

Section 0 – References

<https://bespokeblog.wordpress.com/2011/07/11/basic-data-plotting-with-matplotlib-part-3-histograms/>
<http://www.datarobot.com/blog/ordinary-least-squares-in-python/>
<http://stackoverflow.com/questions/30244251/plt-hist-errors-on-subsetted-data>
<http://stats.stackexchange.com/questions/116315/problem-with-mann-whitney-u-test-in-scipy>
<http://www.weather-and-climate.com/average-monthly-precipitation-Rainfall-inches,New-York,United-States-of-America>

Section 1 – Statistical Test

We want to investigate the relationship between NYC subway ridership and time and weather. Prior to further analysis, descriptive statistics of the dataset was examined to ensure no missing data, obvious data errors, and to get a general sense of the data structure. The time span of our data is the month of May, 2011, recorded at 4-hour time intervals. Sample size is 42649. The mean number of entries hourly is 1886.59, with standard deviation of 2952.39, and max of 32814 (descriptive statistic output at end of report). We'll slice the data further in various ways to examine how, if at all, the weather affects riders' behaviors.

To answer the question whether rain affects ridership in a statistically significant way, Mann-Whitney-U test and was performed (c1) on our two samples, ridership during rainy time versus no rain. We are not assuming equal variance. At critical value of 5%, we reject null hypothesis that the two samples come from same distribution (two-tailed U-test: $p=5.48e-6$). We cannot use t-Test on the dataset because our data is not normally distributed, can be seen clearly with skewed histogram (Figure1). Now we know when it rains, more people ride the subway. Let's consider the amount of rain, defined by precipitation. Light rain is less than or equal to 0.03, and heavy rain is more than 0.03. U-test again, and two-tail p-value is $1.32e-39$. Reject null hypothesis again. So we can see that when it rains lightly (≤ 0.03), the subway stations are busiest.

Sample	Mean of Hourly Entries
No Rain	1845.54
Yes Rain	2028.20
Light Rain	2151.04
Heavy Rain	1246.01

Section 2 -- Linear Regression

Given what we learned in part 1, we want to further investigate the relationship between ridership and various variables in the dataset, and choose predictor(s) to forecast mean hourly entries. I first examined potential variables individually, with gradient decent, I gathered their R^2 values, summarized below. Using combination of the four variables with highest R^2 ('rain', 'temp', 'weekday', 'hour'), our prediction model has R^2 of 0.1045. These four variables make sense that they have higher R^2 than

others. Weather condition as well as rush hours or not, intuitively help explain the mean number of hourly entries. (C2_GradientDescent)

Variable	R^2
'precipi'	0.000766
'meanprecipi'	0.001271
'rain'	0.000667
'tempi'	0.00803
'weekday'	0.02115
'hour'	0.08225

(Gradient Descent) $\text{ENTRIESn_hourly} = 50.294 \cdot \text{rain} + 539.590 \cdot \text{tempi} + 1518.18 \cdot \text{weekday} + 2007.48 \cdot \text{hour}$

I also tried using OLS Regression to fit data. Only using weather-related variables yield R^2 as low as Gradient Descent. After examining various combinations, these variables are used to predict ENTRIESn_hourly : intercept, rain, tempi, weekday and hour. All variables are statistically significant, and R^2 is 0.104. (C2_OLS)

(OLS) $\text{ENTRIESn_hourly} = -391.42 + 80.62 \cdot \text{rain} + 5.87 \cdot \text{tempi} + 951.90 \cdot \text{weekday} + 120.40 \cdot \text{hour}$

With dummy variables, the R^2 for both gradient descent and OLS improve dramatically to 0.377 and 0.481 respectively.

(OLS) $\text{ENTRIESn_hourly} = -316.83 + 56.08 \cdot \text{rain} + 3.52 \cdot \text{tempi} + 975.28 \cdot \text{weekday} + 122.21 \cdot \text{hour} + \text{dummy}$

After trying both regression models, I do not believe linear models are appropriate to model our data. Neither methods have an attractive R^2 , although all variables are statistically significant.

Section 3 – Visualization

Now, let's examine the data graphically. Histograms and scatter plots are useful tools. Histogram helps investigate frequency and skew of data. We know our ENTRIESn_hourly is highly skewed right, histogram can show graphically how the skew is distributed. Scatter plot can help examine relationship between variables. The variable of interest here are ENTRIESn_hourly and hour, only looking at the the station where max ENTRIESn_hourly in the dataset occurs. How does ENTRIESn_houly distribute given hour? Do they peak at rush hour? Do we have an outlier just one day that's skewing the data? A scatter plot can help answer these questions.

Figure 1 shows the histogram of ENTRIESn_hourly , when it rains versus when it doesn't rain. The shape of the two distributions are similar, both are skewed very right. Hence we can only conclude that rain does not have a strong effect on the distribution of ridership. Figure 2 zooms in on the right tail.

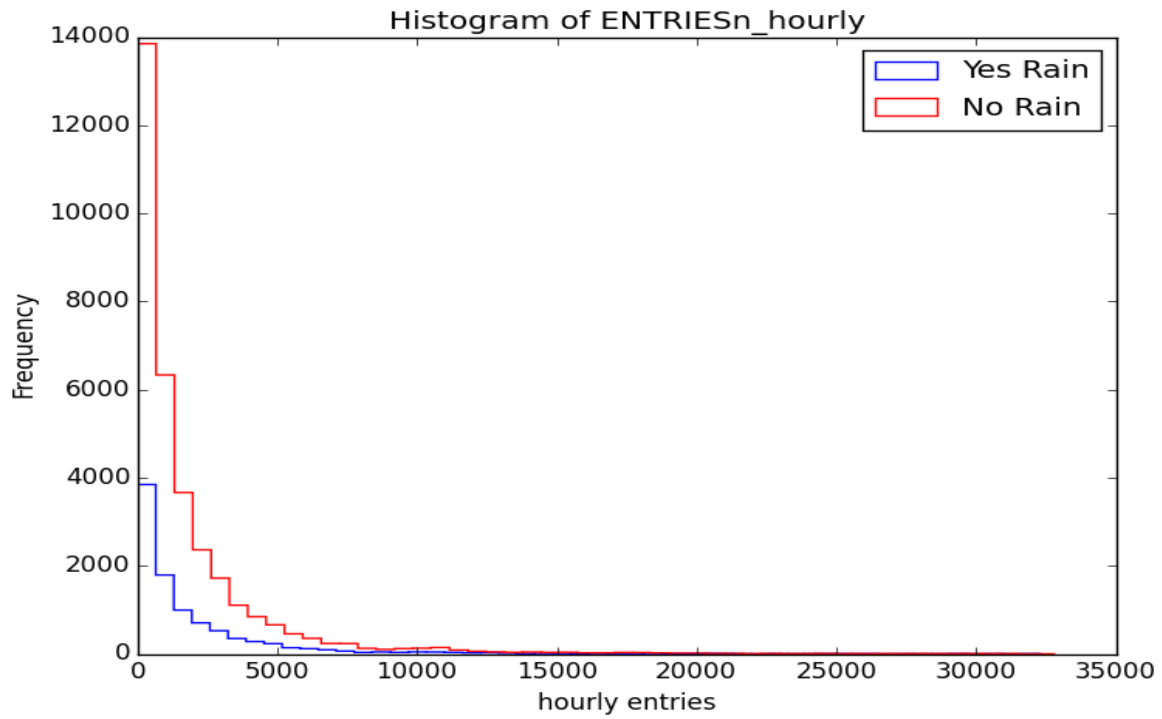


Figure 1

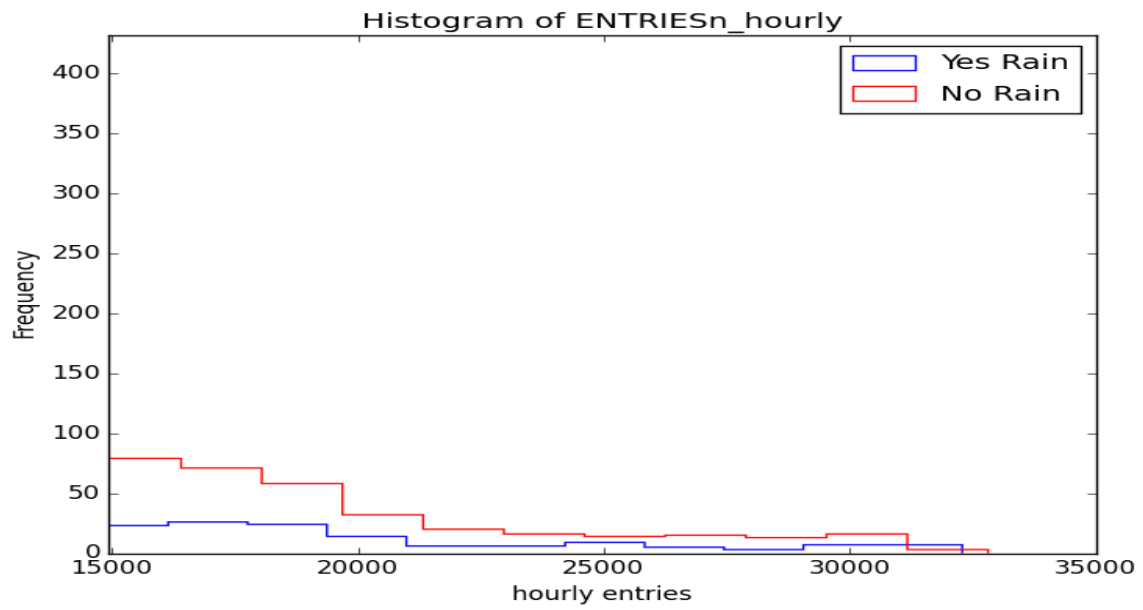


Figure 2

Figure 3 is a scatter plot of ENTRIESn_hourly by hour at the station where the max ENTRIESn_hourly lie, to examine its behavior (R084). The horizontal line indicates the 75th percentile ENTRIESn_hourly of the dataset, and we can see that this station is above the line on almost record. The busiest time at R084 is 8pm, and from noon to midnight, all the records are above the 75th percentile line.

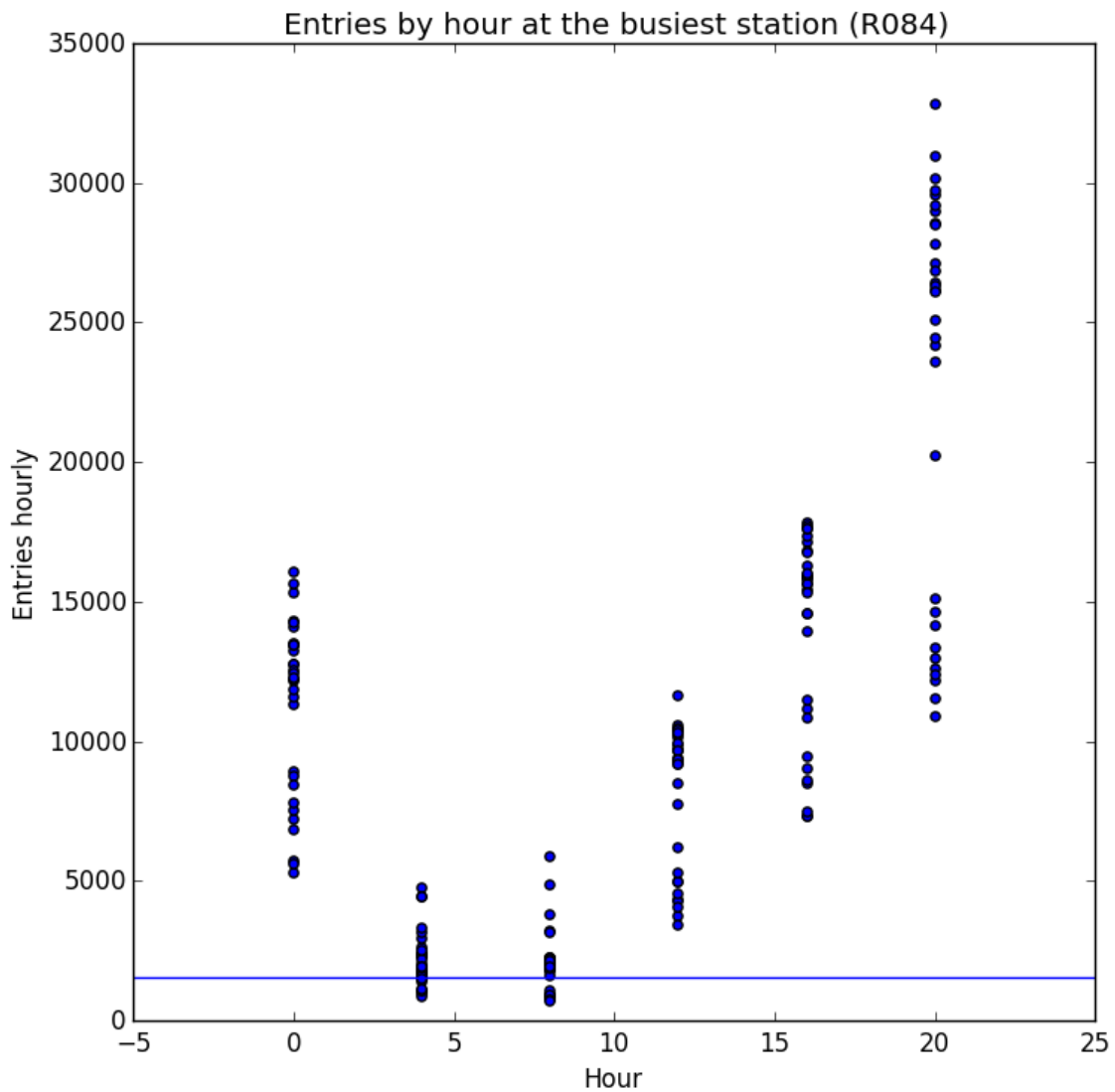


Figure 3

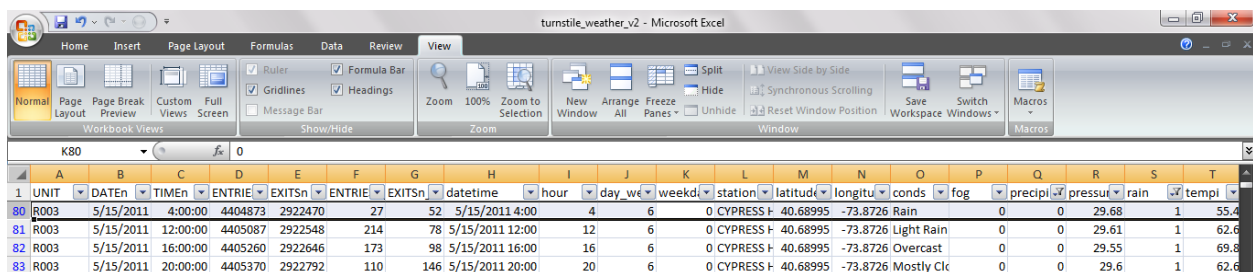
Section 4 – Conclusion

The question we are trying to answer is, do more people ride the NYC subway when it rains. Our statistical test confirms that the answer is yes. We showed that the mean of ridership is statistically different whether it rains, namely higher when it rains. Counter-intuitively, with light rain NYC subway is statistically significantly busier than heavy rain. I believe this is partially due to a much smaller sample size (8 heavy rain observations).

Regression model and statistical tests give conflicting results. Regression shows little effect on ridership. Low R^2 suggests rain does not explain much of the variability in ridership. I tried gradient descent as well as OLS, both yield the same conclusion. However statistical test indicates that rain versus no rain gives statistically significantly different mean. Confounding variables might be the reason for the conflicting results. Regression analysis adjusts out the effects of confounding by holding other variables constant in calculation, but U-test does not take confounding variables into account.

Section 5 – Reflection

Maybe consider focusing data on one station, to minimize effects of geographical impacts of ridership, and have longer time span of data. Our observations might be unusually low in precipitation. The average May precipitation is approximately 4 inches, and the average of our dataset is 0.005 inches. Maybe rain would be a more significant predictor if we had more precipitation in our observed data. Also, the rain column and precipitation seem contradicting. The rain column would indicate 1 (yes rain), and precipitation column would indicate 0 (0 inches of rain observed). Also ridership data is recorded hourly whereas weather data is spanned from daily records. This does not accurately associate hourly spike in ridership with lagged weather observation. If it rains at any time during the day, all the ridership data reflect “rain”.



	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
	UNIT	DATEn	TIMEn	ENTRIEn	EXITSn	ENTRIEn	EXITSn	datetime	hour	day_wé	weekd	station	latitude	longitu	conds	fog	precipi	pressur	rain	tempi
80	R003	5/15/2011	4:00:00	4404873	2922470	27	52	5/15/2011 4:00	4	6	0	CYPRESS	40.68995	-73.8726	Rain	0	0	29.68	1	55.4
81	R003	5/15/2011	12:00:00	4405087	2922548	214	78	5/15/2011 12:00	12	6	0	CYPRESS	40.68995	-73.8726	Light Rain	0	0	29.61	1	62.6
82	R003	5/15/2011	16:00:00	4405260	2922646	173	98	5/15/2011 16:00	16	6	0	CYPRESS	40.68995	-73.8726	Overcast	0	0	29.55	1	69.8
83	R003	5/15/2011	20:00:00	4405370	2922792	110	146	5/15/2011 20:00	20	6	0	CYPRESS	40.68995	-73.8726	Mostly Clc	0	0	29.6	1	62.6

Descriptive Statistic:

	ENTRIESn	EXITSn	ENTRIESn_hourly	EXITSn_hourly	\
count	4.264900e+04	4.264900e+04	42649.000000	42649.000000	
mean	2.812486e+07	1.986993e+07	1886.589955	1361.487866	
std	3.043607e+07	2.028986e+07	2952.385585	2183.845409	
min	0.000000e+00	0.000000e+00	0.000000	0.000000	
25%	1.039762e+07	7.613712e+06	274.000000	237.000000	
50%	1.818389e+07	1.331609e+07	905.000000	664.000000	
75%	3.263049e+07	2.393771e+07	2255.000000	1537.000000	
max	2.357746e+08	1.493782e+08	32814.000000	34828.000000	
	hour	day_week	weekday	latitude	longitude
count	42649.000000	42649.000000	42649.000000	42649.000000	42649.000000
mean	10.046754	2.905719	0.714436	40.724647	-73.940364
std	6.938928	2.079231	0.451688	0.071650	0.059713
min	0.000000	0.000000	0.000000	40.576152	-74.073622
25%	4.000000	1.000000	0.000000	40.677107	-73.987342
50%	12.000000	3.000000	1.000000	40.717241	-73.953459
75%	16.000000	5.000000	1.000000	40.759123	-73.907733
max	20.000000	6.000000	1.000000	40.889185	-73.755383
	fog	...	pressurei	rain	tempi
count	42649.000000	...	42649.000000	42649.000000	42649.000000
mean	0.009824	...	29.971096	0.224741	63.103780
std	0.098631	...	0.137942	0.417417	8.455597
min	0.000000	...	29.550000	0.000000	46.900000
25%	0.000000	...	29.890000	0.000000	57.000000
50%	0.000000	...	29.960000	0.000000	61.000000
75%	0.000000	...	30.060000	0.000000	69.100000
max	1.000000	...	30.320000	1.000000	86.000000
	wspdi	meanprecipi	meanpressurei	meantempi	meanwspdi
count	42649.000000	42649.000000	42649.000000	42649.000000	42649.000000
mean	6.927872	0.004618	29.971096	63.103780	6.927872
std	4.510178	0.016344	0.131158	6.939011	3.179832
min	0.000000	0.000000	29.590000	49.400000	0.000000
25%	4.600000	0.000000	29.913333	58.283333	4.816667
50%	6.900000	0.000000	29.958000	60.950000	6.166667
75%	9.200000	0.000000	30.060000	67.466667	8.850000
max	23.000000	0.157500	30.293333	79.800000	17.083333
	weather_lat	weather_lon			
count	42649.000000	42649.000000			
mean	40.728555	-73.938693			
std	0.065420	0.059582			
min	40.600204	-74.014870			
25%	40.688591	-73.985130			
50%	40.720570	-73.949150			
75%	40.755226	-73.912033			
max	40.862064	-73.694176			