# Uber Data Analysis project
# Report

Academic year: Fall 2021-2022

## Under the guidence of:
## Prof. S.Gopikrishna

## Submitted by:

19BCD7006 - KODURI GOKUL
19BCD7014 - KONGARA SAI HAREESH KUMAR
19BCD7019 - MUPPURI SIVAKUMAR
19BCD7025 - BUSSA SAIBIPIN

# Content:

# Summary:

The working of an Uber dataset includes primary
Data on Uber pickups with details including the date, time of the ride as well as longitude-latitude information.  Uber data analysis project, data storytelling is an important component of Machine Learning through which companies are able to understand the background of various operations. With the help of visualization, companies can avail the benefit of understanding the complex data and gain insights that would help them to craft decisions.The interesting insights that can be derived from a detailed analysis of the dataset. The ggplot2 on the Uber Pickups dataset and at the end, master the art of data visualization in R.the dataset and to know the effect of each field on price with every other field of the dataset. Then The objective is to first explore hidden or previously unknown information by applying exploratory data analytic on y different R models to complete the analysis. The number of pickups is more during weekends. To alleviate the dynamic price surge, we need to manage the 'supply and demand' of cabs through these events of high demand situations. Based on these results, we can expect that the demand will be high as the temperature drops or after business hours or on weekends, so that:  As a Customer: We can plan our trips in advance to avoid paying extra money because of this dynamic price surge.  As an Uber Driver: We can maximize profits by choosing to go on trips when these situations occur.

# Introduction:

Uber launched in NYC in May of 2011, the first city outside of its San Francisco headquarters. NYC is probably the largest and most lucrative rideshare market in the world, with a total demand (for taxis and for-hire vehicles) in 2017 of more than 240 million trips per year.The number of Uber trips per day in NYC is still growing significantly. In 2017 so far, this number has often surpassed 200,000, but the plot below shows that by mid-2015 it was hovering around 120,000.The data also allows us to visualize other interesting trends over time. In the bar charts below, we can see that the demand for Uber is higher . Saturday has the highest demand. Interestingly, Sunday shows a level of demand similar to Wednesday, which is higher than Monday or Tuesday. When looking at the total demand per month along the period of time analyzed, seasonal effects are masked by the consistent month-to-month growth.

# DATA ANALYST:

```r
library(readr)
uber <- read_csv("C:/Users/siva kumar/Downloads/uber-rides-dataset- updated.csv")

## Rows: 678 Columns: 37

##      Column  specification
##   --
## Delimiter: ","
## chr   (20): trip_status, trip_uid, driver_uid, rider_uid, customer, trip_star...
## dbl   (12): pickup_lat, pickup_long, dropoff_lat, dropoff_long,
rub_usd_excha...
## dttm   (2): trip_completed_at, temperature_time ##
time       (3): trip_time, total_time, wait_time

##
## 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.

df <- uber[,c(7:19,22,24:37)]

# clusterring
library(gtools)
library(ClusterR)
library(cluster)
library(ggplot2)


ClusterFunction <- function(clusterDF){

  kmeans.re <- kmeans(clusterDF, centers = 3)
  plot(clusterDF[c(1,2)], col = kmeans.re$cluster)
  points(kmeans.re$centers, col = 1:3, pch = 8, cex = 3)
  y_kmeans <- kmeans.re$cluster clusplot(clusterDF[,c(1,2)],
          y_kmeans, lines
          = 0, shade =
          TRUE, color =
          TRUE, labels =
          2,
```

```
            plotchar = FALSE,
            span = TRUE,
            main = paste("Cluster iris"), xlab =
            'rub_usd_exchange_rate', ylab =
            'total_price')

}

# clustering Prize
ClusterFunction(df[,c(17,18)])
```

## Cluster iris



These two components explain 100 % of the point variabilit

```
#  Clustering  Km
ClusterFunction(df[,c(19,20)])
```

## Cluster iris



These two components explain 100 % of the point variabilit

```
# Hierarchical clustering
Hcluster <- df[,c(22,24:26)]

distance_mat <- dist(Hcluster, method = 'euclidean')
```

```
## Warning in dist(Hcluster, method = "euclidean"): NAs introduced by coercion

Hierar_cl <- hclust(distance_mat, method = "average")
Hierar_cl

##
## Call:
## hclust(d = distance_mat, method = "average") ##
## Cluster method   : average ##
Distance        : euclidean ## Number
of objects: 678

cut_avg <- cutree(Hierar_cl, k = 3)
plot(Hierar_cl)
rect.hclust(Hierar_cl , k = 3, border = 2:6) abline(h = 3,
col = 'red')
```

## Cluster Dendrogram



distance_mat
hclust (*, "average")

```
suppressPackageStartupMessages(library(dendextend))
## Warning: package 'dendextend' was built under R version 4.1.2 avg_dend_obj <-
as.dendrogram(Hierar_cl)
avg_col_dend <- color_branches(avg_dend_obj, h = 3) plot(avg_col_dend)
```

# Data Scientist:

## Dataset:



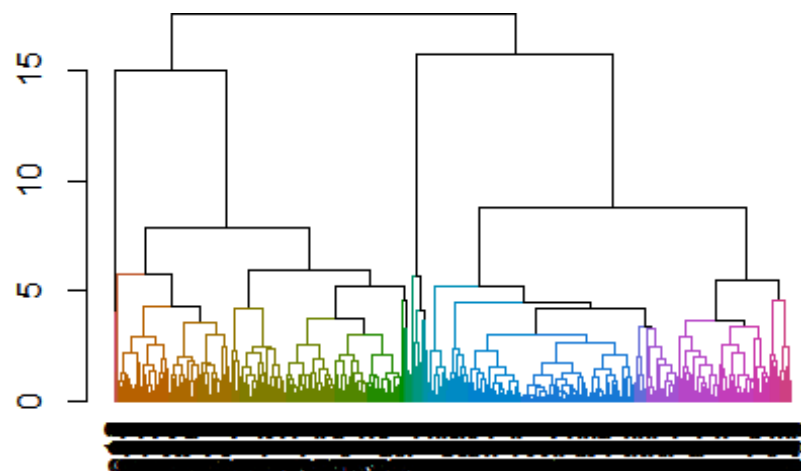| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | trip_comple | trip_status | trip_uid | driver_uid | rider_uid | customer | trip_start_ti | trip_end_tir | trip_time | total_time | wait_time | vehicle_mal | driver_name | vehicle_mal | driver_gend | pickup_lat | pickup_long | dropoff_lat | dropoff_lor | trip_map_ir trip_r |
| 2 | ######### | Completed | ee89076fd9 | 05cfeb269e | 3ffa4a71a5 | stantyan | ######### | ######### | 0:21:33 | 0:29:00 | 0:07:27 | Ford Focus | Maksim | Ford | Male | 60.031438 | 30.329826 | 59.963131 | 30.307655 | [ANONYMIZ [ANO |
| 3 | ######### | Completed | 518be51d4( | 4a4e24874. | 3ffa4a71a5 | stantyan | ######### | ######### | 0:19:27 | 0:26:00 | 0:06:33 | Hyundai Sol | Sergey | Hyundai | Male | 59.963014 | 30.307313 | 60.031351 | 30.329495 | [ANONYMIZ [ANO |
| 4 | ######### | Completed | 6e460cc8a1 | cb249a2bd8 | 3ffa4a71a5 | stantyan | 13-05-2015 | 13-05-2015 | 1:06:53 | 1:23:00 | 0:16:07 | Renault Flue | Oleg | Renault | Male | 60.031529 | 30.329416 | 59.924281 | 30.387561 | [ANONYMIZ [ANO |
| 5 | ######### | Completed | 49613a86a( | d3f73f8151 | 3ffa4a71a5 | stantyan | 16-05-2015 | 16-05-2015 | 0:13:37 | 0:20:00 | 0:06:23 | Mercedes-B | Maksim | Mercedes-B | Male | 59.959883 | 30.311159 | 59.93468 | 30.308489 | [ANONYMIZ [ANO |
| 6 | ######### | Completed | 9896148fde | 1287d21e6 | 3ffa4a71a5 | stantyan | 16-05-2015 | 16-05-2015 | 0:38:54 | 0:49:00 | 0:10:06 | Hyundai Sol | Eduard | Hyundai | Male | 59.934813 | 30.308553 | 60.03147 | 30.329402 | [ANONYMIZ [ANO |
| 7 | ######### | Completed | 5c0312a92f | fc6b151637 | 3ffa4a71a5 | stantyan | 18-05-2015 | 18-05-2015 | 0:16:38 | 0:34:00 | 0:17:22 | Hyundai Sol | Andrey | Hyundai | Male | 59.925603 | 30.321773 | 59.928813 | 30.388147 | [ANONYMIZ [ANO |
| 8 | ######### | Completed | 4ad2e95481 | 1b926e88a | 3ffa4a71a5 | stantyan | 18-05-2015 | 18-05-2015 | 0:40:24 | 0:44:00 | 0:03:36 | Volkswagen | Igor | Volkswagen | Male | 59.9287 | 30.387829 | 60.031336 | 30.329518 | [ANONYMIZ [ANO |
| 9 | ######### | Completed | 1e3935b05. | 439ae2cf8a | 3ffa4a71a5 | stantyan | 19-05-2015 | 19-05-2015 | 0:41:56 | 0:58:00 | 0:16:04 | Hyundai Sol | Muhammec | Hyundai | Male | 60.031592 | 30.330248 | 59.944279 | 30.359076 | [ANONYMIZ [ANO |
| 10 | ######### | Completed | 0eb9a9f7a3 | 75a4c47c32 | 3ffa4a71a5 | stantyan | 19-05-2015 | 19-05-2015 | 0:10:06 | 0:15:00 | 0:04:54 | Hyundai ix3 | Vladimir | Hyundai | Male | 59.945143 | 30.356079 | 59.929122 | 30.388656 | [ANONYMIZ [ANO |
| 11 | ######### | Completed | b56495d14! | 176f50c424 | 3ffa4a71a5 | stantyan | 19-05-2015 | 19-05-2015 | 0:25:30 | 0:33:00 | 0:07:30 | Volkswagen | Roman | Volkswagen | Male | 59.92865 | 30.388166 | 60.031366 | 30.329493 | [ANONYMIZ [ANO |
| 12 | ######### | Completed | 613f3deb51 | e516233991 | 3ffa4a71a5 | stantyan | 20-05-2015 | 20-05-2015 | 0:17:47 | 0:26:00 | 0:08:13 | Hyundai Sol | Aleksey | Hyundai | Male | 59.925556 | 30.321134 | 59.929301 | 30.388775 | [ANONYMIZ [ANO |
| 13 | ######### | Completed | 0d486aced5 | a81952458( | 3ffa4a71a5 | stantyan | 31-05-2015 | 31-05-2015 | 0:14:28 | 0:27:00 | 0:12:32 | Hyundai Sol | Ekaterina | Hyundai | Female | 56.752253 | 60.805972 | 56.794567 | 60.614053 | [ANONYMIZ [ANO |
| 14 | ######### | Completed | d12da2c7ae | f065223fa8 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:15:19 | 0:28:00 | 0:12:41 | Ford Monde | Oleg | Ford | Male | 56.751753 | 60.803867 | 56.796019 | 60.614159 | [ANONYMIZ [ANO |
| 15 | ######### | Completed | 36695e9088 | b897afbe68 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:50:13 | 1:04:00 | 0:13:47 | Ford Focus | Aleksandr | Ford | Male | 56.795387 | 60.612997 | 56.857545 | 60.6063 | [ANONYMIZ [ANO |
| 16 | ######### | Completed | 30bd4a26ca | ef33153fbc | 3ffa4a71a5 | stantyan | ######### | ######### | 0:12:39 | 0:37:00 | 0:24:21 | Hyundai Sol | Sergey | Hyundai | Male | 56.857754 | 60.606077 | 56.857726 | 60.606094 | [ANONYMIZ [ANO |
| 17 | ######### | Completed | 3e95596617 | b6584dd09 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:16:18 | 0:34:00 | 0:17:42 | Kia Rio | Aleksey | Kia | Male | 56.795207 | 60.612907 | 56.857543 | 60.606187 | [ANONYMIZ [ANO |
| 18 | ######### | Completed | 4669e3d96 | c85086781 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:11:27 | 0:25:00 | 0:13:33 | Nissan Alme | Andrey | Nissan | Male | 56.795483 | 60.612903 | 56.8283 | 60.597542 | [ANONYMIZ [ANO |
| 19 | ######### | Completed | a7e321acb: | d505ef8cb3 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:16:39 | 0:30:00 | 0:13:21 | Kia Rio | Ekaterina | Kia | Female | 56.795296 | 60.612885 | 56.817396 | 60.634999 | [ANONYMIZ [ANO |
| 20 | ######### | Completed | 8e9c9603d! | 97e28e41c! | 3ffa4a71a5 | stantyan | ######### | ######### | 0:16:28 | 0:27:00 | 0:10:32 | Hyundai Sol | Aleksandr | Hyundai | Male | 56.795367 | 60.614292 | 56.750943 | 60.801553 | [ANONYMIZ [ANO |
| 21 | ######### | Completed | e01e7067e! | 767a9c87d | 3ffa4a71a5 | stantyan | ######### | ######### | 0:14:13 | 0:17:00 | 0:02:47 | Hyundai Sol | Sergey | Hyundai | Male | 59.925502 | 30.339049 | 59.948717 | 30.303298 | [ANONYMIZ [ANO |
| 22 | ######### | Completed | c9df3a9edb | 151944a8f9 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:21:12 | 0:36:00 | 0:14:48 | Chevrolet C | Oleg | Chevrolet | Male | 59.951232 | 30.310537 | 60.011847 | 30.434048 | [ANONYMIZ [ANO |
| 23 | ######### | Cancelled | b970e99cf7 | a89bc69e8 | 3ffa4a71a5 | stantyan | 14-06-2015 | 14-06-2015 | 0:01:00 | 0:46:00 | 0:45:00 | Mercedes-B | Pavel | Mercedes-B | Male | 59.799903 | 30.273236 | 59.799903 | 30.273236 | [ANONYMIZ [ANO |
| 24 | ######### | Completed | 945269c950 | 2f5f40db71 | 3ffa4a71a5 | stantyan | 14-06-2015 | 14-06-2015 | 0:45:40 | 0:54:00 | 0:08:20 | Kia Cee'd | Maksim | Kia | Male | 59.929035 | 30.356715 | 60.056114 | 30.427891 | [ANONYMIZ [ANO |
| 25 | ######### | Completed | 04f14fe72a | 6b3642ae7! | 3ffa4a71a5 | stantyan | 14-06-2015 | 14-06-2015 | 0:30:44 | 0:40:00 | 0:09:16 | Mitsubishi A | Sergey | Mitsubishi | Male | 60.016182 | 30.409457 | 60.02793 | 30.634919 | [ANONYMIZ [ANO |
| 26 | ######### | Completed | fdbd229175 | 69f4e9a875 | 3ffa4a71a5 | stantyan | 14-06-2015 | 14-06-2015 | 0:25:18 | 0:30:00 | 0:04:42 | Nissan Alme | Artem | Nissan | Male | 59.927837 | 30.337983 | 59.800195 | 30.274495 | [ANONYMIZ [ANO |
| 27 | ######### | Completed | d4984dfe40 | e1992ce90 | 3ffa4a71a5 | stantyan | 14-06-2015 | 14-06-2015 | 0:30:18 | 0:37:00 | 0:06:42 | Honda Civic | Finat | Honda | Male | 59.927697 | 30.338072 | 59.927216 | 30.33929 | [ANONYMIZ [ANO |
| 28 | ######### | Completed | 57bd0828f1 | b2cadbd6e | 3ffa4a71a5 | stantyan | 14-06-2015 | 14-06-2015 | 0:09:33 | 0:19:00 | 0:09:27 | Mercedes-B | Gennadiy | Mercedes-B | Male | 59.927608 | 30.338298 | 59.93286 | 30.345696 | [ANONYMIZ [ANO |
| 29 | ######### | Completed | 268145a88 | a45d6a746 | 3ffa4a71a5 | stantyan | 14-06-2015 | 14-06-2015 | 0:22:31 | 0:29:00 | 0:06:29 | Audi A7 | Ramil | Audi | Male | 59.932907 | 30.345892 | 59.927743 | 30.337948 | [ANONYMIZ [ANO |
| 30 | ######### | Completed | 9622fd4e18 | 78aab4ce1c | 3ffa4a71a5 | stantyan | 15-06-2015 | 15-06-2015 | 0:45:00 | 0:57:00 | 0:12:00 | Mitsubishi C | Evgeniy | Mitsubishi | Male | 60.01593 | 30.407802 | 59.927729 | 30.335885 | [ANONYMIZ [ANO |
| 31 | ######### | Completed | 72801dd19 | 3613586d8( | 3ffa4a71a5 | stantyan | 15-06-2015 | 15-06-2015 | 1:08:04 | 1:22:00 | 0:19:56 | Kia Rio | Kryuchkov | Kia | Male | 59.799339 | 30.274062 | 59.927757 | 30.337976 | [ANONYMIZ [ANO |
| 32 | ######### | Completed | f37af9aad4 | a4c34b4a8 | 3ffa4a71a5 | stantyan | 15-06-2015 | 15-06-2015 | 0:12:28 | 0:27:00 | 0:14:32 | BMW 5-seri | Sergey | BMW | Male | 59.927516 | 30.338067 | 59.918768 | 30.285588 | [ANONYMIZ [ANO |

uber-rides-dataset-updated

Edit



| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 41 | ######### | Completed | cb06523d18 | 66cb6d76e: | 3ffa4a71a5 | stantyan | 16-06-2015 | 16-06-2015 | 0:07:30 | 0:21:00 | 0:13:30 | BMW 5-seri | Yana | BMW | Female | 59.943622 | 30.32664 | 59.927673 | 30.33796 | [ANONYMIZ [ANON |
| 42 | ######### | Completed | 94ff538a9d | 9d9dbe840( | 3ffa4a71a5 | stantyan | 17-06-2015 | 17-06-2015 | 0:10:20 | 0:22:00 | 0:11:40 | Opel Astra | Tatyana | Opel | Female | 59.918799 | 30.285494 | 59.927643 | 30.338161 | [ANONYMIZ [ANON |
| 43 | ######### | Completed | 9e4c25f004 | 1ca0a7c5fb | 3ffa4a71a5 | stantyan | 17-06-2015 | 17-06-2015 | 0:31:09 | 0:38:00 | 0:06:51 | Ford Focus | Vadim | Ford | Male | 59.927612 | 30.338089 | 59.800281 | 30.274314 | [ANONYMIZ [ANON |
| 44 | ######### | Completed | ad25b01455 | c7a97f8bd6 | 3ffa4a71a5 | stantyan | 17-06-2015 | 17-06-2015 | 0:28:17 | 0:30:00 | 0:01:43 | Kia Optima | Nikolay | Kia | Male | 59.92891 | 30.335499 | 59.914004 | 30.341825 | [ANONYMIZ [ANON |
| 45 | ######### | Completed | 03210d1c2a | 6410050a39 | 3ffa4a71a5 | stantyan | 17-06-2015 | 17-06-2015 | 0:29:50 | 0:45:00 | 0:15:10 | Audi Q3 | Andrey | Audi | Male | 59.913822 | 30.341858 | 59.927962 | 30.337645 | [ANONYMIZ [ANON |
| 46 | ######### | Completed | e912743c9 | fc8157aa69 | 3ffa4a71a5 | stantyan | 17-06-2015 | 17-06-2015 | 0:11:39 | 0:27:00 | 0:15:21 | Citroen C4 | Aleksandr | Citroen | Male | 59.929053 | 30.356678 | 59.928041 | 30.336991 | [ANONYMIZ [ANON |
| 47 | ######### | Completed | f981dcaa5c | d3f73f8151 | 3ffa4a71a5 | stantyan | 17-06-2015 | 17-06-2015 | 0:48:13 | 0:57:00 | 0:08:47 | Mercedes-B | Maksim | Mercedes-B | Male | 59.927908 | 30.337908 | 59.800026 | 30.274172 | [ANONYMIZ [ANON |
| 48 | ######### | Completed | 910b27c47f | a7d863211( | 3ffa4a71a5 | stantyan | 29-06-2015 | 29-06-2015 | 0:31:20 | 0:46:00 | 0:14:40 | Mercedes-B | Aleksandr | Mercedes-B | Male | 59.928982 | 30.356505 | 59.800015 | 30.27378 | [ANONYMIZ [ANON |
| 49 | ######### | Completed | e4004a6b4( | bbe2b7c59. | 3ffa4a71a5 | stantyan | ######### | ######### | 1:02:01 | 1:17:00 | 0:14:59 | Chevrolet C | Anton | Chevrolet | Male | 59.938426 | 30.348279 | 60.015673 | 30.409926 | [ANONYMIZ [ANON |
| 50 | ######### | Completed | e49979ce9 | e44cceba3( | 3ffa4a71a5 | stantyan | ######### | ######### | 0:37:41 | 0:45:00 | 0:07:19 | Ford Focus | Valeriy | Ford | Male | 60.016093 | 30.40966 | 59.842317 | 30.319145 | [ANONYMIZ [ANON |
| 51 | ######### | Completed | 4d92ff7d2e | 72d4ed2f0c | 3ffa4a71a5 | stantyan | ######### | ######### | 0:42:26 | 0:55:00 | 0:12:34 | Honda CR-V | Aleksey | Honda | Male | 60.01616 | 30.409194 | 59.84233 | 30.319067 | [ANONYMIZ [ANON |
| 52 | ######### | Completed | 0697b6dd1( | bb10f17be2 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:05:24 | 0:08:00 | 0:02:36 | Hyundai Sol | Denis | Hyundai | Male | 60.015113 | 30.388899 | 60.01593 | 30.407802 | [ANONYMIZ [ANON |
| 53 | ######### | Completed | 54b6a258d( | 6035b1f913 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:23:40 | 2:35:00 | 2:11:20 | Hyundai Sol | Viktor | Hyundai | Male | 56.795427 | 60.612876 | 56.842798 | 60.593533 | [ANONYMIZ [ANON |
| 54 | ######### | Completed | 2a545858dc | b314c4ed5 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:13:29 | 0:33:00 | 0:19:31 | Suzuki Swift | Valentin | Suzuki | Male | 56.835338 | 60.599533 | 56.795247 | 60.612822 | [ANONYMIZ [ANON |
| 55 | ######### | Completed | 2bf08f36d2 | 08e2b4b73( | 3ffa4a71a5 | stantyan | ######### | ######### | 0:30:25 | 0:42:00 | 0:11:35 | Kia Optima | Konstantin | Kia | Male | 56.818351 | 60.538783 | 56.754579 | 60.809916 | [ANONYMIZ [ANON |
| 56 | ######### | Completed | a66a7fde34 | 08e2b4b73( | 3ffa4a71a5 | stantyan | ######### | ######### | 0:24:35 | 0:25:00 | 0:00:25 | Volkswagen | Vladimir | Volkswagen | Male | 56.754657 | 60.809913 | 56.794925 | 60.609017 | [ANONYMIZ [ANON |
| 57 | ######### | Completed | 1ab9e6f3cb | acfdcb8355 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:11:27 | 0:19:00 | 0:07:33 | Hyundai Sol | Dmitriy | Hyundai | Male | 56.795531 | 60.612566 | 56.827389 | 60.598018 | [ANONYMIZ [ANON |
| 58 | ######### | Completed | 328a7ed09: | 2d20c209a5 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:12:23 | 0:28:00 | 0:15:37 | Nissan Alme | Konstantin | Nissan | Male | 56.827501 | 60.597904 | 56.795447 | 60.612977 | [ANONYMIZ [ANON |
| 59 | ######### | Completed | 2f4b8c68a2 | b4617db77f | 3ffa4a71a5 | stantyan | ######### | ######### | 0:14:53 | 0:24:00 | 0:09:07 | Renault Dus | Mihail | Renault | Male | 56.803948 | 60.555167 | 56.795337 | 60.612807 | [ANONYMIZ [ANON |
| 60 | ######### | Completed | 444819db0( | 2e214d1c9f | 3ffa4a71a5 | stantyan | ######### | ######### | 0:15:37 | 0:31:00 | 0:15:23 | Mazda MAZ | Vladimir | Mazda | Male | 56.795439 | 60.612664 | 56.843381 | 60.591195 | [ANONYMIZ [ANON |
| 61 | ######### | Completed | 0fffa07abb | 74042fe65b | 3ffa4a71a5 | stantyan | ######### | ######### | 0:15:41 | 0:22:00 | 0:06:19 | Volkswagen | Anton | Volkswagen | Male | 56.795458 | 60.612796 | 56.842877 | 60.594053 | [ANONYMIZ [ANON |
| 62 | ######### | Completed | 87539006f2 | 51d2fca895 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:13:52 | 0:18:00 | 0:04:08 | Peugeot 20 | Danil | Peugeot | Male | 56.837945 | 60.600501 | 56.796989 | 60.614437 | [ANONYMIZ [ANON |
| 63 | ######### | Completed | 7590c2e055 | 8ddd6d1fc2 | 3ffa4a71a5 | stantyan | 13-09-2015 | 13-09-2015 | 0:18:46 | 0:26:00 | 0:07:14 | Dacia Duste | Nikolay | Dacia | Male | 60.015903 | 30.409211 | 60.004901 | 30.299059 | [ANONYMIZ [ANON |
| 64 | ######### | Completed | 15aaea4b6 | c61a6ab054 | 3ffa4a71a5 | stantyan | 13-09-2015 | 13-09-2015 | 0:20:01 | 0:28:00 | 0:07:59 | Hyundai Sol | Igor | Hyundai | Male | 60.004807 | 30.299834 | 60.016069 | 30.409218 | [ANONYMIZ [ANON |
| 65 | ######### | Completed | 4e1c92354( | 672e11ff3f | 3ffa4a71a5 | stantyan | 13-09-2015 | 13-09-2015 | 0:13:00 | 0:35:00 | 0:22:00 | Chrysler Vo | Nikolay | Chrysler | Male | 56.795383 | 60.611239 | 56.750034 | 60.802219 | [ANONYMIZ [ANON |
| 66 | ######### | Completed | cd2550e2ff | c02023222f | 3ffa4a71a5 | stantyan | 13-09-2015 | 13-09-2015 | 0:29:37 | 2:48:00 | 2:18:23 | Renault Nev | Vladimir | Renault | Male | 56.750988 | 60.799555 | 56.829888 | 60.56663 | [ANONYMIZ [ANON |
| 67 | ######### | Completed | c21b0a0c97 | 0927f61a3f | 3ffa4a71a5 | stantyan | 13-09-2015 | 13-09-2015 | 0:32:39 | 0:39:00 | 0:06:21 | Nissan X-Tra | Vyacheslav | Nissan | Male | 59.79903 | 30.273573 | 60.0161 | 30.409211 | [ANONYMIZ [ANON |
| 68 | ######### | Completed | 5f930c5aa9 | c9f2a70459 | 3ffa4a71a5 | stantyan | 14-09-2015 | 14-09-2015 | 0:16:01 | 0:27:00 | 0:10:59 | Renault Meg | Sergey | Renault | Male | 60.013617 | 30.395937 | 60.03538 | 30.441083 | [ANONYMIZ [ANON |
| 69 | ######### | Completed | c967cf6ce2 | ddb3e70a1 | 3ffa4a71a5 | stantyan | 16-09-2015 | 16-09-2015 | 0:07:45 | 0:19:00 | 0:11:15 | Hyundai Sol | Andrey | Hyundai | Male | 60.016136 | 30.409377 | 60.033 | 30.367926 | [ANONYMIZ [ANON |
| 70 | ######### | Completed | ea8a5c508( | 166cfe40ae | 3ffa4a71a5 | stantyan | 16-09-2015 | 16-09-2015 | 0:11:27 | 0:14:00 | 0:02:33 | Mazda MAZ | Ilya | Mazda | Male | 60.033322 | 30.366692 | 60.016172 | 30.409271 | [ANONYMIZ [ANON |
| 71 | ######### | Completed | bfbadc7fa1 | 83a10afa3b | 3ffa4a71a5 | stantyan | 17-09-2015 | 17-09-2015 | 0:10:47 | 0:18:00 | 0:07:13 | BMW 7-seri | Konstantin | BMW | Male | 60.015875 | 30.409196 | 60.01354 | 30.392946 | [ANONYMIZ [ANON |
| 72 | ######### | Completed | 05ffea32f3 | b563a7c6a8 | 3ffa4a71a5 | stantyan | 17-09-2015 | 17-09-2015 | 0:29:05 | 0:45:00 | 0:15:55 | Opel Astra | Aleksandr | Opel | Male | 60.016183 | 30.409267 | 59.971127 | 30.258715 | [ANONYMIZ [ANON |

uber-rides-dataset-updated

### Sheet 1 (rows 81–111)

| # | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 81 | Completed | 51785bdf5 | 730fbb4bf8 | 3ffa4a71a5 | stantyan | 23-09-2015 | 23-09-2015 | 0:53:13 | 1:00:00 | 0:06:47 | Mitsubishi L | Roman | Mitsubishi | Male | 60.016037 | 30.409295 | 59.934462 | 30.332792 | [ANONYMI | [ANON |
| 82 | Completed | f162b9735f | 8377b615fd | 3ffa4a71a5 | stantyan | 23-09-2015 | 23-09-2015 | 0:22:18 | 0:34:00 | 0:11:42 | Volkswagen | Aleksey | Volkswagen | Male | 60.016139 | 30.409163 | 60.03487 | 30.328099 | [ANONYMI | [ANON |
| 83 | Completed | 920b887efe | 77b9a7fbfc | 3ffa4a71a5 | stantyan | 23-09-2015 | 23-09-2015 | 0:14:33 | 0:21:00 | 0:06:27 | Chevrolet C | Oleg | Chevrolet | Male | 60.034757 | 30.327048 | 60.016137 | 30.409249 | [ANONYMI | [ANON |
| 84 | Completed | b40cb7dc34 | be257214ed | 3ffa4a71a5 | stantyan | 25-09-2015 | 25-09-2015 | 0:41:32 | 0:52:00 | 0:10:28 | Hyundai Sol | Evgeniy | Hyundai | Male | 60.01615 | 30.40959 | 60.092162 | 30.496485 | [ANONYMI | [ANON |
| 85 | Completed | 83fb008944 | 2e0ecf0655 | 3ffa4a71a5 | stantyan | 26-09-2015 | 26-09-2015 | 0:06:44 | 0:18:00 | 0:11:16 | Hyundai Sol | Gennadiy | Hyundai | Male | 60.016327 | 30.409493 | 60.01452 | 30.391047 | [ANONYMI | [ANON |
| 86 | Completed | dc564226eb | 0a66f0b182 | 3ffa4a71a5 | stantyan | 26-09-2015 | 26-09-2015 | 0:16:04 | 0:34:00 | 0:17:56 | Audi A6 | Andrey | Audi | Male | 60.016192 | 30.409451 | 60.058316 | 30.337183 | [ANONYMI | [ANON |
| 87 | Completed | cbaccf2e9b | bd870aa56d | 3ffa4a71a5 | stantyan | 26-09-2015 | 26-09-2015 | 0:10:24 | 0:21:00 | 0:10:36 | Skoda Octav | Albert | Skoda | Male | 56.803946 | 60.555391 | 56.830035 | 60.567308 | [ANONYMI | [ANON |
| 88 | Completed | 8483ac302( | ef33153fbc | 3ffa4a71a5 | stantyan | 26-09-2015 | 26-09-2015 | 0:17:39 | 0:28:00 | 0:10:21 | Opel Astra | Sergey | Opel | Male | 56.830015 | 60.567346 | 56.795512 | 60.614259 | [ANONYMI | [ANON |
| 89 | Completed | 8c4de3cc84 | 95a0a7bfed | 3ffa4a71a5 | stantyan | 26-09-2015 | 26-09-2015 | 0:17:34 | 0:21:00 | 0:03:26 | Hyundai Sol | Mihail | Hyundai | Male | 60.05858 | 30.336323 | 60.016337 | 30.409385 | [ANONYMI | [ANON |
| 90 | Completed | db7b80965( | 1f49c09d58 | 3ffa4a71a5 | stantyan | 27-09-2015 | 27-09-2015 | 0:17:29 | 0:29:00 | 0:11:31 | Toyota RAV | Arkadiy | Toyota | Male | 60.016068 | 30.409528 | 59.95758 | 30.358943 | [ANONYMI | [ANON |
| 91 | Completed | daba5c623t | 2c9a3f60a0 | 3ffa4a71a5 | stantyan | 27-09-2015 | 27-09-2015 | 0:23:19 | 0:30:00 | 0:06:41 | Kia Cee'd Sp | Maksim | Kia | Male | 59.957712 | 30.359077 | 60.0162 | 30.409267 | [ANONYMI | [ANON |
| 92 | Completed | 037a51252c | a65ad3cce5 | 3ffa4a71a5 | stantyan | 27-09-2015 | 27-09-2015 | 0:10:41 | 0:20:00 | 0:09:19 | Hyundai Sol | Oleg | Hyundai | Male | 60.016175 | 30.409491 | 59.986364 | 30.354591 | [ANONYMI | [ANON |
| 93 | Completed | 6b4c37de57 | 161951a22e | 3ffa4a71a5 | stantyan | 27-09-2015 | 27-09-2015 | 0:11:24 | 0:21:00 | 0:09:36 | Lifan Solano | Ivan | Lifan | Male | 59.986255 | 30.355476 | 60.016057 | 30.40921 | [ANONYMI | [ANON |
| 94 | Completed | e73e9ef378 | 5edf276511 | 3ffa4a71a5 | stantyan | 28-09-2015 | 28-09-2015 | 0:05:42 | 0:19:00 | 0:13:18 | Saturn Astra | Vadim | Saturn | Male | 60.0108 | 30.406916 | 60.021694 | 30.372496 | [ANONYMI | [ANON |
| 95 | Completed | b920a8c91( | 290770264( | 3ffa4a71a5 | stantyan | 28-09-2015 | 28-09-2015 | 0:10:55 | 0:19:00 | 0:08:05 | Volkswagen | Aleksandr | Volkswagen | Male | 60.021605 | 30.37288 | 60.016165 | 30.409388 | [ANONYMI | [ANON |
| 96 | Completed | e0578dcecf | 1031ddfed1 | 3ffa4a71a5 | stantyan | 28-09-2015 | 28-09-2015 | 0:05:49 | 0:09:00 | 0:03:11 | Mazda MAZ | Pavel | Mazda | Male | 60.0148 | 30.389778 | 60.016178 | 30.409392 | [ANONYMI | [ANON |
| 97 | Completed | 62a6b43eac | 81ca2b70b( | 3ffa4a71a5 | stantyan | 29-09-2015 | 29-09-2015 | 0:32:15 | 0:44:00 | 0:11:45 | Chevrolet La | Vladimir | Chevrolet | Male | 60.016145 | 30.409429 | 59.957628 | 30.359025 | [ANONYMI | [ANON |
| 98 | Completed | 300357121f | 7744054062 | 3ffa4a71a5 | stantyan | 29-09-2015 | 29-09-2015 | 0:25:52 | 0:36:00 | 0:10:08 | Nissan Alme | Oleg | Nissan | Male | 59.957964 | 30.359245 | 60.016196 | 30.409335 | [ANONYMI | [ANON |
| 99 | Completed | 8bbd296ccf | 4093ae4a2 | 3ffa4a71a5 | stantyan | 29-09-2015 | 29-09-2015 | 0:16:44 | 0:30:00 | 0:13:16 | Hyundai Sol | Pavel | Hyundai | Male | 60.016213 | 30.409406 | 59.989577 | 30.438282 | [ANONYMI | [ANON |
| 100 | Completed | 628e670f26 | ca9d5dd796 | 3ffa4a71a5 | stantyan | 29-09-2015 | 29-09-2015 | 0:22:06 | 0:35:00 | 0:12:54 | Toyota Cord | Aleksandr | Toyota | Male | 59.989583 | 30.438283 | 60.015317 | 30.410871 | [ANONYMI | [ANON |
| 101 | Completed | e65ae2777f | ca9d5dd796 | 3ffa4a71a5 | stantyan | 29-09-2015 | 29-09-2015 | 0:07:21 | 0:16:00 | 0:08:39 | Hyundai Sol | Oleg | Hyundai | Male | 60.016149 | 30.409443 | 60.013977 | 30.393789 | [ANONYMI | [ANON |
| 102 | Completed | 563cd1f3e0 | dbcb06a53c | 3ffa4a71a5 | stantyan | 30-09-2015 | 30-09-2015 | 0:23:53 | 0:32:00 | 0:08:07 | Volkswagen | Evgeniy | Volkswagen | Male | 60.016165 | 30.409247 | 60.053023 | 30.331417 | [ANONYMI | [ANON |
| 103 | Completed | fd7a0b1815 | 05cfeb269e | 3ffa4a71a5 | stantyan | 30-09-2015 | 30-09-2015 | 0:14:20 | 0:19:00 | 0:04:40 | Kia Rio | Maksim | Kia | Male | 60.034892 | 30.327693 | 60.01599 | 30.409209 | [ANONYMI | [ANON |
| 104 | Completed | 0ee664410( | 21444fa76d | 3ffa4a71a5 | stantyan | ######### | ######### | 0:34:17 | 0:41:00 | 0:06:43 | Volkswagen | Dmitriy | Volkswagen | Male | 60.016169 | 30.409301 | 59.981884 | 30.212471 | [ANONYMI | [ANON |
| 105 | Completed | 864c3cefc0 | 89de0fed4e | 3ffa4a71a5 | stantyan | ######### | ######### | 0:02:56 | 0:12:00 | 0:09:04 | Renault Me | Vladimir | Renault | Male | 59.986561 | 30.221763 | 59.987892 | 30.205385 | [ANONYMI | [ANON |
| 106 | Completed | 61c4f1b767 | 20db81eb8 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:33:31 | 0:38:00 | 0:04:29 | Ford Focus | Ilya | Ford | Male | 59.987967 | 30.203975 | 59.94247 | 30.355711 | [ANONYMI | [ANON |
| 107 | Completed | 52091a0c3a | c24e3efee4 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:20:02 | 0:30:00 | 0:09:58 | Volkswagen | Dmitriy | Volkswagen | Male | 60.016196 | 30.409362 | 59.799367 | 30.274014 | [ANONYMI | [ANON |
| 108 | Completed | 176e430aa7 | c24e3efee4 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:51:03 | 0:57:00 | 0:05:57 | Hyundai Sol | Dmitriy | Hyundai | Male | 60.016196 | 30.409362 | 59.799363 | 30.274014 | [ANONYMI | [ANON |
| 109 | Completed | 06dd318a8! | 580a78c6e! | 3ffa4a71a5 | stantyan | 20-10-2015 | 20-10-2015 | 0:38:03 | 0:47:00 | 0:08:57 | Hyundai Sor | Sergey | Hyundai | Male | 59.798797 | 30.273622 | 60.016193 | 30.409293 | [ANONYMI | [ANON |
| 110 | Completed | 22f1e9c384 | 6410050a3! | 3ffa4a71a5 | stantyan | 20-10-2015 | 20-10-2015 | 1:25:19 | 1:43:00 | 0:17:41 | Audi Q3 | Andrey | Audi | Male | 60.015753 | 30.408902 | 59.799858 | 30.27411 | [ANONYMI | [ANON |
| 111 | Completed | 9838c1db78 | 3c07b4aa6? | 3ffa4a71a5 | stantyan | 25-10-2015 | 25-10-2015 | 0:38:12 | 0:46:00 | 0:07:48 | Kia Rio | Andrey | Kia | Male | 59.802778 | 30.32294 | 60.016218 | 30.409376 | [ANONYMI | [ANON |

### Sheet 2 (rows 121–151)

| # | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 121 | Completed | 736ad0b25? | 8c75b2fbca | 3ffa4a71a5 | stantyan | ######### | ######### | 0:34:42 | 0:41:00 | 0:06:18 | Hyundai Sol | Evgeniy | Hyundai | Male | 59.986632 | 30.35402 | 60.016186 | 30.409229 | [ANONYMI | [ANON |
| 122 | Completed | 2be85cdd9( | 093b9e29a! | 3ffa4a71a5 | stantyan | ######### | ######### | 0:10:10 | 0:20:00 | 0:09:50 | Hyundai Sol | Maksim | Hyundai | Male | 60.011793 | 30.378658 | 60.016327 | 30.409227 | [ANONYMI | [ANON |
| 123 | Completed | 6ea41ffb88 | 8e78bc6f87 | 3ffa4a71a5 | stantyan | 26-12-2015 | 26-12-2015 | 0:00:09 | 0:13:00 | 0:12:51 | Ford Kuga | Evgeniy | Ford | Male | 60.016113 | 30.409503 | 60.016113 | 30.409503 | [ANONYMI | [ANON |
| 124 | Cancelled | 7fc047ee4a | 164cdace27 | 3ffa4a71a5 | stantyan | 26-12-2015 | 26-12-2015 | 0:01:00 | 0:22:00 | 0:21:00 | Lifan Solano | Ilya | Lifan | Male | 59.932816 | 30.348415 | 59.932816 | 30.348415 | [ANONYMI | [ANON |
| 125 | Completed | 797457d1ac | 4b3883a81c | 3ffa4a71a5 | stantyan | 29-12-2015 | 29-12-2015 | 0:25:31 | 0:32:00 | 0:06:29 | Volkswagen | Murad | Volkswagen | Male | 56.750887 | 60.800028 | 56.830102 | 60.56737 | [ANONYMI | [ANON |
| 126 | Completed | 27199df0fb | 1f57db561a | 3ffa4a71a5 | stantyan | 29-12-2015 | 29-12-2015 | 0:16:30 | 0:26:00 | 0:09:30 | Hyundai Sol | Aleksey | Hyundai | Male | 56.795715 | 60.612906 | 56.830132 | 60.567294 | [ANONYMI | [ANON |
| 127 | Completed | f7a7e398ba | 8de42f01ee | 3ffa4a71a5 | stantyan | 30-12-2015 | 30-12-2015 | 0:11:31 | 0:27:00 | 0:15:29 | Ford Focus | Dmitriy | Ford | Male | 56.85777 | 60.606151 | 56.830217 | 60.567692 | [ANONYMI | [ANON |
| 128 | Completed | e5a9803b4! | cc78656ac2 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:13:15 | 0:20:00 | 0:06:45 | Chevrolet La | Vladimir | Chevrolet | Male | 56.830087 | 60.567557 | 56.807337 | 60.610955 | [ANONYMI | [ANON |
| 129 | Completed | 4aecc5df22 | 23351ce1b8 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:38:34 | 0:47:00 | 0:08:26 | Volkswagen | Igor | Volkswagen | Male | 59.799157 | 30.273731 | 60.016071 | 30.409215 | [ANONYMI | [ANON |
| 130 | Completed | fd82a418a4 | bbd80a462c | 3ffa4a71a5 | stantyan | 16-01-2016 | 16-01-2016 | 0:18:38 | 0:32:00 | 0:13:22 | Lifan Solano | Evgeniy | Lifan | Male | 59.98653 | 30.354448 | 60.016223 | 30.409012 | [ANONYMI | [ANON |
| 131 | Completed | 4666623a8! | 6e9a080f2b | 3ffa4a71a5 | stantyan | 18-01-2016 | 18-01-2016 | 0:08:10 | 0:19:00 | 0:10:50 | Ford Focus | Svyatoslav | Ford | Male | 59.936137 | 30.31502 | 59.931411 | 30.30366 | [ANONYMI | [ANON |
| 132 | Completed | 9c8a553a7e | 544477de18 | 3ffa4a71a5 | stantyan | 18-01-2016 | 18-01-2016 | 0:06:51 | 0:12:00 | 0:05:09 | Volkswagen | Ramaz | Volkswagen | Male | 60.013349 | 30.396595 | 60.016004 | 30.409274 | [ANONYMI | [ANON |
| 133 | Completed | ba1f368be6 | 50ac5474f3 | 3ffa4a71a5 | stantyan | 27-01-2016 | 27-01-2016 | 0:31:32 | 0:50:00 | 0:18:28 | Hyundai ix3 | Anton | Hyundai | Male | 60.014295 | 30.412211 | 59.938885 | 30.393859 | [ANONYMI | [ANON |
| 134 | Completed | 210315ef68 | 18fa9d0dbb | 3ffa4a71a5 | stantyan | 22-02-2016 | 22-02-2016 | 0:30:06 | 0:37:00 | 0:06:54 | Ford Focus | Aleksey | Ford | Male | 59.989658 | 30.25755 | 60.016073 | 30.409248 | [ANONYMI | [ANON |
| 135 | Completed | f09f20ee43 | 15a6894da4 | 3ffa4a71a5 | stantyan | 25-02-2016 | 25-02-2016 | 0:05:14 | 0:10:00 | 0:04:46 | Volkswagen | Nikolay | Volkswagen | Male | 59.915843 | 30.282688 | 59.910147 | 30.275921 | [ANONYMI | [ANON |
| 136 | Completed | 7dedc1f838 | a7d96f1385 | 3ffa4a71a5 | stantyan | 29-02-2016 | 29-02-2016 | 0:26:25 | 0:33:00 | 0:06:35 | Chevrolet C | Vadim | Chevrolet | Male | 60.016013 | 30.409379 | 59.928975 | 30.356312 | [ANONYMI | [ANON |
| 137 | Completed | 19564cbd0! | 524f217db5 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:18:01 | 0:21:00 | 0:02:59 | Skoda Octav | Aleksey | Skoda | Male | 60.051813 | 30.33642 | 60.015927 | 30.409991 | [ANONYMI | [ANON |
| 138 | Completed | 68e299628! | dce6ffc685? | 3ffa4a71a5 | stantyan | 14-03-2016 | 14-03-2016 | 0:11:03 | 0:14:00 | 0:02:57 | Volkswagen | Maksim | Volkswagen | Male | 59.943463 | 30.323775 | 59.929262 | 30.356958 | [ANONYMI | [ANON |
| 139 | Completed | 2f7ca9c163 | af6d53a1a0 | 3ffa4a71a5 | stantyan | 17-03-2016 | 17-03-2016 | 0:10:50 | 0:19:00 | 0:08:10 | Volkswagen | Dmitriy | Volkswagen | Male | 59.989713 | 30.462224 | 60.015927 | 30.409991 | [ANONYMI | [ANON |
| 140 | Completed | 9fa1d8baa5 | 484892a25? | 3ffa4a71a5 | stantyan | 21-03-2016 | 21-03-2016 | 0:08:10 | 0:14:00 | 0:05:50 | Nissan Alme | Nikolay | Nissan | Male | 59.944748 | 30.35941 | 59.94359 | 30.32399 | [ANONYMI | [ANON |
| 141 | Completed | 5c44c704ce | 5d5611a9bc | 3ffa4a71a5 | stantyan | 25-03-2016 | 25-03-2016 | 0:19:16 | 0:29:00 | 0:09:44 | Audi A3 | Aleksey | Audi | Male | 60.016213 | 30.409305 | 60.042637 | 30.37748 | [ANONYMI | [ANON |
| 142 | Completed | e3ccba6686 | d18e2051d? | 3ffa4a71a5 | stantyan | 25-03-2016 | 25-03-2016 | 0:14:43 | 0:24:00 | 0:09:17 | Toyota Cam | Viktor | Toyota | Male | 60.042747 | 30.37656 | 60.015468 | 30.410383 | [ANONYMI | [ANON |
| 143 | Cancelled | 2cfe7cf24f9 | fbaaa4fe33 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:11:14 | 0:14:00 | 0:02:46 | Lifan Solano | Igor | Lifan | Male | 59.934486 | 30.348009 | 59.938403 | 30.34533 | [ANONYMI | [ANON |
| 144 | Completed | b603f7e2fb | b700e12aef | 3ffa4a71a5 | stantyan | ######### | ######### | 0:04:40 | 0:08:00 | 0:03:20 | Hyundai Sol | Ruslan | Hyundai | Male | 59.938339 | 30.347482 | 59.936145 | 30.320686 | [ANONYMI | [ANON |
| 145 | Completed | 73fd8acb11 | 164cdace27 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:46:02 | 0:51:00 | 0:04:58 | Mercedes-B | Viktor | Mercedes-B | Male | 60.020078 | 30.409435 | 60.103418 | 29.947507 | [ANONYMI | [ANON |
| 146 | Completed | accfabf376« | 5d5611a9bc | 3ffa4a71a5 | stantyan | ######### | ######### | 0:16:04 | 0:26:00 | 0:09:56 | Audi A3 | Aleksey | Audi | Male | 60.015703 | 30.409417 | 60.047282 | 30.35391 | [ANONYMI | [ANON |
| 147 | Completed | df2e21420e | e6d00c4da! | 3ffa4a71a5 | stantyan | ######### | ######### | 0:30:37 | 0:36:00 | 0:05:23 | Volkswagen | Roman | Volkswagen | Male | 59.869218 | 30.33892 | 59.891608 | 30.511607 | [ANONYMI | [ANON |
| 148 | Completed | 6c1c93bc79 | 513f8eb218 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:37:09 | 0:38:00 | 0:00:51 | Ford Focus | Valeriy | Ford | Male | 60.020078 | 30.409435 | 59.800195 | 30.274495 | [ANONYMI | [ANON |
| 149 | Completed | 160f29c8e2 | ef0ef48fb4? | 3ffa4a71a5 | stantyan | ######### | ######### | 0:12:09 | 0:20:00 | 0:07:51 | Kia Rio | Margarita | Kia | Female | 60.016231 | 30.409452 | 60.031687 | 30.428322 | [ANONYMI | [ANON |
| 150 | Completed | 48fb57dd9b | d5d270edb8 | 3ffa4a71a5 | stantyan | 17-04-2016 | 17-04-2016 | 0:16:18 | 0:17:00 | 0:00:42 | Volkswagen | Egor | Volkswagen | Male | 59.986498 | 30.354555 | 60.016138 | 30.409216 | [ANONYMI | [ANON |
| 151 | Completed | a60a08e90« | 47f2730160 | 3ffa4a71a5 | stantyan | 19-04-2016 | 19-04-2016 | 0:17:07 | 0:24:00 | 0:06:53 | Skoda Octav | Viktor | Skoda | Male | 60.058318 | 30.331046 | 60.016136 | 30.409313 | [ANONYMI | [ANON |

### Sheet 3 (rows 241–271)

| # | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 241 | Completed | 706e53114« | 1680202003 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:16:54 | 0:25:00 | 0:08:06 | Hyundai Sol | Mihail | Hyundai | Male | 59.959772 | 30.480846 | 59.91075 | 30.444053 | [ANONYMI | [ANON |
| 242 | Completed | 28f480d141 | 126ea68bb( | 3ffa4a71a5 | stantyan | ######### | ######### | 0:14:52 | 0:21:00 | 0:06:08 | Toyota Cord | Viktor | Toyota | Male | 59.946542 | 30.475798 | 59.959772 | 30.480846 | [ANONYMI | [ANON |
| 243 | Completed | dec533012t | 75c7d46a4! | 3ffa4a71a5 | stantyan | 14-10-2016 | 14-10-2016 | 0:20:04 | 0:28:00 | 0:07:56 | Skoda Octav | Yuriy | Skoda | Male | 59.941415 | 30.366456 | 60.031821 | 30.428298 | [ANONYMI | [ANON |
| 244 | Completed | 8bda2b008? | 70513ea55« | 3ffa4a71a5 | stantyan | 15-10-2016 | 15-10-2016 | 0:22:25 | 0:28:00 | 0:05:35 | Kia Cee'd | Aleksey | Kia | Male | 59.98648 | 30.354853 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 245 | Completed | 95358afd69 | d86f83f26a | 3ffa4a71a5 | stantyan | 18-10-2016 | 18-10-2016 | 0:12:15 | 0:15:00 | 0:02:45 | Volkswagen | Ruslan | Volkswagen | Male | 60.035786 | 30.419227 | 60.061936 | 30.375745 | [ANONYMI | [ANON |
| 246 | Completed | e252856e4( | d86f83f26a | 3ffa4a71a5 | stantyan | 18-10-2016 | 18-10-2016 | 0:11:39 | 0:22:00 | 0:10:21 | Volkswagen | Ruslan | Volkswagen | Male | 60.067314 | 30.377042 | 60.035872 | 30.41767 | [ANONYMI | [ANON |
| 247 | Completed | f4f55df9b6c | 59b001b93? | 3ffa4a71a5 | stantyan | 19-10-2016 | 19-10-2016 | 0:15:45 | 0:23:00 | 0:07:15 | Skoda Octav | Maksim | Skoda | Male | 59.95914 | 30.477307 | 60.031627 | 30.428617 | [ANONYMI | [ANON |
| 248 | Completed | 819ac176fd | caf93b6df8( | 3ffa4a71a5 | stantyan | 23-10-2016 | 23-10-2016 | 0:07:45 | 0:13:00 | 0:05:15 | Ford Focus | Ruslan | Ford | Male | 59.941415 | 30.366456 | 59.944795 | 30.48372 | [ANONYMI | [ANON |
| 249 | Completed | 1df46bf101( | efd973b4b3 | 3ffa4a71a5 | stantyan | 24-10-2016 | 24-10-2016 | 0:18:11 | 0:30:00 | 0:11:49 | Saturn Astra | Dmitriy | Saturn | Male | 59.891657 | 30.511732 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 250 | Completed | 088d605e3« | a22df14b79 | 3ffa4a71a5 | stantyan | 26-10-2016 | 26-10-2016 | 0:16:32 | 0:24:00 | 0:07:28 | Volkswagen | Vyacheslav | Volkswagen | Male | 60.03165 | 30.428316 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 251 | Completed | 961f743e7c | a18931166? | 3ffa4a71a5 | stantyan | 26-10-2016 | 26-10-2016 | 0:35:06 | 0:42:00 | 0:06:54 | Ford Monde | Kirill | Ford | Male | 60.031797 | 30.427917 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 252 | Completed | 855045cc06 | 736f262f69 | 3ffa4a71a5 | stantyan | 27-10-2016 | 27-10-2016 | 0:40:59 | 0:51:00 | 0:10:01 | Hyundai Sol | Valeriy | Hyundai | Male | 60.012595 | 30.399115 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 253 | Completed | bfb9a1a99c | dd7af04184 | 3ffa4a71a5 | stantyan | 28-10-2016 | 28-10-2016 | 0:16:46 | 0:29:00 | 0:12:14 | Renault Dus | Kirill | Renault | Male | 59.941415 | 30.366456 | 60.026073 | 30.42244 | [ANONYMI | [ANON |
| 254 | Completed | 9936a1b72! | 8a31e4d75? | 3ffa4a71a5 | stantyan | ######### | ######### | 0:26:52 | 0:33:00 | 0:06:08 | Ford Monde | Mihail | Ford | Male | 59.941415 | 30.366456 | 60.031642 | 30.428125 | [ANONYMI | [ANON |
| 255 | Completed | 3b47b0f524 | 44c00edf5d | 3ffa4a71a5 | stantyan | ######### | ######### | 1:24:48 | 1:35:00 | 0:10:12 | Ford Focus | Vadim | Ford | Male | 59.982914 | 30.211445 | 59.959646 | 30.480601 | [ANONYMI | [ANON |
| 256 | Completed | db0b34164( | c06e7b818? | 3ffa4a71a5 | stantyan | ######### | ######### | 0:11:04 | 0:19:00 | 0:07:56 | Lifan Solano | Aleksandr | Lifan | Male | 59.93357 | 30.436758 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 257 | Completed | 83fe9c5ed1 | 968aa8f88c | 3ffa4a71a5 | stantyan | ######### | ######### | 0:21:20 | 0:39:00 | 0:17:40 | Kia Rio | Vadim | Kia | Male | 59.95912 | 30.477045 | 60.031602 | 30.428307 | [ANONYMI | [ANON |
| 258 | Completed | 41f664e3d0 | 9a92e4848c | 3ffa4a71a5 | stantyan | ######### | ######### | 0:36:11 | 0:48:00 | 0:11:49 | Toyota Aver | Vladimir | Toyota | Male | 59.959766 | 30.474675 | 60.031728 | 30.427998 | [ANONYMI | [ANON |
| 259 | Completed | 9a41064b2? | dc522dce0c | 3ffa4a71a5 | stantyan | ######### | ######### | 0:18:43 | 0:26:00 | 0:07:17 | Mitsubishi L | Dmitriy | Mitsubishi | Male | 59.959773 | 30.432238 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 260 | Completed | d2d49cd1b( | 64af7d80e5 | 3ffa4a71a5 | stantyan | ######### | ######### | 0:35:40 | 0:48:00 | 0:12:20 | Ford Fiesta | Artem | Ford | Male | 59.91212 | 30.445631 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 261 | Completed | 6b0fea4fab | ef0ef48fb4? | 3ffa4a71a5 | stantyan | 18-11-2016 | 18-11-2016 | 0:20:31 | 0:32:00 | 0:11:29 | Kia Rio | Margarita | Kia | Female | 59.941415 | 30.366456 | 60.031556 | 30.428777 | [ANONYMI | [ANON |
| 262 | Completed | 9a67ecf355 | fdb537e31f | 3ffa4a71a5 | stantyan | 22-11-2016 | 22-11-2016 | 0:09:07 | 0:13:00 | 0:03:53 | Renault Flue | Aleksandr | Renault | Male | 59.946469 | 30.475277 | 59.912907 | 30.475496 | [ANONYMI | [ANON |
| 263 | Completed | e09859a8ea | 3ed895daf4 | 3ffa4a71a5 | stantyan | 22-11-2016 | 22-11-2016 | 0:18:38 | 0:20:00 | 0:01:22 | Peugeot 408 | Maksim | Peugeot | Male | 59.911704 | 30.44503 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 264 | Cancelled | b27215fa3f | 09912ed6d! | 3ffa4a71a5 | stantyan | 25-11-2016 | 25-11-2016 | 0:25:32 | 1:00:00 | 0:34:28 | Hyundai Equ | Sergey | Hyundai | Male | 59.925447 | 30.488825 | 59.925468 | 30.488632 | [ANONYMI | [ANON |
| 265 | Completed | c1a5ebeb74 | 57974b628« | 3ffa4a71a5 | stantyan | 25-11-2016 | 25-11-2016 | 0:08:58 | 0:13:00 | 0:04:02 | Nissan Alme | Dmitriy | Nissan | Male | 59.925447 | 30.488825 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 266 | Completed | d770190f94 | 3847307233 | 3ffa4a71a5 | stantyan | 25-11-2016 | 25-11-2016 | 0:05:56 | 0:09:00 | 0:03:04 | Kia Rio | Boris | Kia | Male | 59.941415 | 30.366456 | 59.946597 | 30.475497 | [ANONYMI | [ANON |
| 267 | Completed | 4e0127f57b | f699025b78 | 3ffa4a71a5 | stantyan | 25-11-2016 | 25-11-2016 | 0:09:52 | 0:16:00 | 0:06:08 | Volkswagen | Sergey | Volkswagen | Male | 59.924753 | 30.489413 | 59.946842 | 30.474938 | [ANONYMI | [ANON |
| 268 | Completed | eb5989591« | a696bd6e3! | 3ffa4a71a5 | stantyan | 25-11-2016 | 25-11-2016 | 0:23:13 | 0:33:00 | 0:09:47 | Ford Focus | Andrey | Ford | Male | 59.941415 | 30.366456 | 59.925485 | 30.488833 | [ANONYMI | [ANON |
| 269 | Completed | ac1490d58? | 7445538d5? | 3ffa4a71a5 | stantyan | 25-11-2016 | 25-11-2016 | 0:47:09 | 0:59:00 | 0:11:51 | Renault San | Aleksandr | Renault | Male | 59.892768 | 30.515795 | 59.941415 | 30.366456 | [ANONYMI | [ANON |
| 270 | Completed | 9a817ac95? | 9e9551459? | 3ffa4a71a5 | stantyan | 26-11-2016 | 26-11-2016 | 0:14:09 | 0:25:00 | 0:10:51 | Audi A7 | Yuriy | Audi | Male | 59.941415 | 30.366456 | 59.925456 | 30.488847 | [ANONYMI | [ANON |
| 271 | Completed | 42e0ca4d2« | cf506ca03d | 3ffa4a71a5 | stantyan | 26-11-2016 | 26-11-2016 | 0:12:05 | 0:15:00 | 0:02:55 | Mitsubishi L | Mugutdin | Mitsubishi | Male | 59.941415 | 30.366456 | 59.488848 | 30.488848 | [ANONYMI | [ANON |

```
A    B         C          D          E        F        G           H           I       J       K        L              M            N            O       P          Q          R          S          T           U
301  ######## Completed  2bc2d7c0f3 c89a86192c 3ffa4a71a5 stantyan  29-12-2016  29-12-2016  0:11:19  0:23:00  0:11:41  Hyundai Sol  Vyacheslav   Hyundai      Male    59.941415  30.366456  59.950601  30.499693  [ANONYMIZ [ANON
302  ######## Completed  8d123ec50a 37e3ee857! 3ffa4a71a5 stantyan  29-12-2016  29-12-2016  0:09:08  0:19:00  0:09:52  Lifan Solanc Sergey       Lifan        Male    59.951163  30.499077  59.941415  30.366456  [ANONYMIZ [ANON
303  ######## Completed  59cf3390ca 63d9ad269c 3ffa4a71a5 stantyan  ########    ########    0:38:12  0:42:00  0:03:48  Mercedes-B  Andrey       Mercedes-B   Male    59.941415  30.366456  59.799797  30.274013  [ANONYMIZ [ANON
304  ######## Completed  37cf2290ca 63d9ad269c 3ffa4a71a5 stantyan  ########    ########    0:38:12  0:42:00  0:03:48  Renault Flue Vsevolod     Renault      Male    59.941415  30.366456  59.799797  30.274013  [ANONYMIZ [ANON
305  ######## Completed  2a5766783( f1c8df3e36 3ffa4a71a5 stantyan  ########    ########    0:11:58  0:21:00  0:09:02  Ford Monde  Yuriy        Ford         Male    59.941415  30.366456  59.95103   30.409437  [ANONYMIZ [ANON
306  ######## Completed  0e123f4309 e8e0cfb4c9 3ffa4a71a5 stantyan  ########    ########    0:11:51  0:17:00  0:05:09  Ford Focus   Aleksey      Ford         Male    59.947577  30.414213  59.941415  30.366456  [ANONYMIZ [ANON
307  ######## Completed  7e72c0b51( 2081b30f5a 3ffa4a71a5 stantyan  ########    ########    0:16:20  0:18:00  0:01:40  Peugeot 20! Aleksandr    Peugeot      Male    59.941415  30.366456  59.950725  30.500285  [ANONYMIZ [ANON
308  ######## Completed  baac5a47b! 3bf0c22123 3ffa4a71a5 stantyan  ########    ########    0:06:04  0:14:00  0:07:56  Hyundai Sol Aleksey      Hyundai      Male    59.949967  30.500365  59.941415  30.366456  [ANONYMIZ [ANON
309  ######## Completed  0af692d4ed dab590332! 3ffa4a71a5 stantyan  ########    ########    0:11:09  0:18:00  0:06:51  Renault Kolc Mihail       Renault      Male    59.941415  30.366456  59.94783   30.411762  [ANONYMIZ [ANON
310  ######## Completed  9929ec769a dc0cb5bbe! 3ffa4a71a5 stantyan  ########    ########    0:14:16  0:18:00  0:03:44  Volkswagen Azad          Volkswagen  Male    59.947695  30.413582  59.941415  30.366456  [ANONYMIZ [ANON
311  ######## Completed  419e92d7b¡ cc5ad05422 3ffa4a71a5 stantyan  ########    ########    0:17:56  0:27:00  0:09:04  Lada Granta Abdula       Lada         Male    59.940385  30.45526   59.941415  30.366456  [ANONYMIZ [ANON
312  ######## Completed  a2e73d35f0 932492a3a! 3ffa4a71a5 stantyan  13-01-2017  13-01-2017  0:18:26  0:26:00  0:07:34  Hyundai Sol Dzhanpolad Hyundai      Male    59.941415  30.366456  60.031686  30.428149  [ANONYMIZ [ANON
313  ######## Completed  b0fa949f94 c7f5f5f831- 3ffa4a71a5 stantyan  13-01-2017  13-01-2017  0:24:05  0:28:00  0:03:55  Volkswagen Ruslan        Volkswagen  Male    60.033742  30.420001  59.941415  30.366456  [ANONYMIZ [ANON
314  ######## Completed  e65d6065e7 90dc4693d¡ 3ffa4a71a5 stantyan  16-01-2017  16-01-2017  0:38:15  0:41:00  0:02:45  Hyundai Sol Aleksandr    Hyundai      Male    59.938962  30.350321  59.941415  30.366456  [ANONYMIZ [ANON
315  ######## Completed  ac5441152! 82a36d768! 3ffa4a71a5 stantyan  19-01-2017  19-01-2017  0:33:12  0:37:00  0:03:48  Renault Me; Dmitriy      Renault      Male    59.799237  30.273759  59.941415  30.366456  [ANONYMIZ [ANON
316  ######## Completed  9e98eddc6; 1d7b33fdd! 3ffa4a71a5 stantyan  25-01-2017  25-01-2017  0:34:54  0:42:00  0:07:06  Kia Rio      Pavel        Kia          Male    59.941415  30.366456  59.957176  30.322328  [ANONYMIZ [ANON
317  ######## Completed  91059eb56; 04ffcd2a2e¡ 3ffa4a71a5 stantyan  25-01-2017  25-01-2017  0:45:09  0:51:00  0:05:51  Ford Focus   Mihail       Ford         Male    59.955558  30.336818  59.959008  30.477597  [ANONYMIZ [ANON
318  ######## Completed  0d0d6d6b9¡ ba5e67041¡ 3ffa4a71a5 stantyan  26-01-2017  26-01-2017  0:35:00  0:38:00  0:03:00  Mazda MAZ Konstantin    Mazda        Male    59.941415  30.366456  59.800062  30.274375  [ANONYMIZ [ANON
319  ######## Completed  aefc739092 1788ef33eb 3ffa4a71a5 stantyan  ########    ########    0:32:18  0:38:00  0:05:42  Hyundai Sol Mariya       Hyundai      Female  59.941415  30.366456  59.918052  30.341649  [ANONYMIZ [ANON
320  ######## Completed  813cd8aeff; b60725f055 3ffa4a71a5 stantyan  ########    ########    0:34:18  0:37:00  0:02:42  Volkswagen Mihail        Volkswagen  Male    59.918062  30.341495  59.941415  30.366456  [ANONYMIZ [ANON
321  ######## Completed  29615b8df1 8107ebb4e! 3ffa4a71a5 stantyan  ########    ########    0:31:21  0:36:00  0:04:39  Nissan Alme Denis        Nissan       Male    59.893612  30.519012  59.941415  30.366456  [ANONYMIZ [ANON
322  ######## Completed  f2103acdfe( d5c9f6dbed 3ffa4a71a5 stantyan  ########    ########    0:15:59  0:26:00  0:10:01  Kia Cee'd    Aleksandr    Kia          Male    59.941415  30.366456  59.940451   30.39729  [ANONYMIZ [ANON
323  ######## Completed  16dc6125ff; 250f251d1d 3ffa4a71a5 stantyan  ########    ########    0:19:00  0:21:00  0:02:00  Kia Sorento  Roman        Kia          Male    59.941224  30.394085  59.941415  30.366456  [ANONYMIZ [ANON
324  ######## Completed  b45aedc53¡ 7124589e8! 3ffa4a71a5 stantyan  14-02-2017  14-02-2017  0:36:58  0:42:00  0:05:02  Ford Fiesta  Igor         Ford         Male    59.799291  30.273972  59.941415  30.366456  [ANONYMIZ [ANON
325  ######## Completed  d3dffbdf95¡ 7f30458dd9 3ffa4a71a5 stantyan  ########    ########    0:23:03  0:33:00  0:09:57  Nissan Alme Viktor       Nissan       Male    59.941415  30.366456  59.916628  30.450753  [ANONYMIZ [ANON
326  ######## Completed  6e99ddbf1a e5c54c2c20 3ffa4a71a5 stantyan  ########    ########    0:20:14  0:26:00  0:05:46  Geely Emgra Narek        Geely        Male    59.916897  30.451212  59.941415  30.366456  [ANONYMIZ [ANON
327  ######## Completed  f8cdcca580 6bce685bc¡ 3ffa4a71a5 stantyan  19-03-2017  19-03-2017  0:19:39  0:26:00  0:06:21  Ford Focus   Viktor       Ford         Male    59.941415  30.366456  59.916841  30.450966  [ANONYMIZ [ANON
328  ######## Completed  5ae364071! 43fb08c568 3ffa4a71a5 stantyan  19-03-2017  19-03-2017  0:14:41  0:20:00  0:05:19  Skoda Fabia Dzhahon      Skoda        Male    59.916893  30.451014  59.941415  30.366456  [ANONYMIZ [ANON
329  ######## Completed  a5a7fdcfd9( 6c1c94c35a 3ffa4a71a5 stantyan  19-03-2017  19-03-2017  0:29:14  0:38:00  0:08:46  Citroen Berl Aleksandr   Citroen      Male    59.941415  30.366456  60.060315   30.29012  [ANONYMIZ [ANON
330  ######## Completed  7807aee5bc 6c1c94c35a 3ffa4a71a5 stantyan  19-03-2017  19-03-2017  0:32:01  0:33:00  0:00:59  Citroen Berl Aleksandr   Citroen      Male    60.060602  30.290772  59.941415  30.366456  [ANONYMIZ [ANON
331  ######## Completed  8c9bdf8550 c336a1898( 3ffa4a71a5 stantyan  23-03-2017  23-03-2017  0:26:59  0:28:00  0:01:01  Nissan Alme Aleksey      Nissan       Male    59.950157  30.409832  59.949298  30.40974   [ANONYMIZ [ANON
332  ######## Completed  4651a3a51/ 0b573ca0d1 3ffa4a71a5 stantyan  ########    ########    0:14:04  0:22:00  0:07:56  Renault Log Georgiy     Renault      Male    59.960016  30.478503  59.91748   30.448984  [ANONYMIZ [ANON
```

uber-rides-dataset-updated

# Cleaning the Dataset:
## Code:

```
#reading dataset
data<-read.csv("uber-rides-dataset.csv")
#number of rows in dataframe
nrow(data)
#number ofcolumns in dataframe
ncol(data)
#structure of dataframe
str(data)
#checking is there any NA values in dataframe
any(is.na(data))
#checking which column contains NA values
for(i in names(data)) {
  if(sum(is.na(data[i]))) {
    print(paste(sum(is.na(data[i]))," NA values in ",i))
  }
}
#replacing NA values of surge_multiplier with mean
data$surge_multiplier[is.na(data$surge_multiplier)]<-mean(data$surge_multiplier,na.rm=TRUE)
#checking NA values in dataframe
any(is.na(data))
#cleaning trip_start_address column
library(stringr)
```

```
b<-c()
for(k in c(1:length(data$trip_start_address))) {
  a<-""
for(i in c(1:nchar(data$trip_start_address[k]))) {
  for(j in c(substr(data$trip_start_address[k],i,i))) {
    if(j %in% letters | j %in% LETTERS | j %in% c(0:9) | j %in% "," | j
%in% ".") {
      a<-paste(a,j)
      a<-str_replace_all(a," ","")


    }
   }
}
  b<-append(b,a)
}
data$trip_start_address<-b

#cleaning trip_end_address column
library(stringr)
b<-c()
for(k in c(1:length(data$trip_end_address))) {
  a<-""
  for(i in c(1:nchar(data$trip_end_address[k]))) {
    for(j in c(substr(data$trip_end_address[k],i,i))) {
      if(j %in% letters | j %in% LETTERS | j %in% c(0:9) | j %in% "," | j
%in% ".") {
        a<-paste(a,j)
        a<-str_replace_all(a," ","")


    }
   }
  }
  b<-append(b,a)
}
data$trip_end_address<-b
#changing trip_completed_at to timestamp
library(lubridate)
b<-c()
for(i in c(1:length(data$trip_completed_at))) {
  a<-strsplit(data$trip_completed_at,"at")
  b<-append(b,mdy_hm(paste(a[[i]][1],a[[i]][2])))
}
```

b
data$trip_completed_at<-b
#suumary of data after cleaning
summary(data)
write.csv(data,"uber-rides-dataset-updated.csv")



```r
1  #reading dataset
2  data<-read.csv("uber-rides-dataset.csv")
3  #number of rows in dataframe
4  nrow(data)
5  #number ofcolumns in dataframe
6  ncol(data)
7  #structure of dataframe
8  str(data)
9  #checking is there any NA values in dataframe
10 any(is.na(data))
11 #checking which column contains NA values
12 for(i in names(data)) {
13   if(sum(is.na(data[i]))) {
14     print(paste(sum(is.na(data[i]))," NA values in ",i))
15   }
16 }
17 #replacing NA values of surge_multiplier with mean
18 data$surge_multiplier[is.na(data$surge_multiplier)]<-mean(data$surge_multiplier,na.rm=TRUE)
19 #checking NA values in dataframe
20 any(is.na(data))
21 #cleaning trip_start_address column
22 library(stringr)
23 b<-c()
24
```

```r
21 #cleaning trip_start_address column
22 library(stringr)
23 b<-c()
24 for(k in c(1:length(data$trip_start_address))) {
25   a<-""
26   for(i in c(1:nchar(data$trip_start_address[k]))) {
27     for(j in c(substr(data$trip_start_address[k],i,i))) {
28       if(j %in% letters | j %in% LETTERS | j %in% c(0:9) | j %in% "," | j %in% ".") {
29         a<-paste(a,j)
30         a<-str_replace_all(a," ","")
31       }
32     }
33   }
34   }
35   b<-append(b,a)
36 }
37 data$trip_start_address<-b
38
39 #cleaning trip_end_address column
40 library(stringr)
41 b<-c()
42 for(k in c(1:length(data$trip_end_address))) {
43   a<-""
44
```

```r
46        if(j %in% letters | j %in% LETTERS | j %in% c(0:9) | j %in% "," | j %in% ".") {
47            a<-paste(a,j)
48            a<-str_replace_all(a," ","")
49
50        }
51      }
52    }
53    b<-append(b,a)
54  }
55  data$trip_end_address<-b
56  #changing trip_completed_at to timestamp
57  library(lubridate)
58  b<-c()
59  for(i in c(1:length(data$trip_completed_at))) {
60    a<-strsplit(data$trip_completed_at,"at")
61    b<-append(b,mdy_hm(paste(a[[i]][1],a[[i]][2])))
62  }
63  b
64  data$trip_completed_at<-b
65  #suumary of data after cleaning
66  summary(data)
67  write.csv(data,"uber-rides-dataset-updated.csv")
68
```

```r
}
return(recursion(v,r+1))
}
}
result <- recursion(a,1)
print(result)
```

# Conclusion:

There were very few clearly erroneous entries in the dataset and a small proportion of suspicious cases or *anomalies* that warrant further internal analysis.All taxi and for-hire-vehicles companies operating in the city, which include Uber, Lyft, and others release their data periodically. An update is published twice a year. The destination data were missing, and an extremely small number of cases had missing trip distance and destination. The imputation method chosen for the latter set was the mean distance and duration of their respective origin-destination pair. The entries with missing destination were left unchanged, although the information from the vast number of complete cases could potentially be used to determine the most probable destination.

The relation between a trip's duration and distance is not entirely linear. Rather, it approximates to a power function because shorter trips, occurring mostly within busy areas of traffic, tend to result in lower average trip speed.Analysis and undoubtedly highlighted the critical importance of a well-defined business problem, which directs all coding efforts to a particular purpose and reveals key details. This business case also attempted to demonstrate the basic use of python in everyday business activities, showing how fun, important, and fun it can  behavioural change of customers raising more complaints about taxi services.

# Reference:

    I.  https://data-flair.training/blogs/r-data-science-project-uber-data-analysis/

    II.https://www.sciencedirect.com/science/article/abs/pii/S2214367X20302027

Thank You