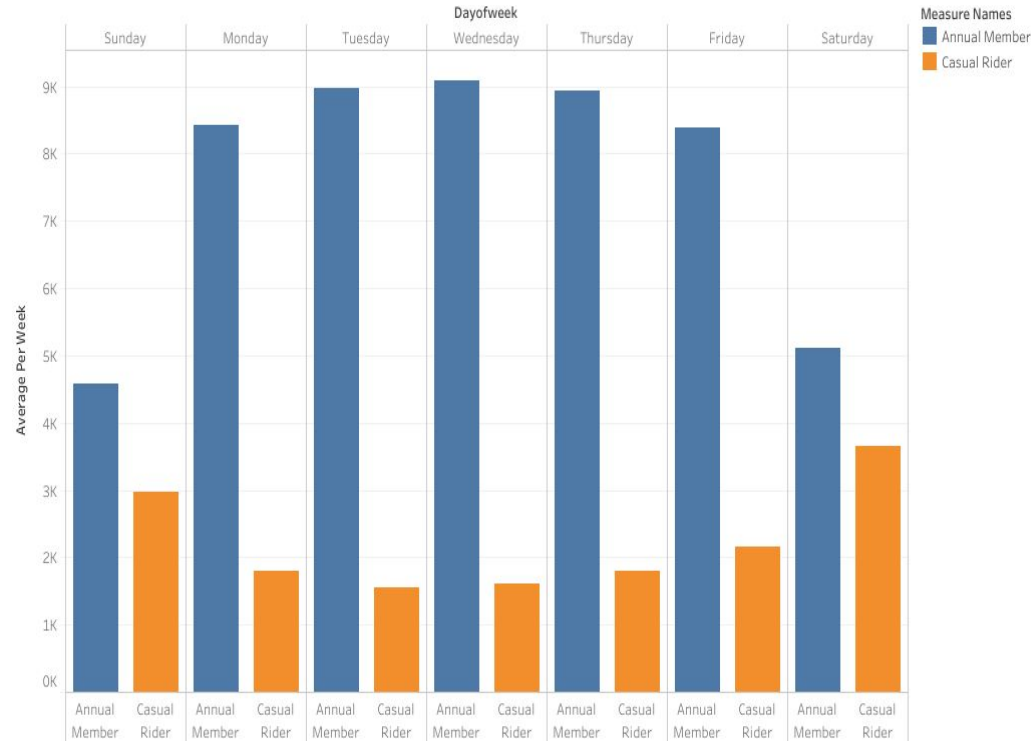
# Code Used for Charts

```
1   select
2   CASE dayofweek
3       WHEN 1 THEN "Sunday"
4       WHEN 2 THEN "Monday"
5       WHEN 3 THEN "Tuesday"
6       WHEN 4 THEN "Wednesday"
7       WHEN 5 THEN "Thursday"
8       WHEN 6 THEN "Friday"
9       WHEN 7 THEN "Saturday"
10      END AS dayofweek,
11  avg(week_frequency) AS Annual_Member
12  from(
13  select
14  weeknumber,
15  dayofweek,
16  count(dayofweek) AS week_frequency
17  from(
18  select
19  start_time,
20  extract(week from start_time) AS weeknumber,
21  extract(dayofweek from start_time) AS dayofweek
22  from bicycle_data.2019_full_clean
23  where usertype = "Subscriber")
24  group by weeknumber, dayofweek)
25  group by dayofweek
26  order by Annual_Member desc
```

2019 Day of Week Trips Comparison



The code used to analyze "Customer" data was the same as "Subscriber". To avoid additional clutter I only showed one code chunk in all my documentations.
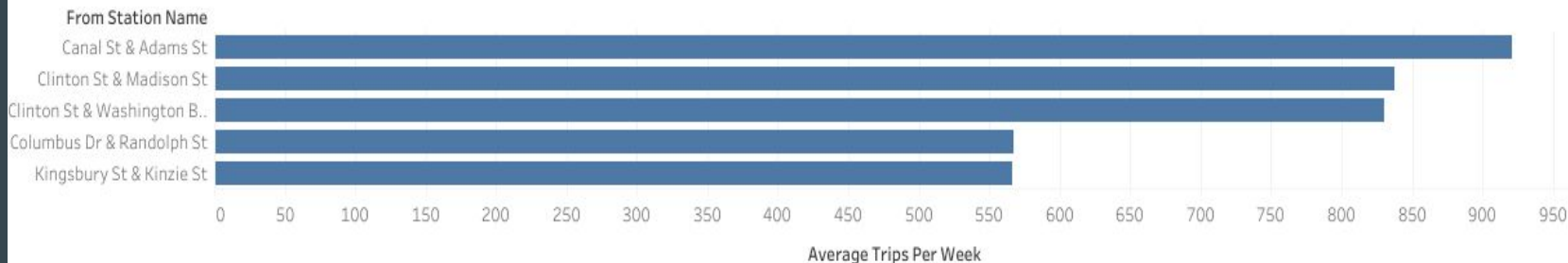
"Subscriber" was manually changed to Annual Member and "Customer" to Casual Rider in all my charts.
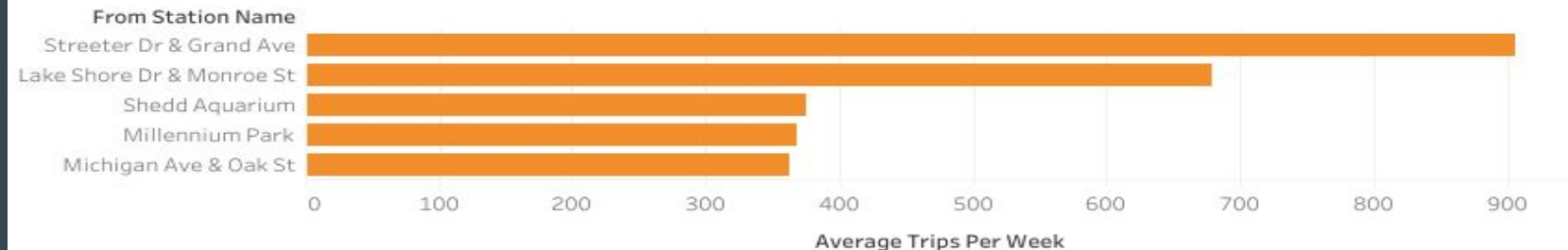
```
 1    select
 2    from_station_name,
 3    avg(frequency) AS avg_freq_weekly
 4    from(
 5    select
 6    from_station_name,
 7    weeknumber,
 8    count(from_station_name) AS frequency
 9    from(
10    select from_station_name,
11    extract(week from start_time) AS weeknumber
12    from bicycle_data.2019_full_clean
13    where usertype = "Customer"
14    group by from_station_name,start_time)
15    group by weeknumber,from_station_name)
16    group by from_station_name
17    order by avg_freq_weekly desc
18    limit 5
19    |
20
```

## Annual Member Top 5 Stations Weekly

From Station Name



Average Trips Per Week

## Casual Rider Top 5 Stations Weekly

From Station Name



Average Trips Per Week

```
1    select
2    weeknumber,
3    count(trip_id) as frequency
4    from(
5    select
6    extract(week from start_time) as weeknumber,
7    trip_id
8    from bicycle_data.2019_full_clean
9    where usertype = "Subscriber"
10   group by start_time,trip_id)
11   group by weeknumber
12   order by weeknumber
```

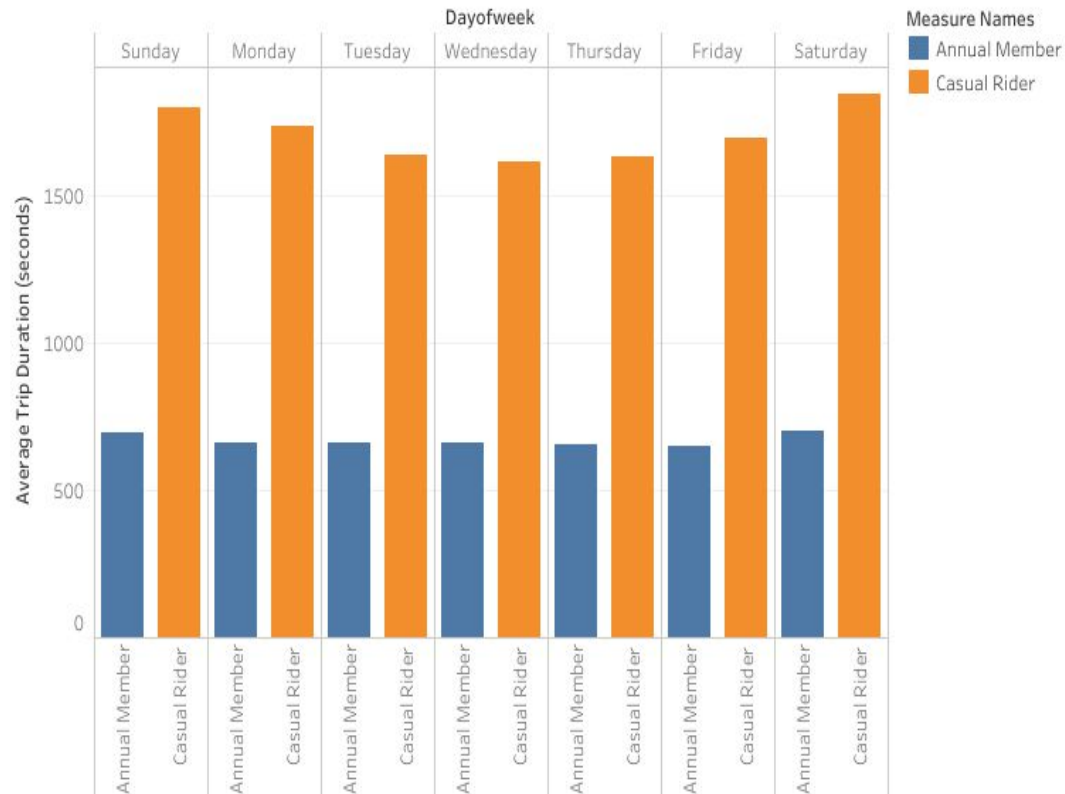## 2019 Total Trips Every Week (Chronological)

```sql
select
avg(tripduration) AS avg_trip,
CASE dayofweek
    WHEN 1 THEN "Sunday"
    WHEN 2 THEN "Monday"
    WHEN 3 THEN "Tuesday"
    WHEN 4 THEN "Wednesday"
    WHEN 5 THEN "Thursday"
    WHEN 6 THEN "Friday"
    WHEN 7 THEN "Saturday"
    END AS dayofweek,
from(
select
extract(dayofweek from start_time) AS dayofweek,
tripduration
from bicycle_data.2019_full_clean
where usertype = "Customer"
group by tripduration, start_time)
group by dayofweek
order by avg_trip desc
```

## 2019 Trip Duration Comparison

Thanks for viewing!