
CZ4015

Final Report

Author:

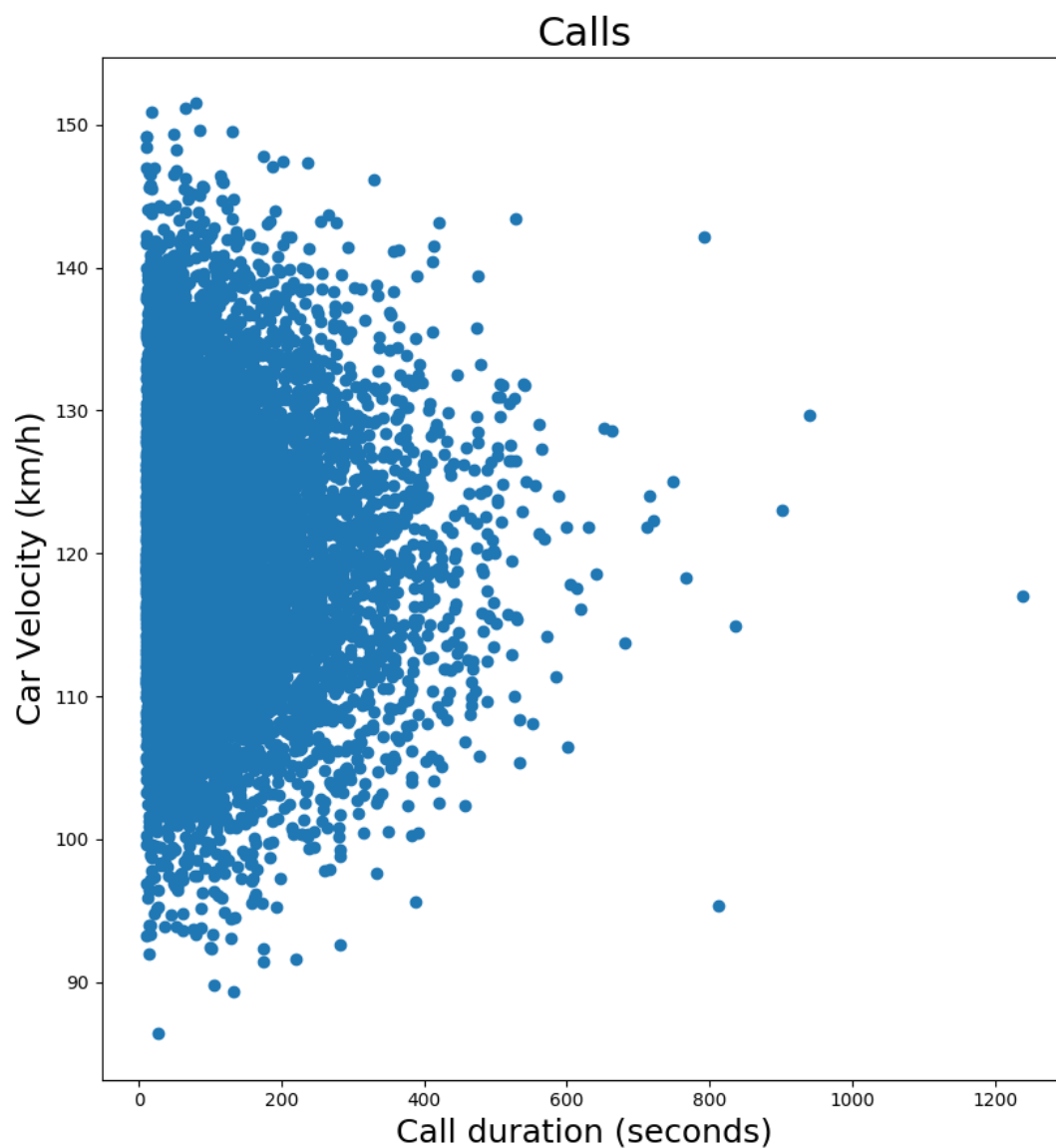
Pavel Jahoda (N1800740K)

1	Input Analysis	1
1.1	Cleaning data	1
1.2	Distribution identification	2
1.3	Hypothesis Testing	3

Input Analysis

Cleaning data

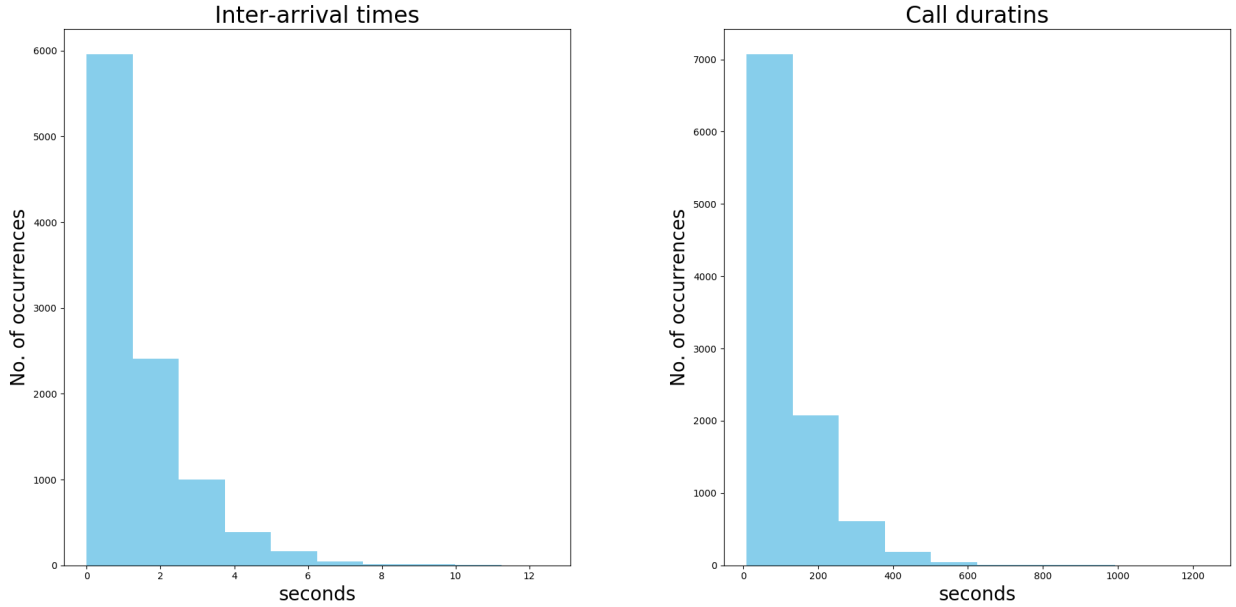
First we need to clean our data from mistakes made during data gathering. To do so, we will visualize our data using scatter plot. The scatter plot below visualizes how long each call lasted and how fast was the person driving when making the call.



The most extreme value we can see is from a call that lasted more than 20 minutes, which is entirely possible, so we won't delete any values. It was also checked that each call started from station between 1-20.

Distribution identification

From the histograms below we can see similarity between the distributions of inter-arrival times and call durations. Both of these histograms resemble probability density function of exponential distribution.



First, we will estimate the parameter of the exponential functions using maximum likelihood estimation. The exponential function is denoted as:

$$\lambda \cdot e^{-\lambda \cdot x} \quad (1)$$

The likelihood function for the parameter lambda given x_1, x_2, \dots, x_n is denoted as:

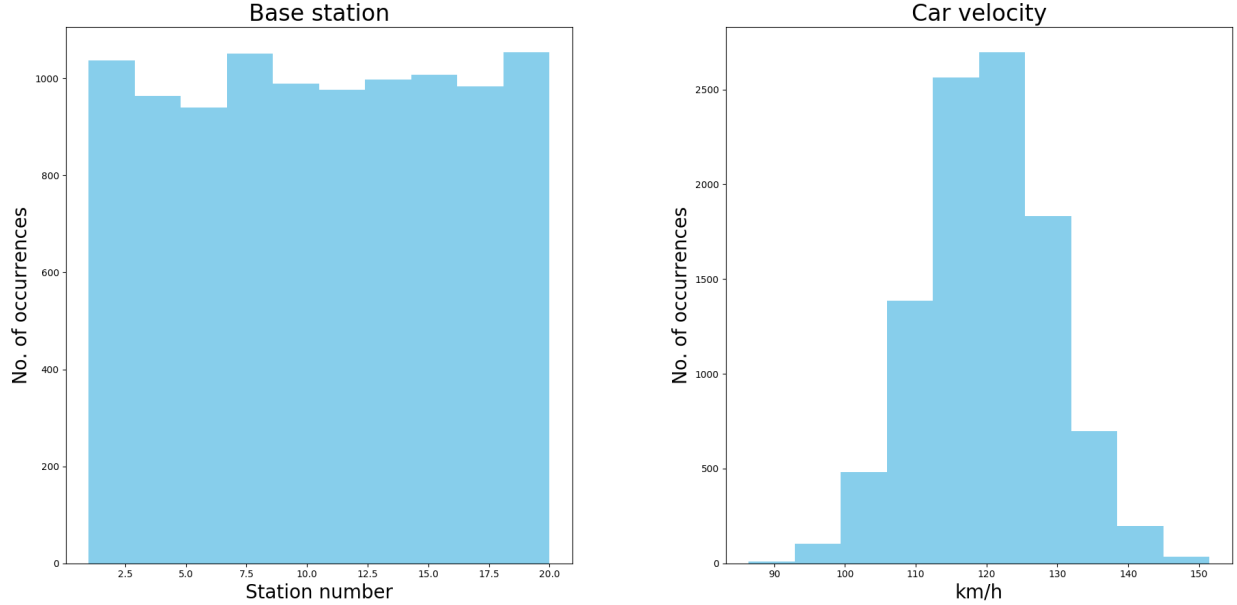
$$\mathcal{L}(\lambda|x_1, x_2, \dots, x_n) = \lambda^n \cdot e^{-\lambda \cdot \sum_{i=1}^n x_i} \quad (2)$$

To find the lambda for which the likelihood function is maximal, we differentiate by lambda and solve for lambda when the derivate is equal to 0, which results in the following formula:

$$\lambda = \frac{n}{\sum_{i=1}^n x_i} \quad (3)$$

Using maximum likelihood function of an exponential distribution we have calculated lambda for the inter-arrival times to be 0.73 and lambda for the call duration times to be 0.009.

On the other hand, the two histograms below show two different distributions. The histogram on the left shows that the stations where the cars are located when the call begins are uniformly distributed from 1 to 20. The histogram of the car velocities (on the right) resembles normal distribution.



Similarly as with the exponential distributions above, we will use maximum likelihood method to calculate the parameters of the normal distribution (car velocities). We get $\mu = 120.07$ and $\sigma^2 = 81.33$.

Hypothesis Testing

After we made distribution identification hypotheses, it is time to do hypothesis testing. First, let's focus on testing our hypothesis about inter-arrival times having exponential distribution with $\lambda = 0.73$ using Pearson's chi-squared test, which has the following formula.

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \quad (4)$$

First we will divide the observations into k mutually exclusive classes that all have the same probability that an observation falls into said class. In the formula above E_i the expected value of observation's that fall into a i th class and O_i is the number of observed values that actually fall into i th class. We define our class as an interval from r_{i-1} to r_i by solving the following formula.

$$\int_0^{r_i} \lambda \cdot e^{-\lambda \cdot x} dx = \frac{i}{k} \quad (5)$$

For $i = 1, 2, \dots, k$. By integrating probability density function we get cumulative distribution function.

$$1 - e^{-\lambda \cdot r_i} = \frac{i}{k} \quad (6)$$

Solving it for r_i , we get

$$r_i = \frac{\ln(-\frac{i}{k} + 1)}{-\lambda} \quad (7)$$

Now that we have our classes we can compare the number of observed values in a class compared to expected number of values in the class to get our χ^2 value. For the inter-arrival times, we got $\chi^2 = 100$, which is less than $\chi_{0.05}^2$, so we cannot reject our null hypothesis that the inter-arrival time has an exponential distribution with $\lambda = 0.73$. If the result of our chi-square test was higher than $\chi_{0.05}^2$, we could reject our null hypothesis and be 95% sure our rejection was correct.

Then, we performed chi-squared test for our other hypotheses. When performing chi-squared test for a hypothesis that call duration times have exponential distribution we got very high χ^2 value, which meant we should reject our hypothesis, but after shifting the values by their lowest value the hypothesis could not have been rejected.