# Analysis of Traffic Delays by Borough Using Bus Delay Data

*Deric Pan, dericpan*

## I.    Introduction:

There are many different means of transportation in large urban areas. For commuters within the New York City metropolitan area, there are many means of transportation, including public buses, charters, taxis, cars, bikes, and of course, the subway system. For many students around the New York City metropolitan area, commuting to school by the subway system is an impossible task. Many students throughout New York City spend hours upon hours commuting to school. Even working people throughout the city struggle with this; they often have to travel through heavy traffic to reach popular economic centers like lower Manhattan and Long Island City. Since the streets and highways of New York City are so congested, commuters often have to schedule their days around traveling to school or work.

The research question that this paper will address is whether traffic delays differ by different times of the day. In order to understand why this question is important, an understanding of the New York City transportation system must be had. New York City consists of five different boroughs. The most densely populated borough, Manhattan, is bound by the East, Hudson, and Harlem rivers. Just south of Manhattan is Brooklyn, a residential and extremely large borough. To the north of Brooklyn and to the east of Manhattan lies Queens, a borough with quickly expanding neighborhoods. The fourth borough, Bronx, is located to the north of Manhattan. The least densely populated borough of New York City is Staten Island, connected to Brooklyn by the Verrazano Bridge. These five boroughs house over 8.6 million residents, with just over 300 square miles of space, making New York City the most densely populated major city in the United States.

The research question at hand is important because it can provide insight on the best time of commute for residents of the five boroughs. Additionally, with New York City's rapidly aging infrastructure, which includes its bridges, subways, streets, and buses, it is becoming increasingly more difficult for New Yorkers to commute through the city. According to C. Rosenzweig et al., massive area wide flooding from a 2007 storm caused system wide outages of the MTA subway system. This phenomenon was not a rare occurrence, and has consistently caused additional delays within New York City's transportation systems (C. Rosenzweig et al., 2011). Potentially understanding when most traffic occurs and how long this traffic occurs for can help lawmakers in New York City to understand where infrastructure funding can best be utilized. More immediately, it can help commuters who live in the five boroughs to better plan their commute.

This paper consists of six sections. Excluding this introduction, there is the *Data* section, which describes the data set, "Bus Breakdowns and Delays". Additionally, it introduces a paper by Yacizi et al., *Breakdown of weather, intersection and recurrent congestion impacts on urban delay in New York City*, which motivated the research question of this paper. The *Methods* section describes the methods that are used for this research. It begins by describing a non parametric bootstrap method to create confidence intervals for the mean differences in delay. Additionally, it describes two permutation tests that will be used, Wilcoxon tests and the Kolmogorov-Smirnov tests. The *Simulations* section provides simulations on testing the accuracy of the methods listed above. *Analysis* directly tests the data through different means, including the Wilcoxon tests and by providing a bootstrapped confidence interval. The *Discussion* section provides insight on the research, how it can be further applied, and how this research paper could be later improved upon A *Bibliography/Works Cited* section is included.
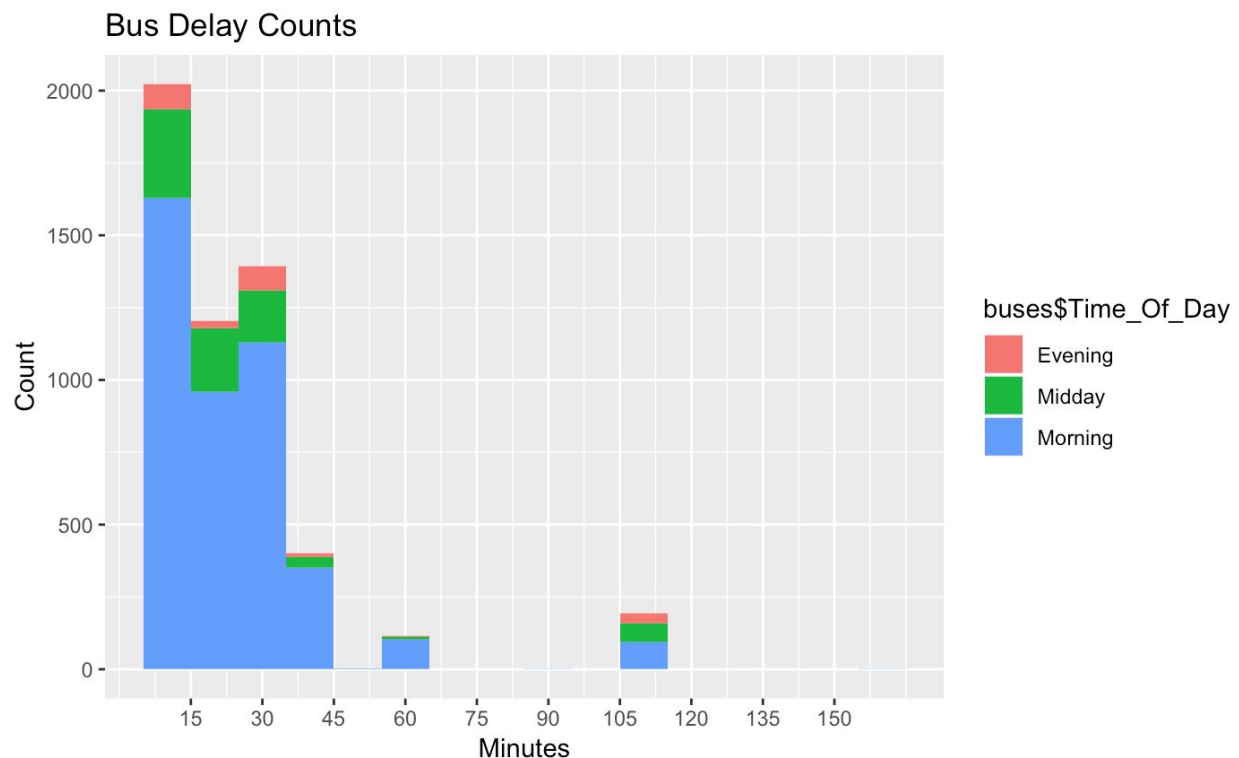
## II.    Data:

The data set used in this research paper is called "Bus Breakdowns and Delays", which was provided by the New York City Department of Education. It was collected from school bus vendors out in the field. The bus vendor staff was instructed to radio the dispatcher at the bus vendor's central office during delays and breakdowns. The staff was instructed to log into the Bus Breakdown and Delay system in order to record the event (City of New York, 2019).

The research in this paper is motivated by the work of Yazici et al. in recording the effects of weather, signaling, and recurrent congestion on urban delay. This paper differs in the fact that it will observe length of delay by borough, comparing delay times by heavy traffic in different neighborhoods. Yazici et al., determined that lower Manhattan experienced high delays during midday, while other boroughs were more delayed during rush hour. This paper will differ in the method by which delay is measured. Instead of measuring delay using taxi data like Yazici et al., this paper aims to do so by using bus delay times.

Since this data set consists of delays of school buses, three important time periods of the day are given. First and foremost, data points regarding rush hour, which in the case of this paper is defined to be from 6:00 AM - 10:00 AM and 4:00 PM to 7:00 PM. Once the data was cleaned, the former was categorized as morning and the latter as evening. Because schools often get out at different times, midday information, which was defined to be from 10:00 AM to 4:00 PM is also provided.
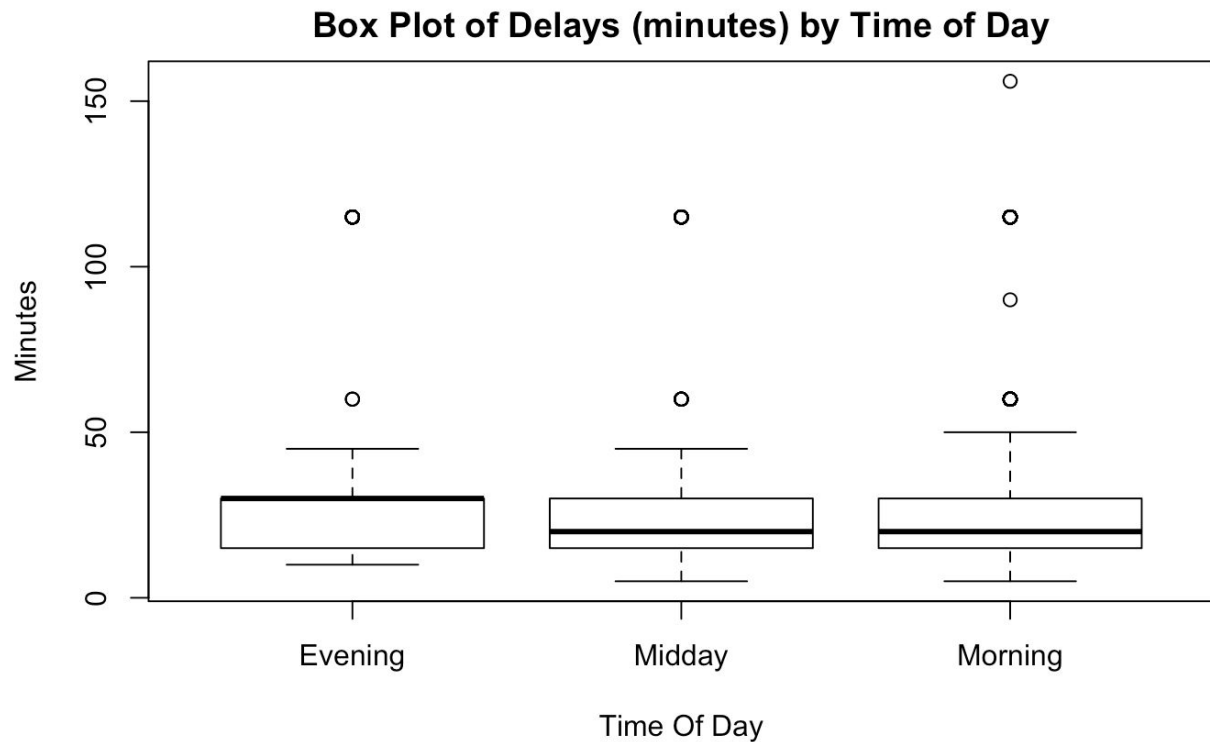
The data set consists of multiple fields, including school year, unique ID of each record, breakdown or delay category, bus number, route number, reason, schools serviced, date and time the event occurred on, date and time the event was recorded, borough, bus company name, delay length, time occurred on, school age, and many more. However, for the analysis in this research paper, only a few categories will be focused on. The categories are as so: school year, reason, route number, borough, length of delay, and length of delay.

The Bus Breakdown and Delay data set consisted of data from the 2015-2016, 2016-2017, and 2017-2018 school years. It consists of data from the five boroughs of New York City, Westchester, Nassau County, New Jersey, and Rockland County. However, we will clean this data in order to remove any places not within the five boroughs of New York. This helps to ensure that we focus mainly on traffic within New York City. Additionally, the reason of bus delay or breakdown will only focus on heavy traffic, as this is the indicator that the research in this paper is focused on. Below is a stacked histogram showing a count of how long delay is. It is split up categorically by different times of the day.
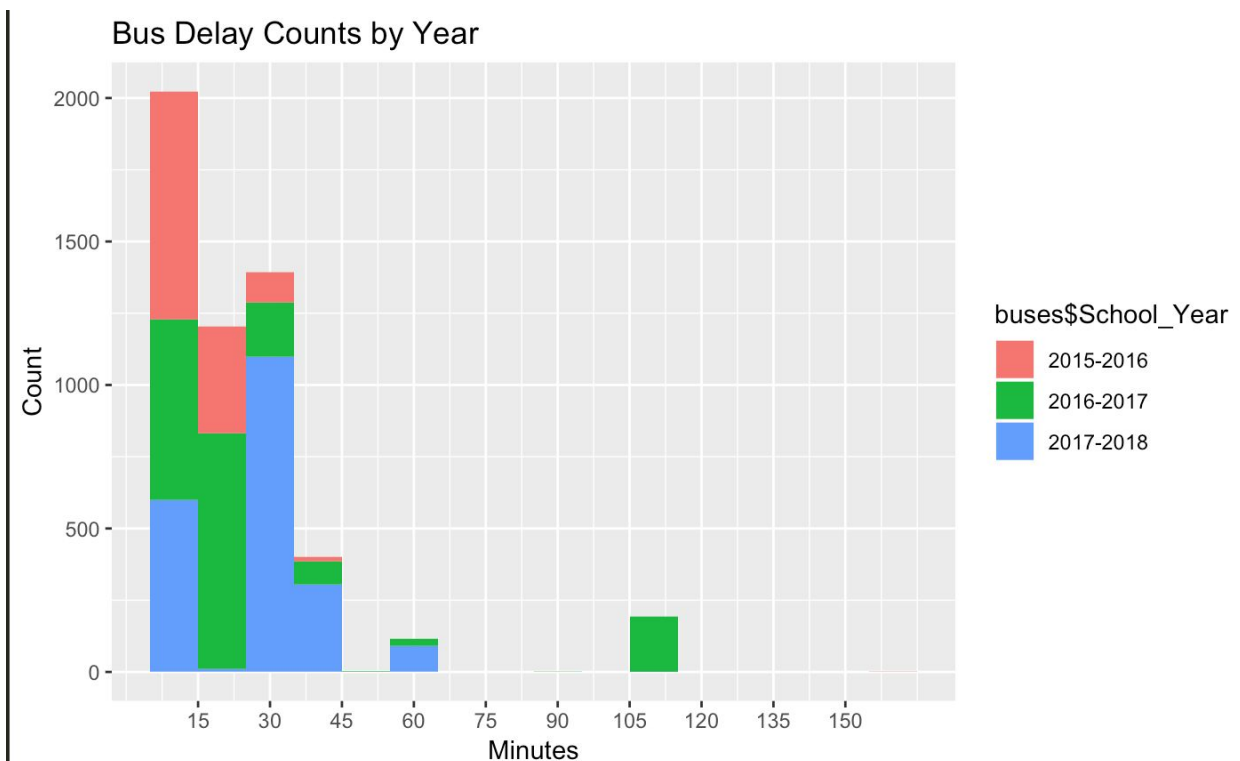
**Bus Delay Counts**



From simply observing the above, it appears that the most common and frequent delays occur during the morning. As of now, the data set is not as complete as one would like. However, research is still performed with the case above, and certain methods, such as the nonparametric bootstrap help to alleviate the lack of data in the evening and midday. This would allow researchers to determine if traffic delay has been increasing annually. This information is further

analyzed with a box plot as shown below:

## Box Plot of Delays (minutes) by Time of Day



The box plot above acts as motivation for the research; particularly by visualizing how medians for midday and morning are the same, but the evening is higher.

Observing the above and comparing it with the data set, this information is not as complete as one would like. With only three years of information, it is hard to determine whether bus delays are trending upward or not. Further research can be conducted once the City of New York receives additional data.

## III.    Methods:

To determine the amount of delay per borough, this paper paper will begin by cleaning the data in order to be able to read in the information and perform tests upon it.  In this paper, time is split accordingly by hour. Morning times are from 6:00 AM to 10:00 AM, midday from 10:00 AM to 4:00 PM, and evening from 7:00 PM. This paper hypothesizes that morning time delays are higher than that of the midday and late afternoon. After cleaning the information, the research in this paper is performed with the help of two permutation tests.

According to Rizzo, permutation tests are based on resampling, but are drawn without replacement. They are often times applied as nonparametric tests of the general hypothesis that

$$H_0: F = G, H_1: F \neq G$$

where F and G are two unspecified distributed. (Maria Rizzo, 2016) To begin the permutation tests, this paper will first assume that the null hypothesis is that distribution of the delays between the different times of the day are the same. This works because Rizzo states that permutation tests can be applied to multi-sample problems, with similar methodology. It can be written like so:

$$H_0: F_1 = F_2 = \ldots = F_k, \qquad H_1: F_i \neq F_j \text{ for some } i, j$$

The first method that will be used is through a simple Wilcoxon test. Wilcoxon tests are used for non parametric data, and since the data used in this paper does not assume known distributions, this would be the appropriate choice. We will run the Wilcoxon test between the distributions of delay from the three times of the day. In this case, the Wilcoxon tests should provide a p-value of less than .05, our significance level that is being tested.
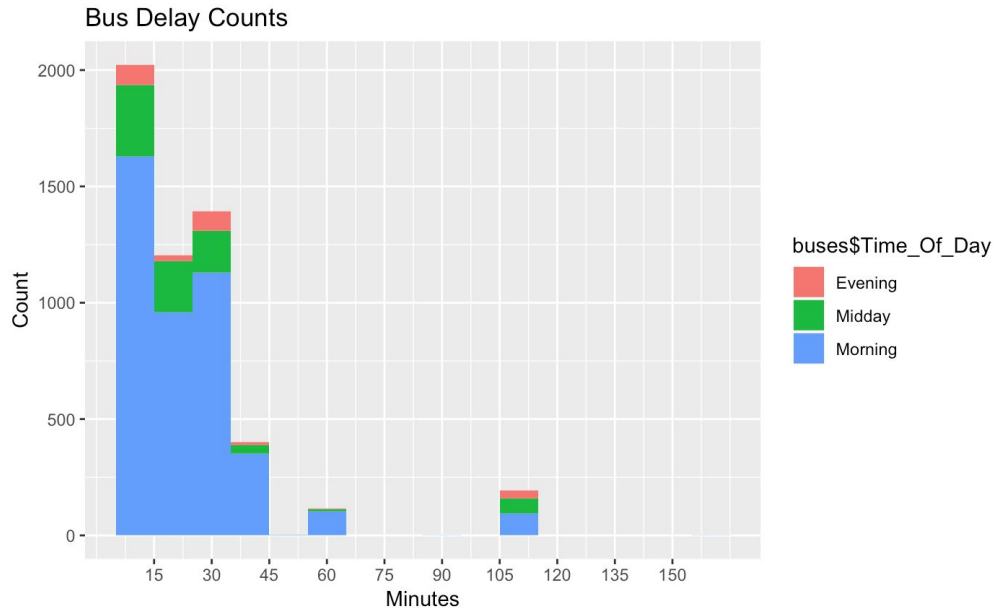
Like the Wilcoxon tests, the Kolmogorov-Smirnov test is also an excellent choice if distributions are not assumed. Kolmogorov-Smirnov tests are used to compare if one variable has identical distribution in two populations. This is relevant to our idea of testing the delay has identical distribution within the population of buses from different boroughs.

When simulations were ran using the K-S test, it was determined that the power for Wilcoxon tests were superior, and thus, information for the K-S test is not shown below.

Furthermore, a nonparametric bootstrap is used to produce the mean delays for different times of day. We will produce a 95% bootstrapped confidence interval for each of the different times of day, allowing us to observe how each one differs.

## IV.    Simulations:

In this section of the paper, the accuracy of the two tests above are analyzed. In order to perform these tests, simulated data for the data set was created. When observing the graph for the count of delay lengths, it is easy to observe that the distribution is Poisson, as shown below.
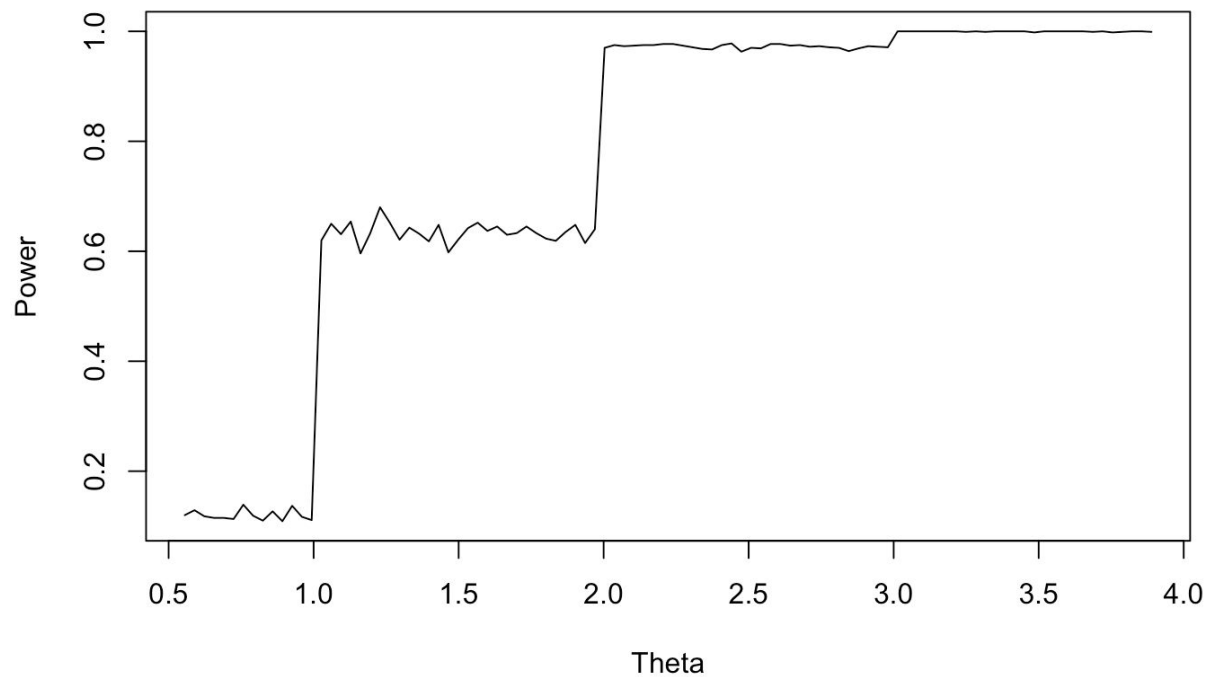


Furthermore, a Poisson distribution would make sense, as the events in the data set are independent and nonnegative with a fixed time period. In this case, we set the distribution to be Poisson with lambda = 20.

This paper then begins the simulations by testing a set of numbers, theta, from -5 to 5, separated by intervals of 0.5. The p-values for each number from this interval is then created into a vector in which the data will be stored. This part of the simulation is responsible for testing however much this additional theta would need to be added to the "mornings" in order to have the same median as the "evenings"

Once this is accomplished the paper works to simply create a 95% confidence interval for these thetas, which is then used to evaluate the power with adding each theta from the set of thetas to the mornings. The results are as shown:

```
[1] 0.5555556 3.8888889
```

Following this, a curve that visualizes the increasing powers for the tests is created. The curve is provided below:

This increasing power curve provides insight on how the increasing thetas are able to increase the accuracy of our tests. Once this power simulation was accomplished, the paper moved onto comparing the chosen process with Wilcoxon tests against those of the K-S tests. As shown below, the power of the simulated data was computed for both Wilcoxon tests and K-S tests, with the additional theta tested to be 1.25. The results are shown below, the former being of the Wilcoxon test and the latter being of the Kolmogorov Smirnov-test:

```
        Exact binomial test

data:  sum(wts.p) and 1000
number of successes = 621, number of trials = 1000, p-value = 1.887e-14
alternative hypothesis: true probability of success is not equal to 0.5
95 percent confidence interval:
 0.5901106 0.6511747
sample estimates:
probability of success
              0.621
```

```
        Exact binomial test

data:  sum(kts.p) and 1000
number of successes = 608, number of trials = 1000, p-value = 8.795e-12
alternative hypothesis: true probability of success is not equal to 0.5
95 percent confidence interval:
 0.5769600 0.6384022
sample estimates:
probability of success
              0.608
```

With the power being higher in the Wilcoxon test, the paper moved forward by using this test, as it was more accurate.

## V.    Analysis:

Median delays for morning, midday, and evenings are respectively shown below.

```
[1] 20
[1] 20
[1] 30
```

As seen above, when observing simple empirical data, without further statistical testing being performed on our data, we are able to see that the median delay time for buses in the morning and midday are both 20 minutes. The median delay time for buses in the afternoon, however, is higher, with a median of 30 minutes.

The means are also provided below in the same order as that of the median:

```
[1] 25.42616
[1] 28.79562
[1] 36.39113
```

In order to provide further insight, the research in this paper provides a 95% bootstrapped interval for means of each delay. With our nonparametric bootstrapped interval, we find that the interval for the mean delays in the morning is between 24.87 and 25.96. The intervals for the mean delays in the midday is between 27.02 and 30.69, and for the evening, it is between 32.26 and 40.53. Nonparametric bootstrapping of the medians provides only 20 as the answer for both the midday and the morning, so it is performed on the means instead.

The bootstrap for the means of the delays in the morning is provided here:

```
BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
Based on 1000 bootstrap replicates

CALL :
boot.ci(boot.out = boot_mean, type = "norm")

Intervals :
Level      Normal
95%   (24.87, 25.96 )
Calculations and Intervals on Original Scale
```

The bootstrap for the means of the delays in the midday is provided here:

```
BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
Based on 1000 bootstrap replicates

CALL :
boot.ci(boot.out = boot_mean, type = "norm")

Intervals :
Level      Normal
95%   (27.02, 30.69 )
Calculations and Intervals on Original Scale
```

The bootstrap for the means of the delays in the evening is provided here:

```
BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
Based on 1000 bootstrap replicates

CALL :
boot.ci(boot.out = boot_mean, type = "norm")

Intervals :
Level      Normal
95%   (32.26, 40.53 )
Calculations and Intervals on Original Scale
```

To provide further evidence that the medians being the same in the morning and midday is significant, we provide a Wilcoxon test, as said in the methods section from above. When the test was run, as shown below, it can be observed that with midday and mornings, we fail to reject the null hypothesis that the distributions have like medians. However, when midday and evening and morning and evening Wilcoxon tests are run, it can be observed that the p-value is less than .05, allowing us to reject the null hypothesis that the distributions have like medians.

Results for the three tests are shown below in the order of Midday/Morning, Midday/Evening, and Morning/Evening:

```
        Wilcoxon rank sum test with continuity correction

data:  How_Long_Delayed by Time_Of_Day
W = 1731300, p-value = 0.5602
alternative hypothesis: true location shift is not equal to 0


        Wilcoxon rank sum test with continuity correction

data:  How_Long_Delayed by Time_Of_Day
W = 116920, p-value = 0.0003389
alternative hypothesis: true location shift is not equal to 0


        Wilcoxon rank sum test with continuity correction

data:  How_Long_Delayed by Time_Of_Day
W = 604440, p-value = 0.0001132
alternative hypothesis: true location shift is not equal to 0
```

As seen above, the p-value when evaluating the medians of the morning and midday is .5602. At the 5% significance level, we do not have enough evidence here to reject the null hypothesis. For the midday and evening comparison, we get a p-value of .0003389, which is low enough to reject the null hypothesis at the 5% significance level. The same is done with the morning and evening comparison, since a p-value of .0001132 is provided.

The tests above are able to provide significant insight to how delay times differ by time of day. Although initially hypothesized that delay times are higher in the mornings than the

midday and evenings, the evaluation of our data provides insight that although delays appear to occur more frequently in the mornings, delay times are higher in the evening.

## VI.    Discussion:

Earlier in the research paper, it was hypothesized that the delays in the morning are the longest. However, with insight provided by the multiple methods used, it can be concluded that the mornings, in fact, do not have the longest time of delay. The Wilcoxon tests allowed us to observe that the medians are the same for the mornings and midday, but the evenings are higher. Additionally, bootstrapped confidence intervals allowed us to observe this data.

Bootstrapped confidence intervals for mean delays during different times of the day revealed that mean delays in the morning were between 24.87 and 25.96 minutes. Mean delays in the midday were between 27.02 and 30.69 minutes, and mean delays in the evening were between 32.26 and 40.53 minutes.

With this in mind, commuters throughout New York City should be cautious when commuting by car in the evening These commuters should allot enough time for travel, or look to using the subway system in order to avoid this evening traffic. Perhaps even leaving their office slightly earlier or later would allow them to avoid evening traffic.

With the knowledge provided from this research paper, government officials in New York City could begin to formulate different ways to offset traffic delays in the city. With further research, statistics could provide insight on where to create new routes, where traffic is highest, and how traffic is changing over time. The data set used in this research paper could be applied to observe changing delay trends on a year-by-year basis.

As this data set expands as years go by, further research should evaluate the hour-by-hour changes in delay trends and furthermore, how delays change by route. By expanding on what was learned in this research paper, data sets that allow urban delay to be tracked should be analyzed and evaluate whether the result is replicable. In the best case scenario, different methods of evaluating urban delay in New York City should come to the same conclusion.

## VII.    Bibliography/Works Cited:

1. City of New York. (2019). Bus Breakdown and Delays (1) [CSV File]. Retrieved from https://catalog.data.gov/dataset/bus-breakdown-and-delays
2. Yazici, M. Anil, et al. "Breakdown of Weather, Intersection and Recurrent Congestion Impacts on Urban Delay in New York City." *Transportation Research Procedia*, vol. 22, 2017, pp. 399–408., doi:10.1016/j.trpro.2017.03.009.

3. Rosenzweig, Cynthia, et al. "Developing Coastal Adaptation to Climate Change in the New York City Infrastructure-Shed: Process, Approach, Tools, and Strategies." *Climatic Change*, vol. 106, no. 1, 2011, pp. 93–127., doi:10.1007/s10584-010-0002-8.
4. Rizzo, Maria L. *Statistical Computing with R*. CRC Press, Taylor & Francis Group, 2016.

# VIII.  Supplemental Materials:

Below, I have my code snippet for which I used the tidyverse library to clean my data. I filtered the data so that it only contained delays from the five boroughs of New York City, omitted any rows where delay time was not given, and finally filtered so that only traffic as a reason for delay would be evaluated. I also converted the given date and time that was occurred on into only times so that it can be used to evaluate the information given. The times were sorted categorically in order to allow us to run hypothesis tests on the data.

```{r}
buses = buses %>% filter(Boro %in% c("Manhattan", "Bronx", "Brooklyn", "Queens", "Staten Island")) %>%
na.omit(How_Long_Delayed) %>% select(School_Year, Route_Number, Reason, Occurred_On, Boro, How_Long_Delayed)
buses = buses %>% filter(Reason == "Heavy Traffic")

buses = buses %>% mutate(Occurred_On = as.POSIXct(buses$Occurred_On, format = "%m/%d/%y %H:%M")) %>% mutate(x =
as.POSIXct(buses$Occurred_On, format = "%m/%d/%y %H:%M")) %>% mutate(Occurred_On = strftime(Occurred_On, format =
"%H:%M:%S")) %>% mutate(Occurred_on = as.POSIXct(buses$Occurred_On, format = "%H:%M:%S")) %>% select(School_Year,
Route_Number, Occurred_On, Boro, How_Long_Delayed, x)

buses$Boro <- droplevels(buses$Boro)
buses$School_Year = droplevels(buses$School_Year, exclude = "2018-2019")


buses$How_Long_Delayed <- gsub("([0-9]*)(.*)","\\1",buses$How_Long_Delayed)
buses$How_Long_Delayed <- gsub("0-15","15",buses$How_Long_Delayed)
buses$How_Long_Delayed <- gsub("16-30","30",buses$How_Long_Delayed)
buses$How_Long_Delayed <- gsub("31-45","45",buses$How_Long_Delayed)
buses$How_Long_Delayed <- gsub("46-60","60",buses$How_Long_Delayed)
buses <- buses[ !(buses$How_Long_Delayed %in% c(1,2,3)), ]

buses$How_Long_Delayed = as.numeric(buses$How_Long_Delayed)
buses = buses %>% na.omit(How_Long_Delayed)

buses = buses %>% mutate(Hour_Of_Day = hour(x))
buses <- buses %>%
        mutate(Time_Of_Day = case_when(
            .$Hour_Of_Day >= 6 & .$Hour_Of_Day < 10 ~ 'Morning',
            .$Hour_Of_Day >= 10 & .$Hour_Of_Day < 16 ~ 'Midday',
               TRUE ~ "Evening"
        ))
worked = buses %>% select(How_Long_Delayed, Time_Of_Day)
```

Below, code snippets used in the research paper are shown

Plot of Delay Counts by Minutes (Further Categorized by Time of Day):

```{r}
ggplot(buses, aes(x=buses$How_Long_Delayed,fill = buses$Time_Of_Day)) +
  geom_histogram(binwidth = 10)+scale_x_continuous(breaks = c(15,30,45,60,75,90,105,120,135,150))+ggtitle("Bus
Delay Length Count") + labs(title = "Bus Delay Counts", x = "Minutes", y = "Count")
```

Plot of Delay Counts by Minutes (Further Categorized by Year):

```r
ggplot(buses, aes(x=buses$How_Long_Delayed,fill = buses$School_Year)) +
  geom_histogram(binwidth = 10)+scale_x_continuous(breaks = c(15,30,45,60,75,90,105,120,135,150)) + labs(title =
"Bus Delay Counts by Year", x = "Minutes", y = "Count")
```

Box Plot of Delays by Time of Day:

```r
ggplot(buses, aes(x=buses$How_Long_Delayed,fill = buses$School_Year)) +
  geom_histogram(binwidth = 10)+scale_x_continuous(breaks = c(15,30,45,60,75,90,105,120,135,150)) + labs(title =
"Bus Delay Counts by Year", x = "Minutes", y = "Count")
```

Wilcox Tests for Delay:

```r
wilcox.test(How_Long_Delayed ~ Time_Of_Day, data = buses, subset = Time_Of_Day %in% c("Midday", "Morning"),
conf.level = .95)

#p-value extremely high, so we know that we fail to reject the null hypothesis that they have equal medians

wilcox.test(How_Long_Delayed ~ Time_Of_Day, data = buses, subset = Time_Of_Day %in% c("Midday", "Evening"),
conf.level = .95)

wilcox.test(How_Long_Delayed ~ Time_Of_Day, data = buses, subset = Time_Of_Day %in% c("Morning", "Evening"),
conf.level = .95)
```

Creating Vector of p-values With Different Thetas

```r
y = buses$How_Long_Delayed
a = buses$How_Long_Delayed[buses$Time_Of_Day == "Morning"]
b = buses$How_Long_Delayed[buses$Time_Of_Day == "Evening"]

thetas = seq(-5,5, length.out = 10)

wp_vec = vector()
power_vec = vector()
for (theta in thetas) {
  new_a = a + theta
  wp = wilcox.test(new_a, b)$p.value
  wp_vec = append(wp_vec, wp)
}
```

Power Simulations with Wilcox vs. K-S, with theta = 1.25

```{r}
wts.p = replicate(1000, {
  wts = wilcox.test(rpois(100, 20),
              rpois(100,20) + 1.25)
  wts$p.value <= 0.05
  })
binom.test(sum(wts.p), 1000, conf.level = .95)

kts.p = replicate(1000, {
  kts = ks.test(rpois(100,20),
              rpois(100,20) + 1.25)
  kts$p.value <= 0.05
})
binom.test(sum(kts.p), 1000, conf.level = .95)
```

Power Simulations:

```{r}
range(thetas[wp_vec >= .05])
```

```{r}
powers_vec = vector()
conf_95_thetas = seq(.5556, 3.8889, length.out = 100)
for (theta in conf_95_thetas){
  wts.p = replicate(1000, {
  wts = wilcox.test(rpois(100, 20),
              rpois(100,20) + theta)
  wts$p.value <= 0.05
  })
powers_vec = append(powers_vec,binom.test(sum(wts.p), 1000, conf.level = .95)$estimate)
}
```

Power Plot:

```{r}
plot(conf_95_thetas, powers_vec, type = "l", xlab = "Theta", ylab = "Power")
```

Median Delays by Different Times of Day:

```{r}
median(buses$How_Long_Delayed[buses$Time_Of_Day == "Morning"])
median(buses$How_Long_Delayed[buses$Time_Of_Day == "Midday"])
median(buses$How_Long_Delayed[buses$Time_Of_Day == "Evening"])
```

Mean Delays by Different Times of Day:

```{r}
mean(buses$How_Long_Delayed[buses$Time_Of_Day == "Morning"])
mean(buses$How_Long_Delayed[buses$Time_Of_Day == "Midday"])
mean(buses$How_Long_Delayed[buses$Time_Of_Day == "Evening"])
```

Bootstrapped Mean Confidence Intervals by Respective Time of Day:

```{r}
mean_boot = function(data, index) {
  data_star = data[index,]
  mean_bus = mean(data_star$`How_Long_Delayed`[data_star$Time_Of_Day == "Midday"])
  return(mean_bus)
}
boot_mean = boot(buses, statistic = mean_boot, R = 1000)
boot.ci(boot_mean, type = "norm")
```

```{r}
mean_boot = function(data, index) {
  data_star = data[index,]
  mean_bus = mean(data_star$`How_Long_Delayed`[data_star$Time_Of_Day == "Morning"])
  return(mean_bus)
}
boot_mean = boot(buses, statistic = mean_boot, R = 1000)
boot.ci(boot_mean, type = "norm")
```

```{r}
mean_boot = function(data, index) {
  data_star = data[index,]
  mean_bus = mean(data_star$`How_Long_Delayed`[data_star$Time_Of_Day == "Evening"])
  return(mean_bus)
}
boot_mean = boot(buses, statistic = mean_boot, R = 1000)
boot.ci(boot_mean, type = "norm")
```