

---

# **Analyzing the NYC Subway Dataset Documentation**

***Release 1.0***

**Ignacio Toledo**

December 29, 2014



## CONTENTS

<b>1</b>	<b>Overview</b>	<b>1</b>
<b>2</b>	<b>Statistical Test</b>	<b>3</b>
2.1	Statistical Test Used . . . . .	3
2.2	Justify the Statistical Test . . . . .	3
2.3	Results . . . . .	3
2.4	Interpretation . . . . .	3
<b>3</b>	<b>Linear Regression</b>	<b>5</b>
3.1	Linear regression algorithm(s) . . . . .	5
3.2	Features used . . . . .	5
3.3	Feature selection: why? . . . . .	5
3.4	Results: R Squared . . . . .	5
3.5	Interpretation and limits . . . . .	5
<b>4</b>	<b>Visualization</b>	<b>7</b>
4.1	Ridership distribution with weather . . . . .	7
4.2	Supporting visualizations . . . . .	7
<b>5</b>	<b>Conclusion</b>	<b>9</b>
<b>6</b>	<b>Reflection</b>	<b>11</b>
6.1	Shortcomings and limitations . . . . .	11
6.2	Insights . . . . .	11



**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.



## STATISTICAL TEST

### 2.1 Statistical Test Used

### 2.2 Justify the Statistical Test

### 2.3 Results

### 2.4 Interpretation





## LINEAR REGRESSION

### 3.1 Linear regression algorithm(s)

#### 3.1.1 Gradient descent

#### 3.1.2 OLS (with statsmodels)

### 3.2 Features used

### 3.3 Feature selection: why?

### 3.4 Results: R Squared

### 3.5 Interpretation and limits



## VISUALIZATION

### 4.1 Ridership distribution with weather

### 4.2 Supporting visualizations



**CONCLUSION**



**REFLECTION**

**6.1 Shortcomings and limitations**

**6.2 Insights**