

# Introduction

## Database *train.csv*

In section B of the term project we were given a data set of Porto taxi trips *train.csv*. The dataset was based off of 442 Taxises in the city of Porto, Portugal. The data was split among 9 categories TRIP\_ID, CALL\_TYPE, ORIGIN\_CALL, ORIGIN\_STAND, TAXI\_ID, TIMESTAMP, DAYTYPE, MISSING\_DATA, POLYTIME and consisted of 1710660 objects sized at 1.94 GB. The TRIP\_ID was a unique identifier for every trip and TAXI\_ID was a unique identifier for the taxi driver who performed the trip. CALL\_TYPE identified the method used to call the taxi: 'A' represented a call from central, 'B' represented if the taxi driver was demanded at a specific stand, and 'C' otherwise. ORIGIN\_CALL represented the phone number of the customer who demanded the taxi ( Sets CALL\_TYPE to 'A'). ORIGIN\_STAND gives each taxi stand a unique identifier (Sets CALL\_TYPE to 'B') .TIMESTAMP shows the unix Timestamp of the start of the trip. DAYTYPE represents the type of day it is at the start of the trip: 'B ' being a holiday or special day, 'C' being the day before a type 'B' day, and 'A' being everything else. Lastly, POLYLINE contains a set of arrays of two elements of global coordination, the longitude and latitude.

## Dealing with a massive dataset

In order to use the dataset, we had to import it into our R data table which at first we did by using the read() function. This created an issue because the data set was taking too long to load in and often times crashing the IDE due to insufficient memory. A temporary fix we decided to use was to only pass in 1000 rows of the full data set. This allowed us to get past the issues of read() and start working on the tasks in part B. In order to fully fix this problem and use the entire data set we decided to use fread() from the data.table library. This allowed us to read the data set faster and use less memory.

## Tasks to solve

- Find density models that accurately represented the trip duration and the time when the driver was busy.
- Explore if the CALL\_TYPE had an effect on the meantime of taxi trips
- Create models that predicted trip times, distances, ... (Add other explicit models)
- Compare how those predictive models relate to one another.

## Language and Libraries used

To solve these tasks our group used R to code the entire project as specified in the instructions. We also used the following libraries throughout the project

- Devtools: Used to assist in the installation of Regtools later in the project.
- Dplyr: Utilized this library to create and edit data frames.

- Data.table: Primarily used for fread() in order to speed up reading of *train.csv*.
- Regtools: Allowed us to utilize many machine learning models

## Exploring the possibility of trip duration

### The trip duration in this database

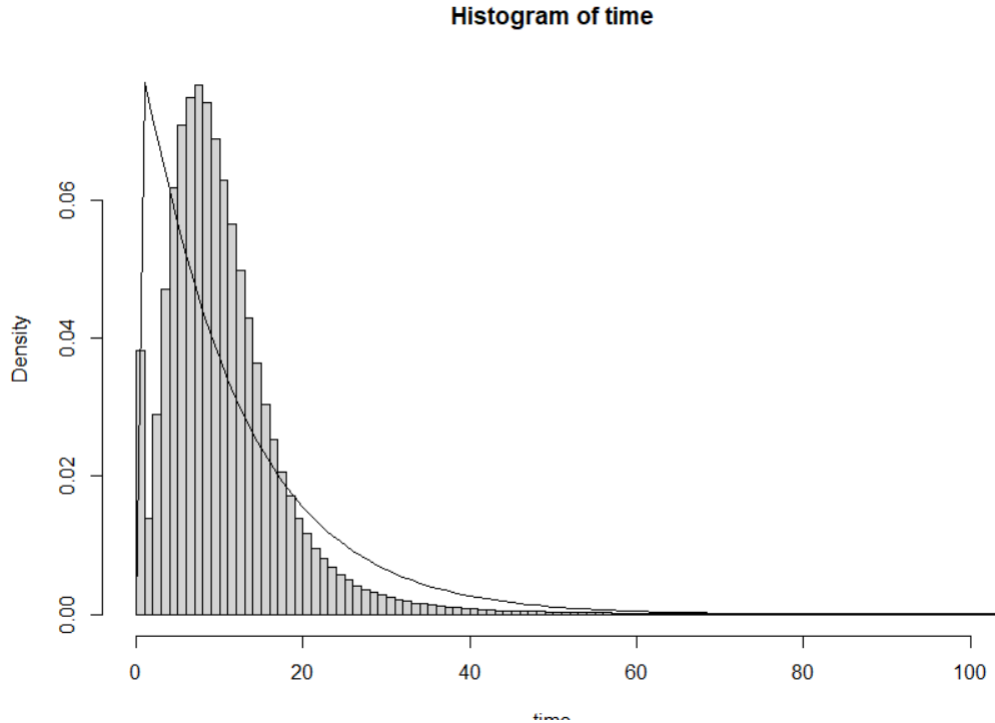
Since POLYLINE represents a set of coordinates recorded every fifteen seconds since the start of the trip, The duration of a trip can be interpreted as the number of intervals existing in the set. In terms of calculation, letting  $N_i$  donates to the number of intervals in  $i^{th}$  set, The duration of  $i^{th}$  trip in second is

$$T_i = 15(N_i - 1) \quad (1)$$

To be more general, we convert the duration into minutes by dividing it by 60.

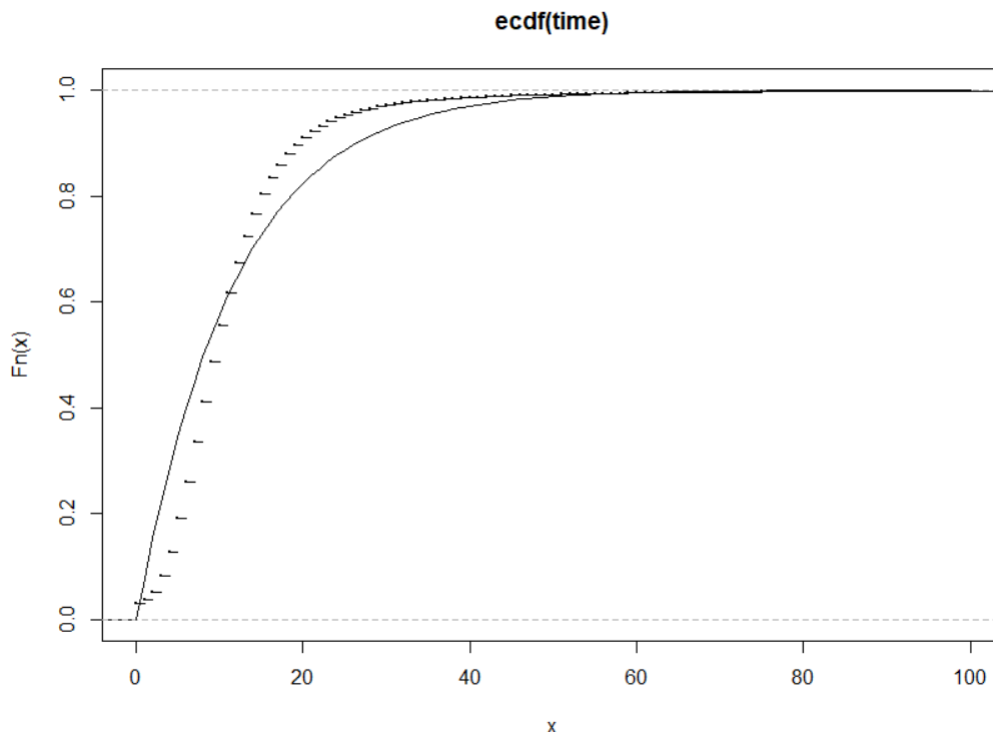
### The probability family of distribution of trip duration:

Now we have a column called trip duration, so we use hist to plot the trip duration into a histogram with the x-axis as the trip duration and the y-axis as the density. The probability of trip duration is found to be a gamma distribution family in the interval  $[0, 970]$ , but because most of the time is distributed within the  $[0, 20]$ , which makes the original histogram hard to be observed, we shrink the range of x-axis to  $[0, 100]$ . However, the gamma distribution seems to only exist in  $[1, 100]$  because the probability that the trip time is under a minute is too high to fit into the gamma distribution.



## Verify the hypothesis:

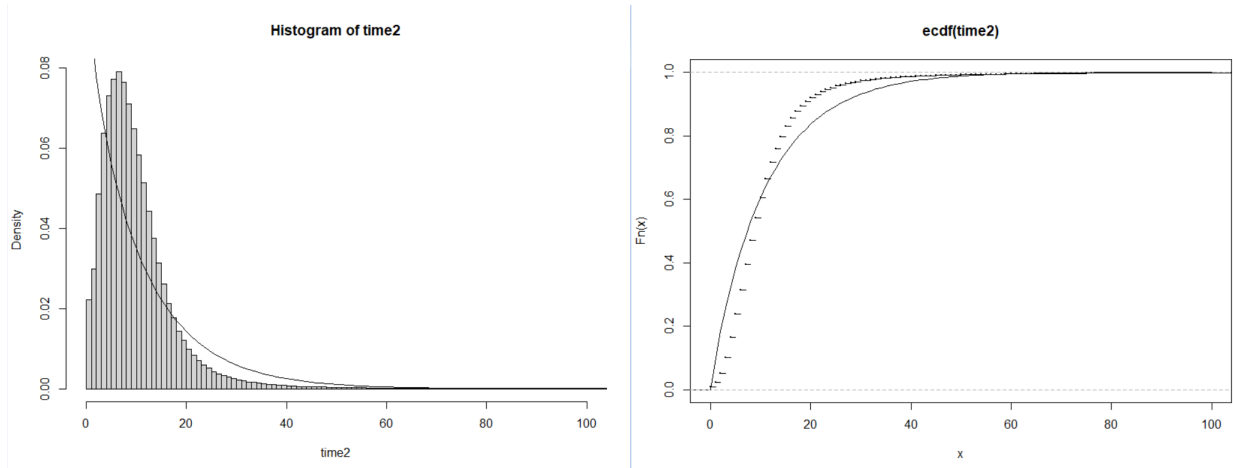
To test whether the probability density of trip duration is the gamma family of distribution, we apply the actual curve of the gamma distribution. Therefore, we need to find the  $r$  and  $\lambda$  value using mean and variance, and the mean and variance of time duration are found to be approximately 11.57 and 130.13, and  $r$  and  $\lambda$  are found to be 1.0294 and 0.0889. However, the gamma distribution curve does not fit the histogram so much,  $r$  value seeming to be too small.



Then, We compare the empirical cdf of  $T$  with the cdf of gamma distribution with  $r$  and  $\lambda$  which we just found. However, they don't match well either.

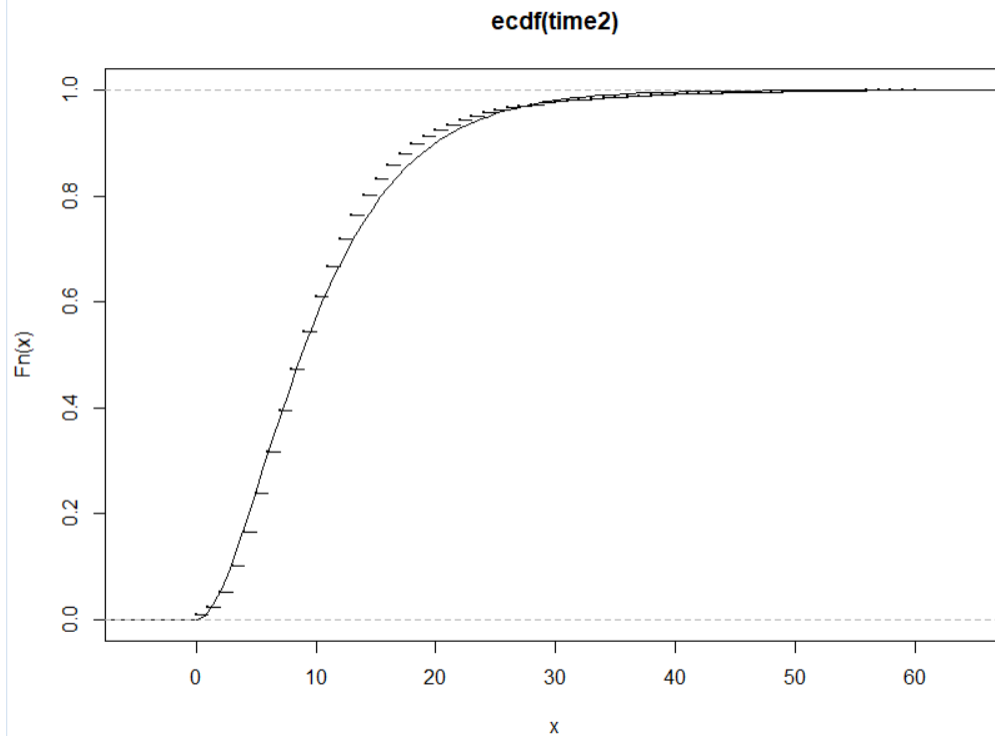
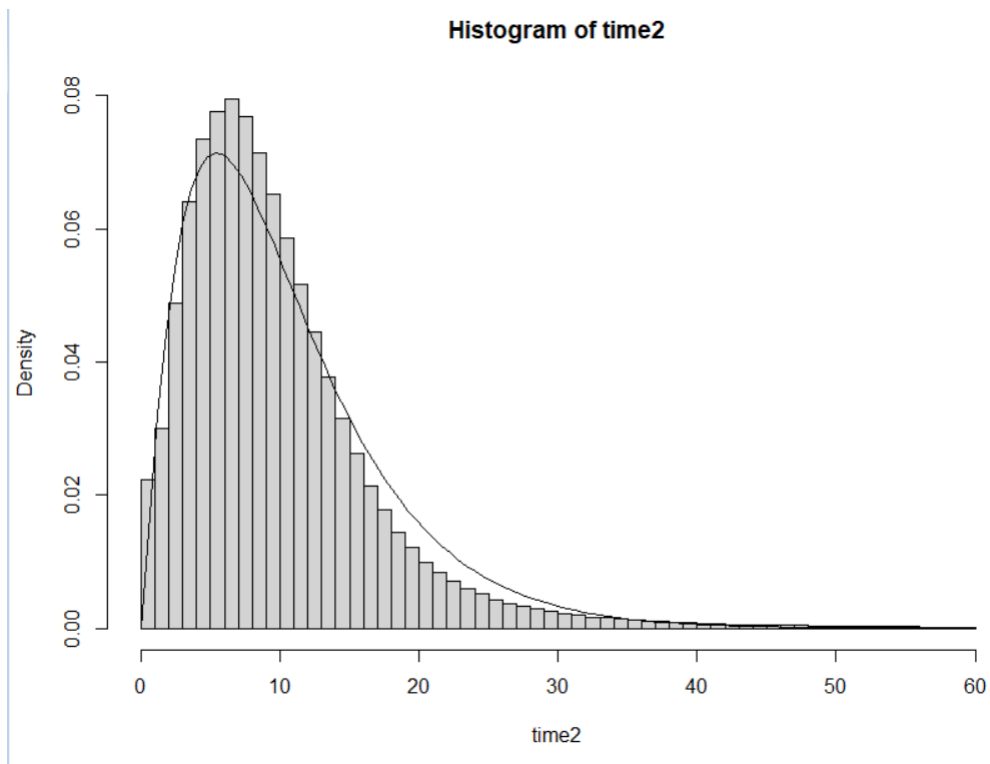
## First Correction

Since the proportion that the time duration is under one minute is too high, we decide to exclude it from the dataset and get a more rational look in the histogram. Then, we get the new mean, variance,  $r$ , and  $\lambda$ , 10.94, 129.88, 1.029, and 0.0889. The new gamma distribution curve turns out to be more inaccurate, but the slope is actually a little bit similar to the histogram's downward slope. When we compare the empirical cdf with the cdf of our new database and new gamma distribution, the difference between them is not smaller than the original difference.



## Second Correction

The unsuccess of fitting a gamma distribution curve to the histogram is due to the high value of variance, which is 130 for the distribution excluding the aberrantly high proportion, that yields a small  $r$ . In fact, the support of  $T$  could be any value less or equal to 970, and that can lead to a very large variance with respect to the concentrated proportion  $[0,20]$ . With a large value of  $r$ , gamma distribution turns out to be steep and inaccurate; therefore. We decide to shrink the range to  $[1,61]$  by excluding every other value from the vector, and that forms a smaller mean and much smaller variance, 10.39 and 51.56, which yields a larger  $r$  and  $\lambda$ , 2.09 and 0.201. Adding the curve, it fits much closer; comparing both empirical cdf and cdf of the gamma distribution, the result, as expected, comes up ideally. Therefore, we can confidently say the density of probability of trip time is the distribution of the gamma family since the proportion of interval  $[1,61]$  occupies 96.4% out of all distribution.



The trip duration shows that people averagely take about 11 minutes to their destination but most likely within 20 minutes. Because most of the probability, about 95%, is distributed in  $[0,25]$ .

The extremely high proportion of trip duration being under one minute may be caused by human errors and instrument errors: drivers mistakenly or accidentally press the in-operation

button, passengers' instantly calling it off, and technique issues. We observe that many POLYLINE don't have value, which explains that the coordination recorder doesn't work correctly in some cases.

### Exploring the distribution of proportion of driver is busy:

#### Definition of Busy:

Literally, drivers being busy means they are working, but in order to calculate the proportion of drivers being busy, we also need to know when they are not busy, waiting for passengers or including they time are off work, We assume that taxi drivers don't have specific time to work and amount of operating time, so we decide to define the total time as ones career life recorded in the database, from the beginning of a driver's first trip to the ending time of his last trip. Since each driver has an unique id number, we can track the sum of trip duration as the total working time of  $i^{th}$  driver, denoting as  $U_i$ , and Let  $T_i$  denotes to the total time of  $i^{th}$  driver,  $R_{i_j}$  denotes to the  $j^{th}$  trip duration of the  $i^{th}$  driver,  $N_i$  denotes to the number of trips  $i^{th}$  driver have completed, and  $B_i$  denotes to the proportion of time  $i^{th}$  driver being busy; the time can be calculated as

$$T_i = TimeStamp_{ilast} - TimeStamp_{ifirst} + R_{iN_i} \quad (2)$$

The busy time is as

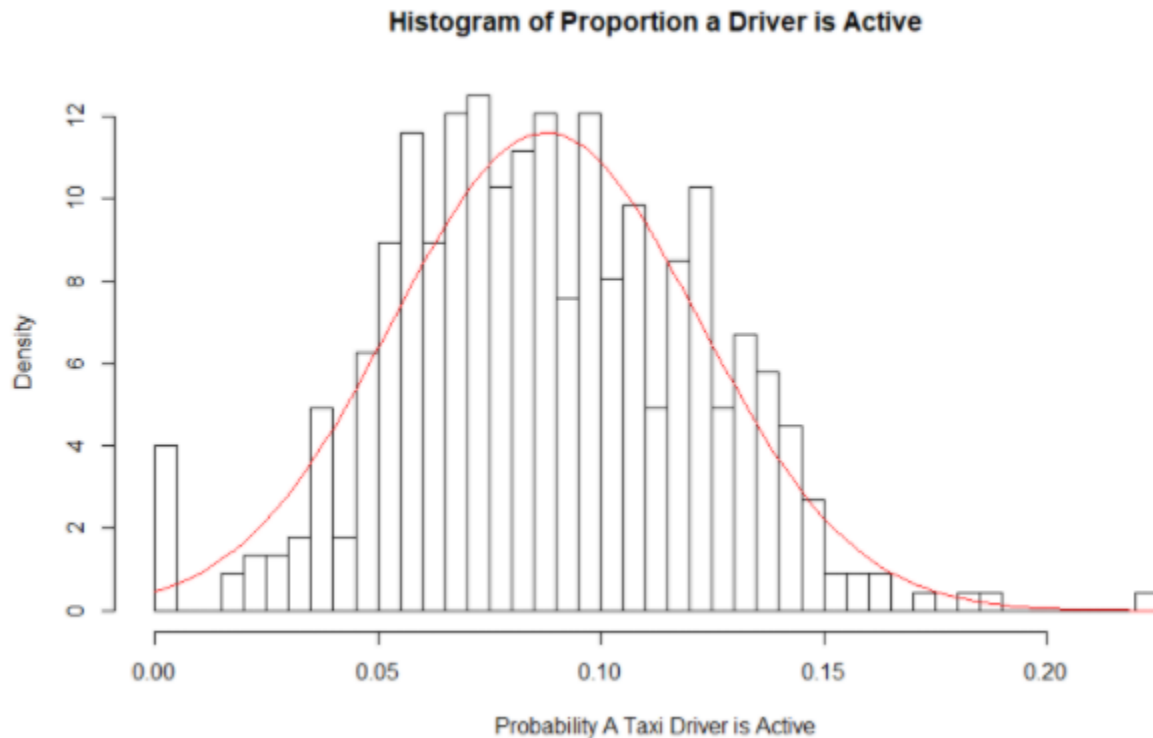
$$U_i = \sum_{k=1}^{N_i} R_{ik} \quad (3)$$

Then, the proportion is

$$B_i = \frac{U_i}{T_i} \quad (4)$$

#### Probability Distribution:

Plotting B into the histogram, we get an approximate normal distribution in the range of (0,0.23]. As we observe there are some extreme cases that drivers only work about 0.005 of their time involved in this career, averagely 0.12 hour a day.



## Investigation of average trip times using CALLTYPE:

```

31 time = time/60 #convert to minutes
32
33 call$time <- time
34
35 mean_time = mean(call[["time"]])
36 var_time = var(call[["time"]])
37
38 a = filter(call, CALL_TYPE == 'A')
39 b = filter(call, CALL_TYPE == 'B')
40 c = filter(call, CALL_TYPE == 'C')
41
42 mean_a = mean(a[["time"]]) #convert to minutes
43 mean_b = mean(b[["time"]]) #convert to minutes
44 mean_c = mean(c[["time"]]) #convert to minutes

```

The trip times used in the data set have been given in terms of UNIX timestamps which make it hard to visualize. In order to get a better understanding of trip times, we decided to convert the timestamps into minutes using the following conversion rates:

- One minute = 60 UNIX units
- One hour = 3,600 UNIX units
- One day = 86,400 UNIX unit

The equation to get minutes from UNIX timestamp:

$$Minutes = \frac{UNIXTIMESTAMP}{60} \quad (5)$$

After we simply applied the `mean()` and `var()` functions to the newly converted `TIMESTAMP` variable to obtain the mean and variance

General Dataset Mean (`TIMESTAMP`):

$$Mean(Overall) = \frac{Minutes(Converted)}{1710670} = 11.940 \quad (6)$$

On the other hand, to get the values with the dependency of `CALLTYPE`, we used `filter()` to sort through the data in order to get the correct `TIMESTAMP` values in regard to the three `CALLTYPE`'s A,B, and C. After we used `mean()` to get a mean of the `TIMESTAMPS` sorted by `CALLTYPE`.

$$Mean(CALLTYPE = A) = \frac{Minutes(Converted)}{1710670} = 12.506 \quad (7)$$

$$Mean(CALLTYPE = B) = \frac{Minutes(Converted)}{1710670} = 11.086 \quad (8)$$

$$Mean(CALLTYPE = C) = \frac{Minutes(Converted)}{1710670} = 12.873 \quad (9)$$

### Analysis:

Difference between trip times using different call types:

$$CallTypeA = 11.940 - 12.506 = -0.566(33.96SecondsSlower) \quad (10)$$

$$CallTypeB = 11.940 - 11.086 = 0.854(51.24SecondsFaster) \quad (11)$$

$$CallTypeC = 11.940 - 12.873 = -0.933(55.98SecondsSlower) \quad (12)$$

$$CallTypeB > CallTypeA > CallTypeC \quad (13)$$

In general, the type of call made to summon a taxi doesn't make much of a difference in regards to trip time, but the fastest would be from a specific stand. Furthermore, we can examine the variance of each call type to determine the spread of the data set compared to the mean.

$$Variance(Overall) = 130.24 \quad (14)$$

$$Variance(CallTypeA) = 72.26 \quad (15)$$

$$Variance(CallTypeB) = 64.76 \quad (16)$$



$$Variance(CallTypeC) = 269.50 \quad (17)$$

The Variance of the 3 Call Types are fairly close with the exception of Call Type C. This indicates that there is a lot of volatility in the trip times for Call Type C. This could be caused by the fact that Call Type C takes riders from various methods instead of a set way of getting riders.

## Moddling the Data

### Models:

In the term project we decided to use the following models to visualize our data:

- Linear
- Linear Polynomial
- k-Nearest Neighbors
- Boosting

### Gathering the parameters:

The first model we decided to take on was modeling the distance vs the speed. In order to plot these two variables on a model we had to obtain the distance and speed. The distance we obtained by utilizing the coordinates from POLYLINE column in the dataset. After we just took the start and end points of the formula and applied the Haversine distance formula. The Haversine (or great circle) distance is the angular distance between two points on the surface of a sphere. It is a very accurate way of computing distances between two points on the surface of a sphere using the latitude and longitude of the two points.

$$(x1, y1) = StartingCoordinates(x1, y1) = EndingCoordinates$$