# COSC 5010 - Data Science for Security Final Project

Libao Jin (ljin1@uwyo.edu)

May 3, 2018

## 1 Motivation

In metropolitan areas, different kind of public transportations such as buses, taxis, metros, etc., play an essential role in people's daily life. However, it is not as efficient as we expected sometimes. Hence, this project mainly looks into the public transportation, i.e., BART, to make some useful suggestion. Here are a bunch questions this project aims to answer:

1. Which BART station is the busiest?
2. What is the least popular BART route?
3. When is the best time to go to SF if you want to find a seat?
4. Which day of the week is the busiest?
5. How many people take the BART late at night?
6. Does the BART ever stop in a station without anyone going off and on?

## 2 Background

BART, short for "Bay Area Rapid Transit", is the transit system severing the San Francisco Bay Area in California. BART operates six routes, 46 stations, and and 112 miles of track. It serves an average weekday ridership of 423,000 people, making it the fifth-busiest rapid transit system in the United States.

This dataset contains daily information on BART ridership for a period covering all of 2016 and part of 2017. Unlike some other rapid transit system datasets, this data includes movements between specific stations (there are just over 2000 station-to-station combinations).

## 3 Project Summary & Major Tasks

### 3.1 Interpretation of the Datasets

We get the data from [1], which provides the following three datasets:

- data-hour-soo-dest-2016.csv: Number of passengers (Throughput) that went between two stations (Origin and Destination) in a given time (DateTime) in 2016 as Table 1.
- data-hour-soo-dest-2017.csv: Number of passengers (Throughput) that went between two stations (Origin and Destination) in a given time (DateTime) in 2017.
- station_info.csv: Information about different BART stations (Location, Description, etc.) as Table 2.

### 3.2 Data Preprocessing/Data Cleaning

Since the data is well-collected, there is not much cleaning required. But in order to make the data analysis process easier, we'd like to preprocess the data such as converting the string read from text files into desired type and format.

Table 1: First five rows of data-hour-soo-dest-2016.csv

| Origin | Destination | Throughput | DateTime |
|--------|-------------|------------|----------|
| 12TH | 12TH | 1 | 2016-01-01 00:00:00 |
| 12TH | 16TH | 1 | 2016-01-01 00:00:00 |
| 12TH | 24TH | 4 | 2016-01-01 00:00:00 |
| 12TH | ASHB | 4 | 2016-01-01 00:00:00 |
| 12TH | BALB | 2 | 2016-01-01 00:00:00 |

### 3.2.1  Dataset Reconstruction

- Convert the string of column `DateTime` of data `data-hour-soo-dest-2016.csv` and `data-hour-soo-dest-2017.csv` to `pandas.Datetime` type as Table 3.

- Extract `Longitude` and `Latitude` from `Location` for the visualization of routes later on, the preprocessed table is as shown in Table 4.

### 3.2.2  Merge the Datasets of 2016 and 2017

Using Pandas to read the `.csv` file, and store each dataset as a `DataFrame`, then using the `pandas.concat` to merge two `DataFrame`s.

## 3.3  Assumptions

In order to find how busy each station is, we consider the throughput of each station as the origin, where riders take the transit. To do so, we made following assumptions:

1. In the dataset, the value of the column `DateTime` is the time that riders arrive at the Origin and take the transit to Destination. Thus, the arriving time is unknown in this case.
2. The throughput of a station is, actually, the number of riders that the train/trainsit carry, rather than the number of people who are waiting in the station.

## 3.4  Calculation of the Daily Throughput of Each Station from 2016 to 2017

We sum up all the throughput of the same origin as the total throughput of that station in the year of 2016 and 2017. Then, we sort the stations by the total throughput and use visualization to find the busiest station.

## 3.5  Calculation of the Total Hourly Throughput of Each Station of Weeks/-Months

We group the stations by weekdays/months, then sum up the hourly throughput of each station and then visualize them to find the best time to take the transit on a certain weekday or in certain month.

## 3.6  Build a Predictor of Throughput for BART Using Machine Learning

We pick the following features to train the model: month of year, day of month, day of week, time, origin, destination, origin longitude, origin latitude, destination latitude, destination longtitude. And then split the datasets into training dataset and test dataset. Using linear regression and decision tree, which provided by the Python module `scikit-learn`. Before we start, we set the criteria to indicate the level of workload of the BART (Table 5).

Table 2: Stations information from stat_info.csv

| Abbreviation | Description | Location | Name |
|---|---|---|---|
| 12TH | 1245 Broadway,... | -122.2714,37.8037,0 | 12th St. Oakland City Center (12TH) |
| 16TH | 2000 Mission S... | -122.4196,37.7650,0 | 16th St. Mission (16TH) |
| 19TH | 1900 Broadway,... | -122.2686,37.8083,0 | 19th St. Oakland (19TH) |
| 24TH | 2800 Mission S... | -122.4181,37.7524,0 | 24th St. Mission (24TH) |
| ASHB | 3100 Adeline S... | -122.2700,37.8528,0 | Ashby (ASHB) |
| BALB | 401 Geneva Ave... | -122.4475,37.7215,0 | Balboa Park (BALB) |
| BAYF | 15242 Hesperia... | -122.1265,37.6969,0 | Bay Fair (BAYF) |
| CAST | 3301 Norbridge... | -122.0756,37.6907,0 | Castro Valley (CAST) |
| CIVC | 1150 Market St... | -122.4141,37.7797,0 | Civic Center/UN Plaza (CIVC) |
| COLS | 7200 San Leand... | -122.1968,37.7536,0 | Coliseum/Oakland Airport (COLS) |
| COLM | 365 D Street, ... | -122.4662,37.6846,0 | Colma (COLM) |
| CONC | 1451 Oakland A... | -122.0290,37.9737,0 | Concord (CONC) |
| DALY | 500 John Daly ... | -122.4690,37.7061,0 | Daly City (DALY) |
| DBRK | 2160 Shattuck ... | -122.2681,37.8701,0 | Downtown Berkeley (DBRK) |
| DELN | 6400 Cutting B... | -122.3167,37.9250,0 | El Cerrito del Norte (DELN) |
| DUBL | 5801 Owens Dr.... | -121.8991,37.7016,0 | Dublin/Pleasanton (DUBL) |
| EMBR | 298 Market Str... | -122.3970,37.7928,0 | Embarcadero (EMBR) |
| FRMT | 2000 BART Way,... | -121.9766,37.5574,0 | Fremont (FRMT) |
| FTVL | 3401 East 12th... | -122.2241,37.7748,0 | Fruitvale (FTVL) |
| GLEN | 2901 Diamond S... | -122.4338,37.7330,0 | Glen Park (GLEN) |
| HAYW | 699 'B' Street... | -122.0870,37.6697,0 | Hayward (HAYW) |
| LAFY | 3601 Deer Hill... | -122.1246,37.8931,0 | Lafayette (LAFY) |
| LAKE | 800 Madison St... | -122.2651,37.7970,0 | Lake Merritt (LAKE) |
| MCAR | 555 40th Stree... | -122.2670,37.8290,0 | MacArthur (MCAR) |
| MLBR | 200 North Roll... | -122.3867,37.6002,0 | Millbrae (MLBR) |
| MONT | 598 Market Str... | -122.4010,37.7894,0 | Montgomery St. (MONT) |
| NBRK | 1750 Sacrament... | -122.2834,37.8739,0 | North Berkeley (NBRK) |
| NCON | 3700 Port Chic... | -122.0246,38.0031,0 | North Concord/Martinez (NCON) |
| OAKL | 1 Airport Driv... | -122.2121,37.7132,0 | Oakland Airport (OAKL) |
| ORIN | 11 Camino Pabl... | -122.1837,37.8783,0 | Orinda (ORIN) |
| PHIL | 1365 Treat Blv... | -122.0560,37.9284,0 | Pleasant Hill/Contra Costa Centre (PHIL) |
| PITT | 1700 West Lela... | -121.9451,38.0189,0 | Pittsburg/Bay Point (PITT) |
| PLZA | 6699 Fairmount... | -122.2989,37.9026,0 | El Cerrito Plaza (PLZA) |
| POWL | 899 Market Str... | -122.4079,37.7844,0 | Powell St. (POWL) |
| RICH | 1700 Nevin Ave... | -122.3530,37.9368,0 | Richmond (RICH) |
| ROCK | 5660 College A... | -122.2513,37.8447,0 | Rockridge (ROCK) |
| SANL | 1401 San Leand... | -122.1608,37.7219,0 | San Leandro (SANL) |
| SBRN | 1151 Huntingto... | -122.4162,37.6377,0 | San Bruno (SBRN) |
| SFIA | International ... | -122.3924,37.6159,0 | San Francisco Int'l Airport (SFIA) |
| SHAY | 28601 Dixon St... | -122.0571,37.6343,0 | South Hayward (SHAY) |
| SSAN | 1333 Mission R... | -122.4439,37.6642,0 | South San Francisco (SSAN) |
| UCTY | 10 Union Squar... | -122.0173,37.5906,0 | Union City (UCTY) |
| WARM | 45193 Warm Spr... | -121.9393,37.5021,0 | Warm Springs/South Fremont (WARM) |
| WCRK | 200 Ygnacio Va... | -122.0675,37.9055,0 | Walnut Creek (WCRK) |
| WDUB | 6501 Golden Ga... | -121.9282,37.6997,0 | West Dublin/Pleasanton (WDUB) |
| WOAK | 1451 7th Stree... | -122.2951,37.8048,0 | West Oakland (WOAK) |

Table 3: First five rows of converted DataFrame of data-hour-soo-dest-2016.csv

| Origin | Destination | Throughput | DateTime | Date | Time | DayOfWeek |
|--------|-------------|-----------|----------|------|------|-----------|
| 12TH | 12TH | 1 | 2016-01-01 00:00:00 | 2016-01-01 | 00:00:00 | Friday |
| 12TH | 16TH | 1 | 2016-01-01 00:00:00 | 2016-01-01 | 00:00:00 | Friday |
| 12TH | 24TH | 4 | 2016-01-01 00:00:00 | 2016-01-01 | 00:00:00 | Friday |
| 12TH | ASHB | 4 | 2016-01-01 00:00:00 | 2016-01-01 | 00:00:00 | Friday |
| 12TH | BALB | 2 | 2016-01-01 00:00:00 | 2016-01-01 | 00:00:00 | Friday |

# 4 Results

## 4.1 Visualization of the Data

### 4.1.1 Visualization of the Routes of BART Lines

To visualize the routes, we'd like to utilize the geographic location provided by Location columns in station_info.csv, which consists of longitude and latitude. The dataset, however, does not provide the exact lines directly. Hence, we obtained the five lines on the website of the BART [2]. Then we obtained the plots of BART lines as Figure 1 shown.

### 4.1.2 Visualization of the Throughput of Each Station

As shown in Figure 2, it is easy to find that the busiest stations (bubbles with large radius or of magenta) are in the downtown of San Francisco, i.e., Montgomery St. Station, Embarcadero Station, Powell St. Station, Civic Center/UN Plaza Station, etc.
Also, we obtained the barplot for the avergae hourly throughput of each station of 2016 and 2017 as shown in 3. Obviously, the throughputs of aforementioned four stations far outweight that of the rest stations.

## 4.2 Daily Throughput of Each Station of 2016 and 2017

The plots of throughput of each station of 2016 and 2017 is as shown in Figure 4, from which we can have a sense of the volume of riderships in each station.

## 4.3 Hourly Throughput of Each Station of a Week

It is shown in Figure 5, 6, 7, 8, 9 that busy hours of each stations on each day of the week. Each row represents the throughput of a certain station while the columns represent the days of week, i.e., the first column stands for Monday, the second for Tuesday, and the like.
In the top four busiest stations, the peak hours are mainly around 17:00-18:00, during which people take the transit to go home. These stations are in downtown San Francisco. Presumably that people come here to work/shopping, etc. As for those stations have peak hours around 9:00-10:00, that indicate those stations are close to the dwelling where people live. And from these plots, we can get the best time (by avoiding the rush hours) to take the BART in order to find a seat.

## 4.4 Linear Regression

Turns the accuracy of the linear regression is pretty low (score: 0.8%). So... we skip this one.

## 4.5 Decision Tree

We used Decision Tree to build a predictor of the throughput of stations. Among features such as month of year, day of month, day of week, time, origin, destination, origin longitude, origin latitude, destination latitude, destination longtitude, we find that the combination of features day of week, time, origin and desination are most effetive. Only use 0.02% of the data as our training dataset (Figure 10), we can reach

Table 4: Preprocessed stations information from stat_info.csv

| Abbreviation | Name | Longitude | Latitude | Description |
|---|---|---|---|---|
| 12TH | 12th St. Oakland City Center (12TH) | -122.271450 | 37.803768 | 1245 Broad... |
| 16TH | 16th St. Mission (16TH) | -122.419694 | 37.765062 | 2000 Missi... |
| 19TH | 19th St. Oakland (19TH) | -122.268602 | 37.808350 | 1900 Broad... |
| 24TH | 24th St. Mission (24TH) | -122.418143 | 37.752470 | 2800 Missi... |
| ASHB | Ashby (ASHB) | -122.270062 | 37.852803 | 3100 Adeli... |
| BALB | Balboa Park (BALB) | -122.447506 | 37.721585 | 401 Geneva... |
| BAYF | Bay Fair (BAYF) | -122.126514 | 37.696924 | 15242 Hesp... |
| CAST | Castro Valley (CAST) | -122.075602 | 37.690746 | 3301 Norbr... |
| CIVC | Civic Center/UN Plaza (CIVC) | -122.414123 | 37.779732 | 1150 Marke... |
| COLS | Coliseum/Oakland Airport (COLS) | -122.196869 | 37.753661 | 7200 San L... |
| COLM | Colma (COLM) | -122.466233 | 37.684638 | 365 D Stre... |
| CONC | Concord (CONC) | -122.029095 | 37.973737 | 1451 Oakla... |
| DALY | Daly City (DALY) | -122.469081 | 37.706121 | 500 John D... |
| DBRK | Downtown Berkeley (DBRK) | -122.268133 | 37.870104 | 2160 Shatt... |
| DELN | El Cerrito del Norte (DELN) | -122.316794 | 37.925086 | 6400 Cutti... |
| DUBL | Dublin/Pleasanton (DUBL) | -121.899179 | 37.701687 | 5801 Owens... |
| EMBR | Embarcadero (EMBR) | -122.397020 | 37.792874 | 298 Market... |
| FRMT | Fremont (FRMT) | -121.976608 | 37.557465 | 2000 BART ... |
| FTVL | Fruitvale (FTVL) | -122.224175 | 37.774836 | 3401 East ... |
| GLEN | Glen Park (GLEN) | -122.433817 | 37.733064 | 2901 Diamo... |
| HAYW | Hayward (HAYW) | -122.087018 | 37.669723 | 699 'B' St... |
| LAFY | Lafayette (LAFY) | -122.124630 | 37.893176 | 3601 Deer ... |
| LAKE | Lake Merritt (LAKE) | -122.265180 | 37.797027 | 800 Madiso... |
| MCAR | MacArthur (MCAR) | -122.267040 | 37.829065 | 555 40th S... |
| MLBR | Millbrae (MLBR) | -122.386702 | 37.600271 | 200 North ... |
| MONT | Montgomery St. (MONT) | -122.401066 | 37.789405 | 598 Market... |
| NBRK | North Berkeley (NBRK) | -122.283440 | 37.873967 | 1750 Sacra... |
| NCON | North Concord/Martinez (NCON) | -122.024653 | 38.003193 | 3700 Port ... |
| OAKL | Oakland Airport (OAKL) | -122.212191 | 37.713238 | 1 Airport ... |
| ORIN | Orinda (ORIN) | -122.183791 | 37.878361 | 11 Camino ... |
| PHIL | Pleasant Hill/Contra Costa Centre (PHIL) | -122.056012 | 37.928468 | 1365 T... ... |
| PITT | Pittsburg/Bay Point (PITT) | -121.945154 | 38.018914 | 1700 West ... |
| PLZA | El Cerrito Plaza (PLZA) | -122.298904 | 37.902632 | 6699 Fairm... |
| POWL | Powell St. (POWL) | -122.407974 | 37.784471 | 899 Market... |
| RICH | Richmond (RICH) | -122.353099 | 37.936853 | 1700 Nevin... |
| ROCK | Rockridge (ROCK) | -122.251371 | 37.844702 | 5660 Colle... |
| SANL | San Leandro (SANL) | -122.160844 | 37.721947 | 1401 San L... |
| SBRN | San Bruno (SBRN) | -122.416287 | 37.637761 | 1151 Hunti... |
| SFIA | San Francisco Int'l Airport (SFIA) | -122.392409 | 37.615966 | Internatio... |
| SHAY | South Hayward (SHAY) | -122.057189 | 37.634375 | 28601 Dixo... |
| SSAN | South San Francisco (SSAN) | -122.443960 | 37.664245 | 1333 Missi... |
| UCTY | Union City (UCTY) | -122.017388 | 37.590630 | 10 Union S... |
| WARM | Warm Springs/South Fremont (WARM) | -121.939313 | 37.502171 | 45193 Warm... |
| WCRK | Walnut Creek (WCRK) | -122.067527 | 37.905522 | 200 Ygnaci... |
| WDUB | West Dublin/Pleasanton (WDUB) | -121.928240 | 37.699756 | 6501 Golde... |
| WOAK | West Oakland (WOAK) | -122.295140 | 37.804872 | 1451 7th S... |

Table 5: Criterion for setting levels of throughput

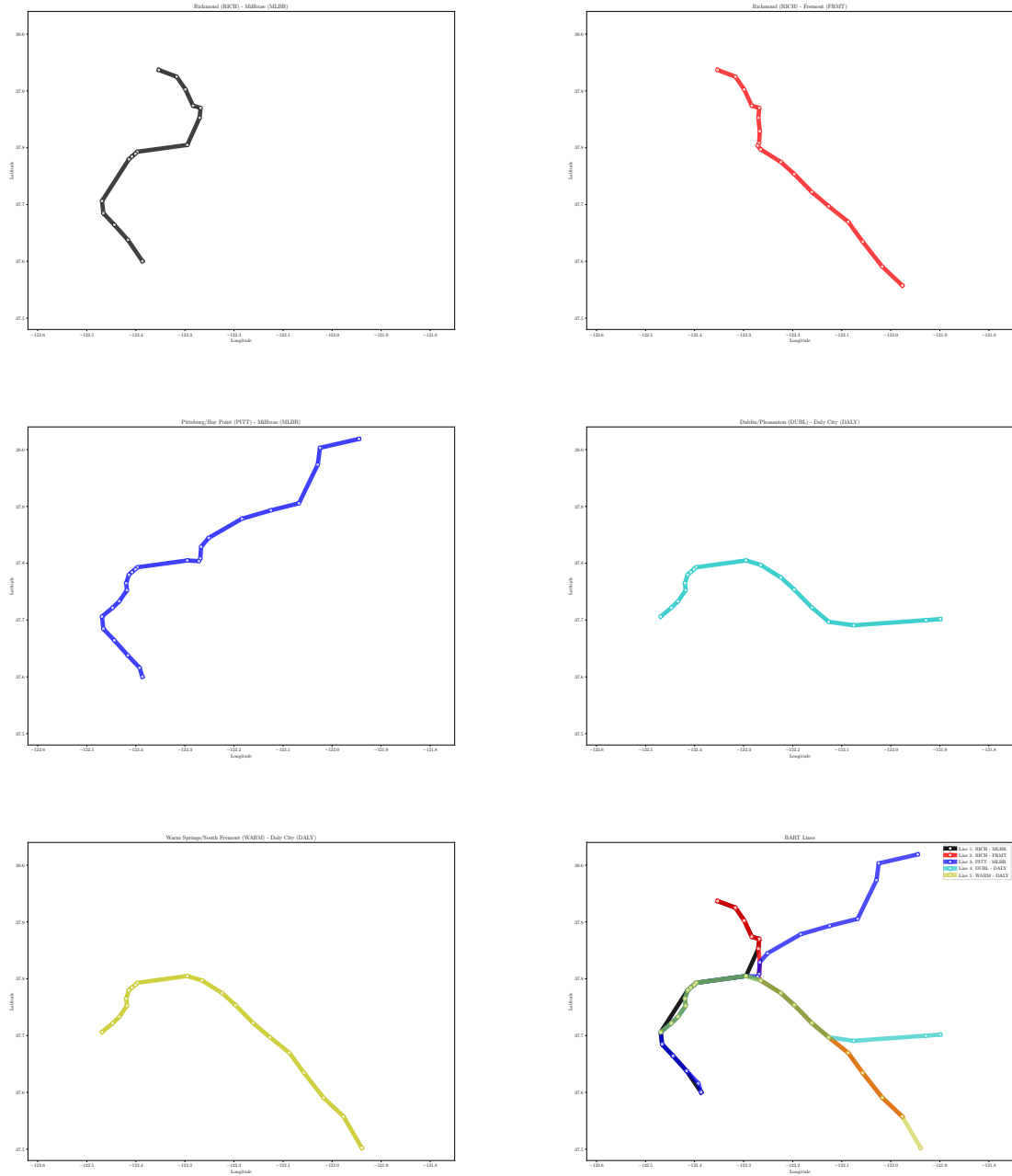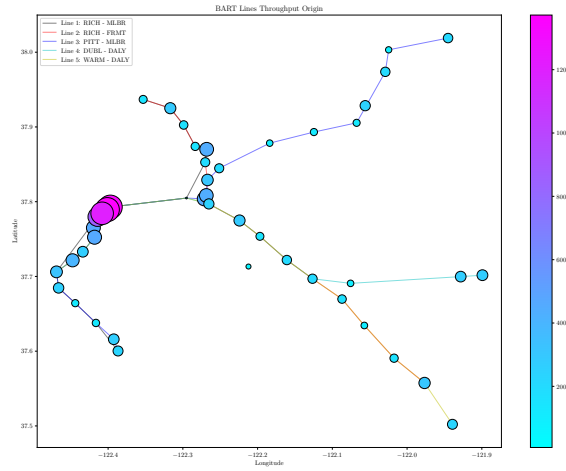| Throughput | Level |
|------------|-------|
| $\leq 5$ | Very low |
| $6 \sim 15$ | Low |
| $16 \sim 30$ | Medium |
| $31 \sim 80$ | High |
| $> 80$ | Very high |



Figure 1: Routes of the BART Lines

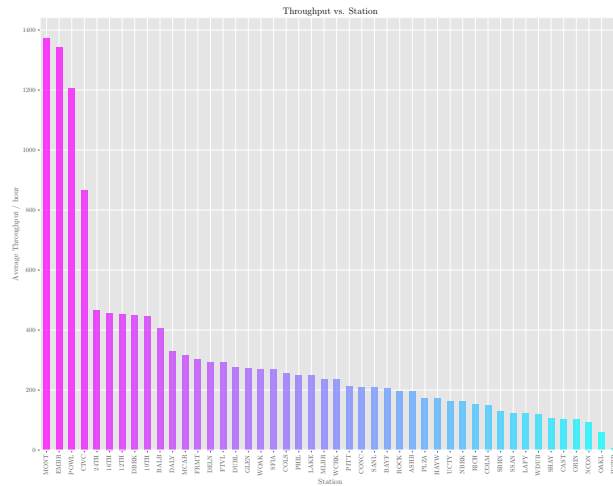Figure 2: Average hourly throughput of each station of 2016 and 2017



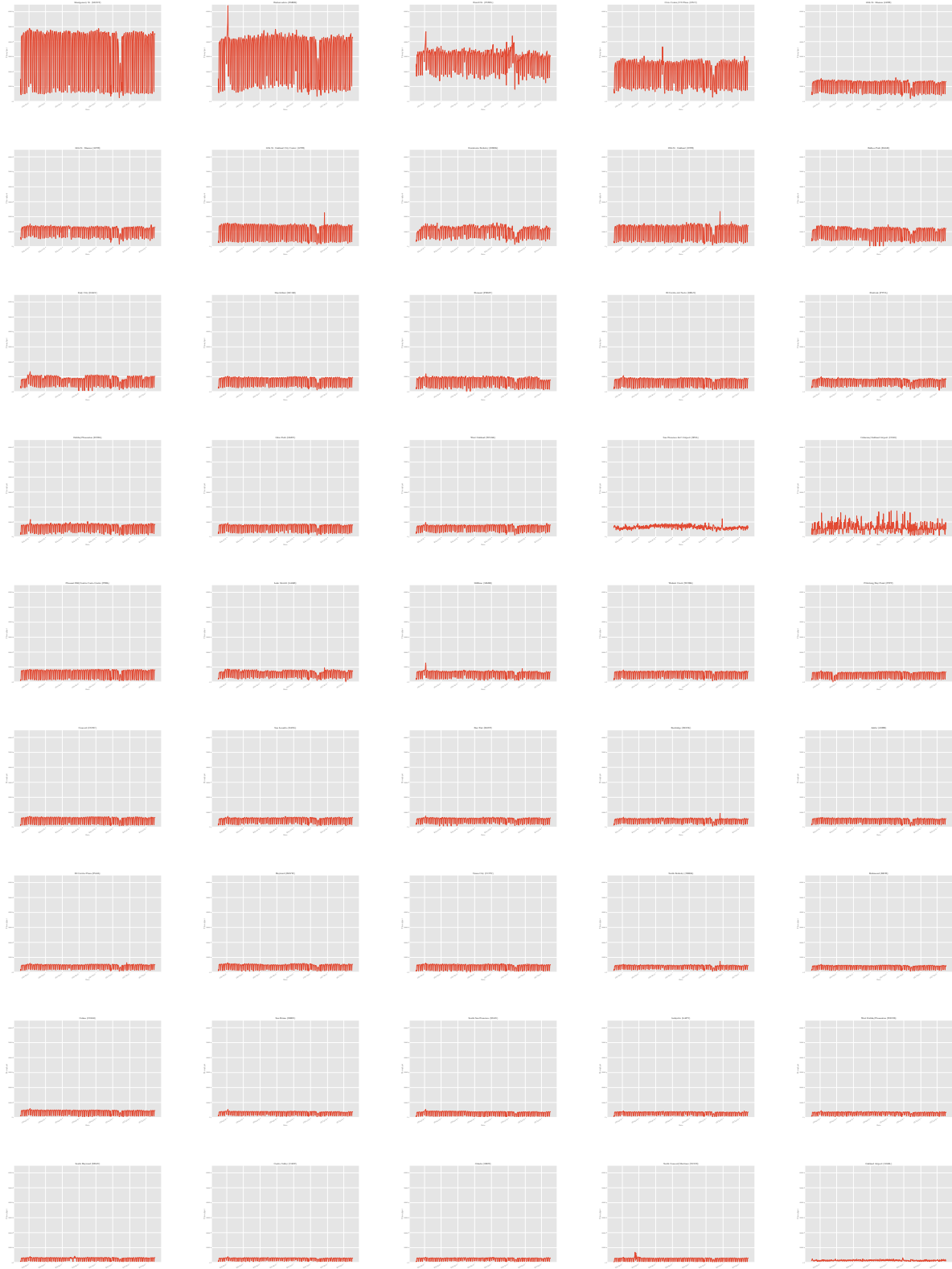Figure 3: Barplot for average hourly throughput of each station of 2016 and 2017

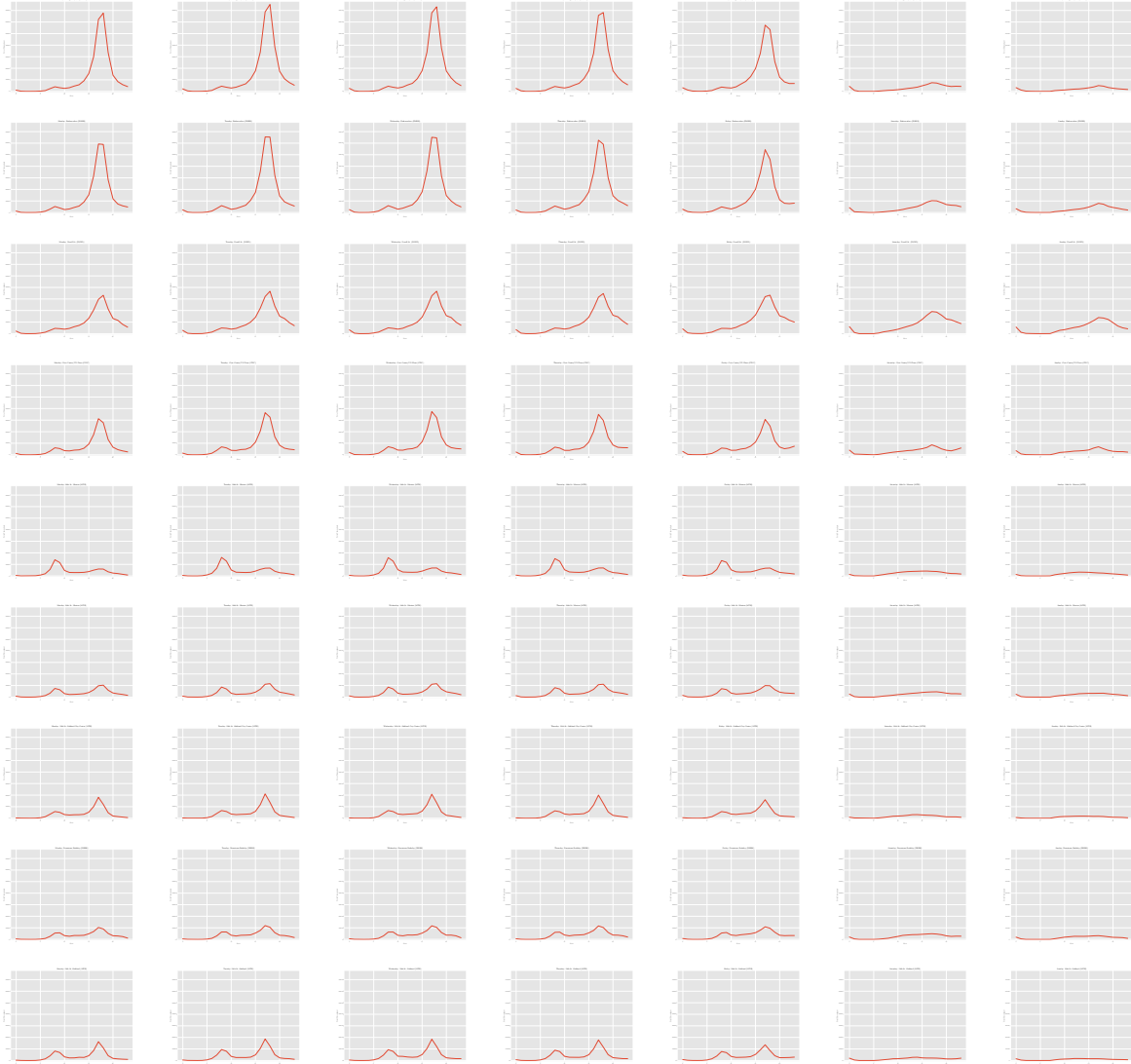Figure 4: Throughput of each station of 2016 and 2017

Figure 5: Total hourly throughput of a week of stations MONT, EMBR, POWL, CIVC, 24TH, 16TH, 12TH, DBRK, 19TH
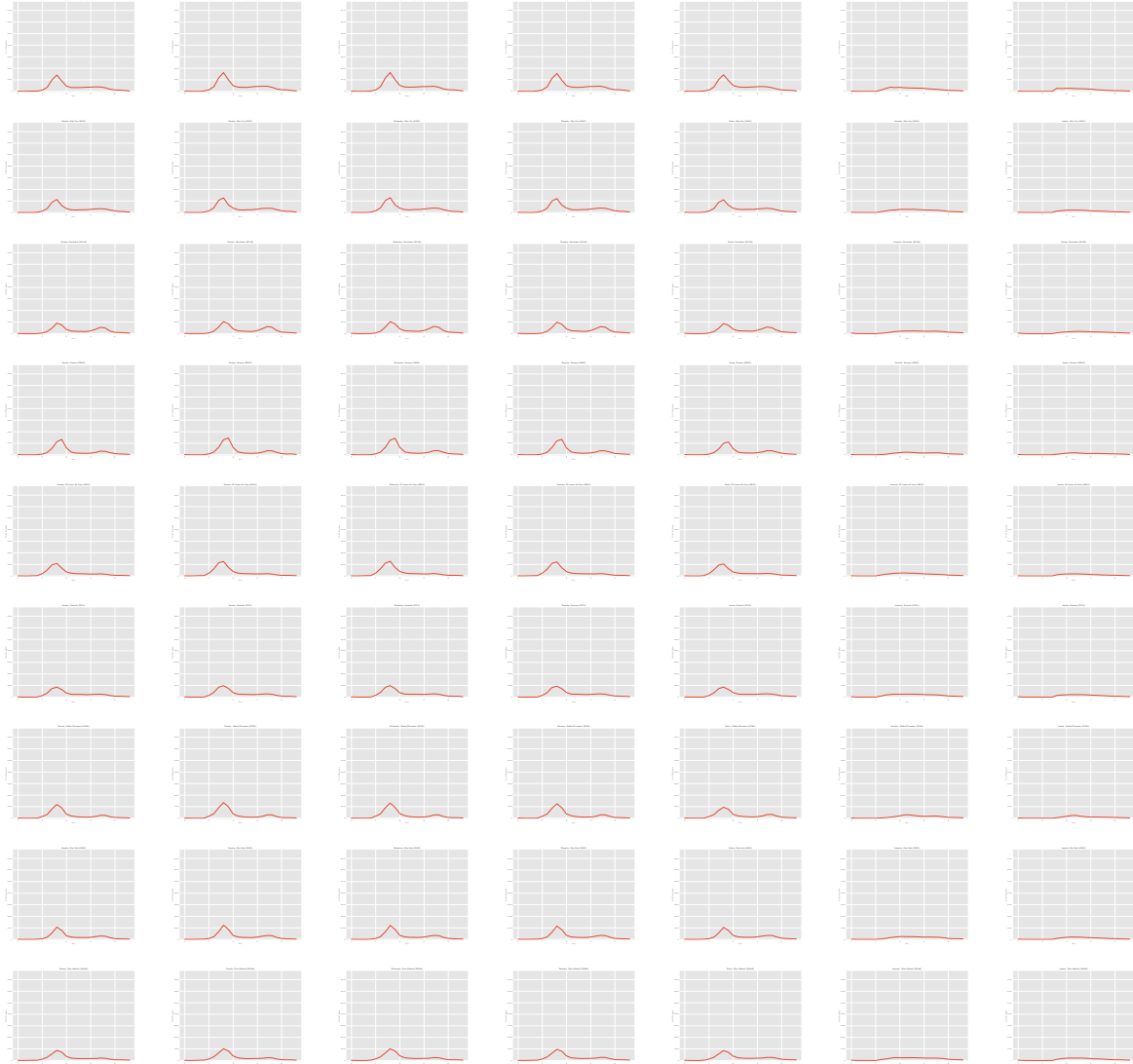
Figure 6: Total hourly throughput of a week of stations BALB, DALY, MCAR, FRMT, DELN, FTVL, DUBL, GLEN, WOAK
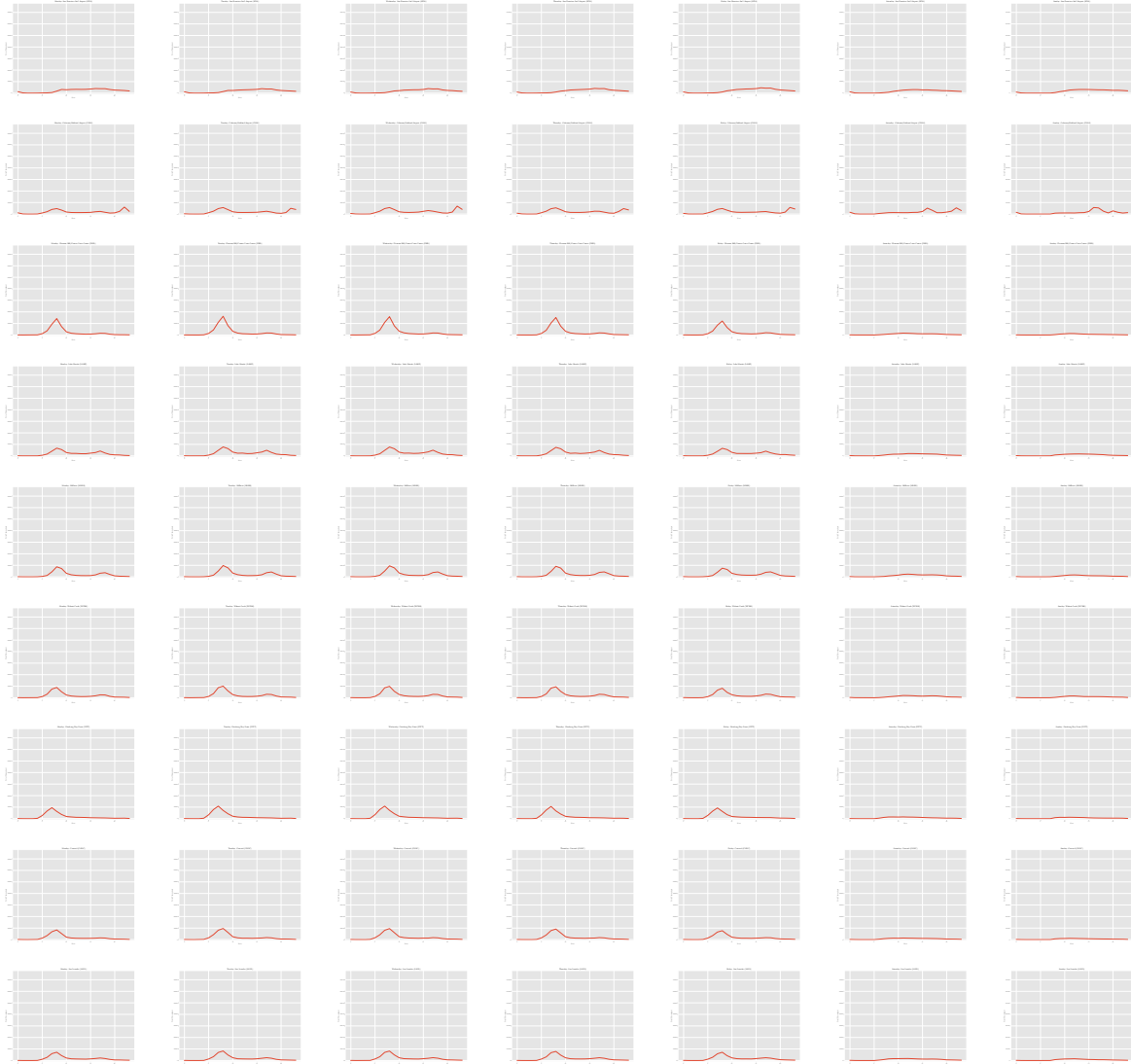
Figure 7: Total hourly throughput of a week of stations SFIA, COLS, PHIL, LAKE, MLBR, WCRK, PITT, CONC, SANL
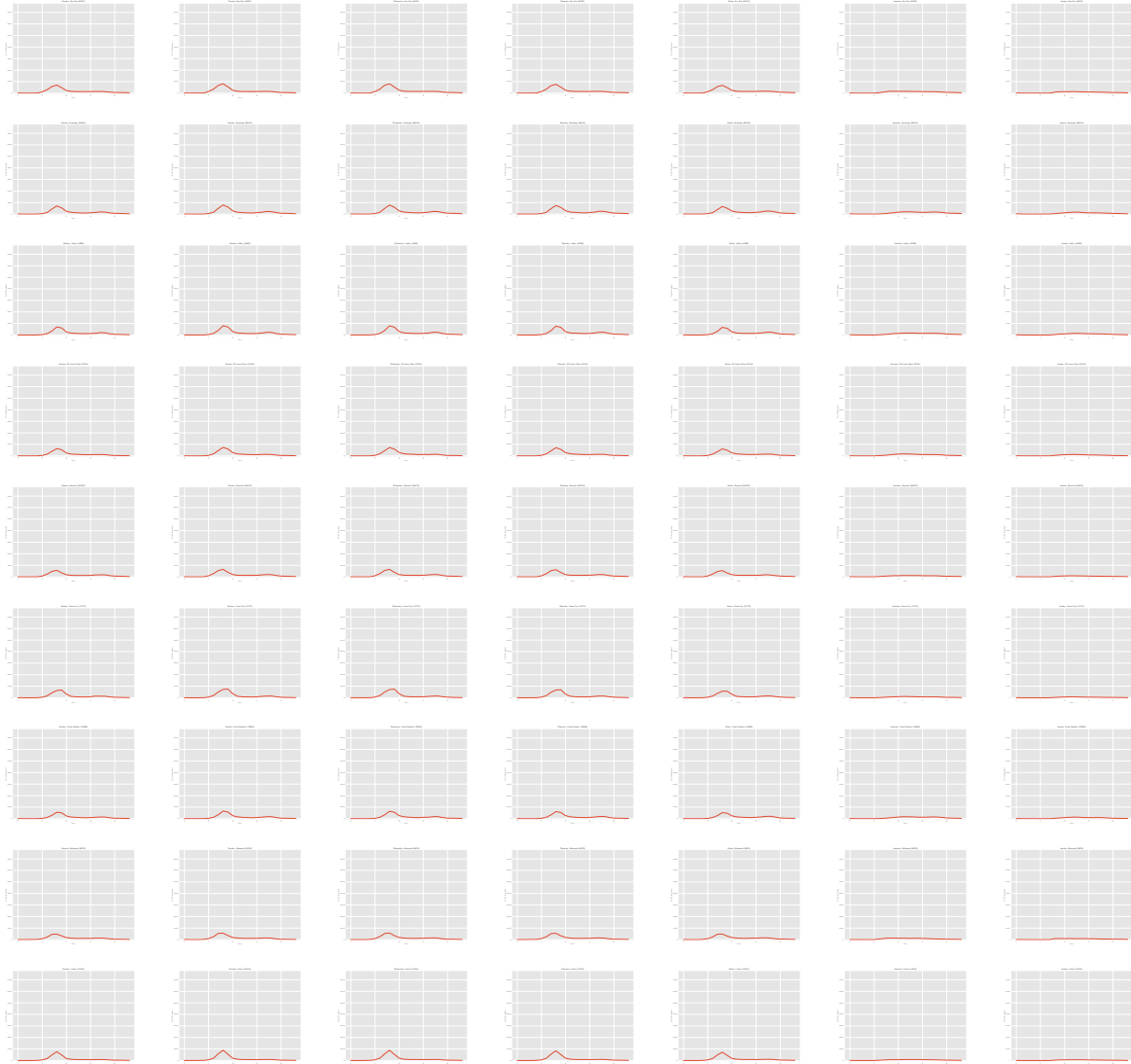
Figure 8: Total hourly throughput of a week of stations BAYF, ROCK, ASHB, PLZA, HAYW, UCTY, NBRK, RICH, COLM
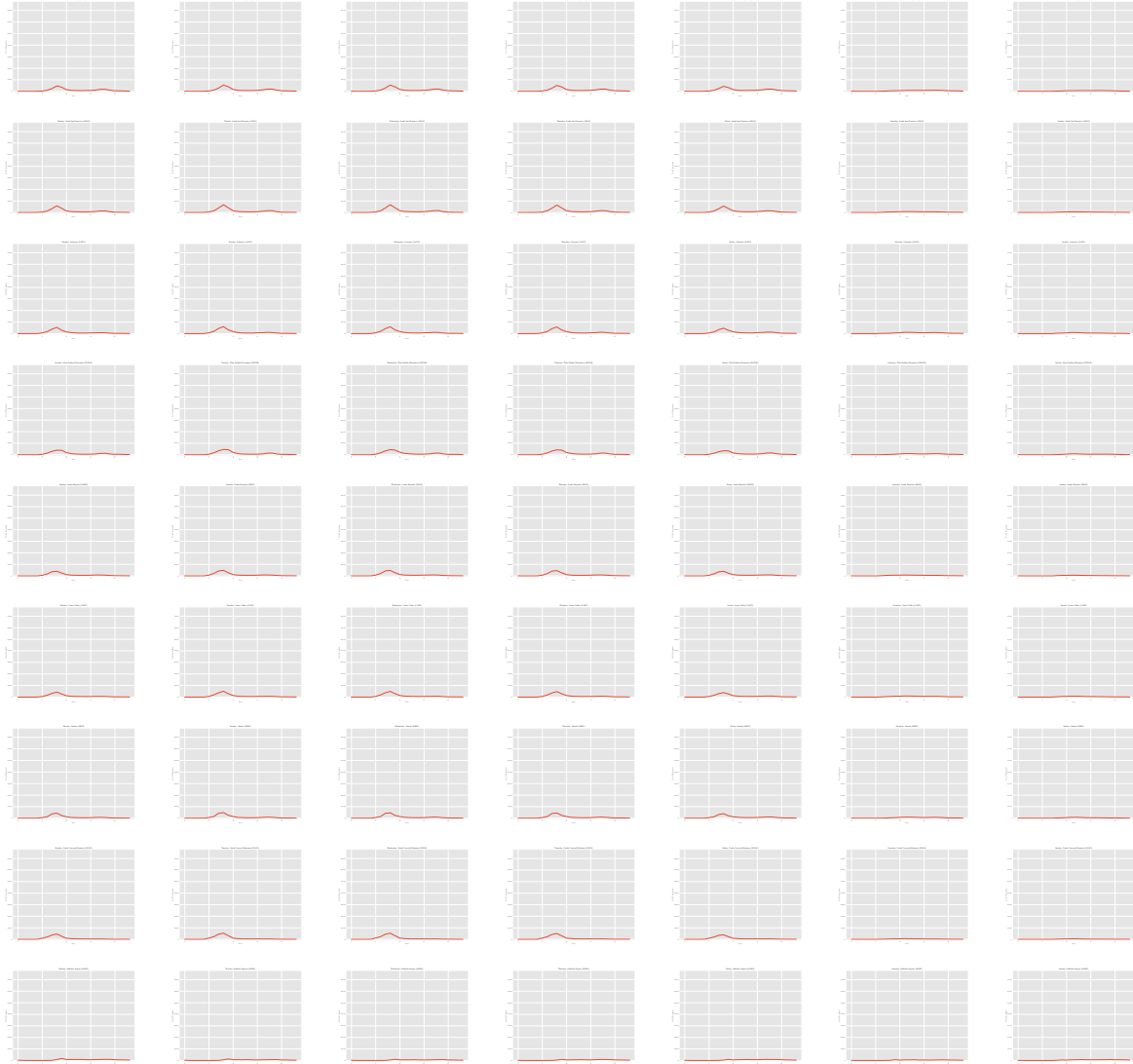
Figure 9: Total hourly throughput of a week of stations SBRN, SSAN, LAFY, WDUB, SHAY, CAST, ORIN, NCON, OAKL

the accuracy rate to 45%, and using 10% of the data as training dataset, the decision tree can reach at 77.59% accuracy rate. As the size of training dataset becomes larger, the accuracy rate goes higher, at the same time the decision tree becomes more complicated (more branches).
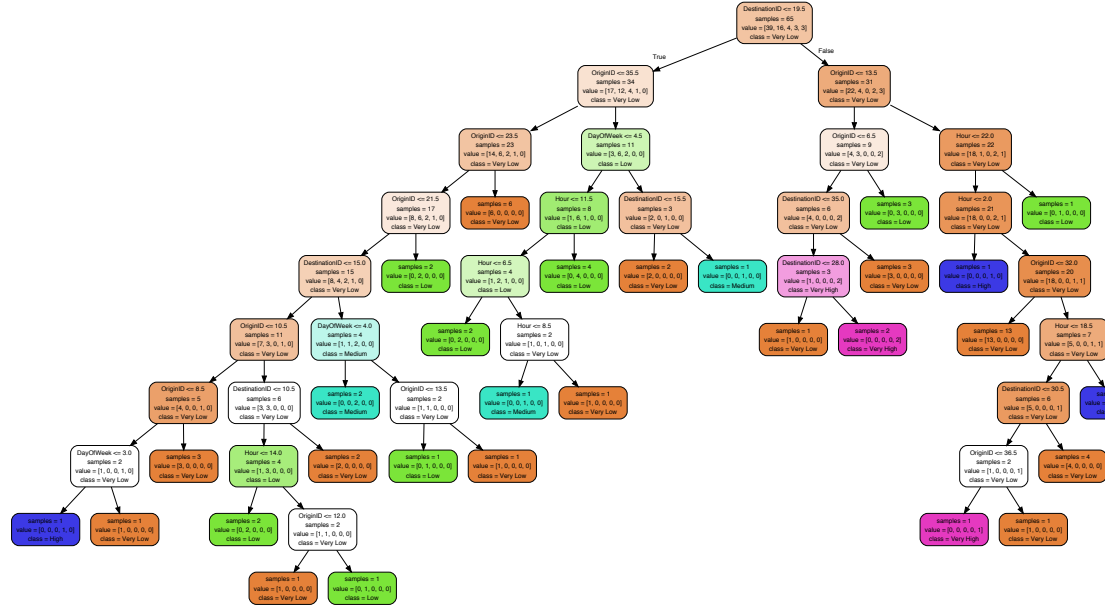


Figure 10: The decision tree obtained by using 0.02% of the data

# 5 Conclusions & Future Work

## 5.1 Conclusions

Now we are able to answer some of the quetions related to the BART as follows:

1. The busiest station is Montgomery St.
2. The least popular BART route would be Richmond – Fremont.
3. The best time to go to downtown San Francisco from, say Dublin/Pleasanton Station, if you want to find a seat is around 11:00 AM.
4. The busiest day of the week is Tuesday.

## 5.2 Future Work

Due to the time limitation, we are unable to build visualize the route of BART lines on a map. So we'd like to realize that and if possible, we'd like to use D3.js to build a interactive interface which enables us to understand the data better by the animation.

# References

[1] Kaggle, *Bay Area Rapid Transit Ridership Datasets*, https://www.kaggle.com/saulfuh/bart-ridership/data, Obtained on 04/30/2018.

[2] BART Official Website, *BART Schedules*, https://www.bart.gov/schedules/bystation, Obtained on 04/30/2018.

[3] Jonathan Bouchet, *BART Transit System*, https://www.kaggle.com/jonathanbouchet/bart-transit-system, Obtained on 04/30/2018.

[4] Xiuhua Han, et al., *Research on Data Mining of Public Transit IC Card and Application*, 2010 International Conference on Intelligent Computation Technology and Automation, https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=arnumber=5522901tag=1, Obtained on 04/30/2018.

[5] Roy Ka-Wei Lee, et al., *Time-Series Data Mining in Transportation: A Case Study on Singapore Public Train Commuter Travel Patterns*, IACSIT International Journal of Engineering and Technology, Vol. 6, No. 5, October 2014, https://pdfs.semanticscholar.org/ffc1/399342b36e178f511af57723f38e37ac7793.pdf, Obtained on 04/30/2018.

# 6 Python Program

```python
#!/usr/bin/env python3
# -*- coding: utf-8 -*-

# import math
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import matplotlib as mpl
from matplotlib.backends.backend_pdf import PdfPages
# from IPython.core.interactiveshell import InteractiveShell
# InteractiveShell.ast_node_interactivity = 'all'

plt.rc('text', usetex=True)
plt.rc('font', family='serif')

__author__ = 'Libao Jin'
__date__ = '04/24/2018'
__email__ = 'ljin1@uwyo.edu'
__copyright__ = 'Copyright (c) 2018 Libao Jin'


def bart_preprocess(bart):
    # Convert string to datetime, split datetime into date and time, and get day of week
    bart['DateTime'] = pd.to_datetime(bart.DateTime)
    bart['Date'] = bart['DateTime'].dt.date
    bart['Time'] = bart['DateTime'].dt.time
    bart['DayOfWeek'] = bart['DateTime'].dt.weekday_name
    # bart.drop(columns='DateTime')
    return bart


def stat_preprocess(stat):
    # Split location into longitude and latitude for visualization later on
    loc = stat.Location.str.split(',', expand=True)
    loc = [pd.to_numeric(loc[i]) for i in loc.columns]
    stat['Longitude'], stat['Latitude'] = loc[0], loc[1]
    columns = ['Abbreviation', 'Name', 'Longitude', 'Latitude', 'Description']
    stat = stat[columns]
    return stat


def generate_bart_routes(bart_lines, stat):
```

```python
43      bart_routes = []
44
45      for line in bart_lines:
46
47          abbr = []
48          name = []
49          fullname = []
50          location = []
51
52          for station in line:
53              tmp = stat[stat['Name'].str.contains(station)]
54              abbr.append(tmp.iloc[0]['Abbreviation'])
55              name.append(station)
56              fullname.append(tmp.iloc[0]['Name'])
57              location.append([tmp.iloc[0]['Longitude'],
58                               tmp.iloc[0]['Latitude']])
59
60          bart_routes.append({
61              'abbr': abbr,
62              'name': name,
63              'fullname': fullname,
64              'location': location
65          })
66      return bart_routes
67
68
69  def visualize_bart_routes(bart_routes, filename='bart_routes.pdf'):
70      color = list('krbcy')
71      alpha = [0.9, 0.8, 0.7, 0.6, 0.5]
72      x_min, x_max, y_min, y_max = -122.62, -121.75, 37.48, 38.04
73
74      with PdfPages(filename) as pdf:
75          plt.style.use('default')
76          plt.rc('text', usetex=True)
77          plt.rc('font', family='serif')
78          for i in range(len(bart_routes)):
79              fig = plt.figure(figsize=(16, 12))
80              loc = np.array(bart_routes[i]['location'])
81              # x, y = loc[:, 0], loc[:, 1]
82              # plt.plot(x, y, '-o', c=color[i], ms=6, lw=4, alpha=0.5)
83              plt.plot(
84                  loc[:, 0],
85                  loc[:, 1],
86                  '-o',
87                  c=color[i],
88                  mfc='w',
89                  ms=8,
90                  lw=9,
91                  alpha=0.75
92              )
93              plt.xlim([x_min, x_max])
94              plt.ylim([y_min, y_max])
95              plt.xlabel('Longitude')
96              plt.ylabel('Latitude')
97              title = '{} - {}'.format(bart_routes[i]['fullname'][0],
98                                       bart_routes[i]['fullname'][-1])
99              plt.title(title)
100             plt.show(block=False)
101             pdf.savefig(fig)
```

```python
102
103          fig = plt.figure(figsize=(16, 12))
104          for i in range(len(bart_routes)):
105              loc = np.array(bart_routes[i]['location'])
106              label = 'Line {:d}: {} - {}'.format(i + 1,
107                                                   bart_routes[i]['abbr'][0],
108                                                   bart_routes[i]['abbr'][-1])
109              # x, y = loc[:, 0], loc[:, 1]
110              # plt.plot(x, y, '-o', c=color[i], ms=6, lw=4, alpha=0.5)
111              plt.plot(
112                  loc[:, 0],
113                  loc[:, 1],
114                  '-o',
115                  c=color[i],
116                  mfc='w',
117                  ms=8,
118                  lw=9,
119                  alpha=alpha[i],
120                  label=label
121              )
122          plt.legend()
123          plt.xlim([x_min, x_max])
124          plt.ylim([y_min, y_max])
125          plt.xlabel('Longitude')
126          plt.ylabel('Latitude')
127          plt.title('BART Lines')
128          plt.show(block=False)
129          # plt.axis('equal')
130          pdf.savefig(fig)
131
132
133  def visualize_throughput(bart_routes, bart_aggregate, column_name, filename='bart_throughput.pdf'):
134      with PdfPages(filename) as pdf:
135          color = list('krbcy')
136          bart_aggregate = bart_aggregate.set_index('Station')
137          x, y = bart_aggregate['Longitude'], bart_aggregate['Latitude']
138          n = len(bart_aggregate['Longitude'])
139          cmap = mpl.cm.cool
140          cs = getattr(bart_aggregate, column_name)
141          x_min, x_max, y_min, y_max = -122.62, -121.75, 37.48, 38.04
142          fig = plt.figure(figsize=(16, 12))
143          plt.style.use('default')
144          plt.rc('text', usetex=True)
145          plt.rc('font', family='serif')
146          for i in range(len(bart_routes)):
147              loc = np.array(bart_routes[i]['location'])
148              label = 'Line {:d}: {} - {}'.format(i + 1,
149                                                   bart_routes[i]['abbr'][0],
150                                                   bart_routes[i]['abbr'][-1])
151              # x, y = loc[:, 0], loc[:, 1]
152              # plt.plot(x, y, '-o', c=color[i], ms=6, lw=4, alpha=0.5)
153              plt.plot(
154                  loc[:, 0],
155                  loc[:, 1],
156                  '-',
157                  c=color[i],
158                  # mfc='w',
159                  # ms=8,
160                  # lw=6,
```

```python
                    alpha=0.5,
                    label=label,
                    zorder=1
                )
            plt.legend()

            plt.scatter(
                x,
                y,
                c=cs,
                s=cs,
                marker='o',
                edgecolors='k',
                cmap=cmap,
                alpha=1,
                zorder=2
            )

            plt.colorbar()
            plt.xlim([x_min, x_max])
            plt.ylim([y_min, y_max])
            plt.xlabel('Longitude')
            plt.ylabel('Latitude')
            plt.axis('equal')
            plt.title('BART Lines ' + cs.name)
            plt.show(block=False)
            pdf.savefig(fig)

            fig = plt.figure(figsize=(16, 12))
            plt.style.use('ggplot')
            plt.rc('text', usetex=True)
            plt.rc('font', family='serif')
            bart_aggregate.sort_values(column_name, inplace=True, ascending=False)
            cs = getattr(bart_aggregate, column_name)
            colors = cmap(np.linspace(1, 0, n))
            cs.plot(kind='bar', color=colors, alpha=0.75)
            # print(list(cs.index))
            plt.ylabel('Average Throughput / hour')
            plt.title('Throughput vs. Station')
            plt.show(block=False)
            pdf.savefig(fig)


def bart_aggregate_throughput(bart, filename):
    number_of_days = len(bart['DateTime'].dt.date.unique())
    number_of_hours = len(bart['DateTime'].dt.time.unique())
    bart_grouped = bart['Throughput'].groupby(bart['Origin']).sum().to_frame()
    bart_grouped['Destination'] = bart['Throughput'].groupby(bart['Destination']).sum()
    bart_grouped.index.names = ['Station']
    bart_grouped.columns = ['Throughput Origin', 'Throughput Destination']
    bart_grouped['Throughput All'] = bart_grouped['Throughput Origin'] + bart_grouped['Throughput Destination']
    bart_grouped.set_index(stat['Abbreviation'])
    bart_grouped.reset_index(level=0, inplace=True)
    bart_grouped['Longitude'] = stat['Longitude']
    bart_grouped['Latitude'] = stat['Latitude']
    bart_grouped[['Throughput Origin', 'Throughput Destination', 'Throughput All']] = \
        bart_grouped[['Throughput Origin', 'Throughput Destination', 'Throughput All']] / (number_of_days * number_of_hours)
    visualize_throughput(bart_routes, bart_grouped, 'Throughput Origin', filename)
    # visualize_throughput(bart_routes, bart_grouped, 'Throughput Destination', filename)
```

```python
220        # visualize_throughput(bart_routes, bart_grouped, 'Throughput All', filename)
221
222
223    def visualize_bart(bart, origins, stat, class_type, group_by, plot_option, filename='bart_overview.pdf'):
224
225        plt.style.use('ggplot')
226        plt.rc('text', usetex=True)
227        plt.rc('font', family='serif')
228
229        with PdfPages(filename) as pdf:
230            for stops in origins:
231                if len(stops) == 1:
232                    data = bart[bart['Origin'] == stops[0]]
233                elif len(stops) == 2:
234                    origin, dest = stops
235                    data = bart[(bart['Origin'] == origin) & (bart['Destination'] == dest)]
236                # Plot the throughput with respect to time (hour) each week/month
237                if plot_option == 1 or group_by == 'hour':
238                    k = len(getattr(data['DateTime'].dt, class_type[0]).unique())
239                    # n = 2                              # number of columns
240                    # m = math.ceil(k / n)  # number of rows
241                    # plt.figure(figsize=(8 * n, 6 * m))
242                    for i in range(k):
243                        fig = plt.figure(figsize=(8, 6))
244                        if class_type[0] == 'weekday':
245                            j = i
246                        else:
247                            j = i + 1
248                        grouped = data[getattr(data['DateTime'].dt, class_type[0]).values == j].groupby(
249                            getattr(data['DateTime'].dt, group_by)).sum()
250                        grouped.sort_index(inplace=True)
251                        plt.plot(grouped['Throughput'])
252                        if len(stops) == 1:
253                            tmp = stat[stat['Abbreviation'] == data.iloc[0]['Origin']]
254                            title = '{}: {}'.format(class_type[1][i],
255                                                    tmp.iloc[0]['Name'])
256                        else:
257                            tmp_1 = stat[stat['Abbreviation'] == data.iloc[0]['Origin']]
258                            tmp_2 = stat[stat['Abbreviation'] == data.iloc[0]['Destination']]
259                            title = '{}: {} - {}'.format(class_type[1][i],
260                                                         tmp_1.iloc[0]['Name'],
261                                                         tmp_2.iloc[0]['Name'])
262                        plt.title(title)
263                        plt.xlabel(group_by.title())
264                        ax = plt.gca()
265                        # fig = plt.gcf()
266                        if group_by == 'date':
267                            xfmt = mpl.dates.DateFormatter('%Y-%m-%d')
268                            ax.xaxis.set_major_formatter(xfmt)
269                            # plt.xticks(rotation=90)
270                            fig.autofmt_xdate()
271                        plt.ylabel('Total Throughput')
272                        plt.ylim([-10000, 780000])
273                        # plt.xticks(grouped.index, list(np.arange(24)))
274                        plt.show(block=False)
275                        pdf.savefig(fig)
276
277                # Plot the throughput with respect to date
278                elif plot_option == 2:
```

```python
                    grouped = data.groupby(getattr(data['DateTime'].dt, group_by)).sum()
                    grouped.sort_index(inplace=True)
                    fig = plt.figure(figsize=(8, 6))
                    plt.plot(grouped['Throughput'])
                    if len(stops) == 1:
                        tmp = stat[stat['Abbreviation'] == data.iloc[0]['Origin']]
                        title = '{}'.format(tmp.iloc[0]['Name'])
                        # title = class_type[1][i] + ': {}'.format(data.iloc[0]['Origin'])
                    else:
                        tmp_1 = stat[stat['Abbreviation'] == data.iloc[0]['Origin']]
                        tmp_2 = stat[stat['Abbreviation'] == data.iloc[0]['Destination']]
                        title = '{} - {}'.format(tmp_1.iloc[0]['Name'],
                                                 tmp_2.iloc[0]['Name'])
                        # title = ': {} - {}'.format(data.iloc[0]['Origin'], data.iloc[0]['Destination'])
                    # plt.legend()
                    plt.xlabel(group_by.title())
                    plt.ylabel('Throughput')
                    plt.xticks(rotation=90)
                    ax = plt.gca()
                    # fig = plt.gcf()
                    xfmt = mpl.dates.DateFormatter('%Y-%m-%d')
                    ax.xaxis.set_major_formatter(xfmt)
                    # plt.xticks(rotation=90)
                    fig.autofmt_xdate()
                    plt.title(title)
                    plt.ylim([0, 65000])
                    plt.show(block=False)
                    pdf.savefig(fig)


if __name__ == '__main__':

    dest_folder = './output'
    plt.rc('text', usetex=True)
    plt.rc('font', family='serif')
    # Load data
    # Data obtained from https://www.kaggle.com/saulfuh/bart-ridership
    date_hour_2016 = '../data/bart-ridership/date-hour-soo-dest-2016.csv'
    date_hour_2017 = '../data/bart-ridership/date-hour-soo-dest-2017.csv'
    stat_info = '../data/bart-ridership/station_info.csv'
    # bart_16 = pd.read_csv(date_hour_2016)
    # bart_17 = pd.read_csv(date_hour_2017)
    # stat = pd.read_csv(stat_info)

    # Data preprocessing
    bart_16 = bart_preprocess(pd.read_csv(date_hour_2016))
    bart_17 = bart_preprocess(pd.read_csv(date_hour_2017))
    bart = pd.concat([bart_16, bart_17], ignore_index=True)
    stat = stat_preprocess(pd.read_csv(stat_info))

    # Visualize the routes according to the BART official website

    line_1 = [
        'Richmond',
        'El Cerrito del Norte',
        'El Cerrito Plaza',
        'North Berkeley',
        'Downtown Berkeley',
        'Ashby',
```

```
338          'West Oakland',
339          'Embarcadero',
340          'Montgomery St.',
341          'Powell St.',
342          'Civic Center/UN Plaza',
343          'Daly City',
344          'Colma',
345          'South San Francisco',
346          'San Bruno',
347          'Millbrae'
348      ]
349
350      line_2 = [
351          'Richmond',
352          'El Cerrito del Norte',
353          'El Cerrito Plaza',
354          'North Berkeley',
355          'Downtown Berkeley',
356          'Ashby',
357          'MacArthur',
358          '19th St. Oakland',
359          '12th St. Oakland City Center',
360          'Lake Merritt',
361          'Fruitvale',
362          'Coliseum/Oakland Airport',
363          'San Leandro',
364          'Bay Fair',
365          'Hayward',
366          'South Hayward',
367          'Union City',
368          'Fremont'
369      ]
370
371      line_3 = [
372          'Pittsburg/Bay Point',
373          'North Concord/Martinez',
374          'Concord',
375          'Walnut Creek',
376          'Lafayette',
377          'Orinda',
378          'Rockridge',
379          'MacArthur',
380          '19th St. Oakland',
381          '12th St. Oakland City Center',
382          'West Oakland',
383          'Embarcadero',
384          'Montgomery St.',
385          'Powell St.',
386          'Civic Center/UN Plaza',
387          '16th St. Mission',
388          '24th St. Mission',
389          'Glen Park',
390          'Balboa Park',
391          'Daly City',
392          'Colma',
393          'South San Francisco',
394          'San Bruno',
395          "San Francisco Int'l Airport",
396          'Millbrae'
```

```python
397        ]
398
399        line_4 = [
400            'Dublin/Pleasanton',
401            'West Dublin/Pleasanton',
402            'Castro Valley',
403            'Bay Fair',
404            'San Leandro',
405            'Coliseum/Oakland Airport',
406            'Fruitvale',
407            'Lake Merritt',
408            'West Oakland',
409            'Embarcadero',
410            'Montgomery St.',
411            'Powell St.',
412            'Civic Center/UN Plaza',
413            '16th St. Mission',
414            '24th St. Mission',
415            'Glen Park',
416            'Balboa Park',
417            'Daly City'
418        ]
419
420        line_5 = [
421            'Warm Springs/South Fremont',
422            'Fremont',
423            'Union City',
424            'South Hayward',
425            'Hayward',
426            'Bay Fair',
427            'San Leandro',
428            'Coliseum/Oakland Airport',
429            'Fruitvale',
430            'Lake Merritt',
431            'West Oakland',
432            'Embarcadero',
433            'Montgomery St.',
434            'Powell St.',
435            'Civic Center/UN Plaza',
436            '16th St. Mission',
437            '24th St. Mission',
438            'Glen Park',
439            'Balboa Park',
440            'Daly City'
441        ]
442
443        bart_lines = [
444            line_1,
445            line_2,
446            line_3,
447            line_4,
448            line_5
449        ]
450
451        bart_routes = generate_bart_routes(bart_lines, stat)
452        visualize_bart_routes(bart_routes, '{}/bart_routes.pdf'.format(dest_folder))
453
454        # bart_aggregate_throughput(bart_16)
455        # bart_aggregate_throughput(bart_17)
```

```
456      bart_aggregate_throughput(bart, '{}/bart_throughput.pdf'.format(dest_folder))
457
458      by_month = [
459          'month',
460          ['January', 'Febrary', 'March', 'April', 'May', 'June', 'July',
461           'August', 'September', 'October', 'November', 'December']
462      ]
463
464      by_weekday = [
465          'weekday',
466          ['Monday', 'Tuesday', 'Wednesday', 'Thursday',
467           'Friday', 'Saturday', 'Sunday']
468      ]
469
470      origins_names = [
471          'MONT', 'EMBR', 'POWL', 'CIVC', '24TH', '16TH', '12TH',
472          'DBRK', '19TH', 'BALB', 'DALY', 'MCAR', 'FRMT', 'DELN',
473          'FTVL', 'DUBL', 'GLEN', 'WOAK', 'SFIA', 'COLS', 'PHIL',
474          'LAKE', 'MLBR', 'WCRK', 'PITT', 'CONC', 'SANL', 'BAYF',
475          'ROCK', 'ASHB', 'PLZA', 'HAYW', 'UCTY', 'NBRK', 'RICH',
476          'COLM', 'SBRN', 'SSAN', 'LAFY', 'WDUB', 'SHAY', 'CAST',
477          'ORIN', 'NCON', 'OAKL'
478      ]
479
480      origins = [[i] for i in origins_names]
481      visualize_bart(bart, origins, stat, by_month, 'date', 2, '{}/bart_overview_1.pdf'.format(dest_folder))
482      visualize_bart(bart, origins, stat, by_weekday, 'hour', 2, '{}/bart_overview_2.pdf'.format(dest_folder))
483      # origins2 = [[origins_names[i], origins_names[i + 1]] for i in range(len(origins_names) - 1)]
484      # visualize_bart(bart, origins2, stat, by_month, 'hour', 1, '{}/bart_overview_3.pdf'.format(dest_folder))
```

```
1   #!/usr/bin/env python3
2   # -*- coding: utf-8 -*-
3
4   __author__ = 'Libao Jin'
5   __date__ = '05/02/2018'
6   __email__ = 'ljin1@uwyo.edu'
7
8
9   import pandas as pd
10  # import numpy as np
11  import matplotlib.pyplot as plt
12  from matplotlib.backends.backend_pdf import PdfPages
13  from sklearn import tree
14  from sklearn import linear_model
15  from sklearn.model_selection import train_test_split
16  from sklearn.externals.six import StringIO
17  import pydot
18
19
20  class BartClassifier(object):
21
22      def __init__(self):
23          plt.style.use('ggplot')
24          plt.rc('text', usetex=True)
25          plt.rc('font', family='serif')
26
27      def load_data(self, filenames):
28          '''Load datasets in batch'''
29          datasets = []
```

```python
30          for filename in filenames:
31              data = pd.read_csv(filename)
32              datasets.append(data)
33          print('Data loaded.')
34          return datasets
35
36      def data_preparation(self, bart, stat):
37          loc = stat['Location'].str.split(',', expand=True)
38          loc = [pd.to_numeric(loc[i]) for i in loc.columns]
39          stat['Longitude'], stat['Latitude'] = loc[0], loc[1]
40          stat_tmp = stat.copy()
41          stat_tmp.index.names = ['Stat_ID']
42          stat_tmp.reset_index(level=0, inplace=True)
43          stat_id = stat_tmp.set_index('Abbreviation')['Stat_ID'].dropna()
44          stat_lon = stat_tmp.set_index('Abbreviation')['Longitude'].dropna()
45          stat_lat = stat_tmp.set_index('Abbreviation')['Latitude'].dropna()
46          bart = bart.drop(bart[(bart['Origin'] == 'WSPR') | (bart['Destination'] == 'WSPR')].index)
47          bart['DateTime'] = pd.to_datetime(bart['DateTime'])
48          bart['DayOfWeek'] = bart['DateTime'].dt.weekday
49          bart['Month'] = bart['DateTime'].dt.month
50          bart['Day'] = bart['DateTime'].dt.day
51          bart['Hour'] = bart['DateTime'].dt.hour
52          print('First session done.')
53          bart['OriginID'] = bart['Origin'].replace(stat_id)
54          bart['OriginLongitude'] = bart['Origin'].replace(stat_lon)
55          bart['OriginLatitude'] = bart['Origin'].replace(stat_lat)
56          print('Second session done.')
57          bart['DestinationID'] = bart['Destination'].replace(stat_id)
58          bart['DestinationLongitude'] = bart['Destination'].replace(stat_lon)
59          bart['DestinationLatitude'] = bart['Destination'].replace(stat_lat)
60          bart['ThroughputLevel'] = bart['Throughput'].apply(self.throughput_level)
61          print('Data prepared.')
62          return (stat, bart)
63
64      def throughput_level(self, throughput):
65          if throughput <= 5:
66              level = 0
67          elif throughput <= 15:
68              level = 1
69          elif throughput <= 30:
70              level = 2
71          elif throughput <= 80:
72              level = 3
73          else:
74              level = 4
75          return level
76
77      def generate_train_test(self, data, feature_keys, target_keys):
78          train_features, test_features, train_labels, test_labels = train_test_split(
79              data[feature_keys],
80              data[target_keys],
81              test_size=0.99998
82          )
83          return (train_features, test_features, train_labels, test_labels)
84
85      def decision_tree(self, data, feature_keys, target_keys, filename='./output/bart_dt.pdf'):
86          train_features, test_features, train_labels, test_labels = self.generate_train_test(
87              data, feature_keys, target_keys)
88          clf = tree.DecisionTreeClassifier()
```

```python
89          clf.fit(train_features, train_labels)
90          test_labels_predict = clf.predict(test_features)
91          hit_rate = sum([1 for i in range(len(test_labels)) if test_labels_predict[i] ==
92                          test_labels.values[i]]) / len(test_labels) * 100
93          print('Decision Tree: Prediction hit/accuracy rate: {:.2f}%'.format(hit_rate))
94          self.visualize(clf, feature_keys, target_keys, filename)
95          train_data, test_data = pd.DataFrame(), pd.DataFrame()
96          train_data[feature_keys], train_data[target_keys], test_data[feature_keys], test_data[target_keys] = \
97              train_features, train_labels, test_features, test_labels
98          test_data = test_data.reset_index(drop=True)
99          test_labels_predict = pd.DataFrame(data={'predict': test_labels_predict})
100         test_data['predict'] = test_labels_predict
101         train_data.to_csv('../data/training_data.csv', encoding='utf-8', index=False)
102         test_data.to_csv('../data/test_data.csv', encoding='utf-8', index=False)
103         return (test_features, test_labels, test_labels_predict, hit_rate)
104
105     def visualize(self, clf, feature_keys, target_keys, filename):
106         feature_names = feature_keys
107         target_names = ['Very Low', 'Low', 'Medium', 'High', 'Very High']
108         dot_data = StringIO()
109         tree.export_graphviz(
110             clf,
111             out_file=dot_data,
112             feature_names=feature_names,
113             class_names=target_names,
114             filled=True,
115             rounded=True,
116             impurity=False
117         )
118         graph = pydot.graph_from_dot_data(dot_data.getvalue())
119         graph[0].write_pdf(filename)
120
121     def linear_regression(self, data, feature_keys, target_keys, filename='./output/bart_lr.pdf'):
122         lm = linear_model.LinearRegression()
123         X = data[feature_keys]
124         y = data[target_keys]
125         lm.fit(X, y)
126         predictions = lm.predict(X)
127         print('Linear Regression: Coefficients: {}\nIntercept: {}\nScore: {}'.format(lm.coef_,
128                                                                     lm.intercept_,
129                                                                     lm.score(X, y)))
130
131         with PdfPages(filename) as pdf:
132             fig, ax = plt.subplots()
133             ax.scatter(y, predictions, edgecolors=(0, 0, 0))
134             ax.plot([min(predictions), max(predictions)], [min(predictions), max(predictions)], 'b--', lw=2)
135             ax.set_xlabel('Measured')
136             ax.set_ylabel('Predicted')
137             plt.show(block=False)
138             pdf.savefig(fig)
139
140     def run(self):
141         # filenames = [
142         #     '../data/bart-ridership/date-hour-soo-dest-2016.csv',
143         #     '../data/bart-ridership/date-hour-soo-dest-2017.csv',
144         #     '../data/bart-ridership/station_info.csv'
145         # ]
146
147         # datasets = self.load_data(filenames)
```

```python
148            # bart = pd.concat(datasets[0:2], ignore_index=True)
149            # stat = datasets[2]
150
151            filenames = [
152                # '../data/bart-ridership/date-hour-soo-dest-2016.csv',
153                '../data/bart-ridership/date-hour-soo-dest-2017.csv',
154                '../data/bart-ridership/station_info.csv'
155            ]
156            datasets = self.load_data(filenames)
157            bart = datasets[0]
158            stat = datasets[1]
159
160            stat, bart = self.data_preparation(bart, stat)
161            print(stat.head())
162            print(bart.head())
163
164            # feature_keys = [
165            #     'OriginID', 'OriginLongitude', 'OriginLatitude',
166            #     'DestinationID', 'DestinationLongitude', 'DestinationLatitude',
167            #     'Month', 'Day', 'Hour', 'DayOfWeek'
168            # ]
169
170            feature_keys = [
171                'OriginID',
172                'DestinationID',
173                'Hour', 'DayOfWeek'
174            ]
175
176            target_keys = ['ThroughputLevel']
177
178            self.decision_tree(bart, feature_keys, target_keys)
179            self.linear_regression(bart, feature_keys, target_keys)
180
181
182 if __name__ == '__main__':
183     bc = BartClassifier()
184     bc.run()
```