

COSC 5010 - Data Science for Security Final Project

Libao Jin (ljin1@uwyo.edu)

May 2, 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 `DataFrames`.

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.

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.

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-des-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

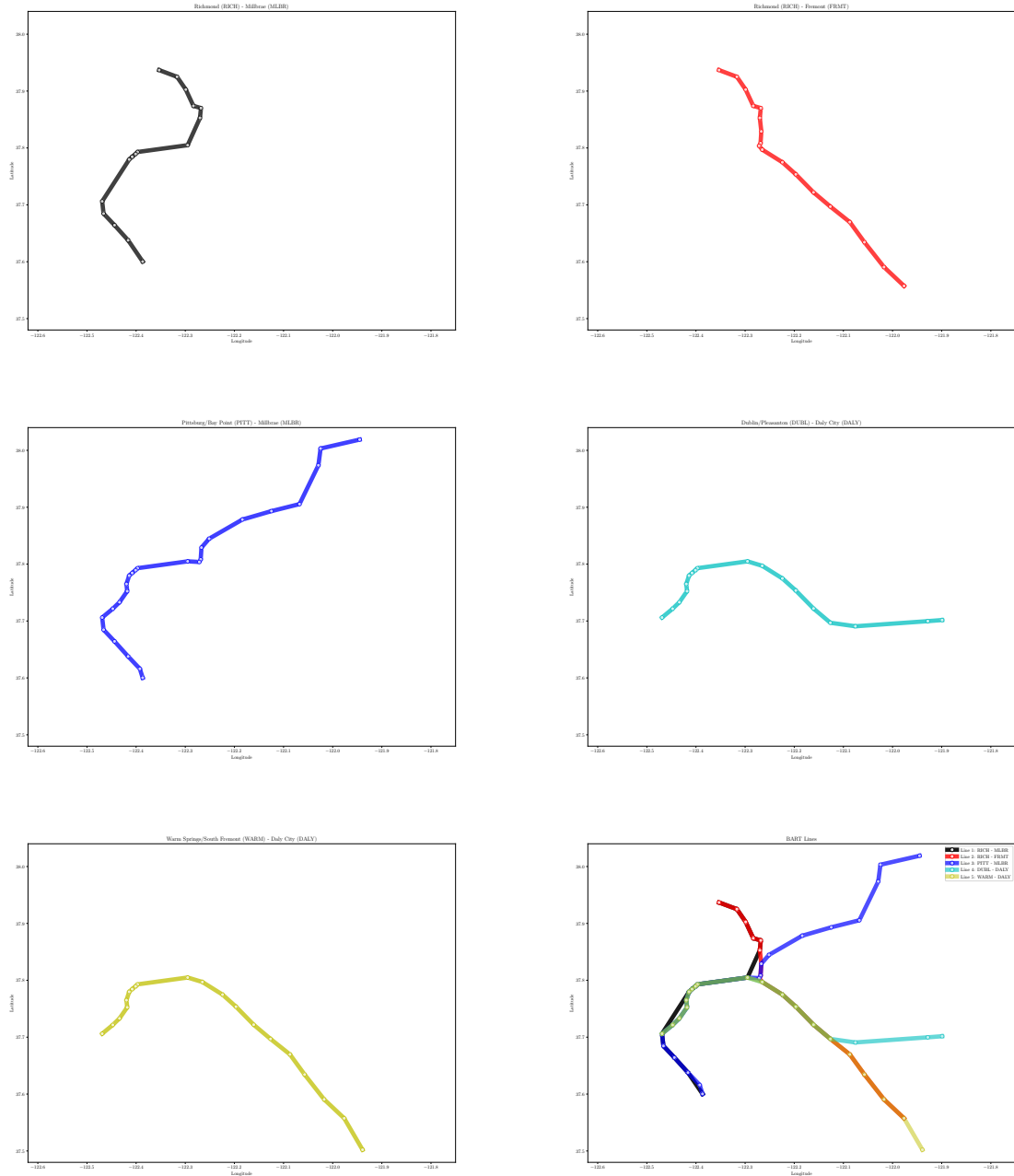


Figure 1: Routes of the BART Lines

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...

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.

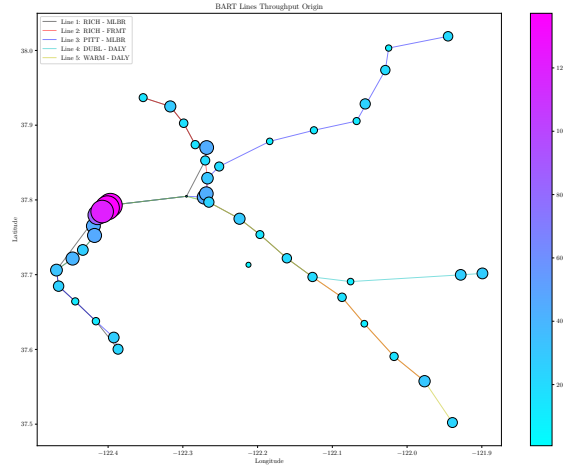


Figure 2: Average hourly throughput of each station of 2016 and 2017

Also, we obtained the barplot for the average hourly throughput of each station of 2016 and 2017 as shown in 3. Obviously, the throughputs of aforementioned four stations far outweigh that of the rest stations.

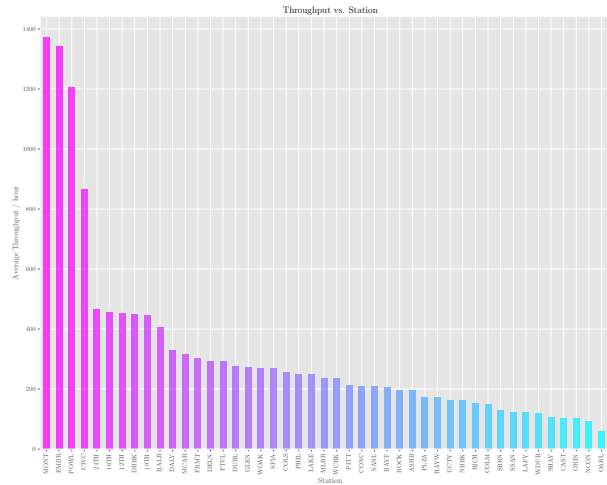


Figure 3: Barplot for average hourly throughput of each station of 2016 and 2017

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.

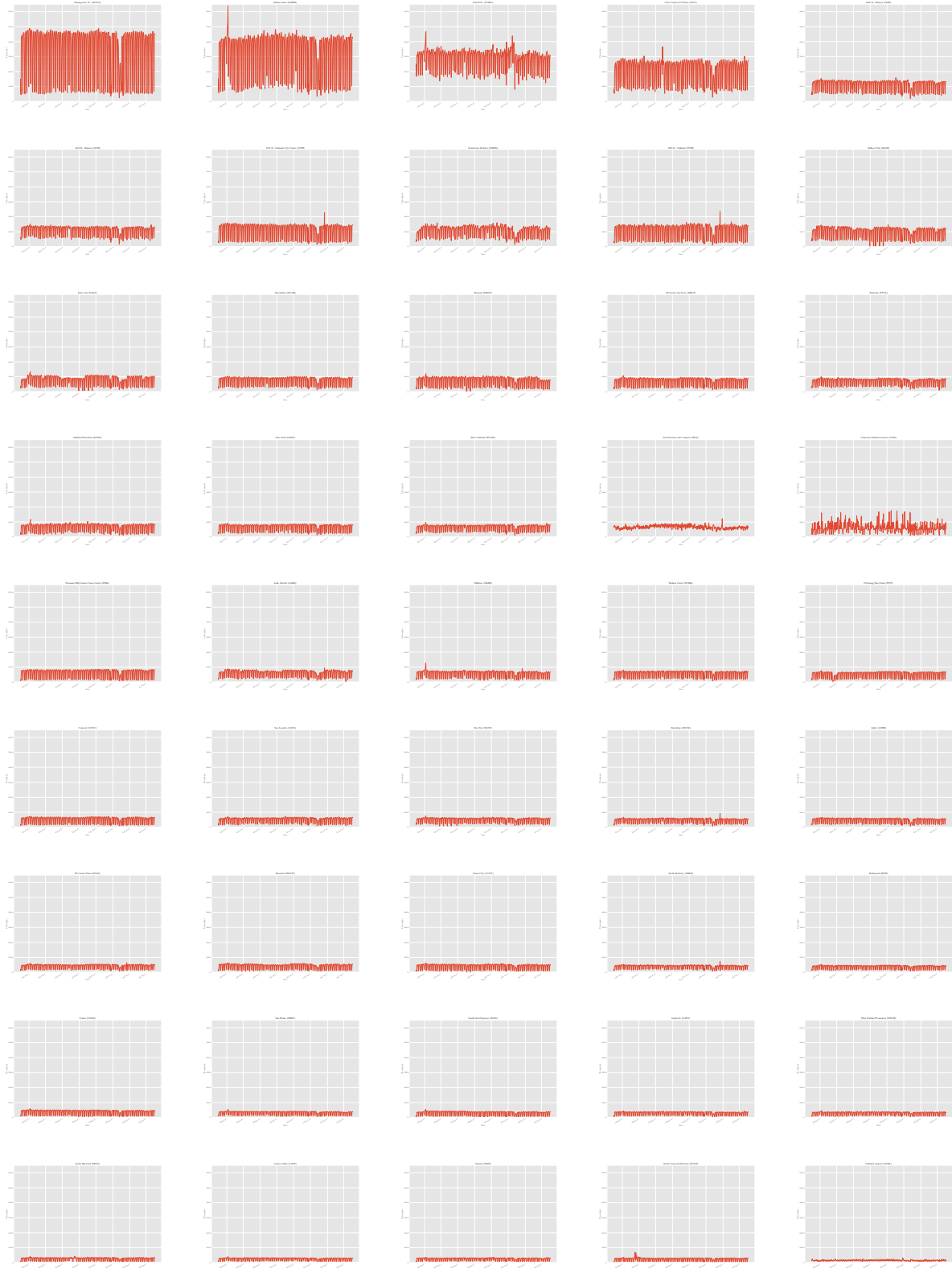


Figure 4: Throughput of each station of 2016 and 2017

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.

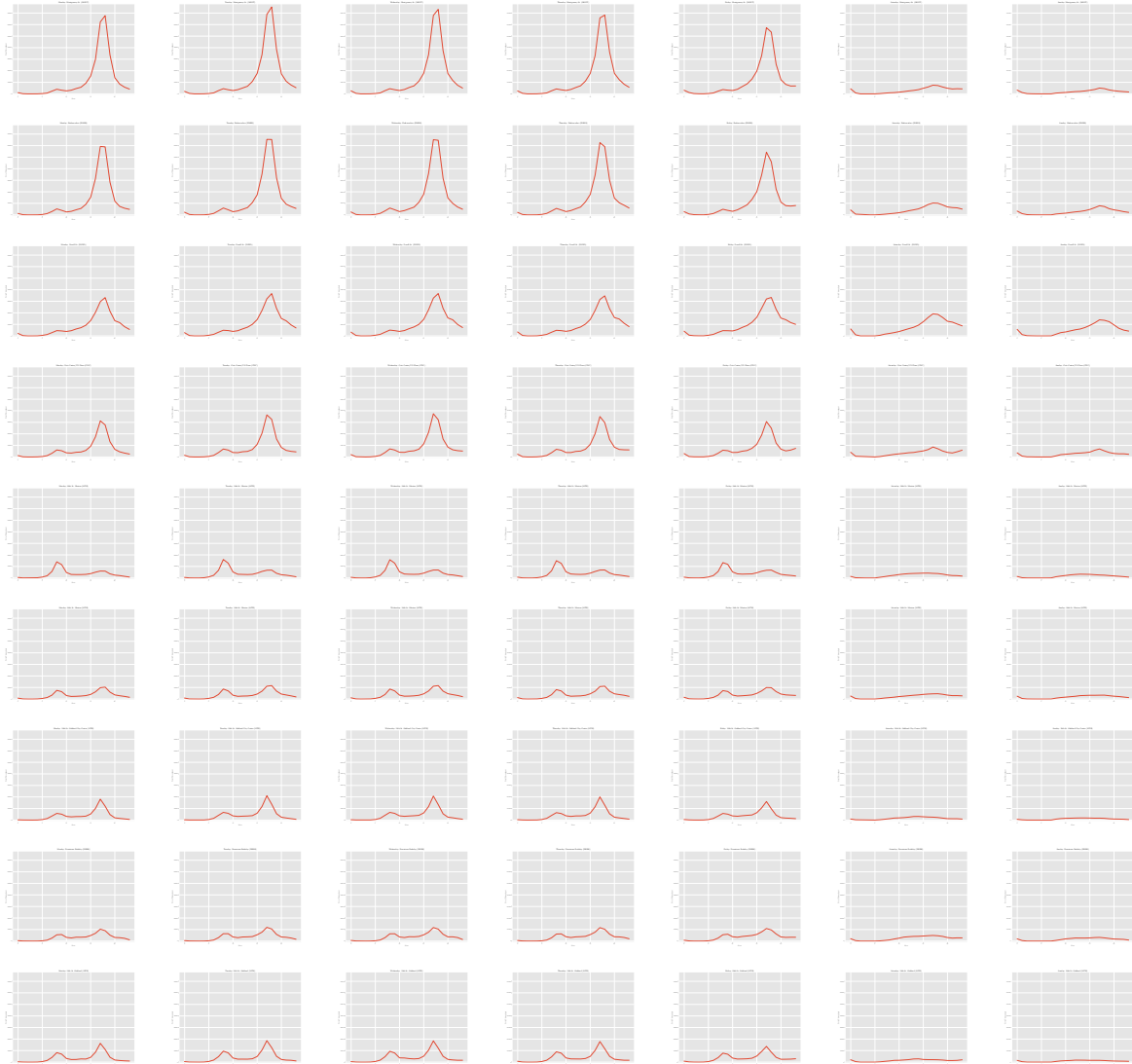


Figure 5: Total hourly throughput of a week of stations MONT, EMBR, POWL, CIVC, 24TH, 16TH, 12TH, DBRK, 19TH

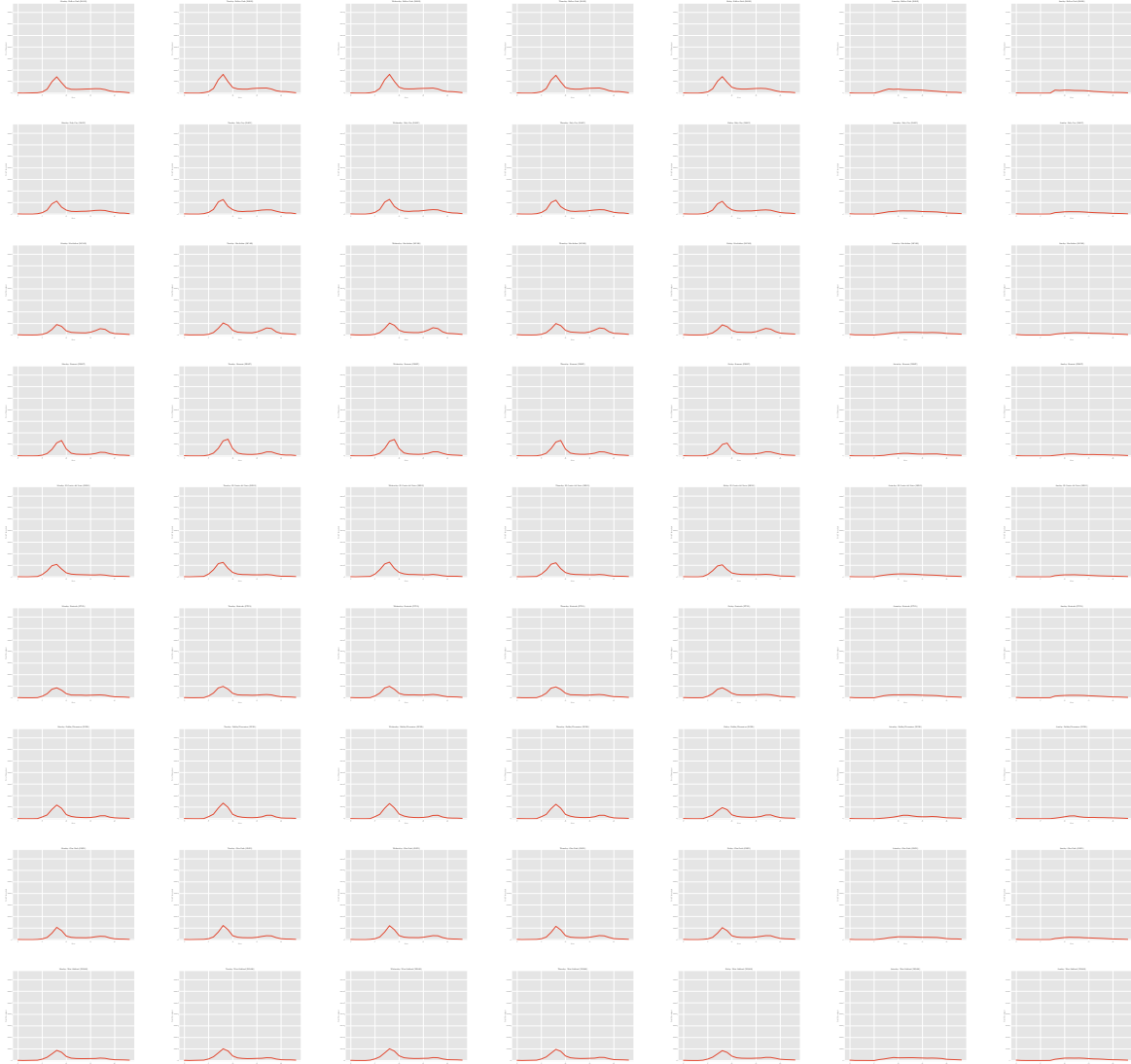


Figure 6: Total hourly throughput of a week of stations BALB, DALY, MCAR, FRMT, DELN, FTVL, DUBL, GLEN, WOAK

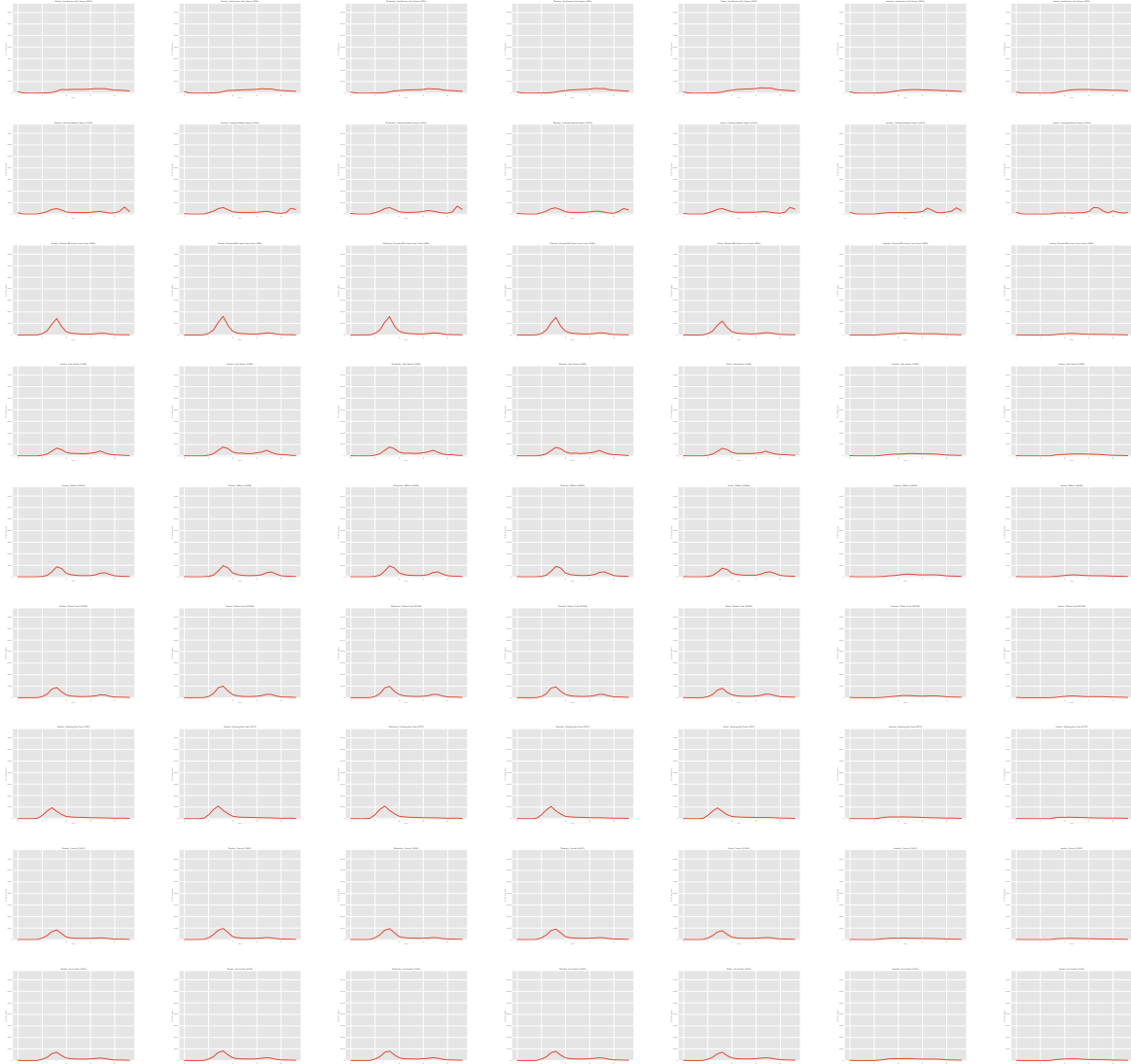


Figure 7: Total hourly throughput of a week of stations SFIA, COLS, PHIL, LAKE, MLBR, WCRK, PITT, CONC, SANL

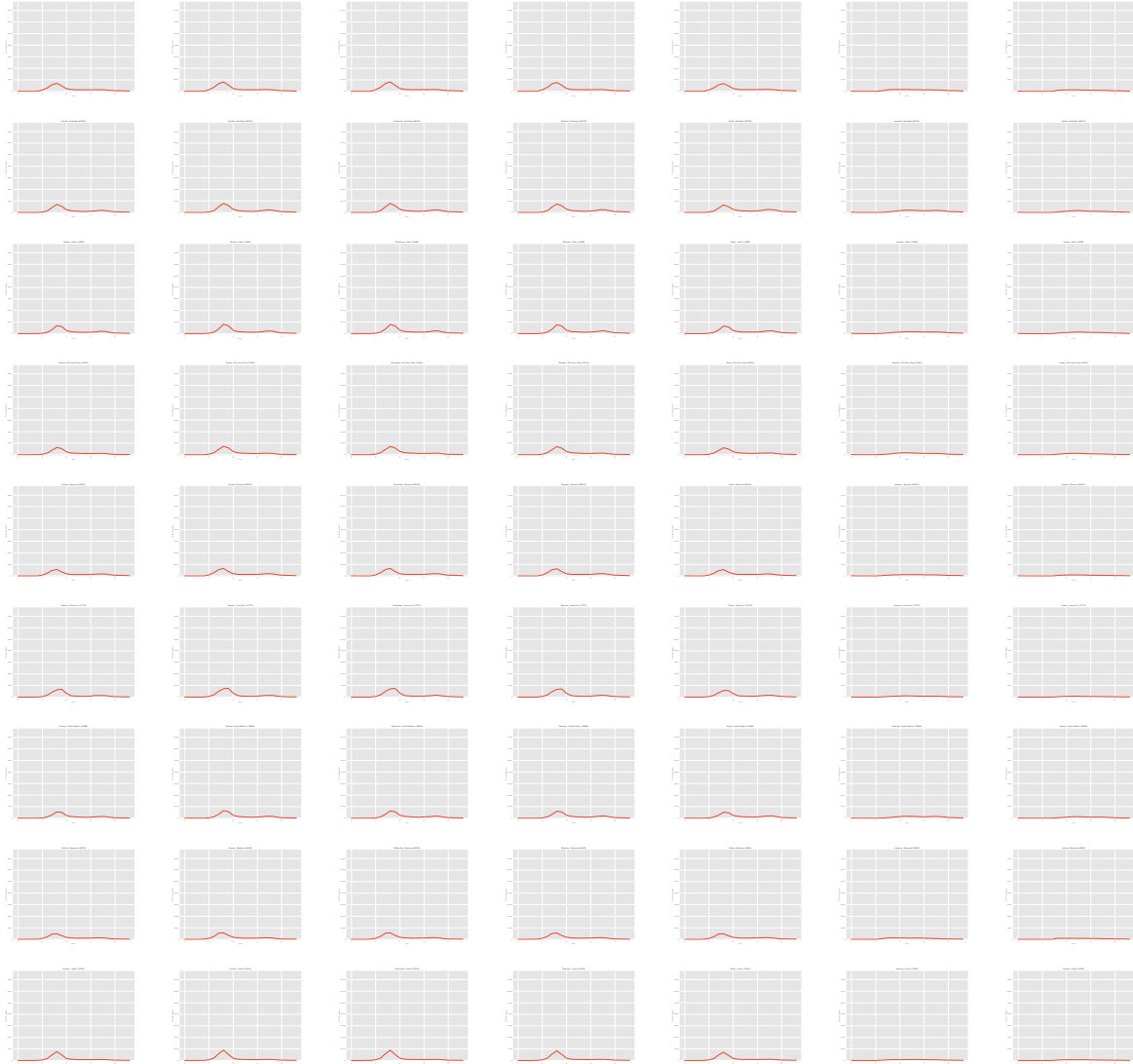


Figure 8: Total hourly throughput of a week of stations BAYF, ROCK, ASHB, PLZA, HAYW, UCTY, NBRK, RICH, COLM

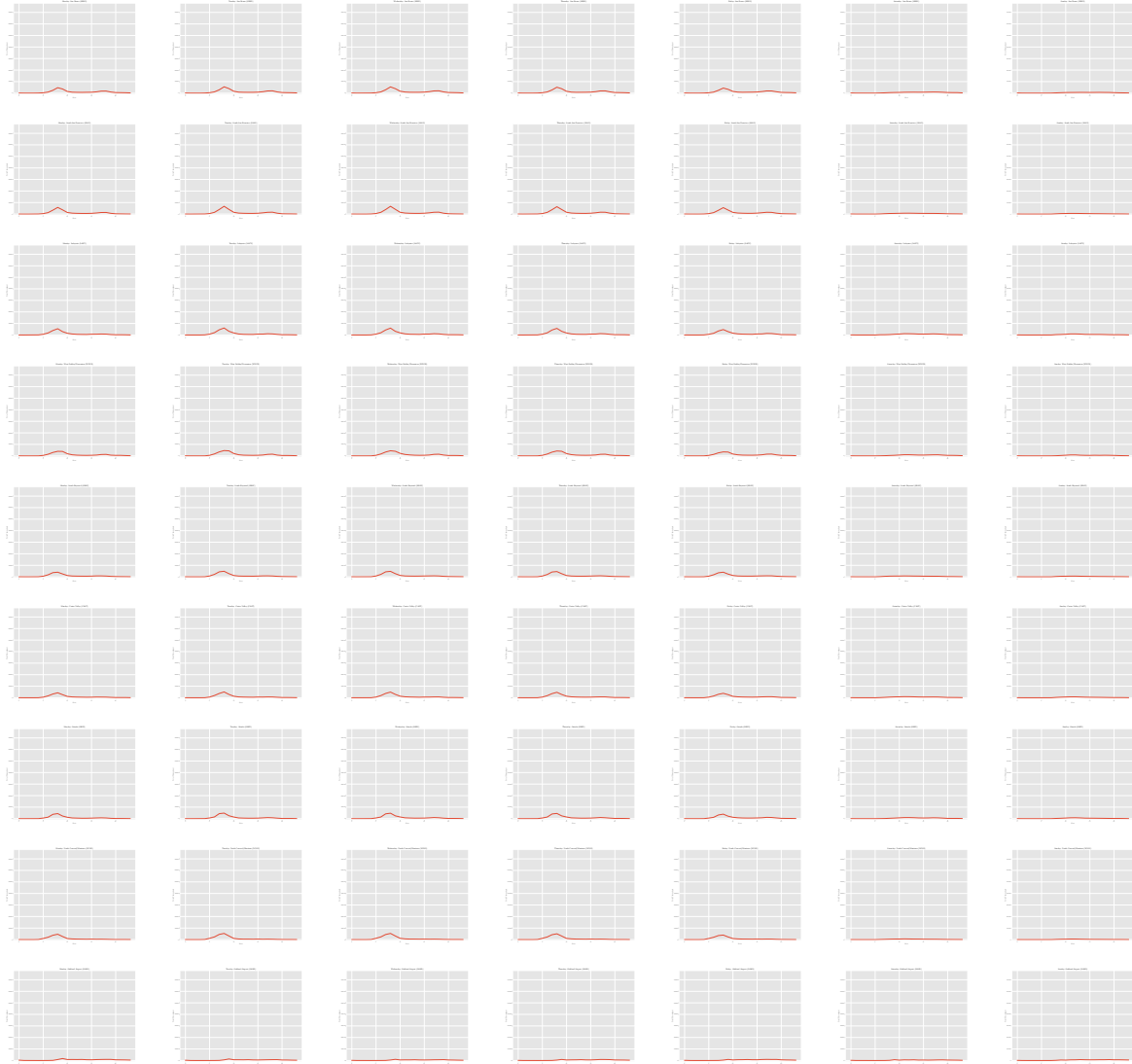


Figure 9: Total hourly throughput of a week of stations SBRN, SSAN, LAFY, WDUB, SHAY, CAST, ORIN, NCON, OAKL

5 Conclusions & Future Work

5.1 Conclusions

Now we are able to answer some of the questions 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

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3
4 # import math
5 import pandas as pd
6 import numpy as np
7 import matplotlib.pyplot as plt
8 import matplotlib as mpl
9 from matplotlib.backends.backend_pdf import PdfPages
10 # from IPython.core.interactiveshell import InteractiveShell
11 # InteractiveShell.ast_node_interactivity = 'all'
12
13 plt.rc('text', usetex=True)
14 plt.rc('font', family='serif')
15
```

```
16 __author__ = 'Libao Jin'
17 __date__ = '04/24/2018'
18 __email__ = 'ljin1@uwyo.edu'
19 __copyright__ = 'Copyright (c) 2018 Libao Jin'
20
21
22 def bart_preprocess(bart):
23     # Convert string to datetime, split datetime into date and time, and get day of week
24     bart['DateTime'] = pd.to_datetime(bart.DateTime)
25     bart['Date'] = bart['DateTime'].dt.date
26     bart['Time'] = bart['DateTime'].dt.time
27     bart['DayOfWeek'] = bart['DateTime'].dt.weekday_name
28     # bart.drop(columns='DateTime')
29     return bart
30
31
32 def stat_preprocess(stat):
33     # Split location into longitude and latitude for visualization later on
34     loc = stat.Location.str.split(',', expand=True)
35     loc = [pd.to_numeric(loc[i]) for i in loc.columns]
36     stat['Longitude'], stat['Latitude'] = loc[0], loc[1]
37     columns = ['Abbreviation', 'Name', 'Longitude', 'Latitude', 'Description']
38     stat = stat[columns]
39     return stat
40
41
42 def generate_bart_routes(bart_lines, stat):
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)
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 {}: {}'.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,
161             alpha=0.5,
162             label=label,
163             zorder=1
164         )
165     plt.legend()
166
167     plt.scatter(
168         x,
169         y,
170         c=cs,
171         s=cs,
172         marker='o',
173         edgecolors='k',
174         cmap=cmap,
175         alpha=1,
176         zorder=2
177     )
178
179     plt.colorbar()
180     plt.xlim([x_min, x_max])
181     plt.ylim([y_min, y_max])
182     plt.xlabel('Longitude')
183     plt.ylabel('Latitude')
184     plt.axis('equal')
185     plt.title('BART Lines ' + cs.name)
186     plt.show(block=False)
187     pdf.savefig(fig)
188
189     fig = plt.figure(figsize=(16, 12))
190     plt.style.use('ggplot')
191     plt.rc('text', usetex=True)
192     plt.rc('font', family='serif')
```



```

193     bart_aggregate.sort_values(column_name, inplace=True, ascending=False)
194     cs = getattr(bart_aggregate, column_name)
195     colors = cmap(np.linspace(1, 0, n))
196     cs.plot(kind='bar', color=colors, alpha=0.75)
197     # print(list(cs.index))
198     plt.ylabel('Average Throughput / hour')
199     plt.title('Throughput vs. Station')
200     plt.show(block=False)
201     pdf.savefig(fig)
202
203
204 def bart_aggregate_throughput(bart, filename):
205     number_of_days = len(bart['DateTime'].dt.date.unique())
206     number_of_hours = len(bart['DateTime'].dt.time.unique())
207     bart_grouped = bart['Throughput'].groupby(bart['Origin']).sum().to_frame()
208     bart_grouped['Destination'] = bart['Throughput'].groupby(bart['Destination']).sum()
209     bart_grouped.index.names = ['Station']
210     bart_grouped.columns = ['Throughput Origin', 'Throughput Destination']
211     bart_grouped['Throughput All'] = bart_grouped['Throughput Origin'] + bart_grouped['Throughput Destination']
212     bart_grouped.set_index(stat['Abbreviation'])
213     bart_grouped.reset_index(level=0, inplace=True)
214     bart_grouped['Longitude'] = stat['Longitude']
215     bart_grouped['Latitude'] = stat['Latitude']
216     bart_grouped[['Throughput Origin', 'Throughput Destination', 'Throughput All']] = \
217         bart_grouped[['Throughput Origin', 'Throughput Destination', 'Throughput All']] / (number_of_days * number_of_hours)
218     visualize_throughput(bart_routes, bart_grouped, 'Throughput Origin', filename)
219     # visualize_throughput(bart_routes, bart_grouped, 'Throughput Destination', filename)
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
263         plt.title(title)
264         plt.xlabel(group_by.title())
265         ax = plt.gca()
266         # fig = plt.gcf()
267         if group_by == 'date':
268             xfmt = mpl.dates.DateFormatter('%Y-%m-%d')
269             ax.xaxis.set_major_formatter(xfmt)
270             # plt.xticks(rotation=90)
271             fig.autofmt_xdate()
272         plt.ylabel('Total Throughput')
273         plt.ylim([-10000, 780000])
274         # plt.xticks(grouped.index, list(np.arange(24)))
275         plt.show(block=False)
276         pdf.savefig(fig)
277
278     # Plot the throughput with respect to date
279     elif plot_option == 2:
280         grouped = data.groupby(getattr(data, 'DateTime').dt, group_by).sum()
281         grouped.sort_index(inplace=True)
282         fig = plt.figure(figsize=(8, 6))
283         plt.plot(grouped['Throughput'])
284         if len(stops) == 1:
285             tmp = stat[stat['Abbreviation'] == data.iloc[0]['Origin']]
286             title = '{}'.format(tmp.iloc[0]['Name'])
287             # title = class_type[1][i] + ': {}'.format(data.iloc[0]['Origin'])
288         else:
289             tmp_1 = stat[stat['Abbreviation'] == data.iloc[0]['Origin']]
290             tmp_2 = stat[stat['Abbreviation'] == data.iloc[0]['Destination']]
291             title = '{} - {}'.format(tmp_1.iloc[0]['Name'],
292                                     tmp_2.iloc[0]['Name'])
293             # title = ': {} - {}'.format(data.iloc[0]['Origin'], data.iloc[0]['Destination'])
294         # plt.legend()
295         plt.xlabel(group_by.title())
296         plt.ylabel('Throughput')
297         plt.xticks(rotation=90)
298         ax = plt.gca()
299         # fig = plt.gcf()
300         xfmt = mpl.dates.DateFormatter('%Y-%m-%d')
301         ax.xaxis.set_major_formatter(xfmt)
302         # plt.xticks(rotation=90)
303         fig.autofmt_xdate()
304         plt.title(title)
305         plt.ylim([0, 65000])
306         plt.show(block=False)
307         pdf.savefig(fig)
308
309 if __name__ == '__main__':
310

```

```
311 dest_folder = './output'
312 plt.rc('text', usetex=True)
313 plt.rc('font', family='serif')
314 # Load data
315 # Data obtained from https://www.kaggle.com/saulfuh/bart-ridership
316 date_hour_2016 = '../data/bart-ridership/date-hour-soo-dest-2016.csv'
317 date_hour_2017 = '../data/bart-ridership/date-hour-soo-dest-2017.csv'
318 stat_info = '../data/bart-ridership/station_info.csv'
319 # bart_16 = pd.read_csv(date_hour_2016)
320 # bart_17 = pd.read_csv(date_hour_2017)
321 # stat = pd.read_csv(stat_info)
322
323 # Data preprocessing
324 bart_16 = bart_preprocess(pd.read_csv(date_hour_2016))
325 bart_17 = bart_preprocess(pd.read_csv(date_hour_2017))
326 bart = pd.concat([bart_16, bart_17], ignore_index=True)
327 stat = stat_preprocess(pd.read_csv(stat_info))
328
329 # Visualize the routes according to the BART official website
330
331 line_1 = [
332     'Richmond',
333     'El Cerrito del Norte',
334     'El Cerrito Plaza',
335     'North Berkeley',
336     'Downtown Berkeley',
337     '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'
```

```
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', 'February', '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))

```