

Bay Area Bike Share Analysis

Introduction

Tip: Quoted sections like this will provide helpful instructions on how to navigate and use an iPython notebook.

Bay Area Bike Share (<http://www.bayareabikeshare.com/>) is a company that provides on-demand bike rentals for customers in San Francisco, Redwood City, Palo Alto, Mountain View, and San Jose. Users can unlock bikes from a variety of stations throughout each city, and return them to any station within the same city. Users pay for the service either through a yearly subscription or by purchasing 3-day or 24-hour passes. Users can make an unlimited number of trips, with trips under thirty minutes in length having no additional charge; longer trips will incur overtime fees.

In this project, you will put yourself in the shoes of a data analyst performing an exploratory analysis on the data. You will take a look at two of the major parts of the data analysis process: data wrangling and exploratory data analysis. But before you even start looking at data, think about some questions you might want to understand about the bike share data. Consider, for example, if you were working for Bay Area Bike Share: what kinds of information would you want to know about in order to make smarter business decisions? Or you might think about if you were a user of the bike share service. What factors might influence how you would want to use the service?

Question 1: Write at least two questions you think could be answered by data.

Answer

Tip: If you double click on this cell, you will see the text change so that all of the formatting is removed. This allows you to edit this block of text. This block of text is written using Markdown (<http://daringfireball.net/projects/markdown/syntax>), which is a way to format text using headers, links, italics, and many other options. You will learn more about Markdown later in the Nanodegree Program. Hit **Shift + Enter** or **Shift + Return**.

Using Visualizations to Communicate Findings in Data

As a data analyst, the ability to effectively communicate findings is a key part of the job. After all, your best analysis is only as good as your ability to communicate it.

In 2014, Bay Area Bike Share held an Open Data Challenge (<http://www.bayareabikeshare.com/datachallenge-2014>) to encourage data analysts to create visualizations based on their open data set. You'll create your own visualizations in this project, but first, take a look at the submission winner for Best Analysis (<http://thfield.github.io/babs/index.html>) from Tyler Field. Read through the entire report to answer the following question:

Question 2: What visualizations do you think provide the most interesting insights? Are you able to answer either of the questions you identified above based on Tyler's analysis? Why or why not?

Answer: Replace this text with your response!

Data Wrangling

Now it's time to explore the data for yourself. Year 1 and Year 2 data from the Bay Area Bike Share's Open Data (<http://www.bayareabikeshare.com/open-data>) page have already been provided with the project materials; you don't need to download anything extra. The data comes in three parts: the first half of Year 1 (files starting 201402), the second half of Year 1 (files starting 201408), and all of Year 2 (files starting 201508). There are three main datafiles associated with each part: trip data showing information about each trip taken in the system (*_trip_data.csv), information about the stations in the system (*_station_data.csv), and daily weather data for each city in the system (*_weather_data.csv).

When dealing with a lot of data, it can be useful to start by working with only a sample of the data. This way, it will be much easier to check that our data wrangling steps are working since our code will take less time to complete. Once we are satisfied with the way things are working, we can then set things up to work on the dataset as a whole.

Since the bulk of the data is contained in the trip information, we should target looking at a subset of the trip data to help us get our bearings. You'll start by looking at only the first month of the bike trip data, from 2013-08-29 to 2013-09-30. The code below will take the data from the first half of the first year, then write the first month's worth of data to an output file. This code exploits the fact that the data is sorted by date (though it should be noted that the first two days are sorted by trip time, rather than being completely chronological).

First, load all of the packages and functions that you'll be using in your analysis by running the first code cell below. Then, run the second code cell to read a subset of the first trip data file, and write a new file containing just the subset we are initially interested in.