

Section 1.

Statistical Test

1.1 Which statistical test did you use to analyze the NYC subway data? Mann Whitney U test

Did you use a one-tail or a two-tail P value? One Tail Test

What is the null hypothesis? The p value is small which mean that the null hypothesis can be rejected - the difference in the data sets is not due to random sampling because the p value is low.

What is your p-critical value? 2%

1.2 Why is this statistical test applicable to the dataset? The data is not normally distributed, and is a large sample. We are investigate one dependant variable (raining or not raining). The two distributions of the sample have the same shape asymptotically downwards.

In particular, consider the assumptions that the test is making about the distribution of ridership in the two samples.

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.

With Rain Mean: 1105.4463767458733

Without Rain Mean: 1090.278780151855

U: 1924409167.0

p: 0.024999912793489721)

1.4 What is the significance and interpretation of these results?

The p value is small which mean that the null hypothesis can be rejected - the difference in the data sets is not due to random sampling because the p value is low. Only 2% of the time will the null hypothesis not be rejected when it is in fact true.

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:

Gradient descent (as implemented in exercise 3.5)

OLS using Statsmodels

Or something different?

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

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^2 value.”

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

2.5 What is your model's R^2 (coefficients of determination) value?

2.6 What does this R^2 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^2 value?

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.

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

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?

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.

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.

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

tion 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?
- 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.
- 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.
- 1.4 What is the significance and interpretation of these results?

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:
 - Gradient descent (as implemented in exercise 3.5)
 - OLS using Statsmodels
 - Or something different?
- 2.2 What features (input variables) did you use in your model? Did you use any dummy variables as part of your features?
- 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^2 value."
- 2.4 What are the coefficients (or weights) of the non-dummy features in your linear regression model?
- 2.5 What is your model's R^2 (coefficients of determination) value?
- 2.6 What does this R^2 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^2 value?

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.

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

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?

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.

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.

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