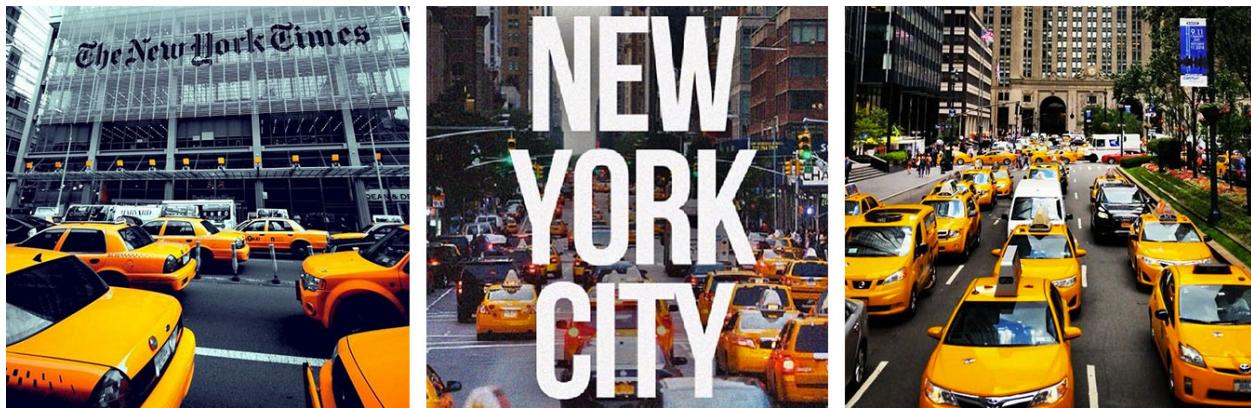


TAXI TAXI - A beginner approach on Exploratory Analysis

Suresh

Date- 2018-02-18



Introduction

This is an basic Exploratory Data Analysis for the NYC Taxi Ride Duration competition.

New York City taxi rides paint a vibrant picture of life in the city. The millions of rides taken each month can provide insight into traffic patterns, road blockage, or large-scale events that attract many New Yorkers. With ridesharing apps gaining popularity, it is increasingly important for taxi companies to provide visibility to their estimated fare and ride duration, since the competing apps provide these metrics upfront. Predicting fare and duration of a ride can help passengers decide when is the optimal time to start their commute.

The primary goal of this project is to predict trip duration of NYC Taxis based on features like trip coordinates, duration date and time. Training dataset has close to 1.4 Million and 630k records in test dataset. Each row contains one taxi trip.

Load required Packages

- Predominantly we have used dplyr, ggplot2 and lubridate packges

```
library(dplyr)
library(tidyr)
library(lubridate)
library(ggplot2)
library(ggmap)
library(stringr)
library(scales)
library(gridExtra)
library(corrplot)
library(RColorBrewer)
library(geosphere)
library(tibble)
library(forcats)
library(xgboost)
library(caret)
library(leaflet)
library(maps)
library(readr)
```

```
library(randomForest)
library(knitr)
```

Load Data

We are using Read CSV function to read train file data into R.

```
taxi <- read_csv("train.csv")

## Parsed with column specification:
## cols(
##   id = col_character(),
##   vendor_id = col_integer(),
##   pickup_datetime = col_datetime(format = ""),
##   dropoff_datetime = col_datetime(format = ""),
##   passenger_count = col_integer(),
##   pickup_longitude = col_double(),
##   pickup_latitude = col_double(),
##   dropoff_longitude = col_double(),
##   dropoff_latitude = col_double(),
##   store_and_fwd_flag = col_character(),
##   trip_duration = col_integer()
## )

set.seed(12345)
taxi1 <- sample_n(taxi, 400000)
```

File Structure & Integrity

In this section we are checking the summary of Taxi Records and also check if any BLANKS or NA values present.

```
glimpse(taxi)

## # Observations: 1,458,644
## # Variables: 11
## # $ id              <chr> "id2875421", "id2377394", "id3858529", "id3...
## # $ vendor_id       <int> 2, 1, 2, 2, 2, 1, 2, 1, 2, 2, 2, 2, 2...
## # $ pickup_datetime <dttm> 2016-03-14 17:24:55, 2016-06-12 00:43:35, ...
## # $ dropoff_datetime <dttm> 2016-03-14 17:32:30, 2016-06-12 00:54:38, ...
## # $ passenger_count <int> 1, 1, 1, 1, 1, 6, 4, 1, 1, 1, 4, 2, 1, 1...
## # $ pickup_longitude <dbl> -73.98215, -73.98042, -73.97903, -74.01004, ...
## # $ pickup_latitude  <dbl> 40.76794, 40.73856, 40.76394, 40.71997, 40....
## # $ dropoff_longitude <dbl> -73.96463, -73.99948, -74.00533, -74.01227, ...
## # $ dropoff_latitude <dbl> 40.76560, 40.73115, 40.71009, 40.70672, 40....
## # $ store_and_fwd_flag <chr> "N", "N", "N", "N", "N", "N", "N", "N"...
## # $ trip_duration    <int> 455, 663, 2124, 429, 435, 443, 341, 1551, 2...
```

```
taxi <- data.frame(taxi)
```

We find below observations on Taxi Records

- Trip *Id* is unique identification of a trip
- *vendor_id* field has only 2 values “1” or “2” assuming two taxi companies

- *pickup/dropoff_datetime* - holds date and time of pickup and dropoff. we need to mutate the fields to get date and time seperately.
- *pickup/dropoff_longitude/latitude* - hold values of geographical coordinates where the meter was activate/deactivated.
- *store_and_fwd_flag* is a flag that indicates whether the trip data was sent immediately to the vendor (“N”) or held in the memory of the taxi because there was no connection to the server (“Y”). Maybe there could be a correlation with certain geographical areas with bad reception?
- *trip_duration* hold the duration in seconds and its our target prediction of ths project.

Lets check the missing values of Taxi and Test Records

```
sum(is.na(taxi))
```

```
## [1] 0
```

We have received Zero count. So we dont have any missing values in Dataset.

Reformatting data for our analysis

For our analysis, we will date fields from characters into date objects. We also change vendor_id as a factor. This makes it easier to visualise relationships that involve these features.

```
taxi <- taxi %>% mutate(
  pickup_datetime = ymd_hms(pickup_datetime),
  dropoff_datetime = ymd_hms(dropoff_datetime),
  vendor_id = factor(vendor_id),
  passenger_count = factor(passenger_count)
)
```

```
## Warning: package 'bindrcpp' was built under R version 3.4.2
```

Consistency check

This code is to check *trip_durations* are consistent with the intervals between *pickup_datetime* and *dropoff_datetime*.

Actual count of Taxi file is 1458644. Count of below check should the same else the records are inconsistent.

```
taxi %>%
  mutate(check = abs(interval(drooff_datetime, pickup_datetime)) + trip_duration) > 0) %>%
  select(check, pickup_datetime, dropoff_datetime, trip_duration) %>%
  group_by(check) %>%
  count()
```

```
## # A tibble: 1 x 2
## # Groups:   check [1]
##   check     n
##   <lgl>    <int>
## 1 FALSE 1458644
```

Feature Visualizations

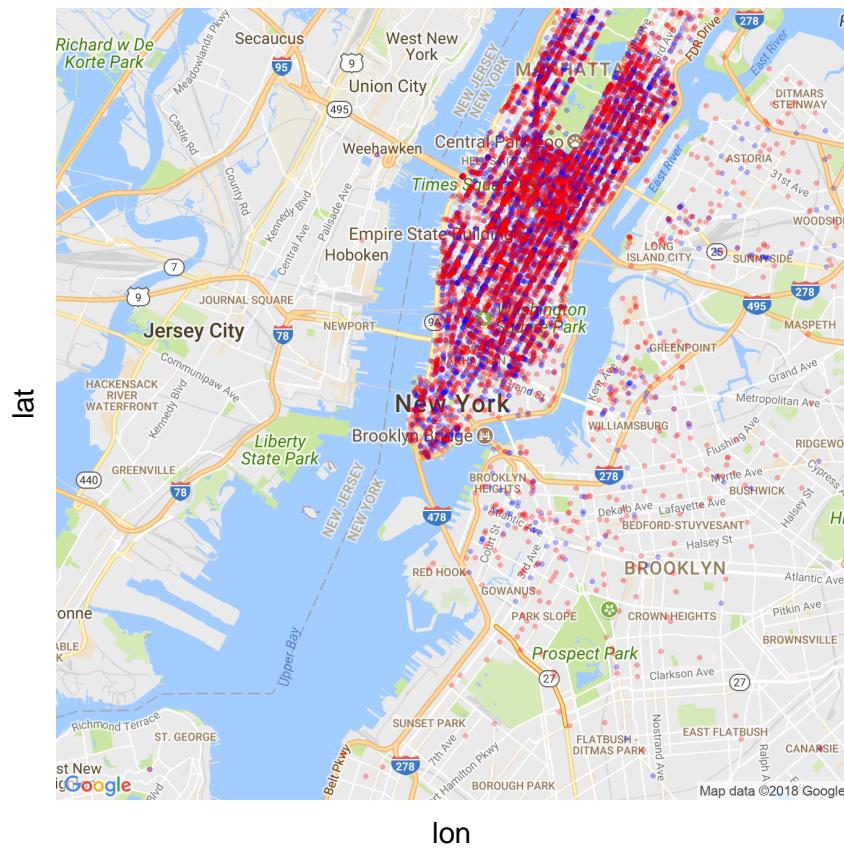
We are starting with NYC Map to check where the most of pickup and dropoff are happening. We took 5000 samples from Taxi file and plotted. It turns out that almost all of our trips were in fact taking place in Manhattan only.

```

if (.exists(".GeocodedInformation")) {
  rm(.GeocodedInformation)
}

my_map <- get_map(location = "New York City", zoom = 12, maptype = "roadmap", source = "google", color =
set.seed(1234)
tax_samp <- sample_n(taxi, 5000)
ggmap(my_map) +
  geom_point(data = tax_samp, aes(x = pickup_longitude, y = pickup_latitude), size = 0.3, alpha = 0.3,
  geom_point(data = tax_samp, aes(x = dropoff_longitude, y = dropoff_latitude), size = 0.3, alpha = 0.3,
  theme(axis.ticks = element_blank(), axis.text = element_blank())

```



Below analysis to check if any abnormal trip duration exists in our data.

```

taxi %>%
  arrange(desc(trip_duration)) %>%
  select(trip_duration, pickup_datetime, dropoff_datetime) %>%
  head(5)

##   trip_duration     pickup_datetime    dropoff_datetime
## 1      3526282 2016-02-13 22:46:52 2016-03-25 18:18:14
## 2      2227612 2016-01-05 06:14:15 2016-01-31 01:01:07
## 3      2049578 2016-02-13 22:38:00 2016-03-08 15:57:38
## 4      1939736 2016-01-05 00:19:42 2016-01-27 11:08:38
## 5      86392   2016-02-15 23:18:06 2016-02-16 23:17:58

```

```

taxi %>%
  arrange(desc(trip_duration)) %>%
  select(trip_duration, pickup_datetime, dropoff_datetime) %>%
  tail(5)

##      trip_duration     pickup_datetime     dropoff_datetime
## 1458640          1 2016-02-22 00:40:25 2016-02-22 00:40:26
## 1458641          1 2016-03-12 02:15:31 2016-03-12 02:15:32
## 1458642          1 2016-01-14 12:33:28 2016-01-14 12:33:29
## 1458643          1 2016-02-06 13:40:27 2016-02-06 13:40:28
## 1458644          1 2016-01-03 16:55:44 2016-01-03 16:55:45

```

Lets check the distributions of pickup_datetime and dropoff_datetime by year.

```

p1 <- taxi %>%
  ggplot(aes(x = pickup_datetime)) +
  geom_histogram(fill = "red", colour = "white", bins = 180) +
  scale_x_datetime(date_breaks = "1 week", date_labels = "%W%b") +
  labs(x = "Pickup dates week/month") +
  theme(axis.text.x = element_text(angle = 90))

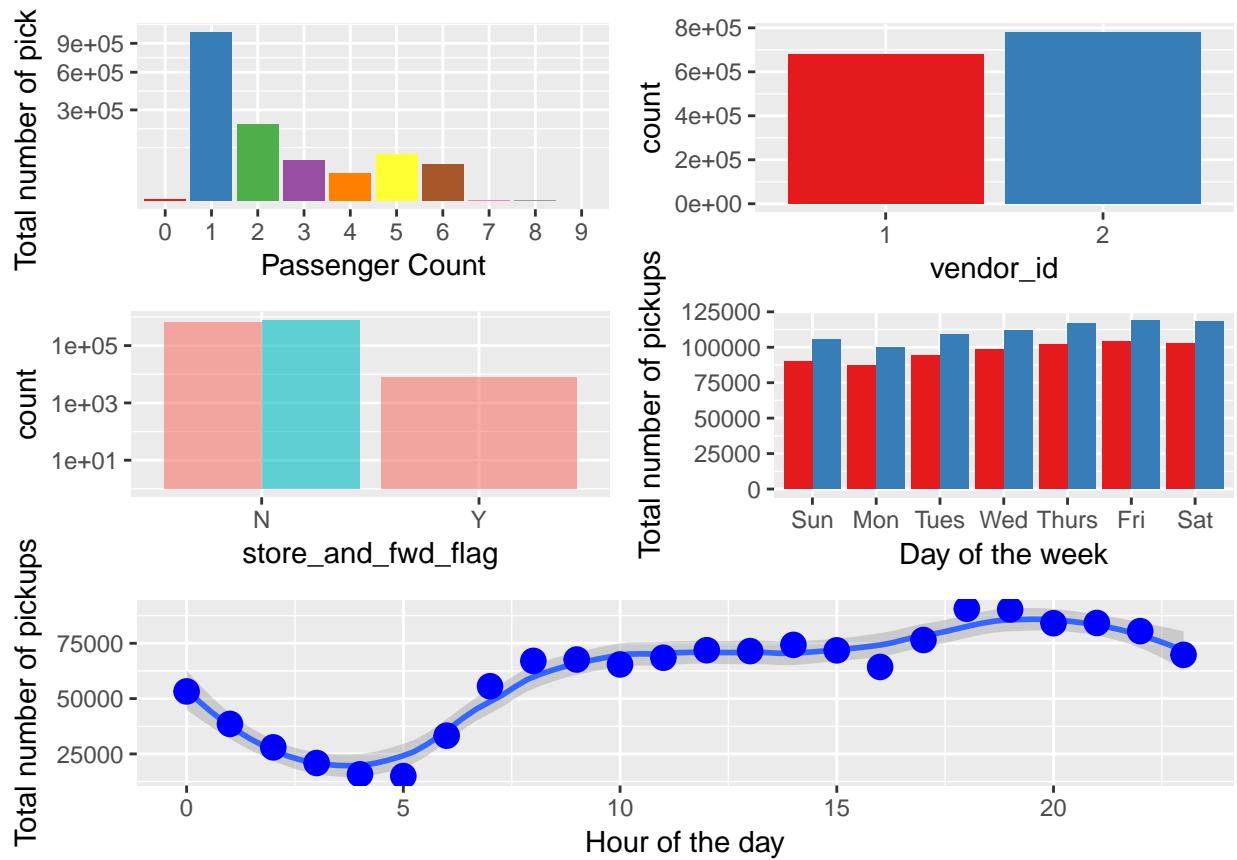
p2 <- taxi %>%
  ggplot(aes(dropoff_datetime)) +
  geom_histogram(fill = "blue", colour = "white", bins = 180) +
  scale_x_datetime(date_breaks = "1 week", date_labels = "%W%b") +
  labs(x = "Dropoff dates week/month") +
  theme(axis.text.x = element_text(angle = 90))

layout <- matrix(c(1, 2), 2, 1, byrow = FALSE)
multiplot(p1, p2, layout = layout)

```



In the below plot we are checking passenger count, vendor_id, total number of pickups on hour/day distribution.



- We found some abnormal trip with zero passenger and more 7 passengers
- We find an interesting pattern with Monday being the quietest day and Friday very busy.
- we find evening hours are busiest hours of the day.

Now lets check how the trends in different vizualization.

- We find Jan and June has less number of trips
- We find During weekends early morning are busy

```
p1 <- taxi %>%
  mutate(
    hpick = hour(pickup_datetime),
    Month = factor(month(pickup_datetime, label = TRUE))
  ) %>%
  group_by(hpick, Month) %>%
  count() %>%
  ggplot(aes(hpick, n, color = Month)) +
  geom_line(size = 1.5) +
  labs(x = "Hour of the day", y = "count")

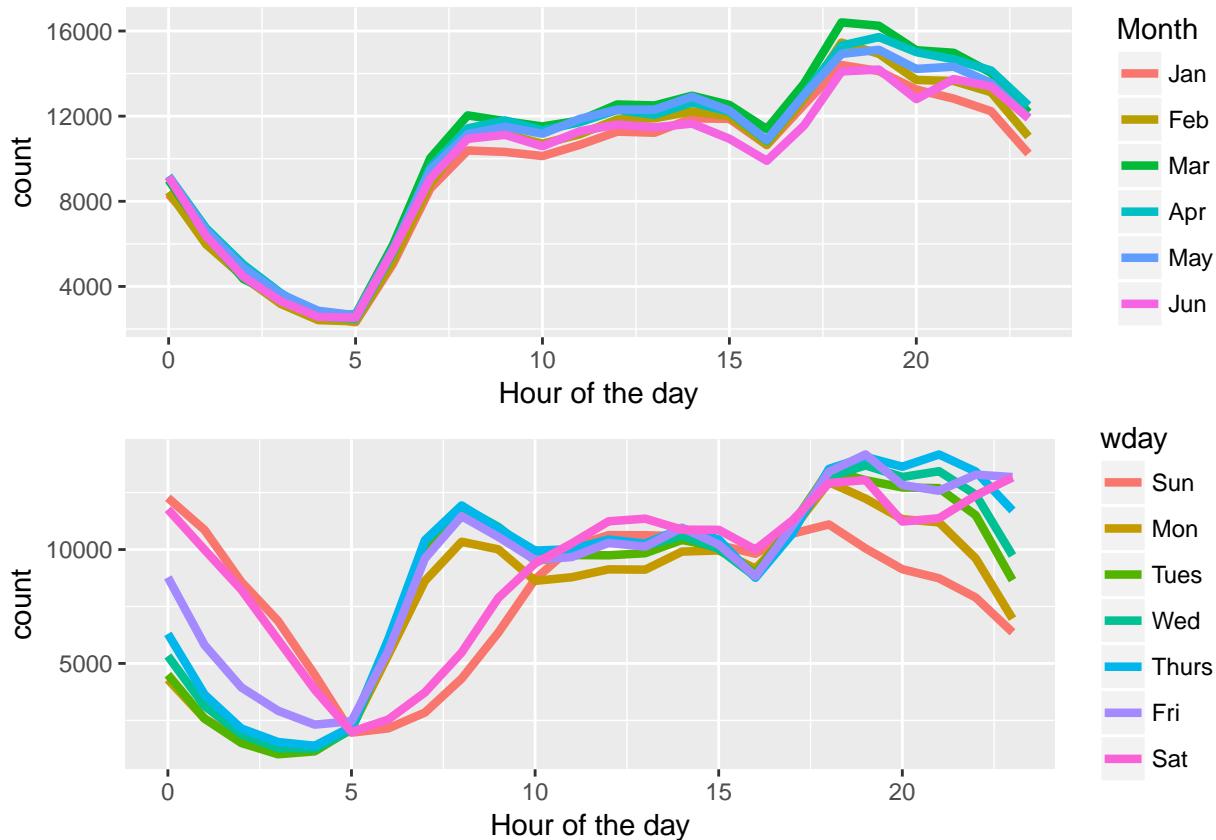
p2 <- taxi %>%
  mutate(
    hpick = hour(pickup_datetime),
    wday = factor(wday(pickup_datetime, label = TRUE))
  ) %>%
  group_by(hpick, wday) %>%
```

```

count() %>%
ggplot(aes(hpick, n, color = wday)) +
geom_line(size = 1.5) +
labs(x = "Hour of the day", y = "count")

layout <- matrix(c(1, 2), 2, 1, byrow = FALSE)
multiplot(p1, p2, layout = layout)

```



In the next slides we trying find the relation between the trip duration and picking up data & time. This will help us identify how strongly correlated to add them in the model.

```

p1 <- taxi %>%
  mutate(wday = wday(pickup_datetime, label = TRUE)) %>%
  group_by(wday, vendor_id) %>%
  summarise(median_duration = median(trip_duration) / 60) %>%
  ggplot(aes(wday, median_duration, color = vendor_id)) +
  geom_point(size = 4) +
  scale_colour_brewer(palette = "Set1") +
  labs(x = "Day of the week", y = "Median duration [min]") +
  theme(panel.background = element_rect(fill = "white"))

p2 <- taxi %>%
  mutate(pickt = hour(pickup_datetime)) %>%
  group_by(pickt, vendor_id) %>%
  summarise(median_duration = median(trip_duration) / 60) %>%

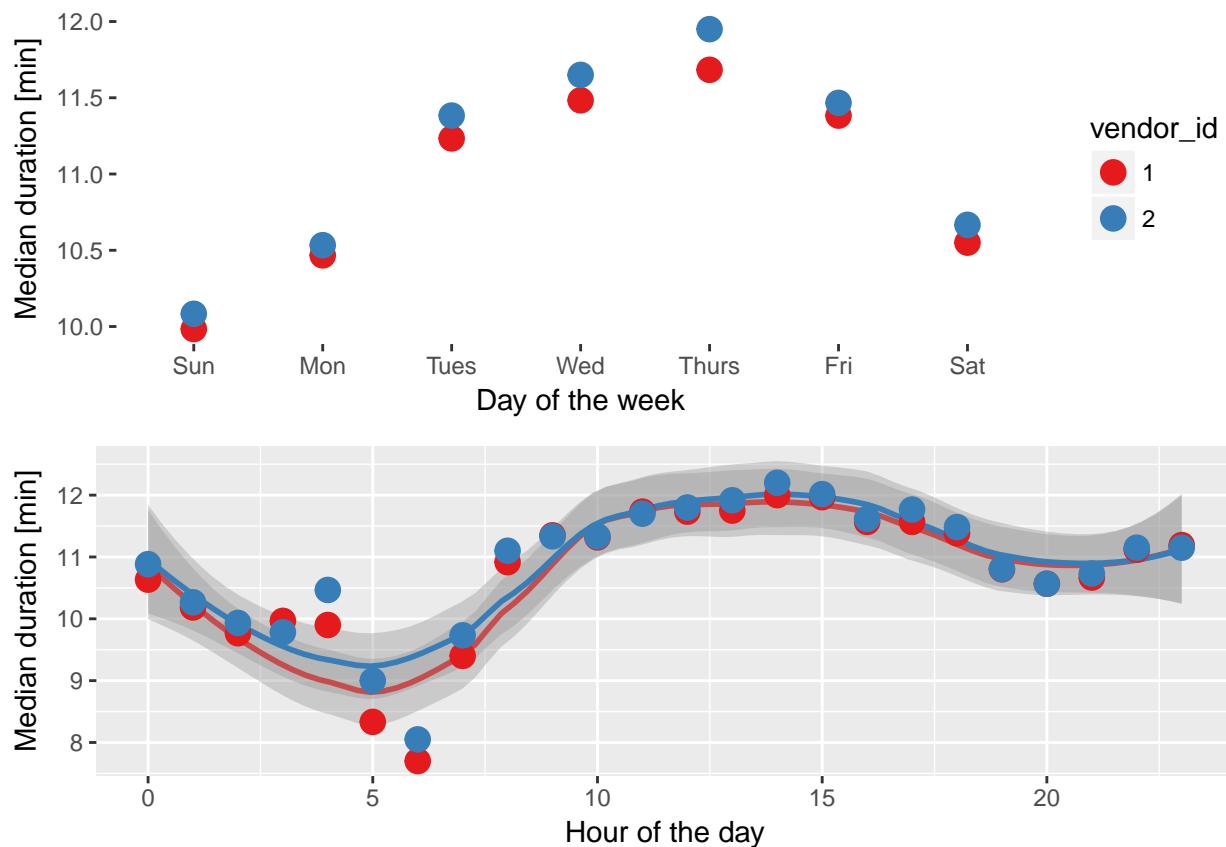
```

```

ggplot(aes(pickt, median_duration, color = vendor_id)) +
  geom_smooth(method = "loess", span = 1 / 2) +
  geom_point(size = 4) +
  scale_colour_brewer(palette = "Set1") +
  labs(x = "Hour of the day", y = "Median duration [min]") +
  theme(legend.position = "none")

layout <- matrix(c(1, 2), 2, 1, byrow = FALSE)
multiplot(p1, p2, layout = layout)

```

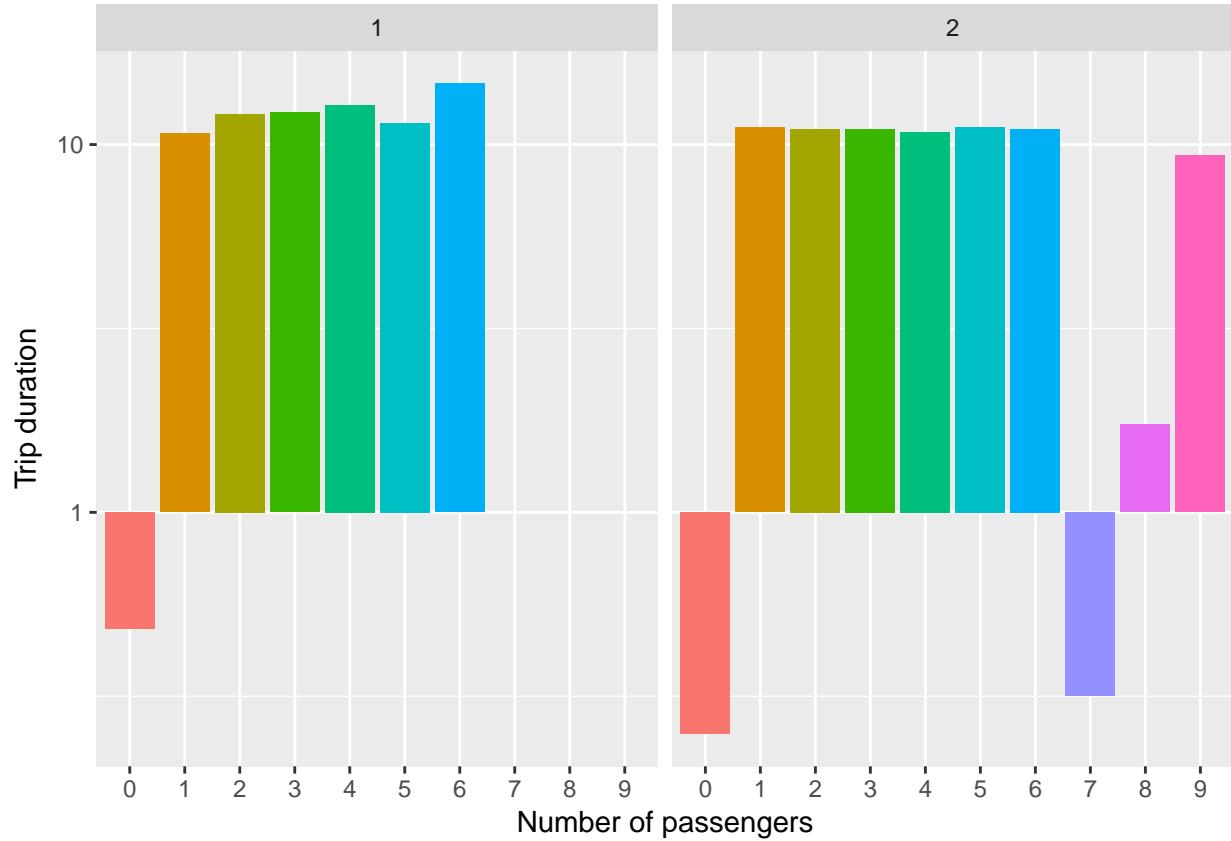


In the our next slide, we are checking for any correlation between passenger count and trip duration.

```

taxi %>%
  group_by(passenger_count, vendor_id) %>%
  summarise(median_duration = median(trip_duration) / 60) %>%
  ggplot(aes(passenger_count, median_duration, fill = passenger_count)) +
  geom_bar(stat = "identity") +
  scale_y_log10() +
  theme(legend.position = "none") +
  facet_wrap(~ vendor_id) +
  labs(y = "Trip duration", x = "Number of passengers")

```

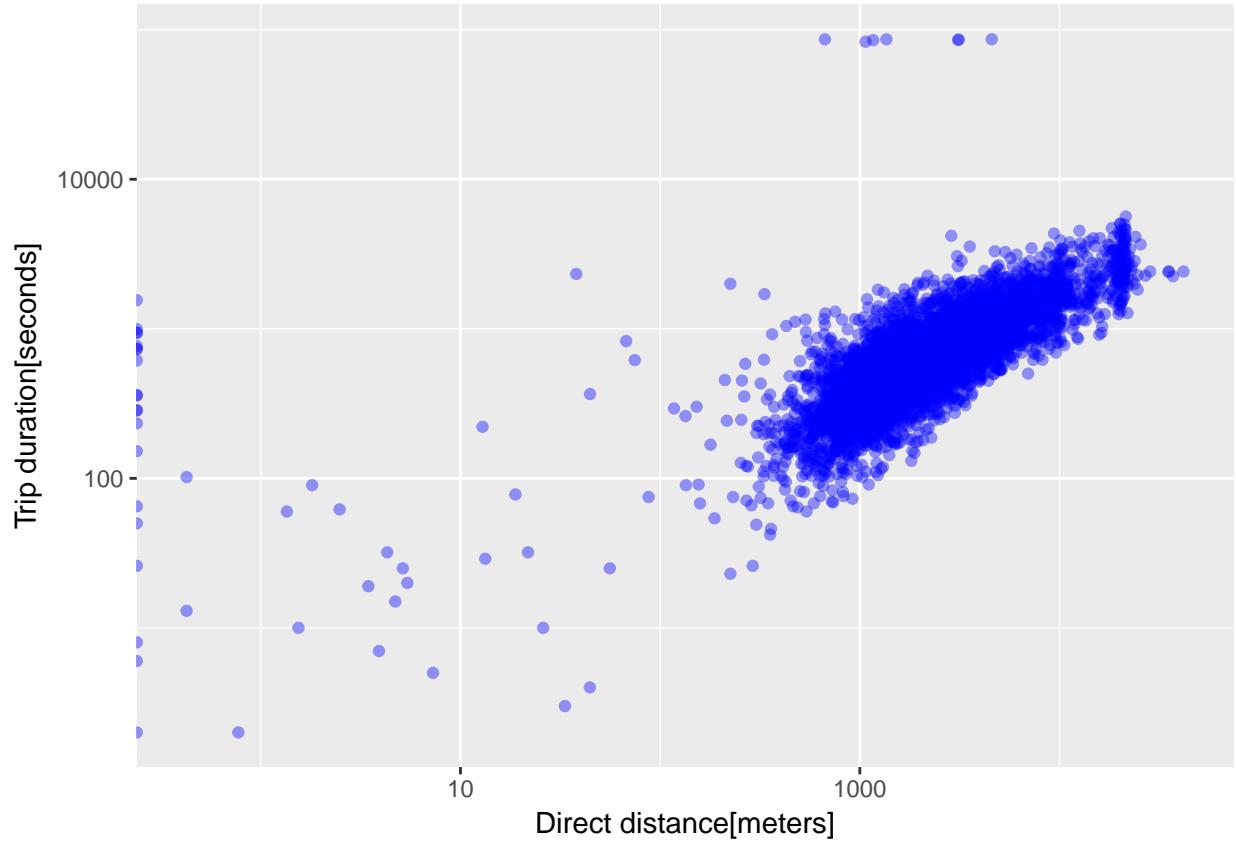


Relation between Time and Direct distance

In this section we are trying to find out the relation between trip duration and direct distance. To derive direct distance, we are using “Geosphere” package. Also, we are trying to find out any significance counts trips made out of Manhattan. Two major airports attract more taxi rides from city. We need to find how significant they are for our modelling.

Lets plot relationship trip duration and distance

```
set.seed(1234)
taxi %>%
  sample_n(5000) %>%
  ggplot(aes(dist, trip_duration)) +
  geom_point(color = "blue", alpha = 0.4) +
  scale_x_log10() +
  scale_y_log10() +
  labs(x = "Direct distance[meters]", y = "Trip duration[seconds]")
```

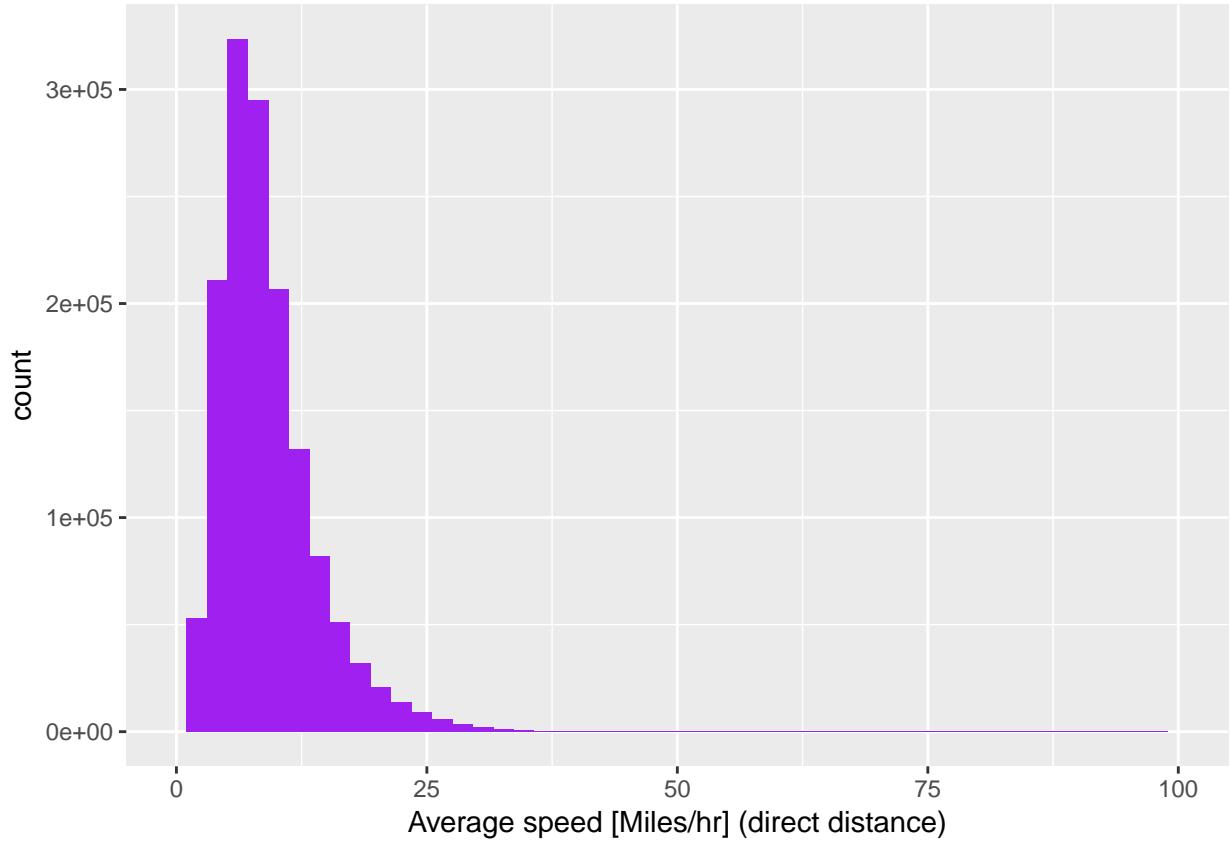


- Its clear observation, that distance increases trip duration.

Lets visualize how speed new york taxis travelling during peak hours and weekends. In order to find bogus values in the datasets, we can find extreme speed records and eliminate them

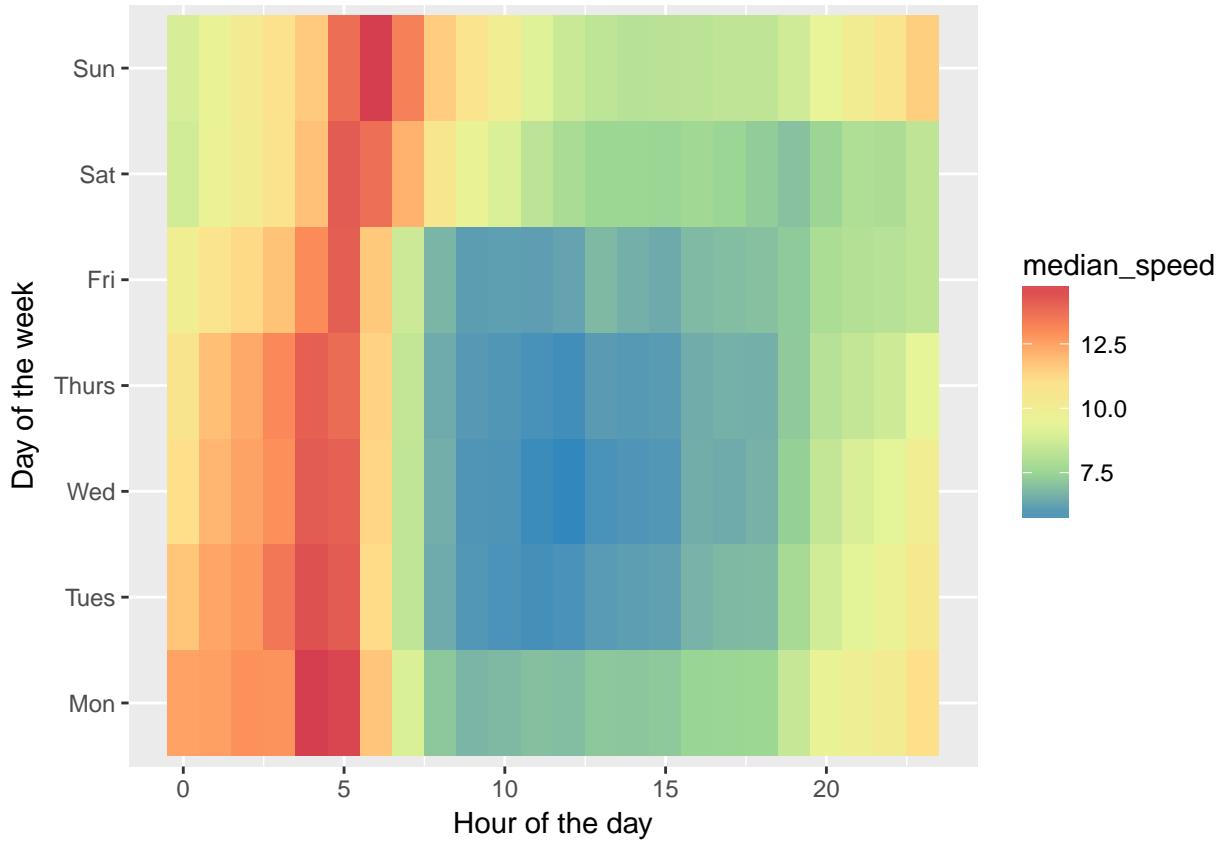
```
taxi %>%
  ggplot(aes(speed)) +
  geom_histogram(fill = "purple", bins = 50) +
  scale_x_continuous(limits = c(0, 100)) +
  labs(x = "Average speed [Miles/hr] (direct distance)")
```

```
## Warning: Removed 89 rows containing non-finite values (stat_bin).
```



Plotting speed on different times using heat map.

```
P1 <- taxi %>%
  group_by(wday, hour) %>%
  summarise(median_speed = median(speed)) %>%
  ggplot(aes(hour, wday, fill = median_speed)) +
  geom_tile() +
  labs(x = "Hour of the day", y = "Day of the week") +
  scale_fill_distiller(palette = "Spectral")
layout <- matrix(c(1, 1), 1, 1, byrow = TRUE)
multiplot(P1, layout = layout)
```



Airport trips has significant number of counts. Since airports are usually not in the city center it is reasonable to assume that the pickup/dropoff distance from the airport could be a useful predictor for longer trip_duration. We have derived a trip is JFK/La Gaurdia based on distance between pickup/dropoff location and airport coordinates. If its near to them, then we assume they are airport trips.

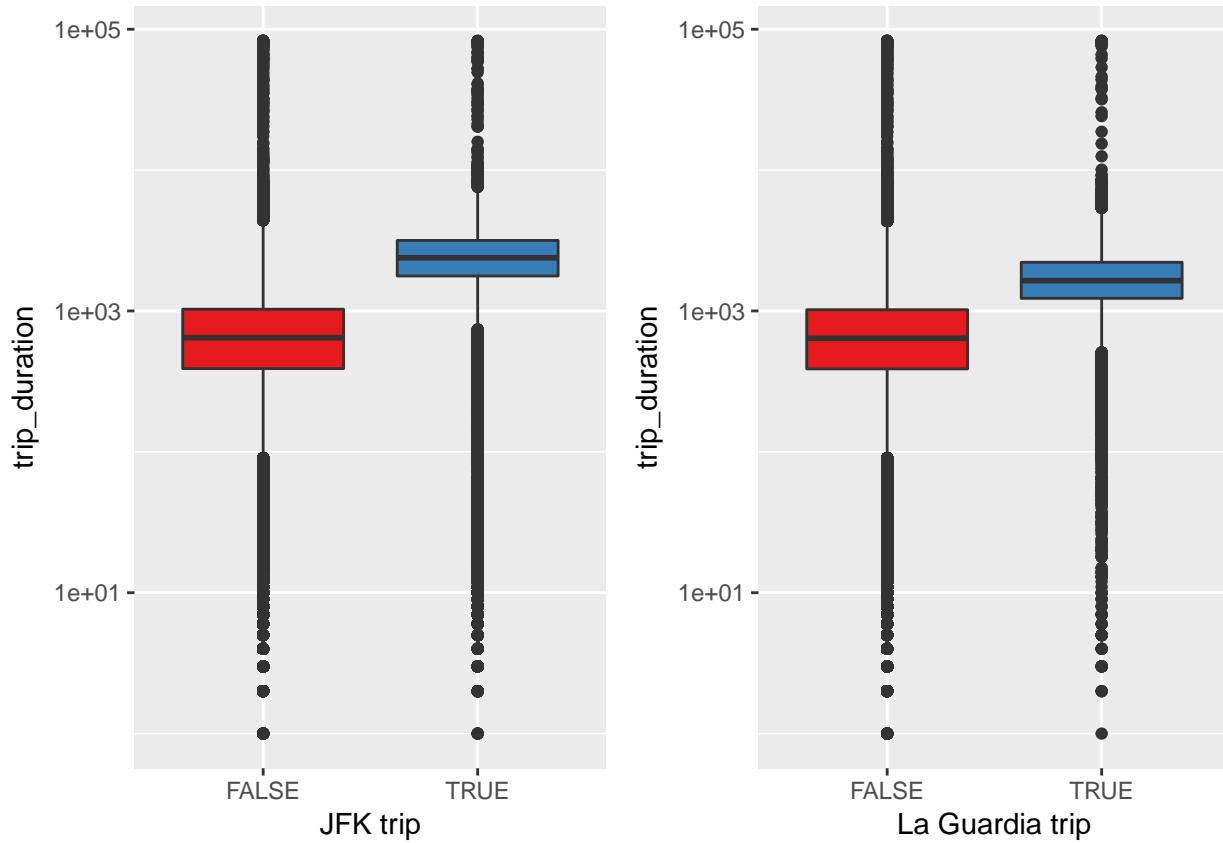
Let's visualize the how significantly trip duration increased for the airport trips. Based on below plot, its evident aiport trips are longer than regular trip to diff parts of city.

```
p1 <- taxi %>%
  filter(trip_duration < 23 * 3600) %>%
  ggplot(aes(jfk_trip, trip_duration, fill = jfk_trip)) +
  geom_boxplot() +
  scale_y_log10() +
  scale_fill_brewer(palette = "Set1") +
  theme(legend.position = "none") +
  labs(x = "JFK trip")

p2 <- taxi %>%
  filter(trip_duration < 23 * 3600) %>%
  ggplot(aes(lg_trip, trip_duration, fill = lg_trip)) +
  geom_boxplot() +
  scale_y_log10() +
  scale_fill_brewer(palette = "Set1") +
  theme(legend.position = "none") +
  labs(x = "La Guardia trip")

layout <- matrix(c(1, 2), 1, 2, byrow = FALSE)
```

```
multiplot(p1, p2, layout = layout)
```



Data Cleaning.

The aim here is to remove trips that have improbable features, such as extreme trip durations or very low average speed.

** Filter trips more than 24 hours.

```
long_hrtrips <- taxi %>%
  filter(trip_duration > 24 * 3600)

long_hrtrips %>%
  arrange(desc(dist)) %>%
  select(pickup_datetime, dropoff_datetime, speed) %>%
  head(05)

##      pickup_datetime    dropoff_datetime      speed
## 1 2016-01-05 00:19:42 2016-01-27 11:08:38 0.023252072
## 2 2016-02-13 22:46:52 2016-03-25 18:18:14 0.012633059
## 3 2016-02-13 22:38:00 2016-03-08 15:57:38 0.006533943
## 4 2016-01-05 06:14:15 2016-01-31 01:01:07 0.001643123
```

** Filter trips shorter than a few minutes

```
min_trips <- taxi %>%
  filter(trip_duration < 5 * 60)
```

```

min_trips %>%
  arrange(dist) %>%
  select(dist, pickup_datetime, dropoff_datetime, speed) %>%
  head(05)

##   dist      pickup_datetime    dropoff_datetime speed
## 1 0 2016-02-29 18:39:12 2016-02-29 18:42:59     0
## 2 0 2016-01-27 22:29:31 2016-01-27 22:29:58     0
## 3 0 2016-01-22 16:13:01 2016-01-22 16:13:20     0
## 4 0 2016-01-18 15:24:43 2016-01-18 15:28:57     0
## 5 0 2016-05-04 22:28:43 2016-05-04 22:32:51     0

** Find trip with zero miles

zero_dist <- taxi %>%
  filter(near(dist, 0))
nrow(zero_dist)

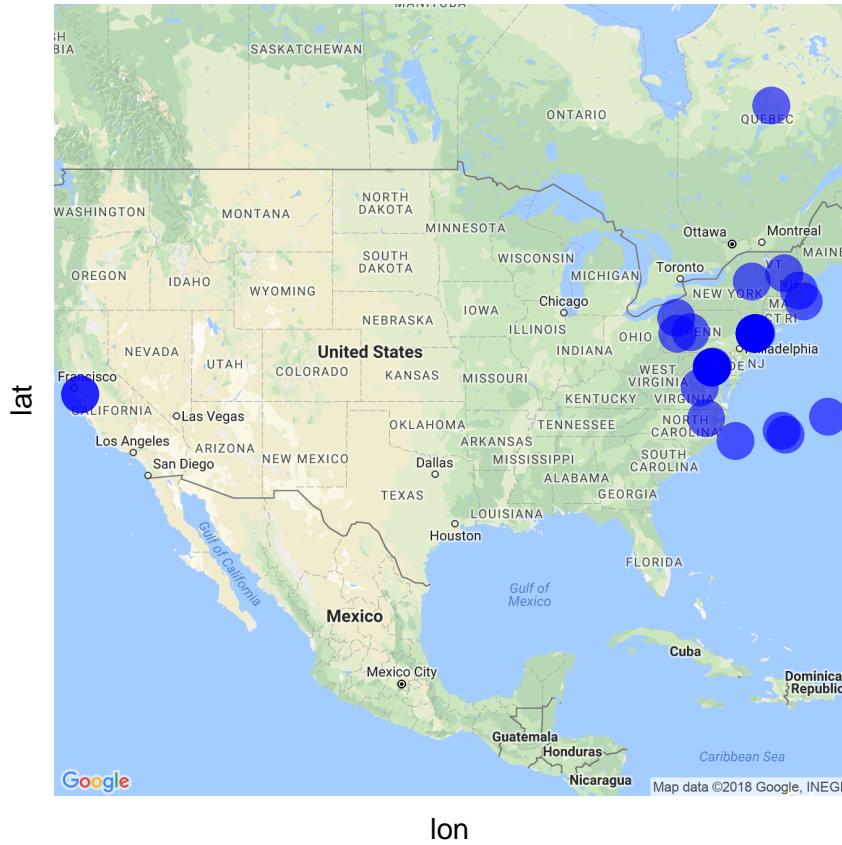
## [1] 5897

** Find trips away from New York

long_dist <- taxi %>%
  filter((jfk_dist_pick > 3e5) | (jfk_dist_drop > 3e5))
long_dist <- data.frame(long_dist)
long_dist_coord <- long_dist %>%
  select(lon = pickup_longitude, lat = pickup_latitude)

my_map1 <- get_map(location = "United States", zoom = 4, maptype = "roadmap", source = "google", color =
gmap(my_map1) +
  geom_point(data = long_dist, aes(x = pickup_longitude, y = pickup_latitude), size = 6, alpha = 0.6, c
  theme(axis.ticks = element_blank(), axis.text = element_blank())

```



Filter all bogus data from taxi dataset

```
taxi <- taxi %>%
  filter(
    trip_duration < 22 * 3600,
    dist > 0 | (near(dist, 0) & trip_duration < 60),
    jfk_dist_pick < 3e5 & jfk_dist_drop < 3e5,
    trip_duration > 10,
    speed < 100
  )
```

External data

We are adding Open Source Routing Machine, OSRM data for each trip. This data is provided by oscarleo and we can download data from here and includes the pickup/dropoff streets and total distance/duration between these two points together with a sequence of travel steps such as turns or entering a highway.

```
fast1 <- read_csv("fastest_train_part_1.csv")
```

```
## Parsed with column specification:
## cols(
##   id = col_character(),
##   starting_street = col_character(),
##   end_street = col_character(),
##   total_distance = col_double(),
##   total_travel_time = col_double(),
```

```

##   number_of_steps = col_integer(),
##   street_for_each_step = col_character(),
##   distance_per_step = col_character(),
##   travel_time_per_step = col_character(),
##   step_maneuvers = col_character(),
##   step_direction = col_character(),
##   step_location_list = col_character()
## )

fast2 <- read_csv("fastest_train_part_2.csv")

## Parsed with column specification:
## cols(
##   id = col_character(),
##   starting_street = col_character(),
##   end_street = col_character(),
##   total_distance = col_double(),
##   total_travel_time = col_double(),
##   number_of_steps = col_integer(),
##   street_for_each_step = col_character(),
##   distance_per_step = col_character(),
##   travel_time_per_step = col_character(),
##   step_maneuvers = col_character(),
##   step_direction = col_character(),
##   step_location_list = col_character()
## )

fast1 <- data.frame(fast1)
fast2 <- data.frame(fast2)

fastest_route <- bind_rows(fast1, fast2)
glimpse(fastest_route)

## # Observations: 1,458,643
## # Variables: 12
## $ id                  <chr> "id2875421", "id2377394", "id3504673", "i...
## $ starting_street     <chr> "Columbus Circle", "2nd Avenue", "Greenwi...
## $ end_street          <chr> "East 65th Street", "Washington Square We...
## $ total_distance       <dbl> 2009.1, 2513.2, 1779.4, 1614.9, 1393.5, 1...
## $ total_travel_time    <dbl> 164.9, 332.0, 235.8, 140.1, 189.4, 138.8, ...
## $ number_of_steps      <int> 5, 6, 4, 5, 5, 5, 2, 13, 6, 9, 4, 4, 10, ...
## $ street_for_each_step <chr> "Columbus Circle|Central Park West|65th S...
## $ distance_per_step    <chr> "0|576.4|885.6|547.1|0", "877.3|836.5|496...
## $ travel_time_per_step <chr> "0|61.1|60.1|43.7|0", "111.7|109|69.9|25....
## $ step_maneuvers       <chr> "depart|rotary|turn|new name|arrive", "de...
## $ step_direction        <chr> "left|straight|right|straight|arrive", "n...
## $ step_location_list    <chr> "-73.982316,40.767869|-73.981997,40.76768...

```

Lets visualize few of the important data points in Histogram to see the distribution.

```

p1 <- fastest_route %>%
  ggplot(aes(total_travel_time)) +
  geom_histogram(bins = 150, fill = "red") +
  scale_x_log10() +
  scale_y_sqrt()

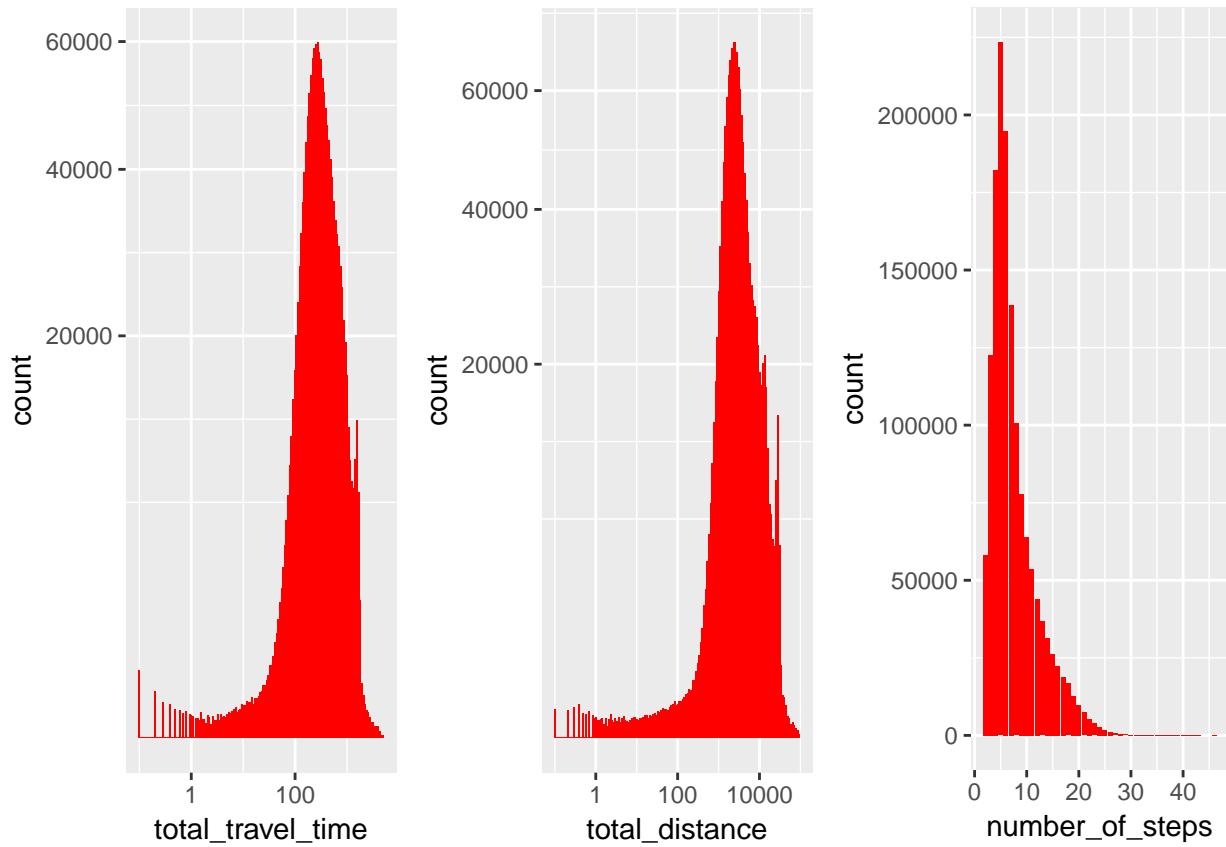
```

```

p2 <- fastest_route %>%
  ggplot(aes(total_distance)) +
  geom_histogram(bins = 150, fill = "red") +
  scale_x_log10() +
  scale_y_sqrt()

p3 <- fastest_route %>%
  ggplot(aes(number_of_steps)) +
  geom_bar(fill = "red")
layout <- matrix(c(1, 2, 3), 1, 3, byrow = TRUE)
multiplot(p1, p2, p3, layout = layout)

```



Joining OSRM file and TAXI file for further analysis.

In a first look at the joined data, we will mainly focus on the total_distance and total_travel_time, when combined with our original training data set:

```

osrm <- fastest_route %>%
  select(
    id, total_distance, total_travel_time, number_of_steps,
    step_direction, step_maneuvers
  ) %>%
  mutate(
    fastest_speed = total_distance / total_travel_time * 2.236,
    left_turns = str_count(step_direction, "left"),
    right_turns = str_count(step_direction, "right"),
  )

```

```

    turns = str_count(step_maneuvers, "turn")
) %>%
select(-step_direction, -step_maneuvers)

taxi <- left_join(taxi, osrm, by = "id")

```

In order to consider fastest route file for modeling, we need to find out total_travel_time in OSRM vs drip_duration recorded in taxi file.

```

p1 <- taxi %>%
  ggplot(aes(trip_duration)) +
  geom_density(fill = "red", alpha = 0.5) +
  geom_density(aes(total_travel_time), fill = "blue", alpha = 0.5) +
  scale_x_log10() +
  coord_cartesian(xlim = c(5e1, 8e3))

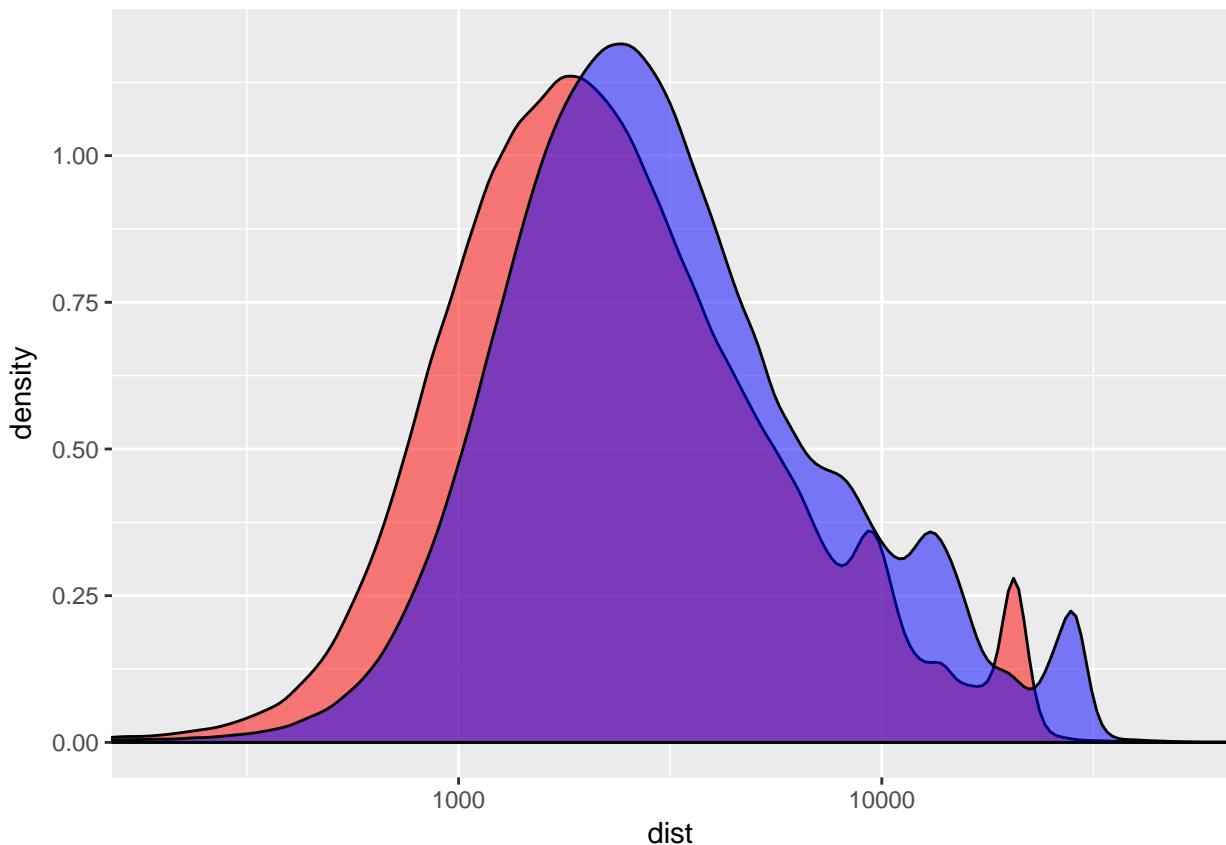
```

Check for trip distance aswell.

```

taxi %>%
  ggplot(aes(dist)) +
  geom_density(fill = "red", alpha = 0.5) +
  geom_density(aes(total_distance), fill = "blue", alpha = 0.5) +
  scale_x_log10() +
  coord_cartesian(xlim = c(2e2, 5e4))

```

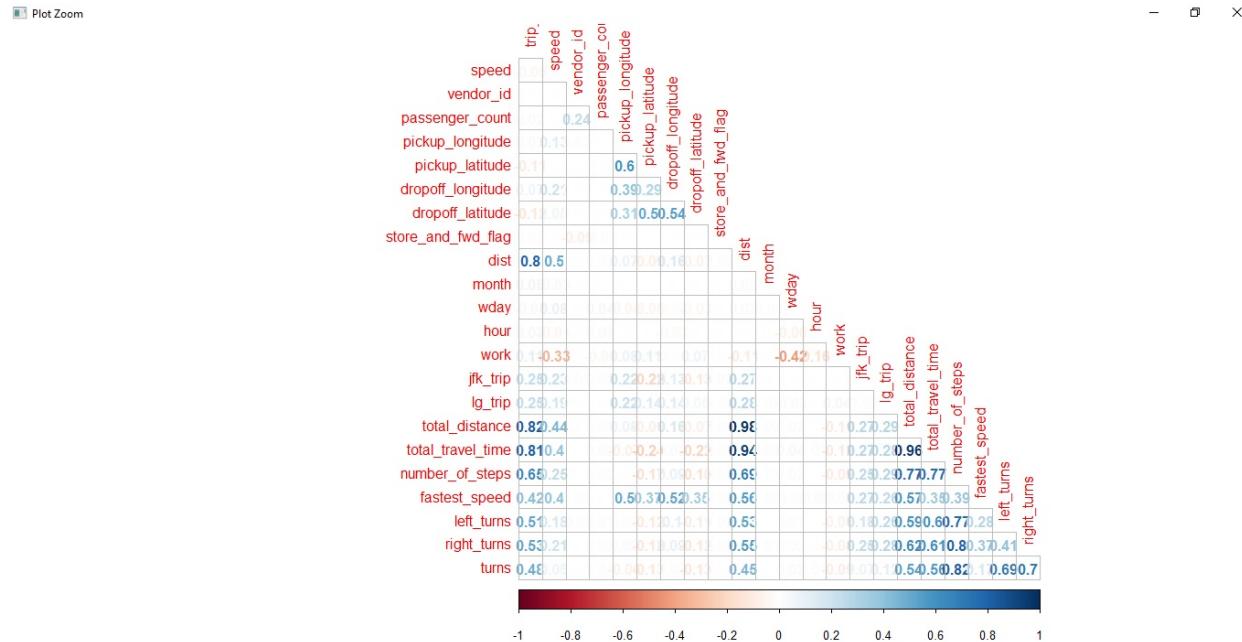


Correlation Matrix

Lets vizualize the relations between our variables using correlation matrix. We are using corrplot package for plotting.

NOTE: To resolve memory issue, we have run the corrplot outside markdown and attached screenshot below.

```
# taxi %>%
#   select(
#     -id, -pickup_datetime, -dropoff_datetime, -jfk_dist_pick,
#     -jfk_dist_drop, -lg_dist_pick, -lg_dist_drop, -date
#   ) %>%
#   mutate(
#     passenger_count = as.integer(passenger_count),
#     vendor_id = as.integer(vendor_id),
#     store_and_fwd_flag = as.integer(as.factor(store_and_fwd_flag)),
#     jfk_trip = as.integer(jfk_trip),
#     wday = as.integer(wday),
#     month = as.integer(month),
#     work = as.integer(work),
#     lg_trip = as.integer(lg_trip)
#   ) %>%
#   select(trip_duration, speed, everything()) %>%
#   cor(use = "complete.obs", method = "spearman") %>%
#   corrplot(type = "lower", method = "number", diag = FALSE)
```



- The strongest correlations with the trip.duration are seen for the direct distance as well as the total_distance and total_travel_time derived from the OSRM data.
- Number of turns, presumably mostly via the number_of_steps, are having an impact on the trip.duration.
- Another effect on the trip.duration can be seen for our engineered features jfk_trip and lg_trip indicating journeys to either airport.
- Features like store_and_fwd_flag, vendor_id, passenger_count, speed are having little to no correlation, So we are removing from our model.

Model Selection and Prediction

Next is our most important step to find best method to predict the duration. For modelling we are taking a subset of 400K trips and run our model on the subset data. In this report, we discuss three models: Linear Regression, Random Forest, Gradient Boosting. We evaluate these models based on the Root Mean Square Logarithmic Error (RMSLE) and Pseudo R². We also discuss the importance of various features in our prediction algorithms.

To assess the mode performance we will run a cross-validation step and also split our training subset data into a *train*, *validation*, *test* data set. Thereby, the model performance can be evaluated on a sample that the algorithm has not seen. We split our data into 70/15/15 fractions using a tool from the caret package.

```
set.seed(4321)
trainIndex <- createDataPartition(taxi1$trip_duration, p = 0.7, list = FALSE, times = 1)

train <- taxi1[trainIndex, ]
valid <- taxi1[-trainIndex, ]

trainIndex <- createDataPartition(valid$trip_duration, p = 0.5, list = FALSE, times = 1)
valid <- valid[trainIndex, ]
test <- valid[-trainIndex, ]
```

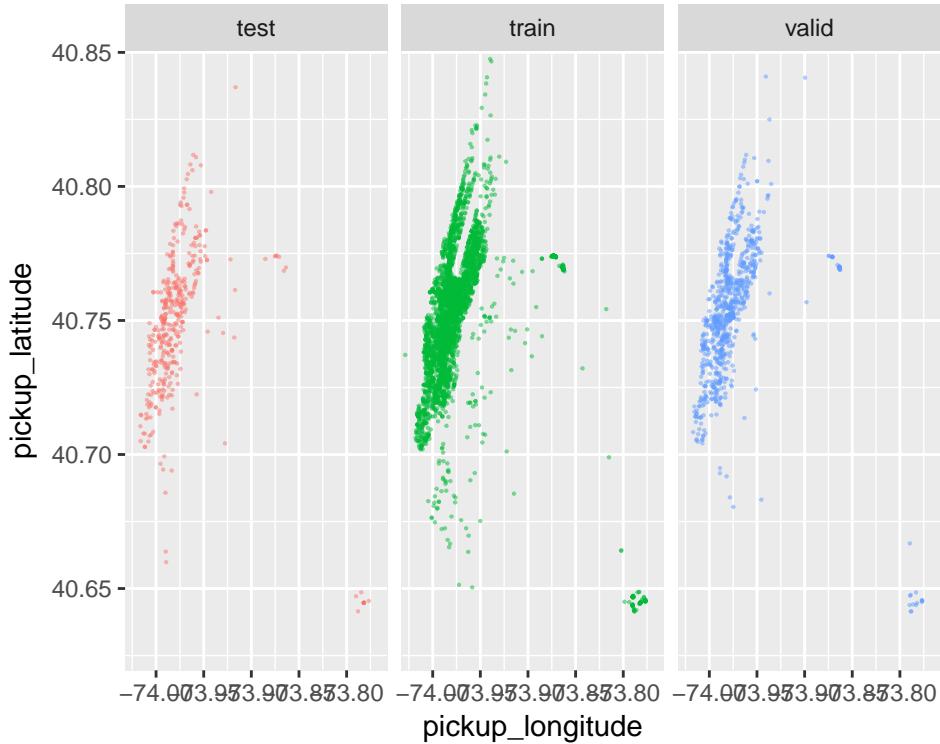
We will start with XGBoost model now. Inorder to make sure we are training a feature which is relevant to test data. Plotting Test and Training Data overlap to see they are relevant.

```
combine <- bind_rows(
  train %>% mutate(dset = "train"),
  test %>% mutate(dset = "test"),
  valid %>% mutate(dset = "valid")
)

combine <- combine %>% mutate(dset = factor(dset))

pick_good <- combine %>%
  filter(pickup_longitude > -75 & pickup_longitude < -73) %>%
  filter(pickup_latitude > 40 & pickup_latitude < 42)
pick_good <- sample_n(pick_good, 5000)

pick_good %>%
  ggplot(aes(pickup_longitude, pickup_latitude, color = dset)) +
  geom_point(size = 0.1, alpha = 0.5) +
  coord_cartesian(xlim = c(-74.02, -73.77), ylim = c(40.63, 40.84)) +
  facet_wrap(~ dset) +
  guides(color = guide_legend(override.aes = list(alpha = 1, size = 4))) +
  theme(legend.position = "none")
```



*** Data formatting

Here we will format the selected features to turn them into integer columns, since many classifiers cannot deal with categorical values. We have formatted Taxi data already. However, we need to redo same formating with combined records which intially joined both train and test data.

```
# airport coordinates again
jfk_coord <- tibble(lon = -73.778889, lat = 40.639722)
la_guardia_coord <- tibble(lon = -73.872611, lat = 40.77725)

# derive distances
pick_coord <- combine %>%
  select(pickup_longitude, pickup_latitude)
drop_coord <- combine %>%
  select(dropoff_longitude, dropoff_latitude)
combine$dist <- distCosine(pick_coord, drop_coord)
combine$bearing <- bearing(pick_coord, drop_coord)

combine$jfk_dist_pick <- distCosine(pick_coord, jfk_coord)
combine$jfk_dist_drop <- distCosine(drop_coord, jfk_coord)
combine$lg_dist_pick <- distCosine(pick_coord, la_guardia_coord)
combine$lg_dist_drop <- distCosine(drop_coord, la_guardia_coord)

combine <- combine %>%
  mutate(
    pickup_datetime = ymd_hms(pickup_datetime),
    dropoff_datetime = ymd_hms(dropoff_datetime),
    date = date(pickup_datetime)
  )
# join OSRM data with our data.
```

```
combine <- left_join(combine, osrm, by = "id")
```

Reformat data in to integer format.

```
combine <- combine %>%
  mutate(
    store_and_fwd_flag = as.integer(factor(store_and_fwd_flag)),
    vendor_id = as.integer(vendor_id),
    month = as.integer(month(pickup_datetime)),
    wday = wday(pickup_datetime, label = TRUE),
    wday = as.integer(fct_relevel(wday, c("Sun", "Sat", "Mon", "Tues", "Wed", "Thurs", "Fri"))),
    hour = hour(pickup_datetime),
    work = as.integer((hour %in% seq(8, 18)) & (wday %in% c("Mon", "Tues", "Fri", "Wed", "Thurs"))),
    jfk_trip = as.integer((jfk_dist_pick < 2e3) | (jfk_dist_drop < 2e3)),
    lg_trip = as.integer((lg_dist_pick < 2e3) | (lg_dist_drop < 2e3))
  )
```

Just to make sure our data is in proper format.

```
glimpse(combine)
```

```
## Observations: 369,955
## Variables: 32
## $ id              <chr> "id3963323", "id3456697", "id1742665", "id3...
## $ vendor_id       <int> 1, 2, 2, 1, 1, 2, 2, 2, 2, 2, 1, 2, 2...
## $ pickup_datetime <dttm> 2016-05-10 18:48:50, 2016-05-23 06:26:11, ...
## $ dropoff_datetime <dttm> 2016-05-10 18:52:13, 2016-05-23 06:37:39, ...
## $ passenger_count <int> 1, 1, 1, 1, 1, 6, 1, 1, 2, 1, 1, 6, 6...
## $ pickup_longitude <dbl> -73.95911, -73.99098, -73.96302, -73.99679, ...
## $ pickup_latitude  <dbl> 40.77194, 40.69990, 40.77079, 40.73747, 40....
## $ dropoff_longitude <dbl> -73.94935, -73.97195, -73.96848, -73.95267, ...
## $ dropoff_latitude <dbl> 40.77263, 40.74971, 40.76645, 40.77478, 40....
## $ store_and_fwd_flag <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
## $ trip_duration    <int> 203, 688, 435, 1493, 232, 400, 1983, 399, 1...
## $ dset             <fctr> train, train, train, train, train, train, ...
## $ dist              <dbl> 826.8169, 5772.1227, 666.8075, 5576.1335, 9...
## $ bearing           <dbl> 84.71270, 16.20674, -136.28083, 41.94746, 1...
## $ jfk_dist_pick    <dbl> 21164.74, 19119.76, 21314.66, 21370.35, 219...
## $ jfk_dist_drop    <dbl> 20635.22, 20381.57, 21330.58, 21002.03, 210...
## $ lg_dist_pick     <dbl> 7316.155, 13184.157, 7655.511, 11368.855, 7...
## $ lg_dist_drop     <dbl> 6489.398, 8919.056, 8170.702, 6754.238, 645...
## $ date              <date> 2016-05-10, 2016-05-23, 2016-02-14, 2016-0...
## $ total_distance   <dbl> 1569.4, 9400.6, 1094.2, 6694.7, 1154.2, 236...
## $ total_travel_time <dbl> 146.0, 688.6, 125.3, 663.3, 163.3, 218.8, 1...
## $ number_of_steps   <int> 5, 13, 6, 6, 4, 6, 19, 3, 12, 3, 5, 15, 4, ...
## $ fastest_speed     <dbl> 24.03547, 30.52533, 19.52619, 22.56799, 15....
## $ left_turns        <int> 3, 3, 1, 1, 0, 7, 2, 4, 0, 2, 2, 1, 2, 6...
## $ right_turns       <int> 1, 7, 3, 3, 2, 3, 8, 0, 3, 1, 2, 10, 2, 2, ...
## $ turns              <int> 3, 4, 4, 3, 2, 2, 11, 1, 4, 1, 3, 4, 2, 4, ...
## $ month              <int> 5, 5, 2, 1, 3, 2, 6, 5, 1, 5, 4, 5, 5, 1, 5...
## $ wday               <int> 4, 3, 1, 5, 7, 2, 2, 3, 3, 2, 1, 4, 6, 6...
## $ hour               <int> 18, 6, 12, 22, 22, 8, 4, 1, 21, 12, 1, 12, ...
## $ work               <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
## $ jfk_trip            <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
## $ lg_trip              <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
```

Next, we are selecting Feature to be in the model. I have run the feature selection outside this document, to select optimal features which gives less RMSLE values. Based on that below are the selected Features will be used across all algorithms. Our evaluation metric is Root Mean Squared Logarithmic Error. In order to easily simulate the evaluation metric in our model fitting we replace the trip_duration with its logarithm. We adding +1 to duration to avoid infinity value incase of zeros.

```
# convert trip_duration to log

combine <- combine %>%
  mutate(trip_duration = log(trip_duration + 1))

# predictor variables
train_cols <- c(
  "total_travel_time", "total_distance", "hour", "dist",
  "vendor_id", "jfk_trip", "lg_trip", "wday", "month",
  "pickup_longitude", "pickup_latitude", "lg_dist_drop", "jfk_dist_drop"
)
# target variable
y_col <- c("trip_duration")
# auxilliary variable
aux_cols <- c("dset")
# cleaning variable
clean_cols <- c("jfk_dist_pick")
# ----

# all relevant columns for train/test
cols <- c(train_cols, y_col, aux_cols, clean_cols)
combine <- combine %>%
  select_.dots = cols

# split train/test
train <- combine %>%
  filter(dset == "train") %>%
  select_.dots = str_c("-", c(aux_cols)))
test <- combine %>%
  filter(dset == "test") %>%
  select_.dots = str_c("-", c(aux_cols)))

valid <- combine %>%
  filter(dset == "valid") %>%
  select_.dots = str_c("-", c(aux_cols)))
```

Next we clean the date before we start processng the date. We will do celaning only on Training dataset inorder to identify overfitting if we see any variance with Validation dataset.

```
valid <- valid %>%
  select_.dots = str_c("-", c(clean_cols))

train <- train %>%
  filter(
    trip_duration < 24 * 3600,
    jfk_dist_pick < 3e5 & jfk_dist_drop < 3e5
  ) %>%
  select_.dots = str_c("-", c(clean_cols)))
```

XGBoost parameters and fitting

In order for *XGBoost* to properly ingest our data samples we need to re-format them slightly:

```
# convert to XGB matrix
temptrn <- train %>% select(-trip_duration)
tempval <- valid %>% select(-trip_duration)

dtrain <- xgb.DMatrix(as.matrix(temptrn), label = train$trip_duration)
dvalid <- xgb.DMatrix(as.matrix(tempval), label = valid$trip_duration)
dtest <- xgb.DMatrix(as.matrix(test))
```

Now we define the meta-parameters that govern how *XGBoost* operates. See here for more details. The *watchlist* parameter tells the algorithm to keep an eye on both the *training* and *validation* sample metrics.

```
xgb_params <- list(
  colsample_bytree = 0.7, # variables per tree
  subsample = 0.7, # data subset per tree
  booster = "gbtree",
  max_depth = 5, # tree levels
  eta = 0.3, # shrinkage
  eval_metric = "rmse",
  objective = "reg:linear",
  seed = 4321
)

watchlist <- list(train = dtrain, valid = dvalid)
```

And here we *train* our classifier, i.e. fit it to the *training* data. To ensure reproducability we set an R seed here.

NOTE: Recuing number of rounds to 60 due to memory and cpu issues. Actual submission model was trined with 100 rounds.

```
set.seed(4321)
gb_dt <- xgb.train(
  params = xgb_params,
  data = dtrain,
  print_every_n = 10,
  watchlist = watchlist,
  nrounds = 60
)

## [1]  train-rmse:4.228354 valid-rmse:4.230548
## [11] train-rmse:0.446516 valid-rmse:0.451522
## [21] train-rmse:0.414842 valid-rmse:0.422537
## [31] train-rmse:0.408781 valid-rmse:0.419076
## [41] train-rmse:0.403583 valid-rmse:0.416451
## [51] train-rmse:0.399969 valid-rmse:0.414858
## [60] train-rmse:0.396142 valid-rmse:0.413225
```

After the fitting we are running a 5-fold cross-validation (CV) to estimate our model's performance.

NOTE: Reducing number of rounds to 60 due to memory and cpu issue.

```
xgb_cv <- xgb.cv(xgb_params,
                    dtrain,
                    early_stopping_rounds = 10,
```

```

    print_every_n = 10,
    nfold = 5,
    nrounds = 60)

## [1] train-rmse:4.227942+0.000540    test-rmse:4.227885+0.001478
## Multiple eval metrics are present. Will use test_rmse for early stopping.
## Will train until test_rmse hasn't improved in 10 rounds.
##
## [11] train-rmse:0.446897+0.001701    test-rmse:0.449664+0.006800
## [21] train-rmse:0.414190+0.002299    test-rmse:0.420438+0.006133
## [31] train-rmse:0.406372+0.001270    test-rmse:0.415853+0.007264
## [41] train-rmse:0.400154+0.001346    test-rmse:0.412874+0.006854
## [51] train-rmse:0.395460+0.001353    test-rmse:0.411433+0.006853
## [60] train-rmse:0.391992+0.001228    test-rmse:0.410073+0.006527

```

Feature importance

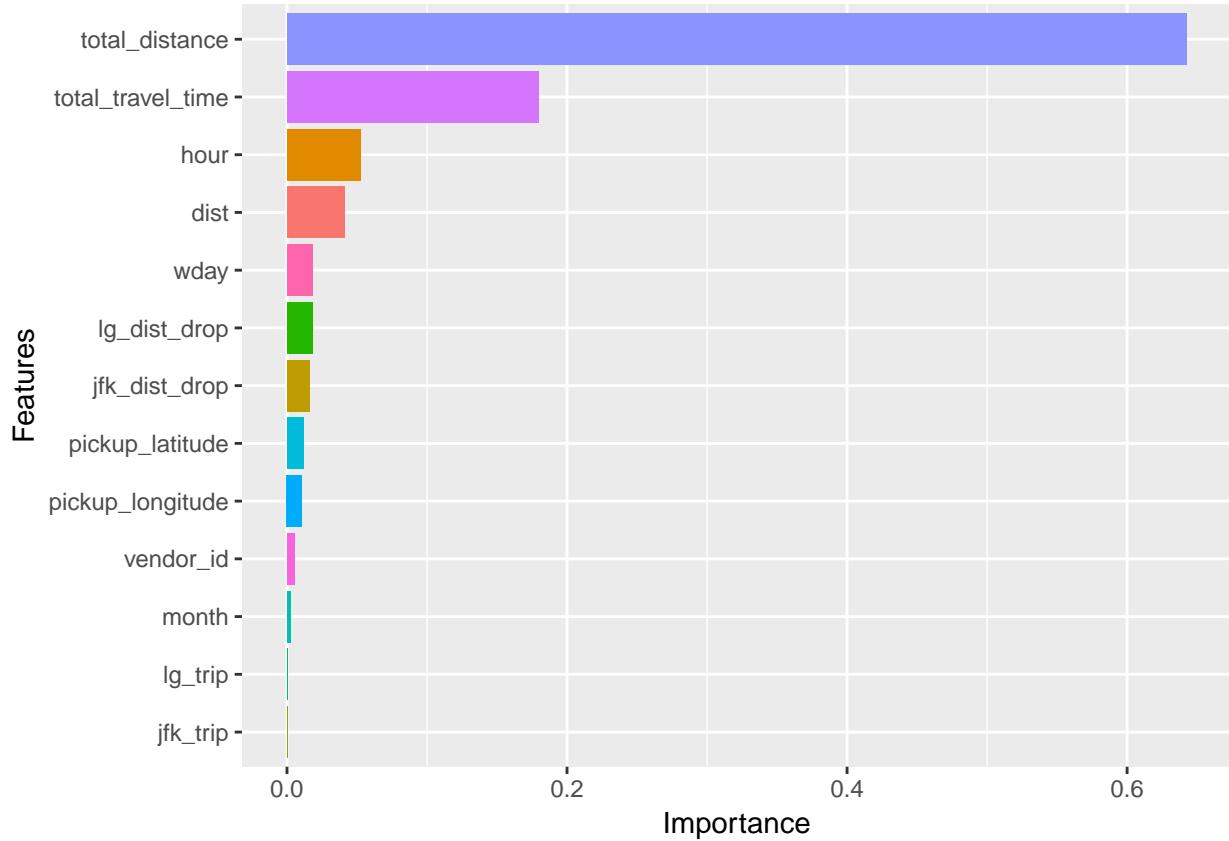
After training we will check which features are the most important for our model. This can provide the starting point for an iterative process where we identify, step by step, the significant features for a model. Here we will simply visualise these features:

```

imp_matrix <- as.tibble(xgb.importance(feature_names = colnames(train %>% select(-trip_duration))), mode = "wide")

imp_matrix %>%
  ggplot(aes(reorder(Feature, Gain, FUN = max), Gain, fill = Feature)) +
  geom_col() +
  coord_flip() +
  theme(legend.position = "none") +
  labs(x = "Features", y = "Importance")

```



Write the XGBoost prediction to csv file and later will submit in Kaggle to know the actual score. Please refer below screenshot to see submission score.

```
test_preds <- predict(gb_dt, dtest)
xgb_test <- test %>%
  mutate(pred_trip_time = test_preds)

xgb_test %>% write_csv("XGBOOST_TEST_PREDICTION.csv")

xgb_RMSLE <- RMSE(xgb_test$pred_trip_time, xgb_test$trip_duration)
xg_rsq <- 1 - (mean((xgb_test$pred_trip_time - xgb_test$trip_duration) ^ 2) / var(xgb_test$trip_duration))
```

Results of XG Boost predictions are RMSE -> 0.4038426 and Pseudo R^2 -> 0.7397317

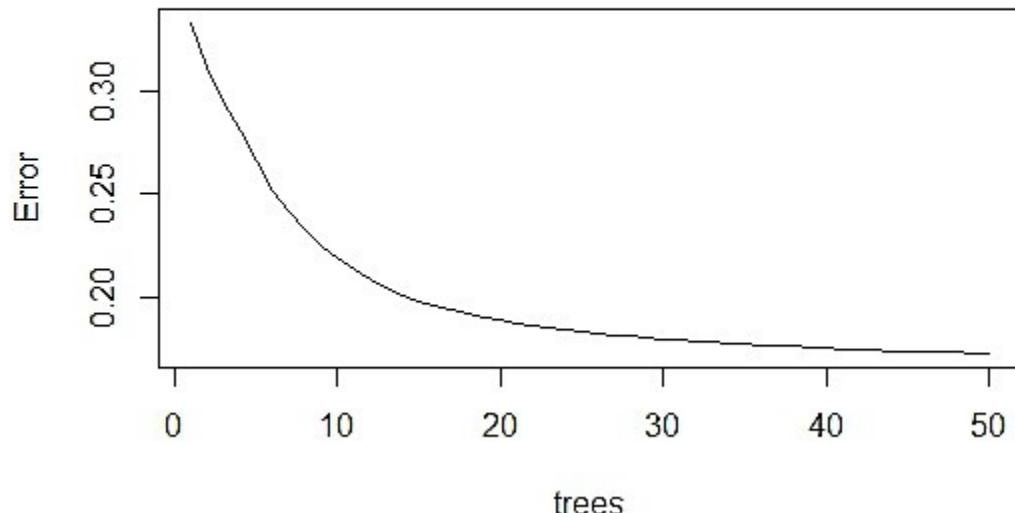
Next we start with Random Forest modelling. Random Forest modelling takes more CPU time to run in local machine. So we have runt he model outside Markdown and attached results for the same. Below is the model and its error plots. Overall Random Forest algorithm performed well than XG boost for our feature selection and data.

```
rf_model <- randomForest(trip_duration ~ total_travel_time + total_distance + hour + dist + vendor_id +
+ jfk_trip + lg_trip + wday + month + pickup_longitude + pickup_latitude + lg_dist_drop, data = train,
na.action = na.omit, importance = TRUE, ntree = 50, do.trace=TRUE)
```

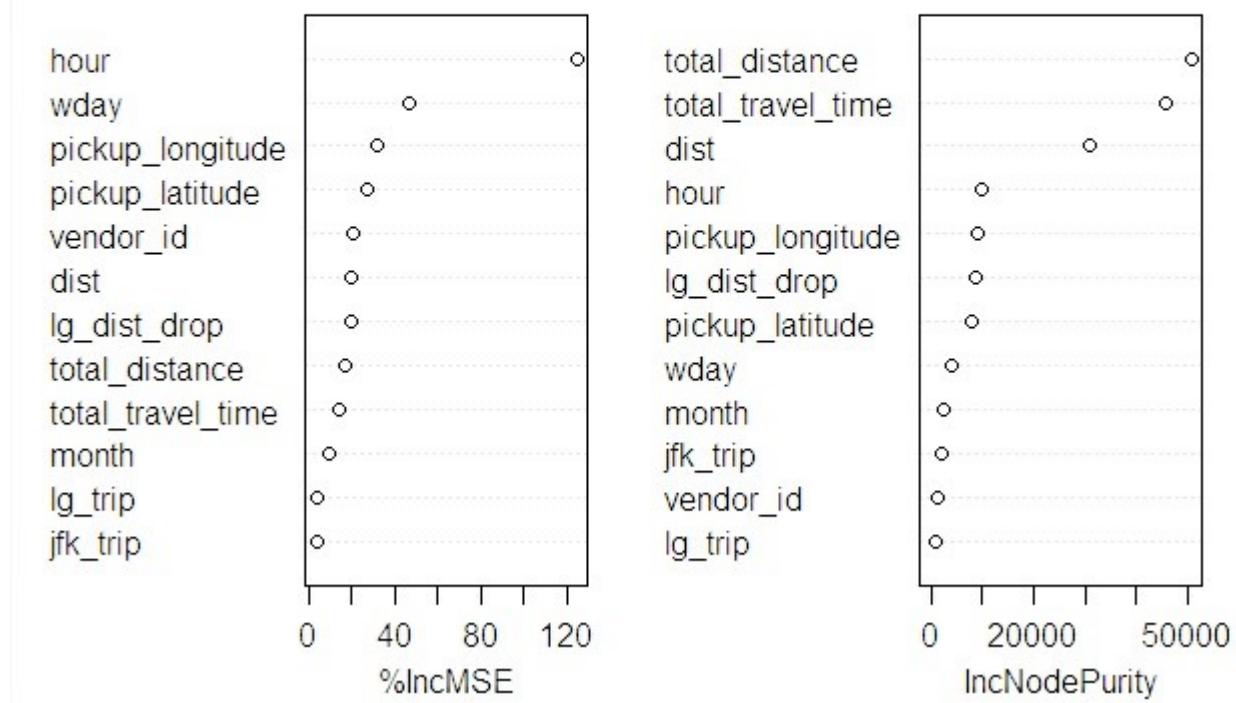
Results of XG Boost predictions are RMSE -> 0.3995 and Pseudo R^2 -> 0.7452

Listed below are important features predicted by Random Forest model and error trend as tree grows.

RandomForest



RandomForest Var Imp



In order to review the results of Random Forest we are reading the results of prediction test file and plot down the line for comparison.

```
rf_test <- read_csv("RF_TEST_PREDICTION.csv")
```

Next is our base line model Linear Regression.

```
lr_md1 <- lm(
  trip_duration ~ total_travel_time + total_distance + hour + dist +
  vendor_id + jfk_trip + lg_trip + wday + month +
  pickup_longitude + pickup_latitude + lg_dist_drop,
  data = train,
  na.action = na.omit
)
summary(lr_md1)

##
## Call:
## lm(formula = trip_duration ~ total_travel_time + total_distance +
##     hour + dist + vendor_id + jfk_trip + lg_trip + wday + month +
##     pickup_longitude + pickup_latitude + lg_dist_drop, data = train,
##     na.action = na.omit)
##
## Residuals:
##    Min      1Q  Median      3Q     Max 
## -5.1588 -0.2841  0.0093  0.3015  5.7972 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) -7.460e+00  3.741e+00 -1.995   0.0461 *  
## total_travel_time 3.030e-03  1.316e-05 230.181  < 2e-16 *** 
## total_distance -5.029e-05  1.295e-06 -38.822  < 2e-16 *** 
## hour          5.767e-03  1.610e-04  35.814  < 2e-16 *** 
## dist          -1.011e-05  1.392e-06 -7.260  3.88e-13 *** 
## vendor_id      2.952e-02  2.059e-03 14.338  < 2e-16 *** 
## jfk_trip       -6.136e-01  1.147e-02 -53.507  < 2e-16 *** 
## lg_trip        -4.410e-01  6.994e-03 -63.057  < 2e-16 *** 
## wday          2.986e-02  5.119e-04  58.323  < 2e-16 *** 
## month         1.683e-02  6.099e-04  27.597  < 2e-16 *** 
## pickup_longitude 6.927e-01  4.123e-02  16.802  < 2e-16 *** 
## pickup_latitude  1.575e+00  4.555e-02  34.576  < 2e-16 *** 
## lg_dist_drop   -2.031e-05  4.060e-07 -50.032  < 2e-16 *** 
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
##
## Residual standard error: 0.5432 on 279981 degrees of freedom
## Multiple R-squared:  0.5377, Adjusted R-squared:  0.5377 
## F-statistic:  2.713e+04 on 12 and 279981 DF,  p-value: < 2.2e-16 

lr_pred <- predict(lr_md1, test)
lr_test <- test %>%
  mutate(pred_trip_time = lr_pred)

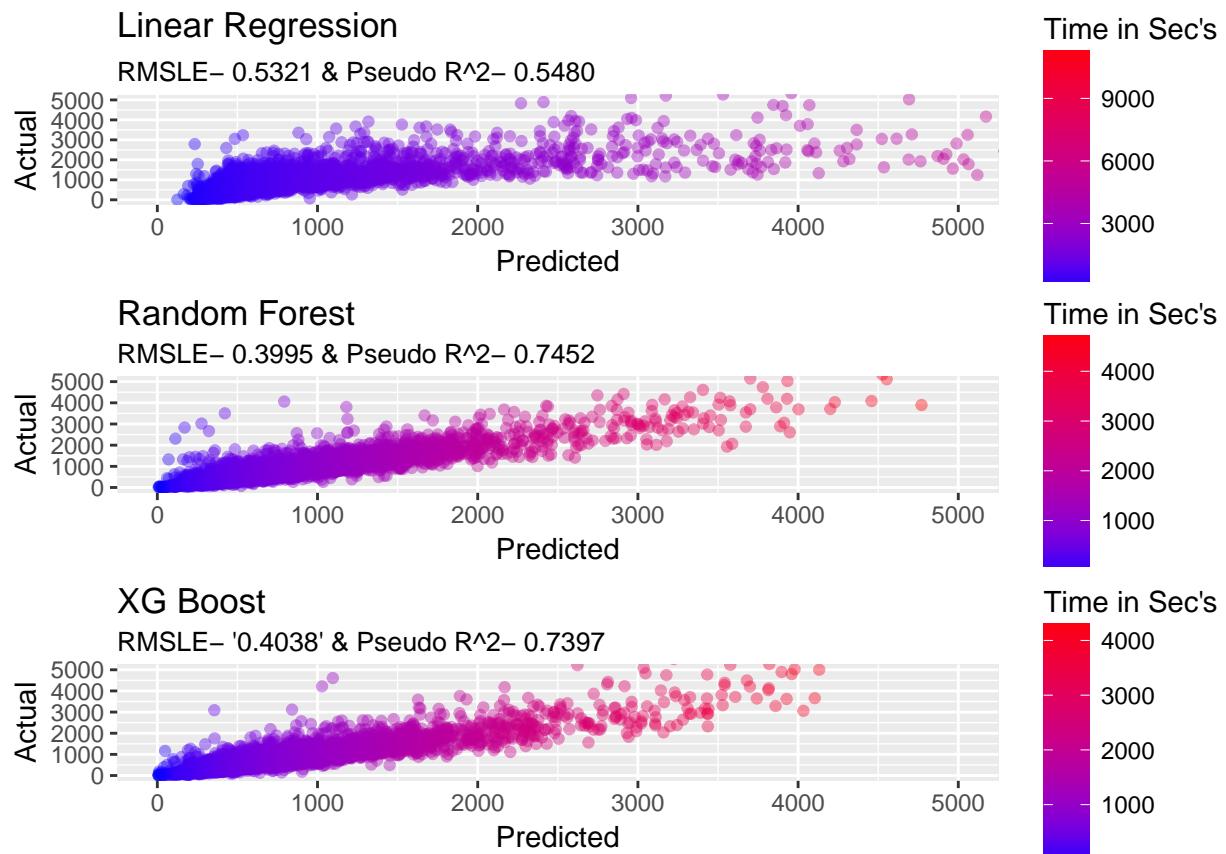
lr_test %>% write_csv("LREG_TEST_PREDICTION.csv")

lr_RMSLE <- RMSE(lr_test$pred_trip_time, lr_test$trip_duration)
lr_rsq <- 1 - (mean((lr_test$pred_trip_time - lr_test$trip_duration) ^ 2) / var(lr_test$trip_duration))
```

The results of all the models are compared in the below table

Algorithm	RMLSE	Pseudo R ²
Linear Regression	0.5321394	0.5480943
Random Forest	0.3995000	0.7452000
XG Boost	0.4038426	0.7397317

In order to visualize how well the models (randomforest, xgboost, Linear Regression) perform, we plot the actual versus predicted trip duration scatter plot. The plots in the figure suggest that both Random Forest and XG Boost models perform well on the test set. Most predictions are almost close to the true values.



References:

- EDA - NYC TAXI EDA The fast & the curious
- EDA - From EDA to the Top (LB 0.367)
- Mapping techniques
- GGPLOT
- Machine Learning
- XGBOOST Basics