

## 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:

Histograms with matplotlib

<http://pandas.pydata.org/pandas-docs/stable/visualization.html#histograms>

Working with text in matplotlib

[http://matplotlib.org/users/pyplot\\_tutorial.html](http://matplotlib.org/users/pyplot_tutorial.html)

Remove non-letters from a string

[http://www.codecademy.com/wiki/Remove\\_non-letters\\_from\\_a\\_string#Python](http://www.codecademy.com/wiki/Remove_non-letters_from_a_string#Python)

SQL AVG() Function

[http://www.w3schools.com/sql/sql\\_func\\_avg.asp](http://www.w3schools.com/sql/sql_func_avg.asp)

Python Dictionaries

<http://learnpythonthehardway.org/book/ex39.html>

Arithmetics with dates

<https://docs.python.org/2/library/datetime.html>

Python ggplot examples

<https://pypi.python.org/pypi/ggplot/0.4.7>

Reading from stdin and files

[http://www.tutorialspoint.com/python/file\\_next.htm](http://www.tutorialspoint.com/python/file_next.htm)

Seaborn plotting library

[http://stanford.edu/~mwaskom/software/seaborn/tutorial/quantitative\\_linear\\_models.html#plotting-simple-regression-with-regplot](http://stanford.edu/~mwaskom/software/seaborn/tutorial/quantitative_linear_models.html#plotting-simple-regression-with-regplot)

Mann-Whitney U Test :

<http://docs.scipy.org/doc/scipy-0.15.1/reference/generated/scipy.stats.mannwhitneyu.html>

<http://wikiofscience.wikidot.com/technology1:menn-whitney-u-test>

Critical values and p values

<http://www.itl.nist.gov/div898/handbook/prc/section1/prc131.htm>

Interpretation of  $r^2$  and assessment of goodness of fit

<http://blog.minitab.com/blog/adventures-in-statistics/regression-analysis-how-do-i-interpret-r-squared-and-assess-the-goodness-of-fit>

## 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:

I used a Mann-Whitney U test, since the distributions of entries are not normal from visual check. See Fig. 1.

The null and alternative hypothesis are considered as follows:

H0 : samples of entries of rain and no\_rain scenarios come from the same distribution

HA : samples of entries of rain and no\_rain scenarios come from different distributions.

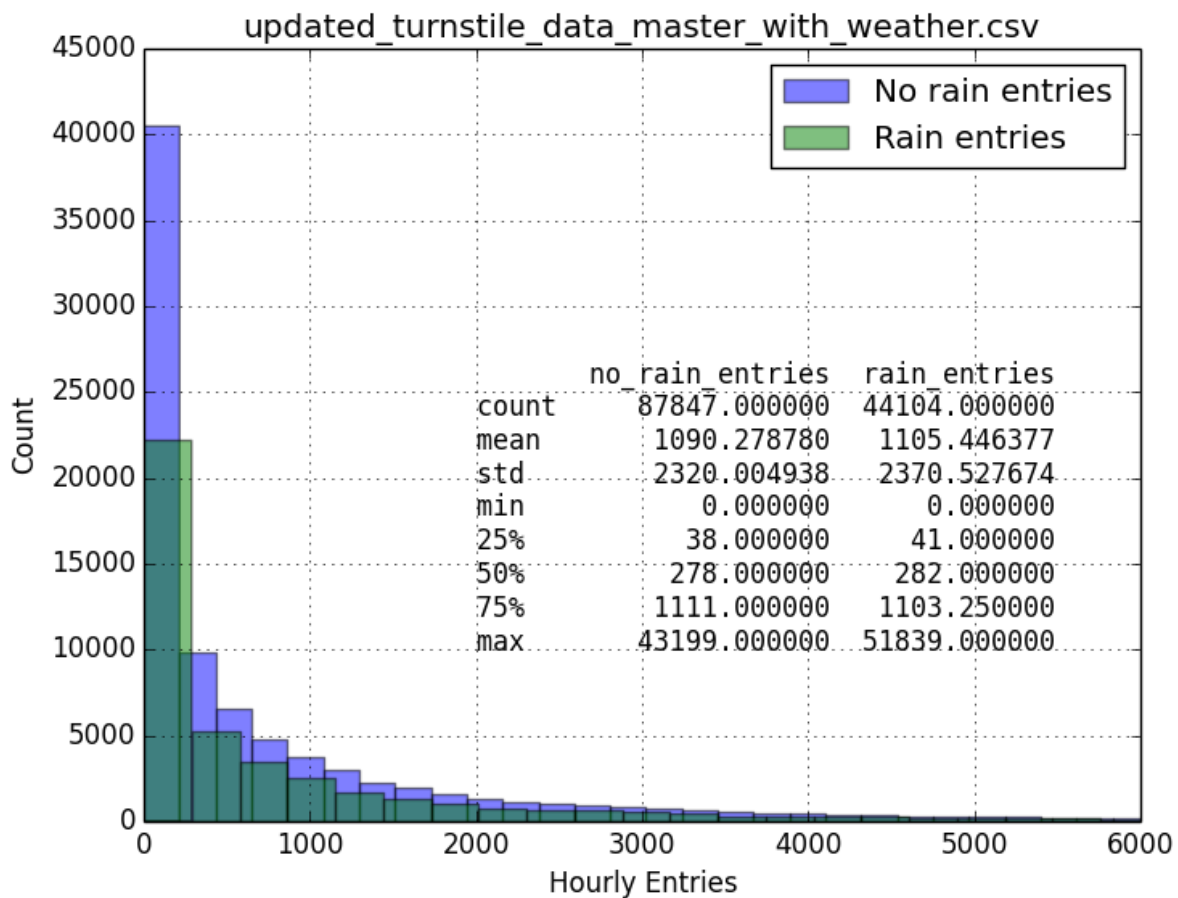


Fig.1: Histogram of turnstile entries - separated by Rain and No rain

A two-tailed test for  $\alpha = 0.05$  was considered, so p-critical is 0.025.

The scipy function `scipy.stats.mannwhitneyu()` gives p-values for one-sided, it is required to multiply the obtained p-value by 2.

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:

The Mann Whitney u test does not assume our data is drawn from any particular underlying probability distribution, and should work for this case where we cannot assume dataset follow normal distribution.

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:

Mann Whitney u test gives the following results:

(rain\_mean, no\_rain\_mean, U, p) =

(1105.4463767458733, 1090.278780151855, 1924409167.0, 0.024999912793489721)

1.4 What is the significance and interpretation of these results?

Answer:

Since the p-value is lower than 0.025 I reject the null hypothesis, that is, we conclude that entries when rain are higher than entries without rain, with a 0.05 chance that I am making a mistake.

Note this a result for a generic day, regardless if it is a working day or weekend.

Interestingly, if I restrict analysis only to working days by filtering by weekday (weekday dummy column was created and indicates whether a row has a date which is a weekend or not. See answer 2.2). the results are:

(rain\_mean, no\_rain\_mean, U, p) =

(1198.1502978290941, 1304.536235167018, 985531035.5, 3.1874387801245787e-15)

So, for a working day, non-rainy days show higher average entries. And with a very small p-value indicating that  $H_0$  can be rejected with a very low chance of making a Type I error.

On the other hand, the results for weekends, that is for (weekday == 0) are

(rain\_mean, no\_rain\_mean, U, p) =

(727.16610989517335, 686.57462735570289, 129024195.5, 0.00030968316696356743)

and the conclusions here are similar to the generic day, which is that on weekends people ride more often when it is raining.

The Mann-Whitney results for these three scenarios are summarized in Table 1.

Also, assuming normality of the distributions, the corresponding Welch two sample t-test results were calculated and are presented also in Table 1. In this case, the only null hypothesis rejection would be for working days.



	Mean Entries		p - value (Mann Whitney u test)	p-value (Welch's two sample t-test)
				
<b>generic day</b>	1105.44	1090.27	0.0249999	0.269506
<b>working day</b>	1198.15	1304.53	3.187438E-15	9.036605E-10
<b>weekend</b>	727.16	686.57	0.00030968	0.02866

Table 1: mean entries in rainy and non-rainy days

## Section 2. Linear Regression

2.1 What approach did you use to compute the coefficients theta and produce prediction for  $ENTRIES_{n\_hourly}$  in your regression model:

Gradient descent (as implemented in exercise 3.5)

OLS using Statsmodels

Or something different?

Answer:

I used Gradient descent as implemented in exercise 3.5.

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

In addition to the dummy variables related to the turnstile UNIT, I wrote a program to include day\_week & weekday to “turnstile\_data\_master\_with\_weather.csv” and saved this new file as “updated\_turnstile\_data\_master\_with\_weather.csv”. Adding weekday and day\_week to features list, increased  $r^2$  from 0.458554208399 to 0.468759744238.

From Fig. 2 that was created for Problem Set 4, I see that 30.may.2011 presented low entries. In fact, it was memorial day in the US, a national holiday. So, I added an additional dummy variable holiday,

This adds an additional column for holiday (0 = no, 1 = yes).  $r^2$  is improves to 0.470541486962. The number of iterations was kept at 75.

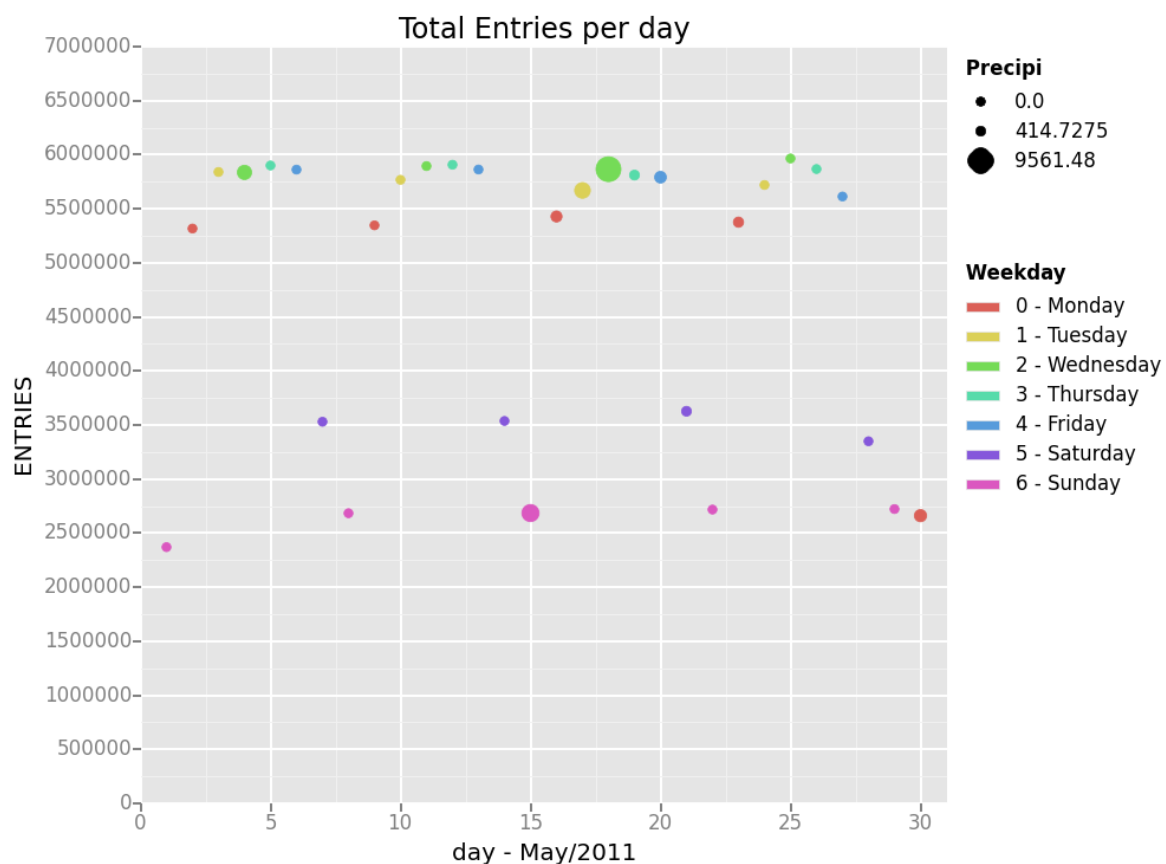


Fig.2: Entries and precipitation per day in May/2011

Finally, per the Fig.3 I see the number of entries do not show a linear relationship with the hour of day. It increases and decreases showing peaks at rush hours. So, forcing a relationship such as  $Y = \dots + \theta_H \cdot X_H + \dots$ , can lead to high errors and not help in predicting the entries. So, I decided to create additional dummy variables for the hours that were added to the feature set. Removing 'Hour' from the features and with these hour dummy variables,  $r^2$  increased to 0.514958381474.

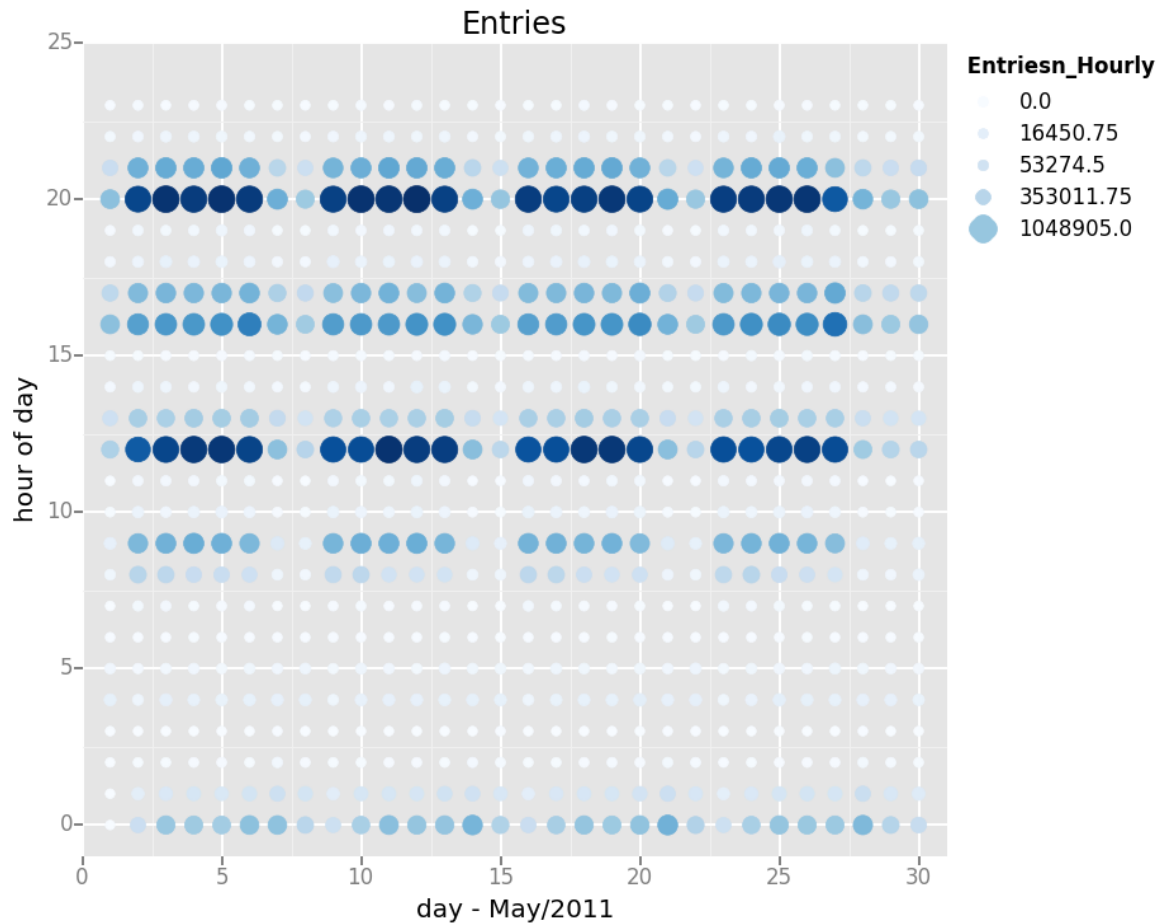


Fig.3: Entries per hour and day in May/2011

In summary, the cumulative improvements in  $r^2$  are as follows:

	$r^2$
original feature set	0.4585
add weekday and day_week	0.4687
add holiday	0.4705
add hour dummy variables	0.5150

Table 2:  $r^2$  improvements for different feature sets

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 R<sup>2</sup> value."

Answer:

The choice of weekday, day\_week and holiday and are based on observation of the graph in Fig.2 explained in 2.2 and justified by improvement of r<sup>2</sup>.

The dummy variables for UNITS came already chosen with the code already set by default. Removing them decreases the r<sup>2</sup> a lot, to around 0.0458122001164, so entries are highly dependent on the unit.

The dummy variables for hour of day were chosen, with reasons already explained in 2.2.

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

Answer:

The coefficients for the non-dummy variables are:

```
rain:    6.0959
precipi: 1.6881
mintempi: -6.9026
maxtemp: 23.0811
day_week: 19.4411
weekday: 281.2628
holiday: -119.0045
ones: 1095.34848
```

The positive value for rain agrees with the first line 'generic day' of Table 1: More rain, more riders.

Now, to check the negative effect of rain on ridership for working days, I ran regression only with samples coming from those days by restricting `dataframe = dataframe[dataframe.weekday == 1]` and removing weekday from the features list. As expected, the coefficient for rain became -1.30550547, a value consistent with line 'working day' of Table 1: More rain, less riders.

2.5 What is your model's R<sup>2</sup> (coefficients of determination) value?

Answer:

The value of r<sup>2</sup> is 0.514958381474.

2.6 What does this R<sup>2</sup> value mean for the goodness of fit for your regression model? Do you think this linear model to predict ridership is appropriate for this dataset, given this R<sup>2</sup> value?

Answer:

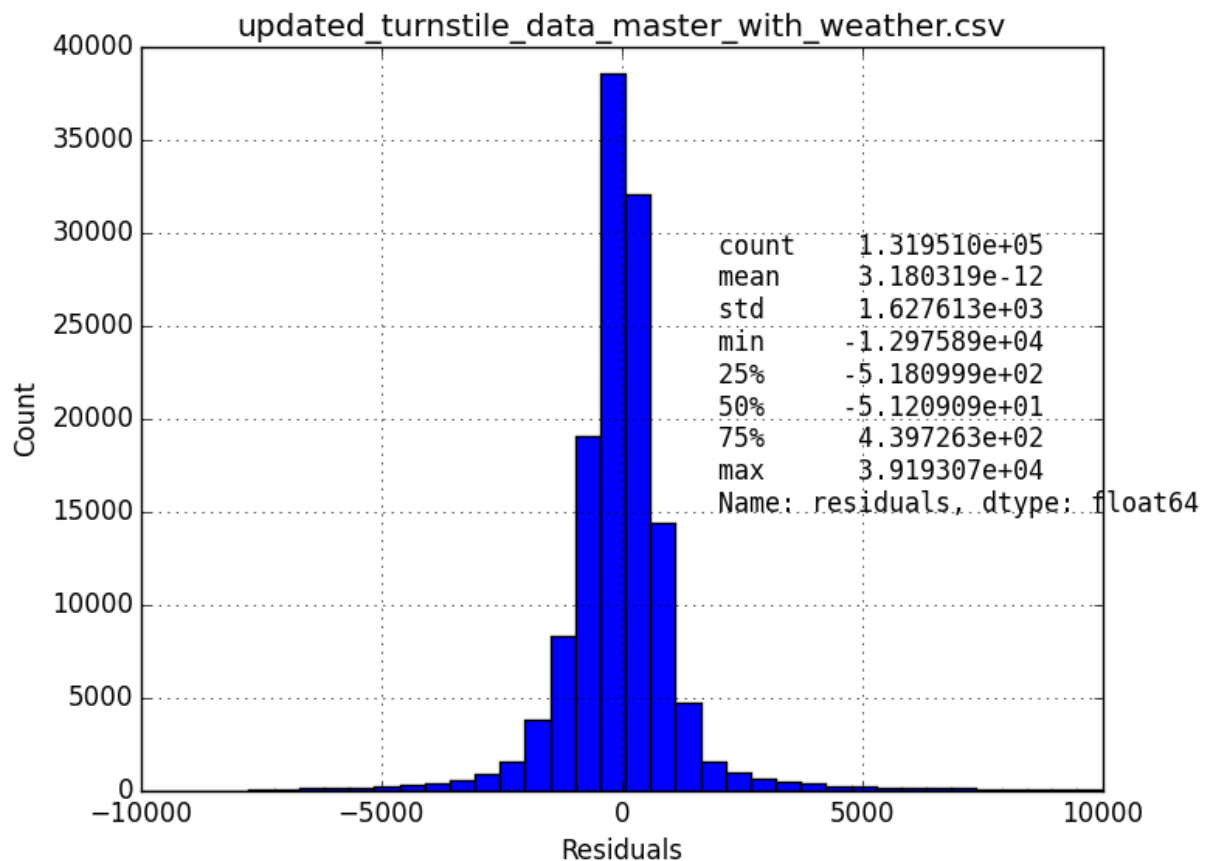


Fig.4: Residuals histogram

Residuals, calculated as the difference between the predicted and the actual values, show a normal shaped distribution, with zero mean, but long tails. These large values, for example below -518 (1st quartile) and above 440 (3rd quartile) represent 50% of the samples.

Now, if we think that for infrastructure dimensioning, capacity planning and cost estimation, we should be more concerned about high demands on the subway system rather than idle times, we should check the residuals when entries are high. I arbitrarily define “high” and “very high” as the 75th and 90th percentile of a no-rain working day distribution, which are 1373 and 3473, respectively. This residual distributions are depicted in Fig. 5. First thing we see is that mean is positive and increases as passenger traffic increases. We see our model runs the risk of underestimating ridership more often. Secondly, standard deviation also increases, achieving more than twice of the distribution of Fig. 4.

So, I conclude the linear model is not too reliable, which is also supported by the low value of  $r^2$  at 51%. More bad news: at busier riderships, just when we would need most an accurate prediction, our model gets less reliable.

### 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.

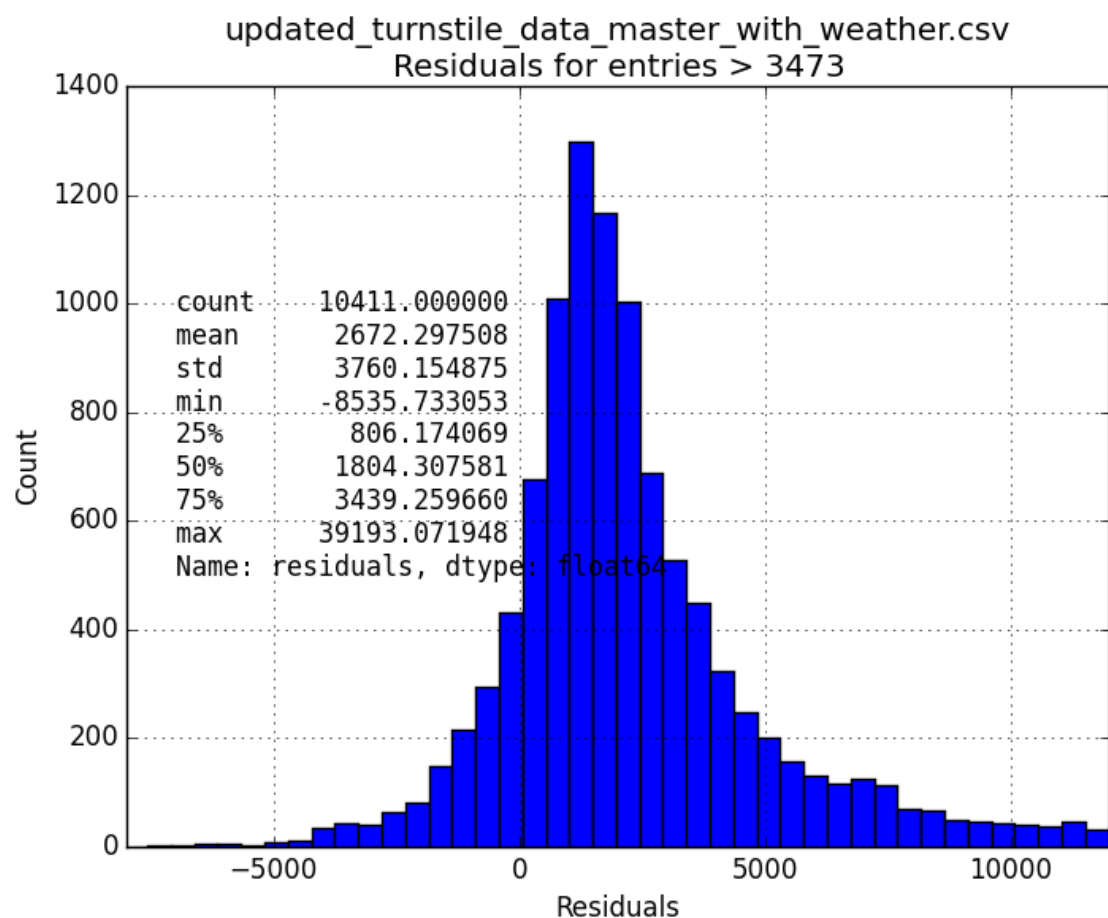
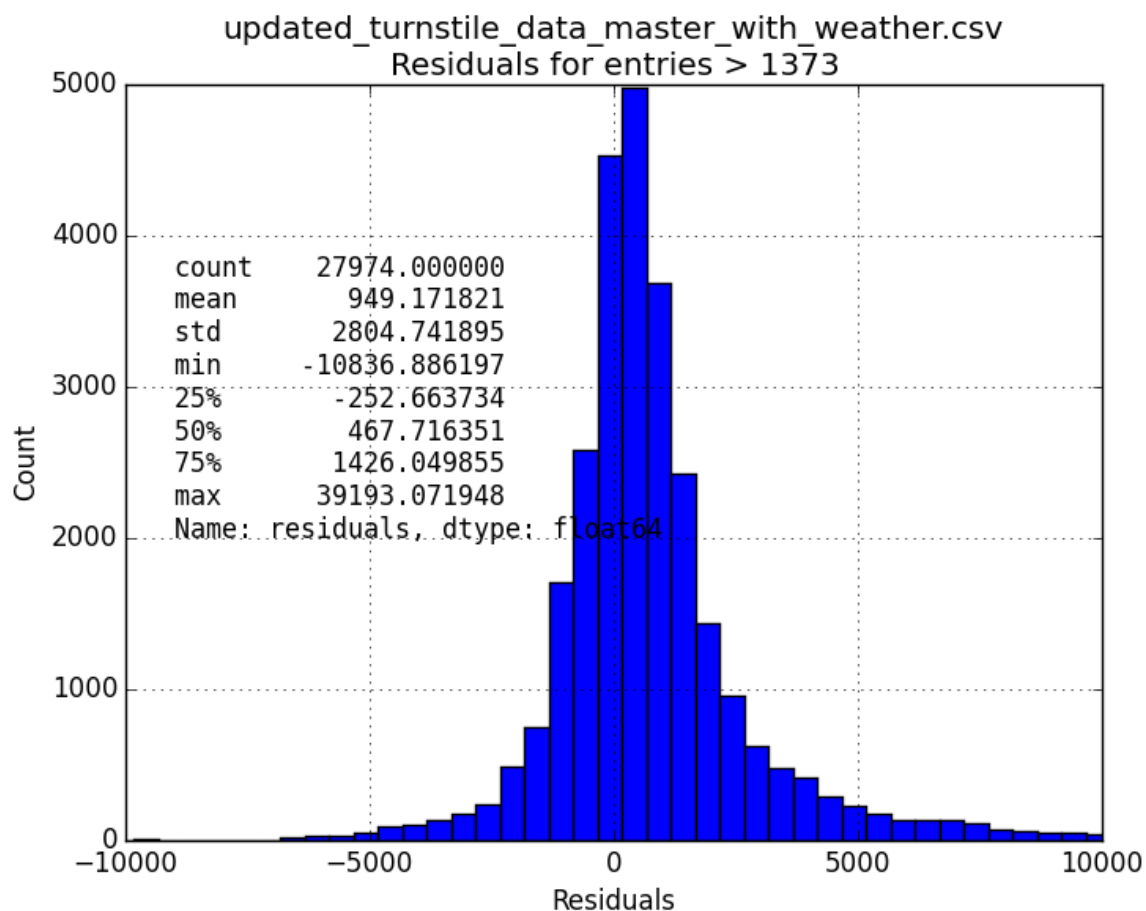


Fig.5: Residuals histograms for entries above 1373 and 3473



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:

The histogram in Fig.1 provides the visualization requested.

This visualization shows the counts of entries corresponding to rain and no-rain in defined ranges (bin sizes). I see that it is not shaped as a normal distribution, and that no-rain entries dominates at lower values.

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:

Examples of plots with conclusions were presented in Fig. 2 and Fig. 3.

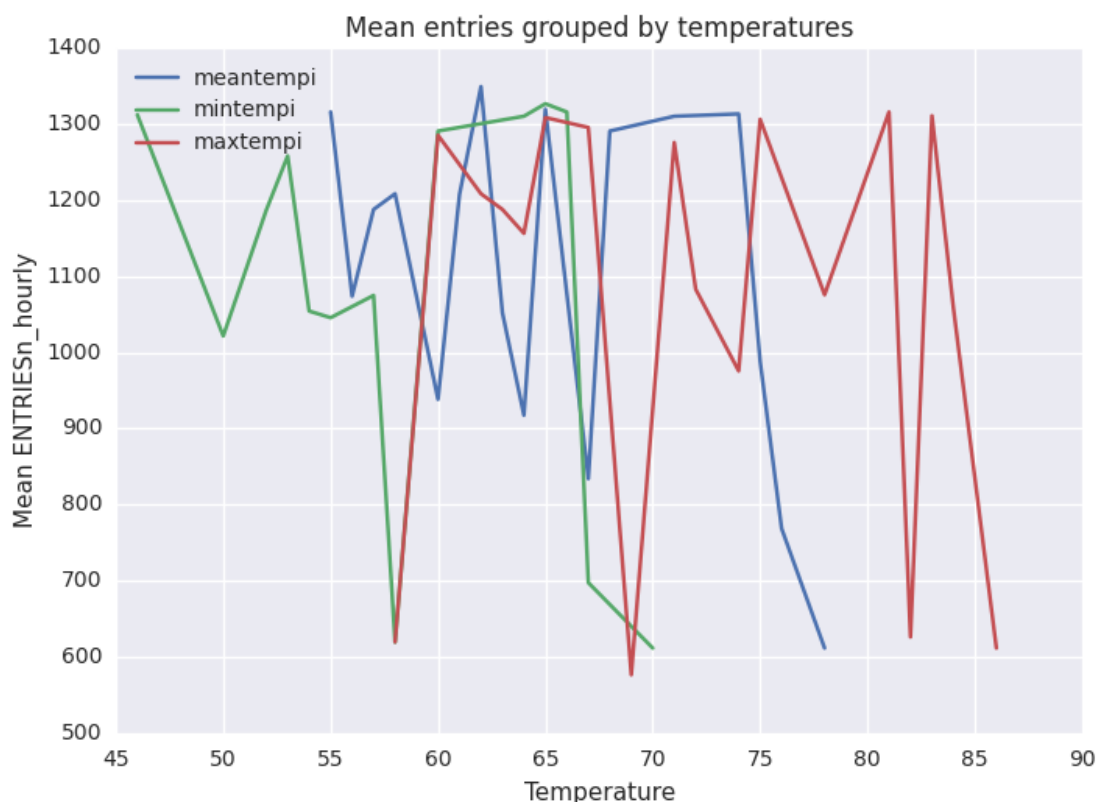


Fig.4: Mean entries grouped by temperatures

In addition, I tried to check whether the temperatures and barometric pressure had correlation with the entries and the kind of relationship.

From Fig. 4 I can see the mean entries grouped by three temperature columns: meantempi, mintempi, maxtempi. By choosing only meantempi and mintempi, only a very small improvement is observed:  $r^2 = 0.514964513027$ , when compared to the last  $r^2$  value obtained in the last features choice of answer to 2.2.

Similarly, no linearity observation can be obtained by Fig. 5, indicating that barometric pressure would not bring a significant improvement to  $r^2$ , found out to be  $r^2 = 0.515053815292$ .

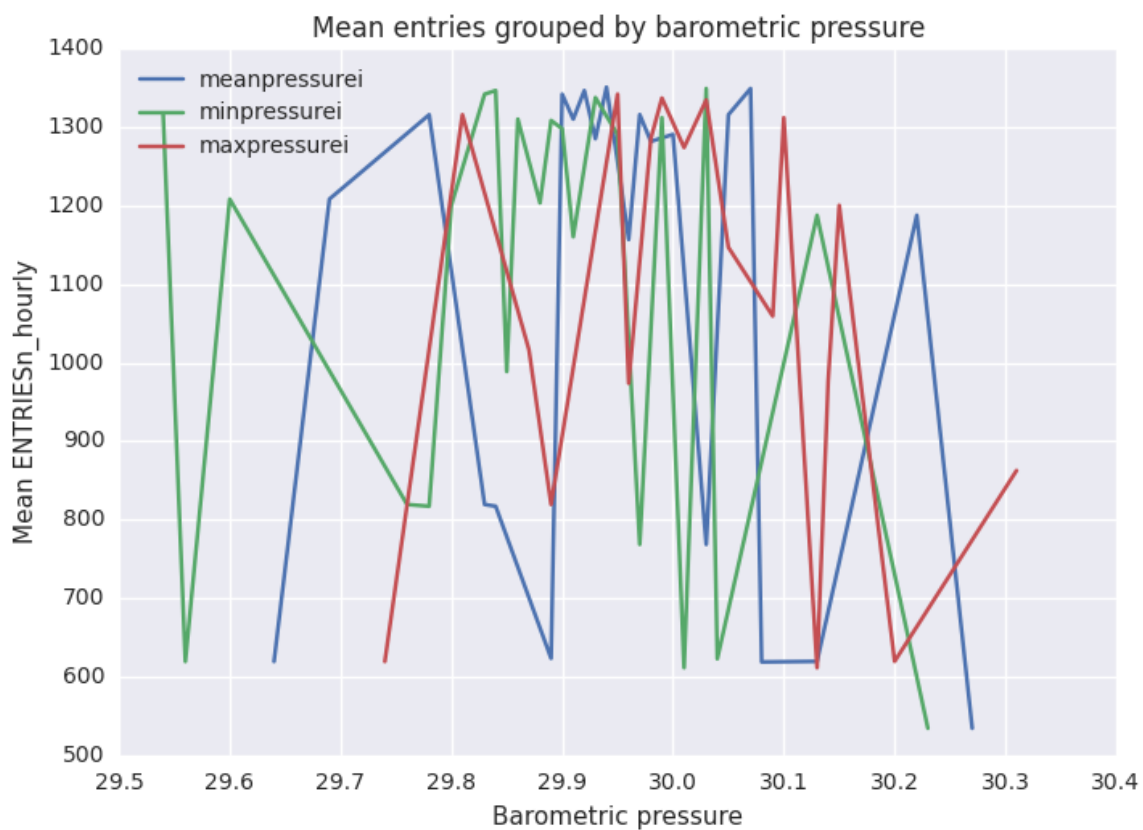


Fig.5: Mean entries grouped by barometric pressure

## 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:

The results of Table 1 of Question 1.4 along with the analysis of coefficients of Question 2.4, support this answer.

If a generic day is randomly chosen, we have more riders when it is raining. This conclusion considers  $\alpha = 0.05$ . If I had chosen  $\alpha = 0.01$  this would not be valid and I would conclude that there is no difference between ridership volume between rainy and non-rainy days.

However, when I analyze working days and weekends separately, the conclusion is:

- a) for a working day, non-rainy days have more riders than rainy days
- b) for weekends, rainy days have more riders than non-rainy days

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:

The Mann-Whitney u tests confirm that the conclusions in 4.1 are valid since p-values are very low. In fact, for both weekends and working days, we can safely reject the null hypothesis that both come from same distribution, since  $p < 0.1\%$ .

Also, from answer to question 2.4, the coefficient of rain feature is positive, so for a generic day the more rain, the more riders.

It can also be noted that coefficient for weekday (281.2628) is positive and its magnitude higher than other features. It is an indication that among the chosen features, this is the most important non-dummy factor on the prediction of the number of riders, more than other features such as rain or temperature. The values of different coefficients  $\theta$  can be compared because of the normalization.

Dummy variables show high coefficients, but since these variables are of 'yes'/'no' type, with most of the time being a 'no' (which means close to zero after normalization), they only affect ridership when the focus is in one particular turnstile unit. It almost seems like I am 'cheating', trying to conform a prediction formula to a bunch of discrete 0/1 variables, increasing the computational time because of quadratic nature of matrices product, instead of trying to achieve a more elegant and compact polynomial.

Other weather factors, such as barometric pressure, temperature show only marginal contributions to the prediction accuracy of ridership.

## 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:

Dataset,

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

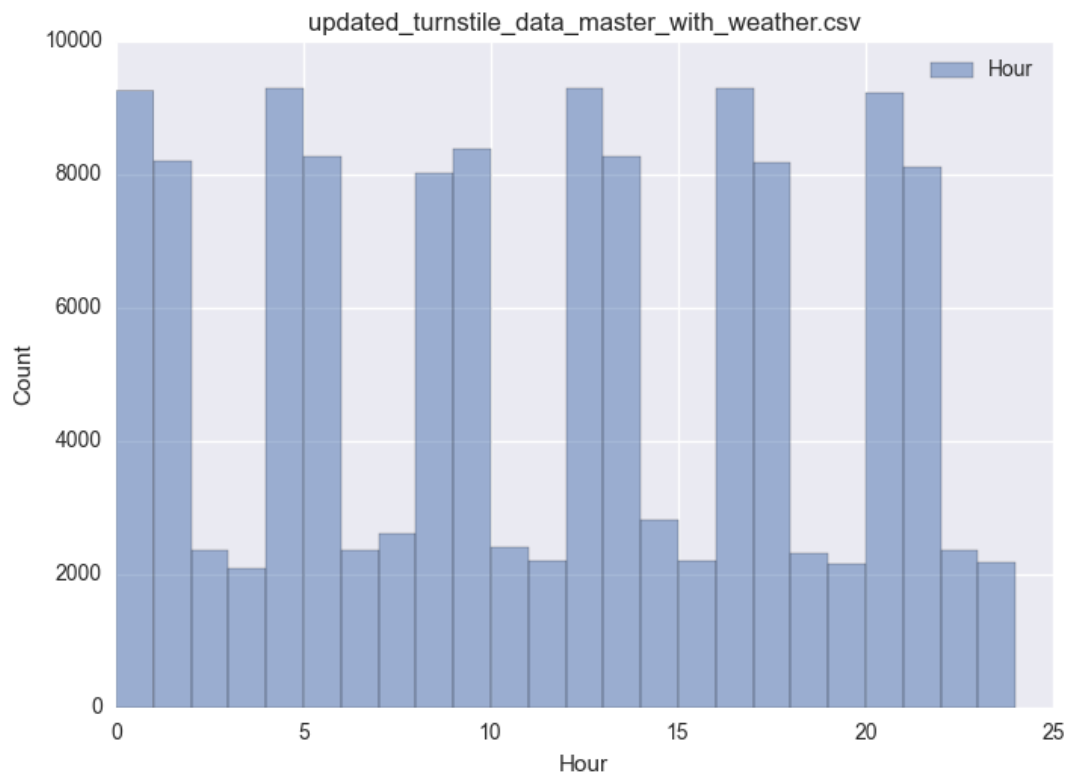


Fig.6: Hour value histogram

Answer:

the original data set showed inconsistent entries for hour. While most of the rows fit into the (0, 4, 8, 12, 16, 20) window, with a comparable number in the (1, 5, 9, 13, 17, 21) window, and the rest distributed among the remaining time slots. Initially, I was thinking this inconsistency would privilege low Hour values, but it was not the case. In any case, more checks are needed to see if this affects unevenly rain vs no-rain entries or working day versus weekend entries. In quick glance on the data I suspect these are more related to UNIT than to other factors.

Three shortcomings, that is lack of day\_week, weekday and holiday, were corrected by a script that added those columns to the file. Using these columns as features of the regression model proved to be effecting in increasing the  $r^2$ .

Finally, I think that one month of data tells us only what could happen in that month in different years, that is, from data of May/11, we could predict ridership of May/12, May/13,.... . Seasonal increases or decreases on ridership would be captured by this model.

About the analysis, from histograms in Fig. 4 and Fig. 5 I see that the model does not do a good job in predicting ridership with some non negligible long tails. Reinforcing this conclusion, Fig. 7

shows residuals plotted against the sample index, with two different views. On the top figure, the line presents high and low peaks in cyclical pattern, indicating that the residuals correlate with the row sequence of the input data. The column hierarchy is: date, turnstile unit and time, that is, there are 30 blocks of rows, one block for each day, and within each block the rows are organised by unit first and then by hour. So, the each of the 30 peaks correspond to highest residuals (busiest stations) of each day of May/2011.

At the bottom figure, we plot residuals with an area of the bubble proportional to the entry. Supporting the conclusions of 2.6 we see that high residuals correspond to instances of higher entries, a situation where the model underestimates the actual results. Also, we observe the small dots, corresponding to more idle periods. We see they correlate with negative residuals (i.e.,  $\text{prediction} > \text{sample}$ ), here a situation where the model overestimates the actual results.

Also, Fig. 8 shows there is relationship between entries and residuals, which is approximately linear.

A final comment about the limitation of the linear model. It looks very difficult to model all the ups and downs of the entries, such as observed in Fig. 7, via a linear equation with finite number of terms.

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

From the improved dataset, I see that there are more columns available for analysis. For example, 'conds' column (expanded into dummy columns) could be an interesting feature to check - to see if improves the accuracy of prediction.

Also, in the improved dataset, a test I would try is to check if stations in downtown have the same rain/no-rain ridership as stations away from downtown. My intuition is that in downtown, riders would not have choice of not taking the subway, while in less dense areas, people may use more cars when it is raining - but I may be wrong. To check if a station is in downtown or not, I would use latitude and longitude along with some criteria to define which are downtown stations.

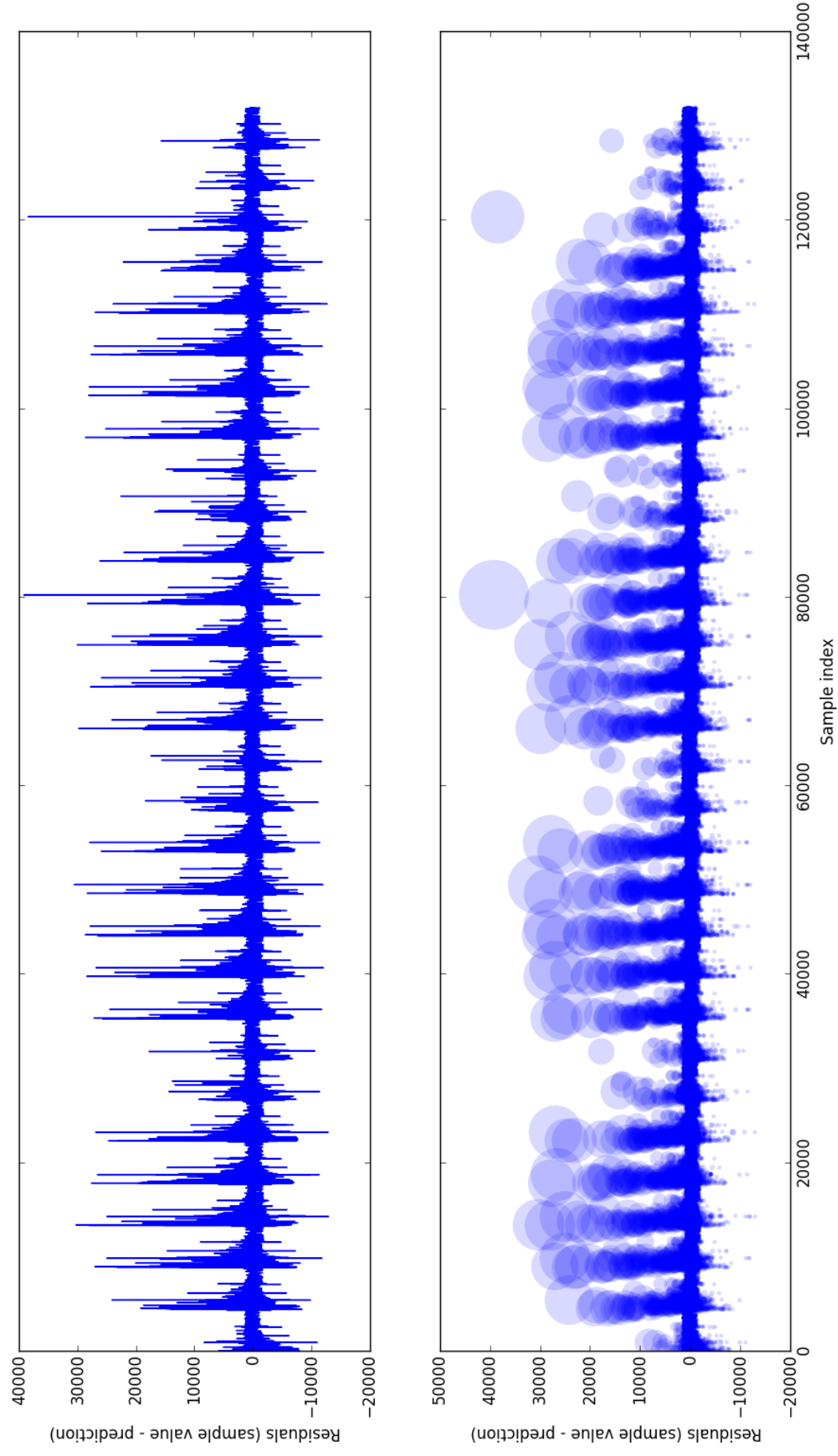


Fig.7: Residuals plot and bubble plot (size = entries) against sample index

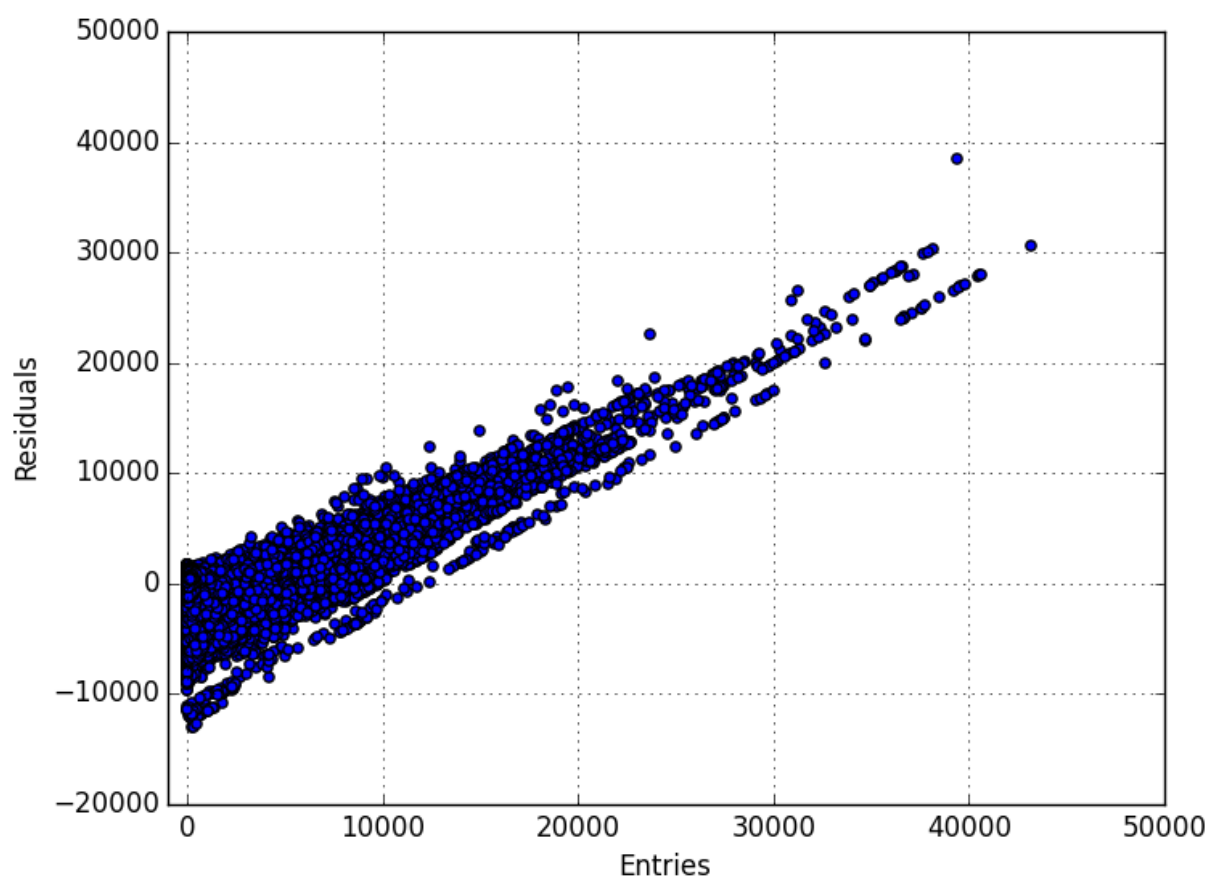


Fig.8: Relationship between residuals and entries