## Google Data Analytics Capstone Project. Case Study 1: How Does a Bike-Share Navigate Speedy Success?

Paul W January 29, 2024

# **Business Task**

Company Background

memberships. Customers who purchase single-ride or full day passes are referred to as casual riders. Customers who purchase annual memberships are Cyclistic members. 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, the marketing director believes that maximizing the number of annual members will be key to future growth. The

approach that helped make these things possible was the flexibility of its pricing plans: single-ride passes, full-day passes, and annual

In 2016, Cyclistic launched a successful bike-share offering. Since then, the program has grown to a fleet of 5,824 bicycles that are geotracked 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

however, the marketing team needs to better understand how annual members and casual riders differ, why casual riders would buy a bike trip data to identify trends.

marketing director has a clear goal: Design marketing strategies aimed at converting casual riders into annual members. In order to do that, membership, and how digital media could affect their marketing tactics. The director and her team are interested in analyzing the Cyclistic historical The purpose of this case study and report is to do an exploratory data analysis of the provided data and identify any data points, descriptive statistics or trends that may help the marketing team achieve their goals. The analysis is broken up into these sections that are steps of the data analysis process from the course:

 Prepare Process Analyze

Share

Act

Ask

Ask - Three Questions

Cyclistic Data 2022 ↑ Sort → ■ View → Date modified A Home Microsoft Excel C... 202201-divvy-tripdata 10/4/2023 2:42 PM 18,567 KB 🔼 Gallery 202202-divvy-tripdata 10/4/2023 2:42 PM Microsoft Excel C... 20,638 KB Paul - Personal 202203-divvy-tripdata 10/4/2023 2:42 PM Microsoft Excel C... 50,533 KB 202204-divvy-tripdata 10/4/2023 2:44 PM Microsoft Excel C... 65,341 KB Desktop

42.01280 -87.66591 42.01256 -87.67437 casual

42.01276 -87.66597 42.01256 -87.67437 casual

202207-divvy-tripdata 10/4/2023 2:46 PM Microsoft Excel C... Pictures 202208-divvy-tripdata 10/4/2023 2:46 PM Microsoft Excel C... Music 202209-divvy-publictripdata 10/4/2023 2:46 PM Microsoft Excel C... Videos 202210-divvy-tripdata 10/4/2023 2:46 PM Microsoft Excel C... Google Data Analyst Cert 202211-divvy-tripdata 10/4/2023 2:46 PM Microsoft Excel C... 202212-divvy-tripdata 10/4/2023 2:47 PM Microsoft Excel C... Example of table structure found in csv files: † rideable\_type † started\_at † ended\_at start\_station\_name start\_station\_id end\_station\_name end\_station\_id end\_station\_id end\_station\_id start\_lat start\_lng end\_lat end\_stat member\_casual 1 C2F7DD78E82EC875 electric\_bike 2022-01-13 11:59:47 | 2022-01-13 12:02:44 | Glenwood Ave & Touhy Ave | 525 Clark St & Touhy Ave RP-007 2 A6CF8980A652D272 electric bike 2022-01-10 08:41:56 2022-01-10 08:46:17 Glenwood Ave & Touhy Ave 525 Clark St & Touhy Ave RP-007 2022-01-25 04:53:40 2022-01-25 04:58:01 Sheffield Ave & Fullerton Ave TA1306000016 Greenview Ave & Fullerton Ave TA1307000001 4 CBB80ED419105406 classic\_bike 2022-01-04 00:18:04 2022-01-04 00:33:00 Clark St & Bryn Mawr Ave KA1504000151 Paulina St & Montrose Ave TA1309000021 5 DDC963BFDDA51EEA classic\_bike 2022-01-20 01:31:10 2022-01-20 01:37:12 Michigan Ave & Jackson Blvd TA1309000002 State St & Randolph St TA1305000029 Install and Load Packages # Load libraries library("tidyverse") library("leaflet") library("dplyr") library("magrittr") library("readr")

library("ggplot2") library("stringr") library("geosphere") library("sf") library("maps") library("lubridate")

map\_df(~read\_csv(.))

#Data\_Combined\_2022 <- Data\_Combined\_2022 %>% slice(1:100000)

Combine csv files into one dataframe. Check csv files for issues that would prevent combining them Instead of inspecting each csv file individually, this method combines all of the csv files for 2022 into one dataframe. It also checks for anything that can produce errors such as differences found in column names, column types, missing or extra columns etc. The console will give a report for each file showing if there were errors or if the file successfully passed inspection. No errors were found and files were combined successfully: Data\_Combined\_2022 <- list.files(path = "C:/Users/paulw/Desktop/Google Data Analyst Cert/Cyclistic Data 2022",

Process - Data Processing and Data Validation Initial Review of Tables and Any Documents For many projects, a data dictionary/layout is normally reviewed by the analyst. This project did not include one which would be very helpful when inspecting the fields and checking for errors. Also, it is good to know exactly what is being reported in a field and to know what lengths, data types etc. are to be expected. We will have to look at the fields and try to get an idea of what data types should be in those columns and if they contain facets, see if the correct facet values are there. Here we can start with str() to get a summary of what is going on in the dataframe and get names of columns, the data types found in the columns and a few sample values of what is found in each column: str(Data\_Combined\_2022) ##  $spc_tbl_[5,667,717 \times 13]$  (S3:  $spec_tbl_df/tbl_df/tbl/data.frame$ ) ## \$ ride\_id : chr [1:5667717] "C2F7DD78E82EC875" "A6CF8980A652D272" "BD0F91DFF741C66D" "CBB80ED419105 406" ... ## \$ rideable\_type : chr [1:5667717] "electric\_bike" "electric\_bike" "classic\_bike" "classic\_bike" ...

## \$ start\_station\_name: chr [1:5667717] "Glenwood Ave & Touhy Ave" "Glenwood Ave & Touhy Ave" "Sheffield Ave &

## \$ end\_station\_name : chr [1:5667717] "Clark St & Touhy Ave" "Clark St & Touhy Ave" "Greenview Ave & Fullerto

## \$ started\_at : POSIXct[1:5667717], format: "2022-01-13 11:59:47" "2022-01-10 08:41:56" ... ## \$ ended\_at : POSIXct[1:5667717], format: "2022-01-13 12:02:44" "2022-01-10 08:46:17" ...

## \$ start\_station\_id : chr [1:5667717] "525" "525" "TA1306000016" "KA1504000151" ...

#### .. rideable\_type = col\_character(), .. started\_at = col\_datetime(format = ""), .. ended\_at = col\_datetime(format = ""), .. start\_station\_name = col\_character(), .. start\_station\_id = col\_character(), .. end\_station\_name = col\_character(),

Total Rows in Data\_Combined\_2022 = 5,667,717 Total Rows where a column has NA = 1,298,357. Roughly 23% of the rows for 2022 have an AnyRowsNA <- Data\_Combined\_2022[rowSums(is.na(Data\_Combined\_2022))>1,] #Return data frame with rows containing N

Cleaned\_Data\_2022\$hour\_am\_pm <- ifelse(</pre> Cleaned\_Data\_2022\$started\_at\_hour == 0, "12 AM", Cleaned\_Data\_2022\$started\_at\_hour == 12, "12 PM", ifelse( Cleaned\_Data\_2022\$started\_at\_hour < 12,</pre> paste0(Cleaned\_Data\_2022\$started\_at\_hour, " AM"), paste0(Cleaned\_Data\_2022\$started\_at\_hour - 12, " PM")

Cleaned\_Data\_2022 <- Cleaned\_Data\_2022 %>% mutate(ride\_distance\_feet = distGeo(matrix(c(start\_lng, start\_lat),

Cleaned\_Data\_2022 <- Cleaned\_Data\_2022 %>% mutate(ride\_time\_minutes = difftime(ended\_at,started\_at,units="mins")) Cleaned\_Data\_2022\$ride\_time\_minutes <- as.numeric(gsub(" mins", "", Cleaned\_Data\_2022\$ride\_time\_minutes)) #reform

Add new columns: ride time minutes, started at hour, start day of week, start month, ride distance feet

# Add a column for the week day the ride started on and a column for the month it was in

Cleaned\_Data\_2022\$start\_month <- month(Cleaned\_Data\_2022\$started\_at, label = TRUE, abbr = FALSE)

# Add a ride\_distance\_feet column that is the distance between starting and ending stations.

Cleaned\_Data\_2022\$start\_day\_of\_week <- weekdays(Cleaned\_Data\_2022\$started\_at)</pre>

Cleaned\_Data\_2022\$started\_at\_hour <- hour(Cleaned\_Data\_2022\$started\_at)</pre>

# Add a column ride\_time\_minutes that is the difference in start and end times, convert to minutes

The ride\_id column is used to identify individual tirps. Checking if all values are unique by comparing the row count to the distinct count of ride\_id's. Distinct\_ride\_ids <- n\_distinct(Data\_Combined\_2022\$ride\_id)</pre> Distinct\_ride\_ids ## [1] 5667717 Total\_Rows\_2022 <- nrow(Data\_Combined\_2022) Total\_Rows\_2022 ## [1] 5667717 An example of a facet column would be a column like member\_casual that has repeated values that stratify the data into categories. Even though we do not have a data dictionary it appears that this column should have only two values - "member" or "casual". The rideable\_types column should also have values "electric\_bike", "classic\_bike" and "docked bike". The station names and ids should be checked as well but that will be mentioned in the next subsection:

### start station id 📤 13109 13109

13109

13109

13109

Find any extreme latitude and longitude data points that do not belong in the data set

check\_member\_casual <- unique(Cleaned\_Data\_2022\$member\_casual)</pre>

check\_rideable\_type <- unique(Cleaned\_Data\_2022\$rideable\_type)</pre>

## [1] "electric\_bike" "classic\_bike" "docked\_bike"

-84 -

-25 -

end\_Ing

-50 -

-75 -

from the data

Analyze

Trips per Weekday

100K

400K

300K

Total Trips

100K

150K

100K ·

50K

200K

100K

## Total Minutes Spent 42,182,990 32,512,947 ## Max. Minutes Spent 34,354.067 1,493.233 ## Avg. Minutes Spent 23.99316 12.45175 ## Median Minutes Spent 13.850000 8.983333

Casual riders spend more time on averge and in total riding their bikes than Member riders.

member

casual

Total Trips

Time Spent Summary.

##

## 1

## 2

## 3

## 6

## 7

## 8

## 9

## 10

##

## 1

## 2

## 3

## 6

## 7

## 8

## 9

## 10

Share

Tableau Public Report

ood Dale

dison

Bensenville

occur more between 12pm and 5pm.

Values found in member\_casual column are correct:

check\_member\_casual

check\_rideable\_type

## [1] "casual" "member"

Values found in rideable-types are correct:

may affect the map negatively when trying to load or display it. They may also throw off any distance and travel time calculations. Starting and Ending geolocations should be clustered together when using a scatter plot that has longituted on the y axis and latitude on the x axis. Extraneous points will appear away from the cluster of Chicago points. In the images below, the clustered points belonging to Chicago are in the -75 to -84 range for longitude points and latitude points are in the 40 to 43 range. The data that does not belong to Chicago is away from these clusters and outside of these ranges: Starting Station Latitude and Longitude points:

#This next line can be used to generate the plot but will have to be made separately and then inserted into the d

ocument as an image. R Studio freezes when trying to execute this while knitting the document

 $\#(checking\_start\_plot <- ggplot(Cleaned\_Data\_2022, aes(x = start\_lat, y = start\_lng)) + geom\_point())$ 

knitr::include\_graphics("C:/Users/paulw/Desktop/Google Data Analyst Cert/R Documents/Not\_In\_Chicago\_1.png")

The latitude and longitude points in the data should all be located in the Chicago area. Any extreme points that do not belong in the Chicago area

Clark St & Winnemac Ave

Clark St & Winnemac Ave

Clark St & Winnemac Ave

43 44 45 start lat Ending Station Latitude and Longitude points: #This next line can be used to generate the plot but will have to be made separately and then inserted as an imag e into the document. R Studio freezes when trying to execute this while knitting the document  $\#(checking\_end\_plot <- ggplot(Cleaned\_Data\_2022, aes(x = end\_lat, y = end\_lng)) + geom\_point())$ knitr::include\_graphics("C:/Users/paulw/Desktop/Google Data Analyst Cert/R Documents/Not\_In\_Chicago\_2.png")

10 30 end\_lat Identiyfing and removing the bad Latitude and Longitude points. Renaming the data frame: #Going online and looking up each of these points in this dataframe will show they are not in Chicago  $Not\_Chicago\_Locations <- subset(Cleaned\_Data\_2022, start\_lat < 40 \mid start\_lng > -75 \mid end\_lat < 40 \mid end\_lng > -70 \mid end\_lat < 40 \mid end\_l$ #Removing the rows that have ride data outside the Chicago area

Tuesday Wednesday Thursday Friday Monday Saturday Sunday Trips by the Hour

Member trips peak both in the morning and later in the evening at 5pm. Casual trips peak at 5 pm but do not peak in the morning. Casual trips

Trips By The Hour

casual member

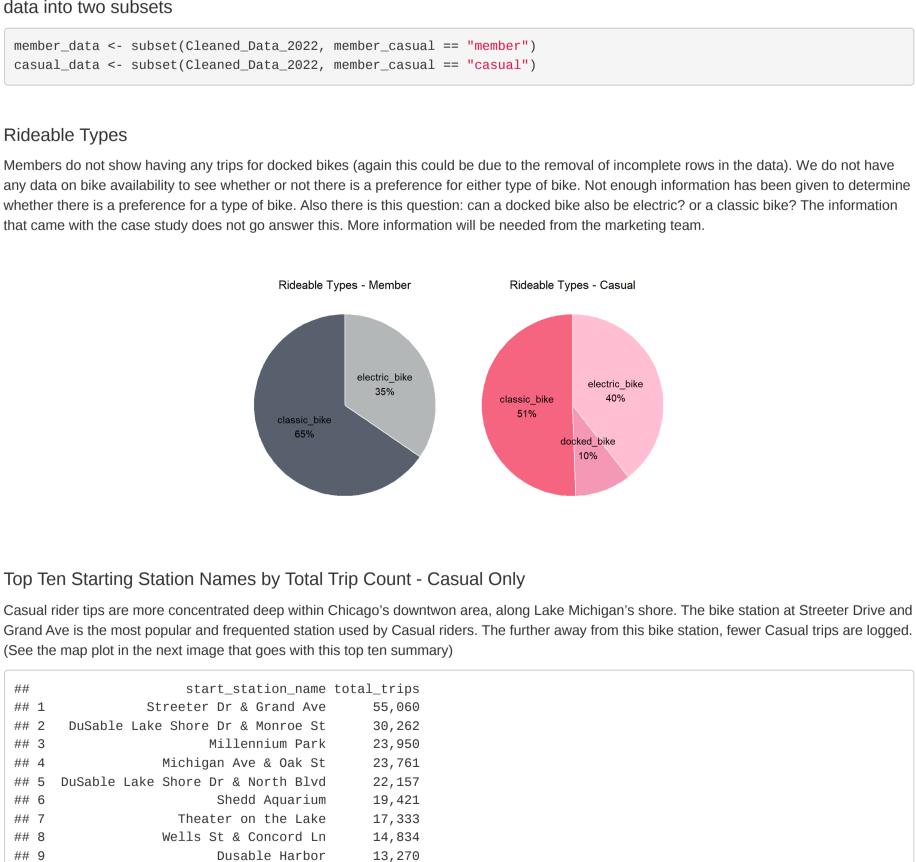
Hour of the Day

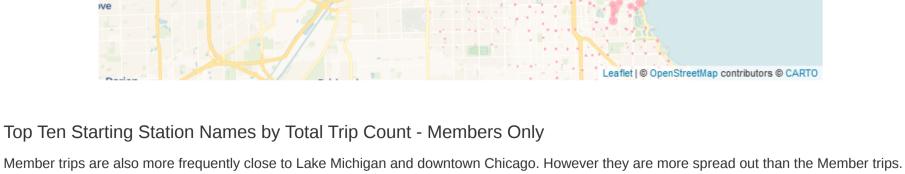
Month

Weekday Trips

member\_casual casual

Casual riders are more likely to rent a bike on the weekend on Saturday and Sunday than during the work week.





20,581

19,674

18,828

18,251 18,007

17,581

17,504

17,496

(See the map plot in the next image that goes with this top ten summary).

Kingsbury St & Kinzie St

Loomis St & Lexington St

Clinton St & Madison St

University Ave & 57th St

Ellis Ave & 60th St

Wells St & Elm St

also a link to the code documentation listed below the Tableau Public link:

## 4 Clinton St & Washington Blvd

+

Franklin Park

Clark St & Elm St

Wells St & Concord Ln

start\_station\_name total\_trips

10,000 hlake 12,000 Melrose Park 14,000 Mavwood Cicero Brookfield

The Member trips also appear to be more evenly spread out among the popular bike stations that are near the lake shore and downtown area.

More documentation for this project can also be found on this Github repository page https://github.com/pwindsor37115/Paul-W-s-Google-Data-Analytics-Capstone-Project — Act Top three recommendations given to the marketing department:

Besides this R markdown document, below is a link to a report and dashboard made in Tableau Public that can be sent to the marketing

https://public.tableau.com/app/profile/paul.windsor/viz/GoogleDataAnalyticsProject-CyclisticCaseStudy/Story1

department as well. The Tableau Public report is set up as a presentation that goes over the analyis findings and the recommendations. There is

not be examined b/c of this. • At the stations, offer a discount to casual riders if they fill out a survey to get more information on how and why they purchase day passes. The data suggests that member riders may be using the bikes more for traveling to work or nearby their home, while the casual riders appear to be using the bikes more for recreational purposes. Maybe there are customers who live downtown but are not members that may be better persuaded into purchasing a membership if we can gather more information from them as well. • Offer seasonal passes, instead of just day passes, for bikes while offering other services. Begin advertising the seasonal passes during the spring months in March and April.

data issues such as these that need to be reviewed and also will need a better description of the business processes and how data is recorded. Also another tool other than Tableau Public needs to be used due to the limit on data that can be uploaded - yearly trends could

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 Prepare - Preparation of Data Description of all data sources used and technical information The data for this project was provided by Coursera in the form of many csv files. Each csv file has bike ride data with a unique id to represent each bike trip made and each file has data for a particular month and year. These files can be made to create one large fact table of bike ride data. No files were offered to create dimension tables. Also sql was not needed to pull the files - they were made already available in several folders. Instead of using sql to prepare the files, R was used to prepare the data and perform an analysis. R Markdown was used to create the report document and further analysis and presentation was done with Tableau Public. Only one year's worth of data could be used for this report because of the limits to how much data can be uploaded to Tableau Public. The documentation from the project also notes that we will not be able to identify individual users and purchases because there are no credit card numbers or identification numbers given to identify them. There were no data dictionaries or layouts provided, further complicating the analysis performed. Here are some screenshots of the csv files and a sneak peak view of how the tables with their columns appear: Files used: 202205-divvy-tripdata 10/4/2023 2:44 PM Microsoft Excel C... 114,780 KB Downloads 202206-divvy-tripdata 10/4/2023 2:44 PM Microsoft Excel C... 140,228 KB 149,306 KB 142,148 KB 138,135 KB 109,293 KB 66,348 KB 35.612 KB

41.98359 -87.66915 41.96151 -87.67139 casual 41.87785 -87.62408 41.88462 -87.62783 member library("scales") library("viridis") library("gridExtra") pattern="\*.csv", full.names = T) %>%

## ## \$ end\_station\_id : chr [1:5667717] "RP-007" "RP-007" "TA1307000001" "TA1309000021" ... ## \$ start\_lat : num [1:5667717] 42 42 41.9 42 41.9 ... ## \$ start\_lng : num [1:5667717] -87.7 -87.7 -87.7 -87.7 -87.6 ... ## \$ end\_lat : num [1:5667717] 42 42 41.9 42 41.9 ... ## \$ end\_lng : num [1:5667717] -87.7 -87.7 -87.7 -87.7 -87.6 ... ## \$ member\_casual : chr [1:5667717] "casual" "casual" "member" "casual" ... ## - attr(\*, "spec")= ## .. cols( .. ride\_id = col\_character(),

Fullerton Ave" "Clark St & Bryn Mawr Ave" ...

.. end\_station\_id = col\_character(),

.. start\_lat = col\_double(), .. start\_lng = col\_double(), .. end\_lat = col\_double(), .. end\_lng = col\_double(),

Add new columns

at as number

#Add a column that has the hour the ride started

ncol=2), matrix(c(end\_lng, end\_lat), ncol=2))\*3.2808399)

Data Validation and More Cleaning

n Ave" "Paulina St & Montrose Ave" ...

.. member\_casual = col\_character() ## - attr(\*, "problems")=<externalptr> Initial Cleaning of Columns with Strings Clean leading and trailing spaces that may be before or after text in columns that have strings Columns\_with\_strings = c("ride\_id", "rideable\_type", "start\_station\_name", "start\_station\_id", "end\_station\_name") e", "end\_station\_id", "member\_casual") Data\_Combined\_2022[Columns\_with\_strings] <- lapply(Data\_Combined\_2022[Columns\_with\_strings], trimws)</pre> Identify the columns in the dataframe that have a NA (null) value. Are these expected to occasionally have an NA value? NA in at least one column. Normally would get with someone to figure out why almost a quarter of the data is unusable.... AnyColsNA <- colnames(Data\_Combined\_2022[colSums(is.na(Data\_Combined\_2022))>1]) A value. Total\_Rows\_2022 <- nrow(Data\_Combined\_2022) #Total Rows in the data frame so far Total\_Rows\_wNA\_2022 <- nrow(AnyRowsNA) #Total Rows with one or more fields with NA Cleaned\_Data\_2022 <- Data\_Combined\_2022[rowSums(is.na(Data\_Combined\_2022))==0,] #Dropping rows from the datafra me that have NA values and renaming the dataframe

This is where I have my own steps inserted into the process step used by this course. This is where the Process step overlaps with the Analyze step. Normally, the data is validated before further processing. The data is checked against a data dictionary or any notes we have and some testing occurs. If errors are found, then the data may have to be sent back or you may have to go and pull a new query. While this is going on , the rest of the project can be continued while waiting for feedback on the data quality and any questions about the data issues to be answered. Checking columns that contain unique identifiers - does the column have unique values? Both the row count for the data frame and the distinct number of ride id's match Columns with facet data - do these have any unexpected values?

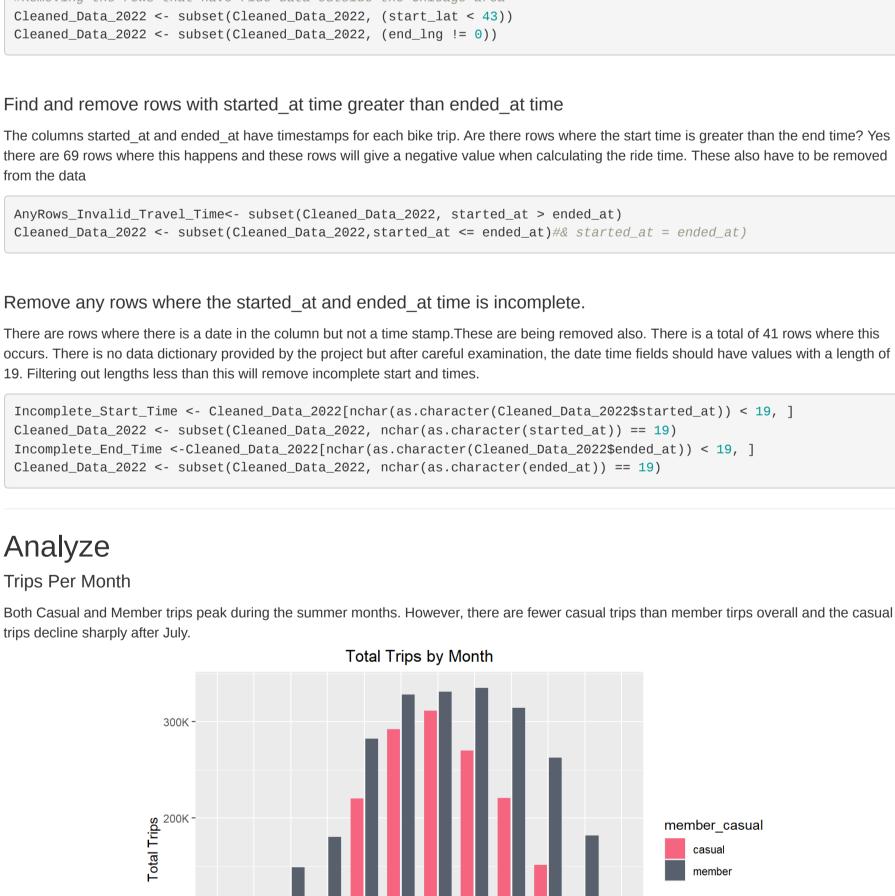
#### An example of multiple station names listed with a station id: end\_station\_name Broadway & Thorndale Ave Broadway & Thorndale Ave 13109 Broadway & Thorndale Ave 13109 Broadway & Thorndale Ave 13109 Clark St & Winnemac Ave

Several station id's were found where they had several station names attached to them. Should there not be one station id per station name?

Columns with facet data continued - why are there several station names found per station id?

Again there is no data dictionary or layout to refer to and there is not any more information provided in the project handout.

-76 -



## Min. Minutes Spent 0 **Distance Traveled Summary** Casual riders spend more time on average and in total on their bikes than Members, but Members' travel distances are slightly more than what Casual riders travel. casual member ## Total Distance Traveled 12,397,480,273 17,532,883,606 98,941.86 99,375.86 ## Max. Ride Distance Feet 7,051.531 6,714.711 ## Avg. Ride Distance Feet ## Median Ride Distance Feet 5,415.157 4,918.963 ## Min. Ride Distance Feet

Before creating visuals such as top ten stations by total rides for members and casual riders, will need to split the

Elmwood Park **ELMHURST** Melrose Park 'illa Park Maywood 4GO Wester Springa Grange

12,779

Starting Stations and Total Rides, Casual Only

5,000

-10,000 15,000

20,000

-25,000

-30,000 35,000

40,000

Clark St & Armitage Ave

Map plot of Starting Stations and Total Trips - Casual Riders only.

Franklin Park

Streeter Dr & Grand Ave 16,208 Map plot of Starting Stations and Total Trips - Members only. Starting Stations and Total Rides, Member Only 2,000 4,000 6,000 8,000

Cyclistic Data Analysis Code Documentation https://github.com/pwindsor37115/Paul-W-s-Google-Data-Analytics-Capstone-Project/blob/main/Cyclistic%20Data%20Analysis.R More data and information is needed from the rest of the team to properly use this data. The large amount of data that had to be removed because it was incomplete is a big concern. This affects the report negatively and the ability to come to any conclusions. Also there are the issues of multiple station names for each station id, missing start station data on rows but end station data is populated, etc. Examples of