

Name - Punit Mehndiratta

Analyzing the NYC Subway Dataset

Questions

Overview

This project consists of two parts. In Part 1 of the project, you should have completed the questions in Problem Sets 2, 3, 4, and 5 in the Introduction to Data Science course.

This document addresses part 2 of the project. Please use this document as a template and answer the following questions to explain your reasoning and conclusion behind your work in the problem sets. You will attach a document with your answers to these questions as part of your final project submission.

Section 0. References

Please include a list of references you have used for this project. Please be specific - for example, instead of including a general website such as stackoverflow.com, try to include a specific topic from Stackoverflow that you have found useful.

Answer - References used -

Reference 1 - Reindexing timeseries from dtype to datetime-dtype

stackoverflow.com/questions/13654699/reindexing-pandas-timeseries-from-object-dtype-to-datetime-dtype

Reference 2- Python strptime and strftime behaviour

<https://docs.python.org/2/library/datetime.html#strftime-strptime-behavior>

Reference 3 - String formatting and strip punctuation , formatting in data frames

stackoverflow.com/questions/265960/best-way-to-strip-punctuation-from-a-string-in-python

stackoverflow.com/questions/16729483/converting-strings-to-floats-in-a-dataframe

Reference 4 - ggplot documentation

docs.ggplot2.org/current/geom_histogram.html

ggplot.yhathq.com/docs/geom_histogram.html

Reference 5- Logging tutorial Python

<https://docs.python.org/2/howto/logging.html#logging-basic-tutorial>

Reference 6 - OLS in Python

statsmodels.sourceforge.net/devel/examples/notebooks/generated/ols.html

Reference no. 7- Graphpad Statistical Guide- Using One value Vs two values P test

http://graphpad.com/guides/prism/6/statistics/index.htm?one-tail_vs_two-tail_p_values.htm

Reference no. 8- - Statsoft information on Multiple Regression

<http://www.statsoft.com/Textbook/Multiple-Regression#residual>

Section 1. Statistical Test

1.1 Which statistical test did you use to analyze the NYC subway data? Did you use a one-tail or a two-tail P value? What is the null hypothesis? What is your p-critical value?

Answer -

We will be using the two tailed p value , as we don't know the direction of our p values in advance of test.

In order to analyze the number of entries at subway Vs weather(rain or no rain) , we performed Mann-Whitney U test, since the distribution was not normal.

Mann-Whitney U test provided one tailed p value of $p = 0.025$, so the two tailed p value will be $p = +0.05$ or -0.05 which is lower than or equal to our p-critical value ($p \leq 0.05$ two sided p critical) and number of entries were significant i.e. > 20 hence the null hypothesis is considered not to be true.

The Null hypothesis suggest that there is no statistical difference between Number of Hourly entries when it rain or it doesn't rain.

1.2 Why is this statistical test applicable to the dataset? In particular, consider the assumptions that the test is making about the distribution of ridership in the two samples.

Answer-

Since the distribution is not normal (based on Histogram plot for the data as well as results of Shapiro-Wilk test shows it not to be a normal or probability density function), we can't perform

the Welch's t-test, hence we performed Mann-Whitney U test which is applicable for testing two populations with unknown distribution.

Here we assume that the number of Hourly entries at one particular time interval were independent of number of Hourly entries at another time interval (i.e. the two data did not influence each other hence allowing us to Use Mann-Whitney Test)

Also the number of entries are significant (> 20) hence the Mann-Whitney U test provides a good statistical test. It showed that there is statistical difference between subway ridership on Rainy Vs Non rainy days.

1.3 What results did you get from this statistical test? These should include the following numerical values: p-values, as well as the means for each of the two samples under test.

Answer -

Mean of Hourly entries when Rain =1105.4463767458733

Mean of Hourly entries when not raining = 1090.278780151855,

u value for the data =1924409167.0

p value = 0.05 (which is \leq p critical value of 0.05 for two sided test)

1.4 What is the significance and interpretation of these results?

Answer -

The Test showed that there is statistical and significant difference between subway ridership on Rainy Vs Non rainy days. Hence null hypothesis is not true.

On average there is more ridership during rainy days rather than non rainy days.

Section 2. Linear Regression

2.1 What approach did you use to compute the coefficients theta and produce prediction for ENTRIESn_hourly in your regression model:

1. Gradient descent (as implemented in exercise 3.5)
2. OLS using Statsmodels
3. Or something different?

Answer - Gradient Descent

Also i have used OLS using statsmodel

2.2 What features (input variables) did you use in your model? Did you use any dummy variables as part of your features?

Answer -

The features(input variables) used for the analysis are

'rain' (whether it rained or not)

'percipi' (precipitation at that time)

'Hour' (hour of the day)

'mintempi' (minimum temperature)

'fog' (foggy or not)

'meanwindspd' (mean wind speed)

The dummy variables used is 'UNIT' (representing different stations).

2.3 Why did you select these features in your model? We are looking for specific reasons that lead you to believe that

the selected features will contribute to the predictive power of your model.

- Your reasons might be based on intuition. For example, response for fog might be: "I decided to use fog because I thought that when it is very foggy outside people might decide to use the subway more often."
- Your reasons might also be based on data exploration and experimentation, for example: "I used feature X because as soon as I included it in my model, it drastically improved my R2 value."

Answer -

Based on Intuition (ex- fog and rain as well as cold temperature cause people to use subway more often generally) and based on data experimentation, i decided to try certain variables and check whether it improved my R2 value which it did in many cases. Adding 'fog' . Hour, mintempi and meanwindspd to rain, percipi and Hour variables improved the R2 value to 0.4655. (using gradient descent)

Using OLS the value is 0.485 (somewhat increased) .

2.4 What are the coefficients (or weights) of the non-dummy features in your linear regression model?

Answer -

The weights of non-dummy features were represented as below obtained using OLS (by removing unit from the feature and calculating the coefficients)

Section 3. Visualization

Please include two visualizations that show the relationships between two or more variables in the NYC subway data.

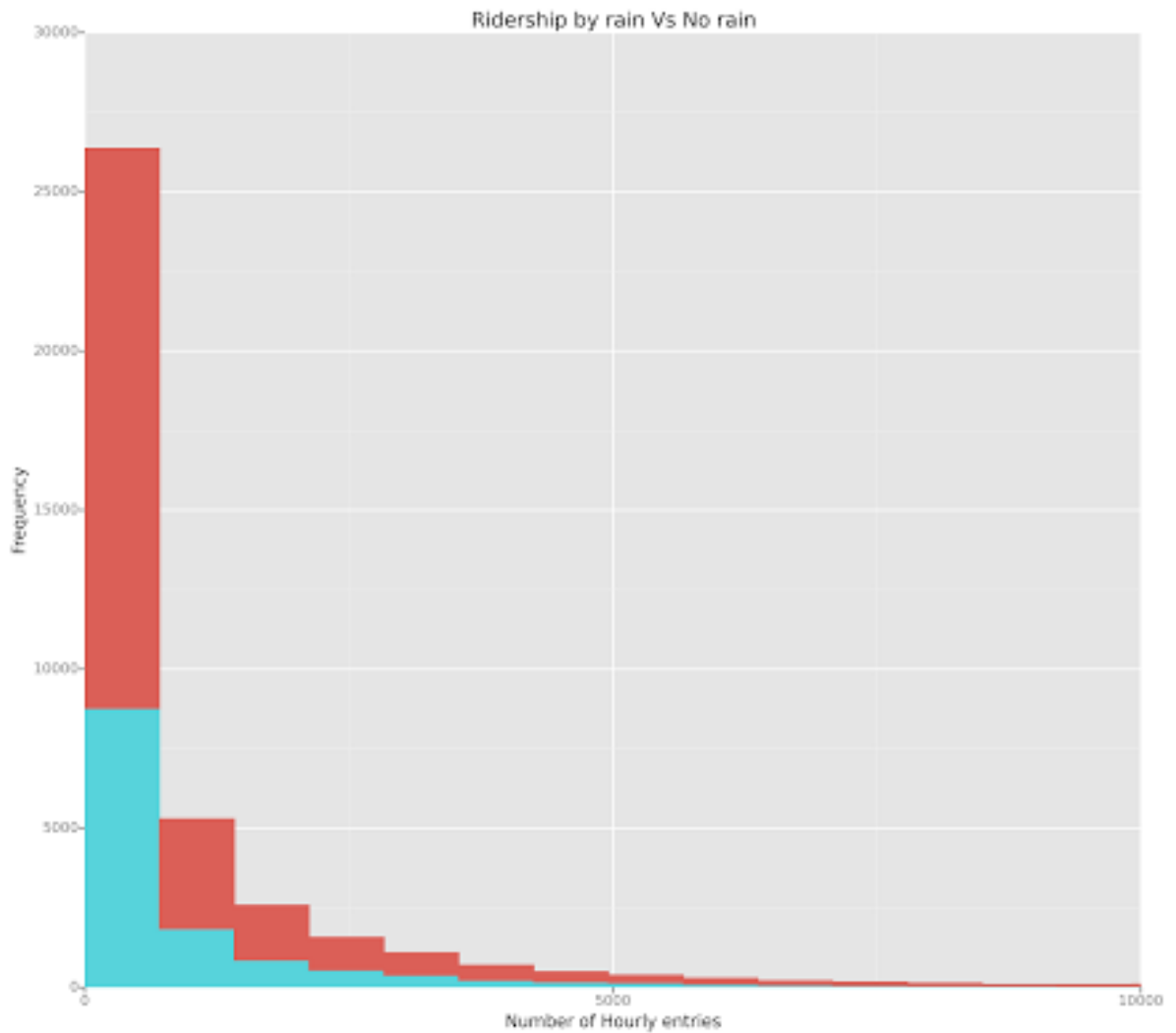
Remember to add appropriate titles and axes labels to your plots. Also, please add a short description below each figure commenting on the key insights depicted in the figure.

3.1 One visualization should contain two histograms: one of `ENTRIESn_hourly` for rainy days and one of `ENTRIESn_hourly` for non-rainy days.

- You can combine the two histograms in a single plot or you can use two separate plots.
- If you decide to use two separate plots for the two histograms, please ensure that the x-axis limits for both of the plots are identical. It is much easier to compare the two in that case.
- For the histograms, you should have intervals representing the volume of ridership (value of `ENTRIESn_hourly`) on the x-axis and the frequency of occurrence on the y-axis. For example, each interval (along the x-axis), the height of the bar for this interval will represent the number of records (rows in our data) that have `ENTRIESn_hourly` that falls in this interval.
- Remember to increase the number of bins in the histogram (by having larger number of bars). The default bin width is not sufficient to capture the variability in the two samples.

Answer -

Histogram for Ridership when it Rains Vs Not raining (binwidth =700)



ENTRIESn_hourly - **X axis**

Frequency - **Y axis**

Ridership with 'rain' - **Blue**

Ridership when its not raining - **Red**

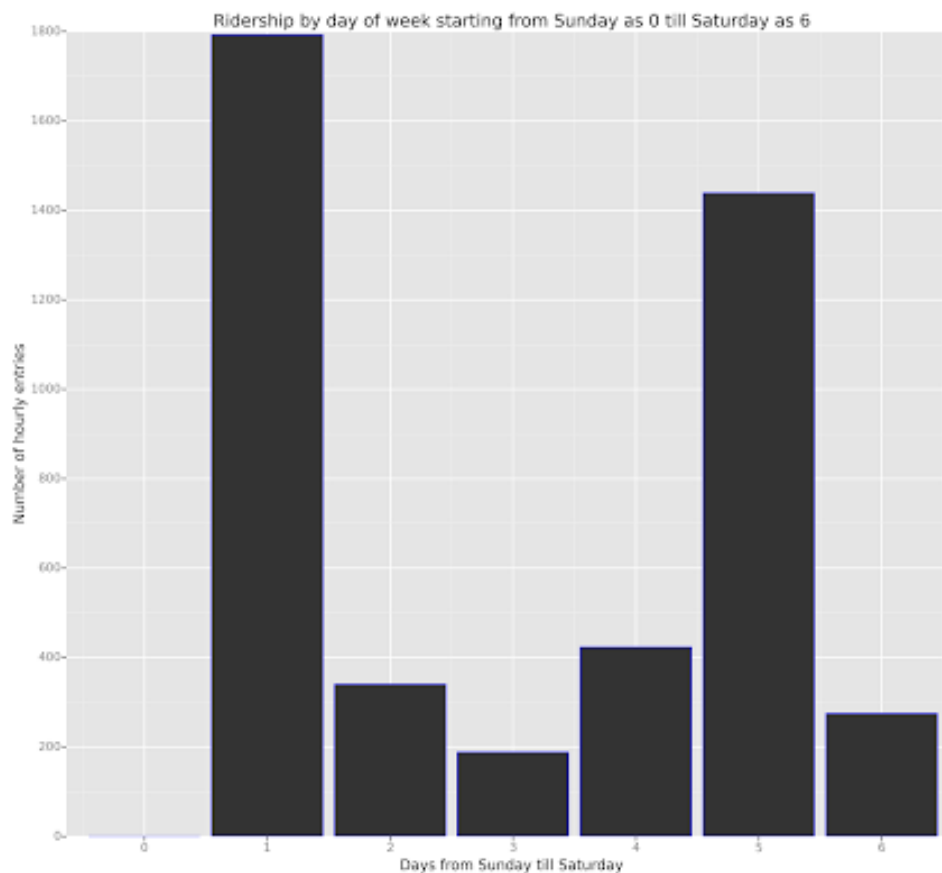
The above figure shows that more people were riding subway on 'Non rainy' time in general than when it was raining. (I rather mention it as Non rainy time rather than 'Non rainy days') Also it shows that the ridership plot is not a normal distribution.

3.2 One visualization can be more freeform. You should feel free to implement something that we discussed in class (e.g., scatter plots, line plots) or attempt to implement something more advanced if you'd like. Some suggestions are:

- Ridership by time-of-day
- Ridership by day-of-week

Answer - Below are many visualization for ridership by day of the week , or day of the week with variability of Rain , or ridership by Hour of the day (Line plot and Histogram)

Visualization no. 1 - Ridership by Day of the week - Histogram



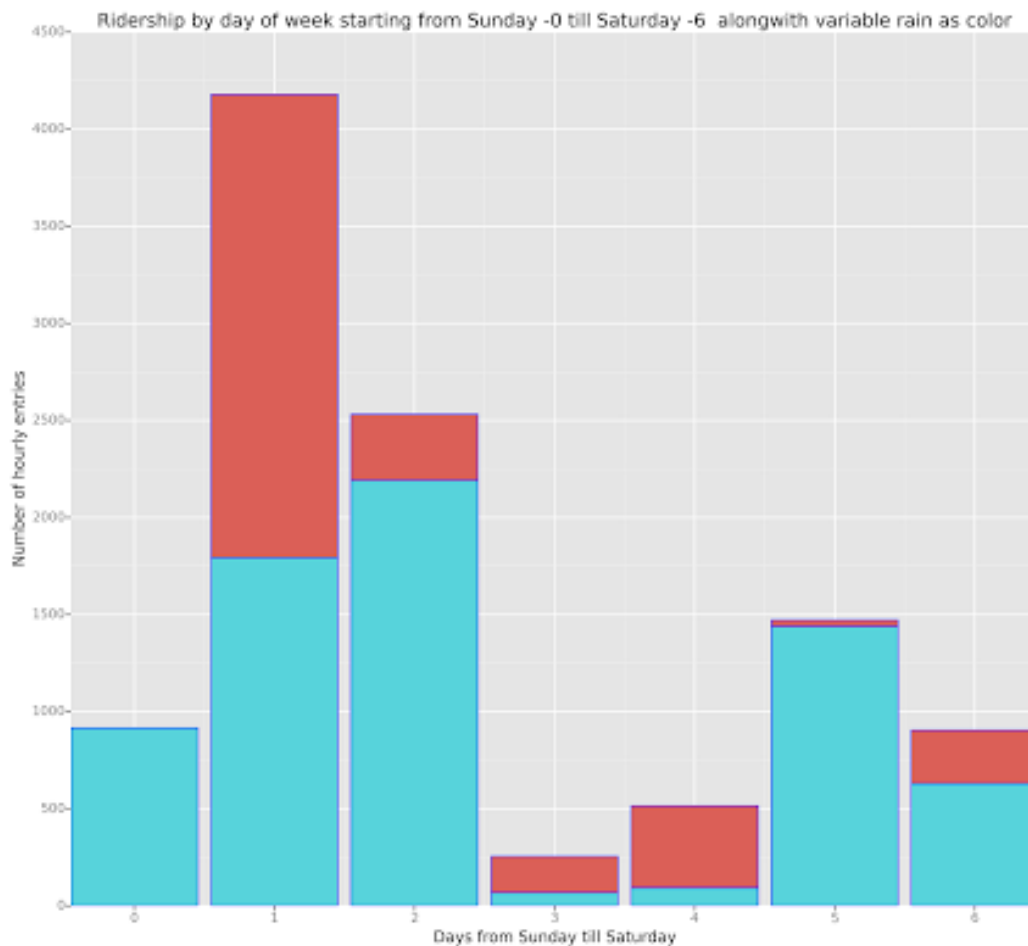
Day of the week starting from Sunday as '0, Monday as '1' till Saturday as '6
- X axis

ENTRIESn_hourly - **Y axis**

As we can see based on above visualization , the busiest days are Monday , while the least busy day is Sunday.

In the above plot we have not used 'rain' as a variable.

Visualization no. 2- Ridership by Day of the week - with filling variable 'rain'



Day of the week starting from Sunday as '0, Monday as '1' till Saturday as '6

- **X axis**

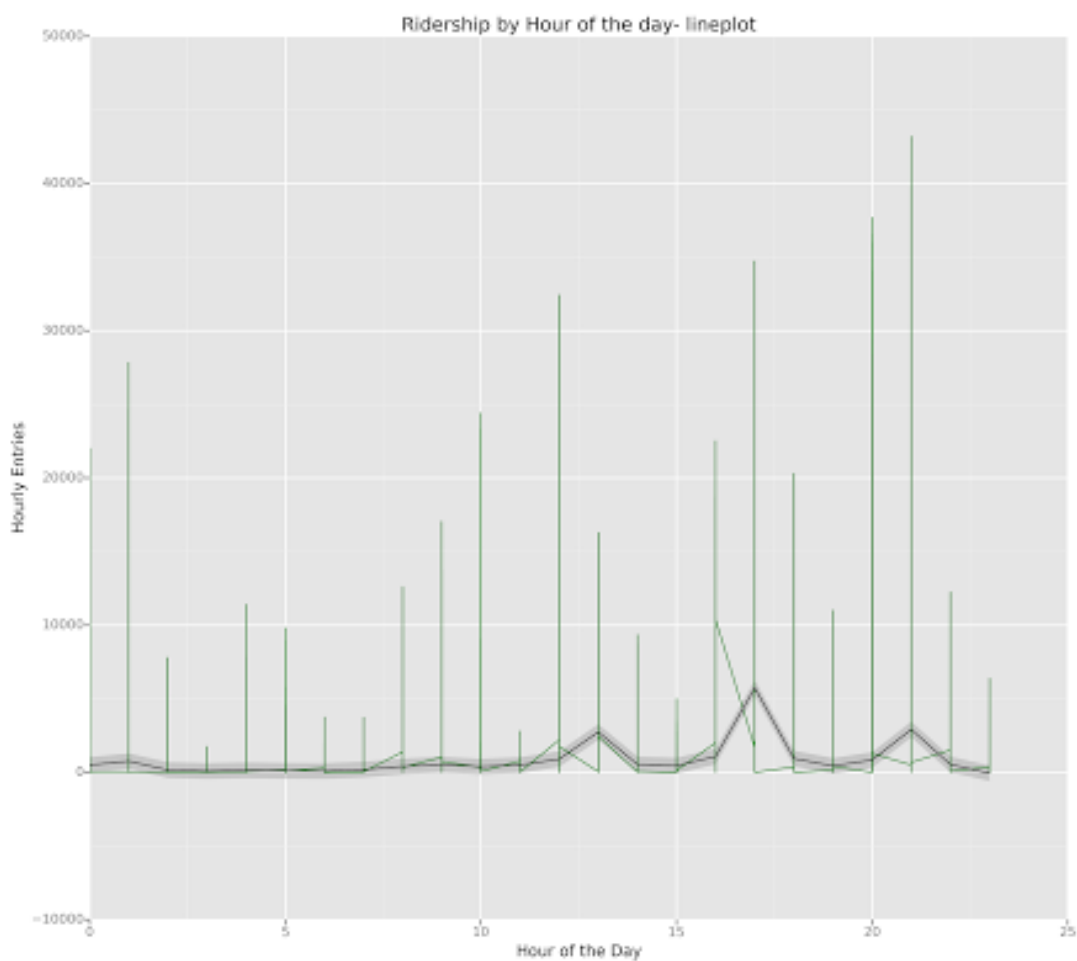
ENTRIESn_hourly - **Y axis**

Ridership with 'rain' - **Blue**

Ridership when its not raining - **Red**

Here we see that though more people take subway during rain' - Blue color, still on Monday we have people taking subway more (even when it was not raining). While on weekend , we have less people taking subway and mostly when its raining.

Visualization no. 3- Ridership by Hour of the Day - Lineplot

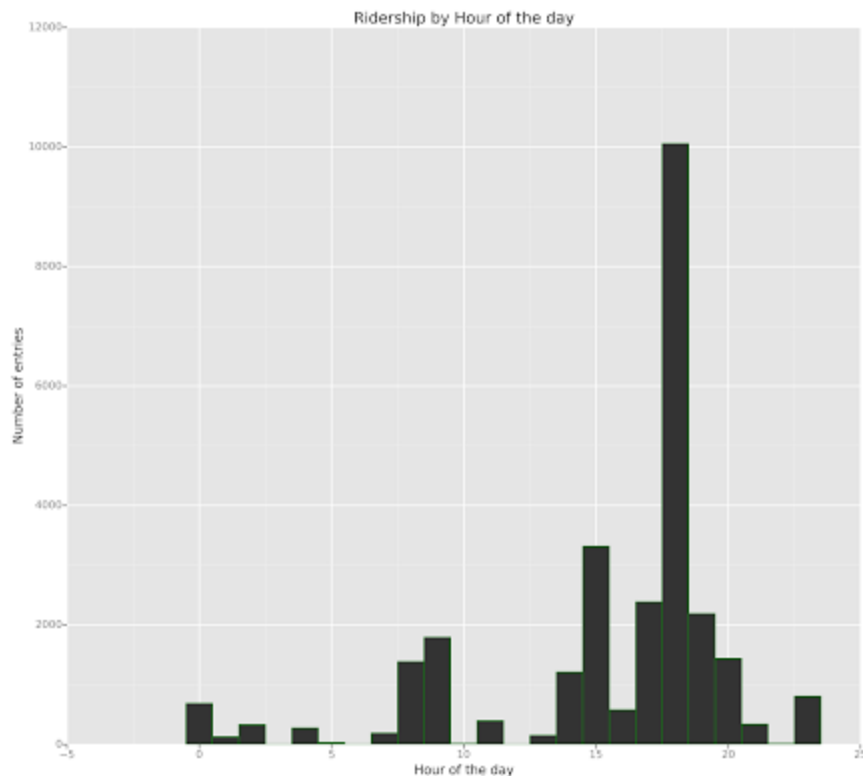


Hour of the day starting from midnight
ENTRIESn_hourly -

X axis
Y axis

Based on ridership plot above , we can see that more people use subway in the evening hours in between 5:30 - 6:30 pm time frame. While the subway ridership is minimum during certain time intervals (late night , early morning ,noon)

Visualization no. 4- Ridership by Hour of the Day -Histogram (same variables as Visualization no. 3 except this one is a Histogram)



Hour of the day starting from midnight **X axis**
ENTRIESn_hourly - **Y axis**

Based on ridership plot above , we can see that more people use subway in the evening hours in between 5:30 - 6:30 pm time frame. While the subway ridership is minimum during certain time intervals (late night , early morning ,noon)

Section 4. Conclusion

Please address the following questions in detail. Your answers should be 1-2 paragraphs long.

4.1 From your analysis and interpretation of the data, do more people ride the NYC subway when it is raining or when it is not raining?

Answer -

More people ride the subway when its raining than when its not.

4.2 What analyses lead you to this conclusion? You should use results from both your statistical tests and your linear regression to support your analysis.

Answer -

Here we have used the standard set of data.

Based on the statistical analysis(Mann-Whitney U test) of hourly entries (ridership) for the hourly interval when its raining versus when its not raining we can see that the Mann Whitney test provides us with a p value (standard data set) of 0.0249 which is less than the p critical value. Because of this we reject the null hypothesis that there is no statistical difference between ridership when its raining versus not raining. Also we can see that the average ridership (U mean value) is higher(1105.446) when its raining versus when its not raining. (1090.27).

Based on regression analysis(using gradient descent method) with normalized features and using the features which impact the ridership the most - i.e. rain , fog , time of the day , unit , minimum temperature , mean wind speed , we can see that our coefficient of determination is quite good (almost 0.5) which shows that the above factors which includes rain can be used to predict fairly the ridership in the NYC subway.

The above analysis shows that NYC subway ridership is higher when its raining Vs when its not raining.

Section 5. Reflection

Please address the following questions in detail. Your answers should be 1-2 paragraphs long.

5.1 Please discuss potential shortcomings of the methods of your analysis, including:

1. Dataset,
2. Analysis, such as the linear regression model or statistical test.

Answer -

The shortcoming of dataset (using original dataset , as that was used in the analysis),

- The dataset used contain information of data in May 2011 which is one of the heaviest rainy season in NYC.
- Also there are some units which have very heavy ridership in NYC than the others. (ex- R008 is much less busy unit (station) than R170 irrespective of whether it rained or not rained

The dataset would have been better if we checked the data for a month also when it does not rain that much or if we have data from another year for same month. Thats the only shortcoming i can see.

The shortcoming of analysis method (using original dataset)-

The linear regression method assumed linear relationship between the dependent and the independent variables, but as we can see from the plots , there are variables such as time of the day which have non linear relationship with the ridership.

The linear regression method used by me (gradient descent) provides a fairly good model but i am unable to improve the R2 value more than 0.46.

Through using OLS i find the R2 value increased to 0.485

5.2 (Optional) Do you have any other insight about the dataset that you would like to share with us?

During the OLS model analysis of subway data , i find the warnings related to multicollinearity issues . While checking the correlation between various variables used by me in the model , i found out that the correlation between variables 'rain' , 'fog' and 'precipi' are around 0.43(between rain and fog) and 0.56 (between rain and precipi).