

Visualization, Time-series Analysis and Community Detection on CitiBike stations

Yun Yan (yy1533@nyu.edu)

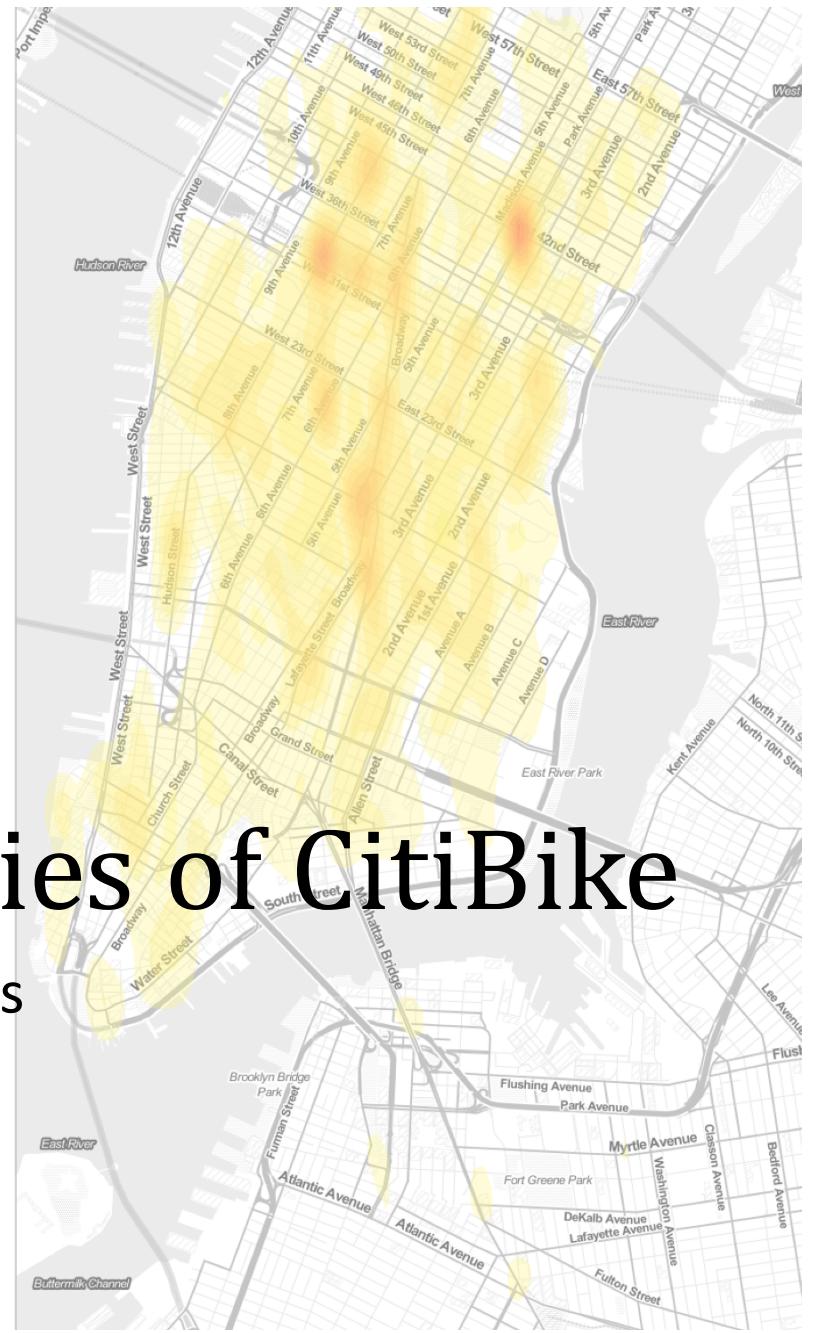
Willian Zhang (willian@nyu.edu)

Outline

- Introduction to Citibike dataset and possible tasks
- What does the Citibike activity looks like within 24-hour?
- Is there any specific purpose of using Citibike depending on stations?
- Can we model / forecast the daily usage?
- Is there any community of stations?

Visualizing Activities of CitiBike

Picking during weekdays of all users



Trips of CitiBike

1. Activity: Picking / Docking
2. Day: all days / Weekdays / Weekends
3. Type: all users / Member / One-time Customer

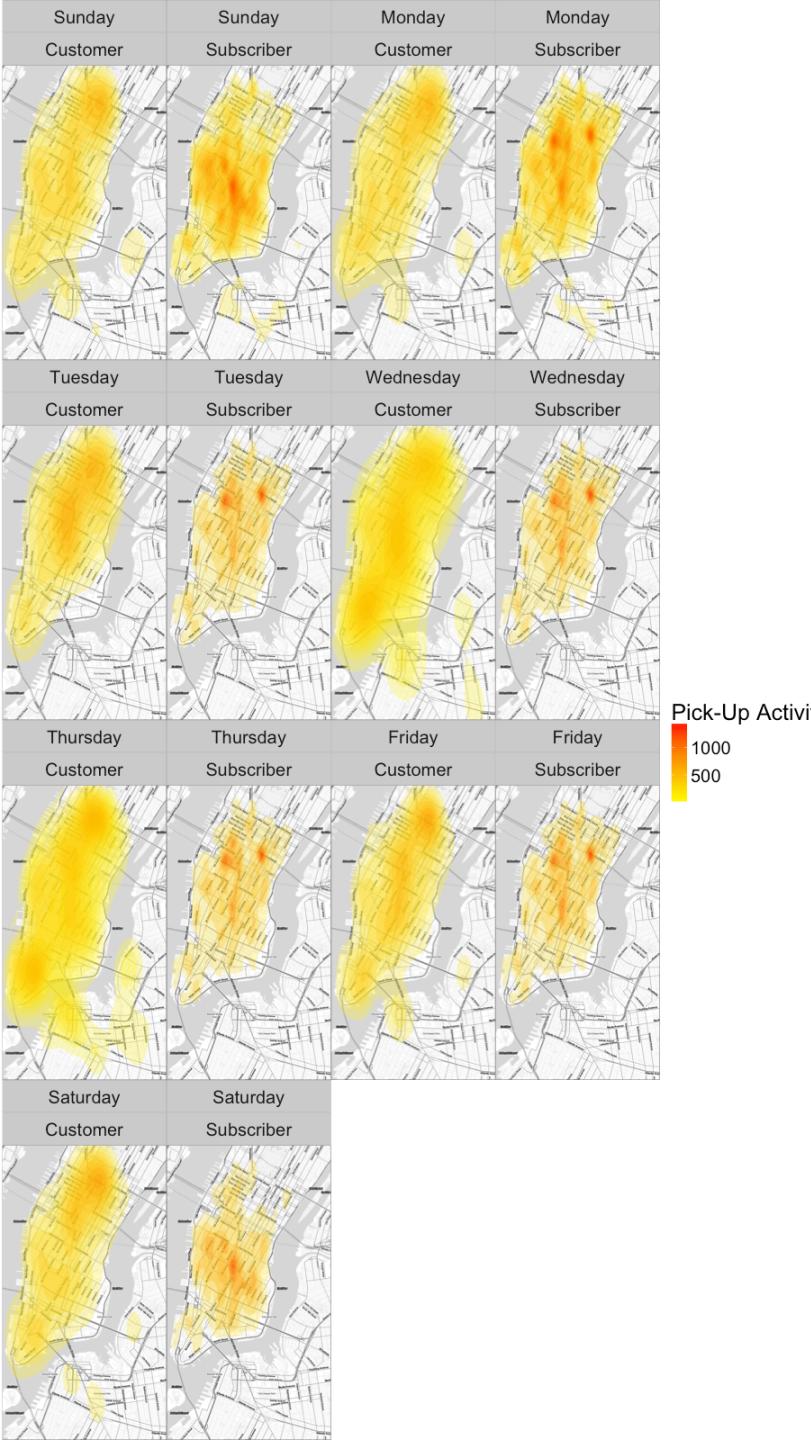
All Activities of CitiBike

- Picking during all days
- One-time Customer v.s. Members



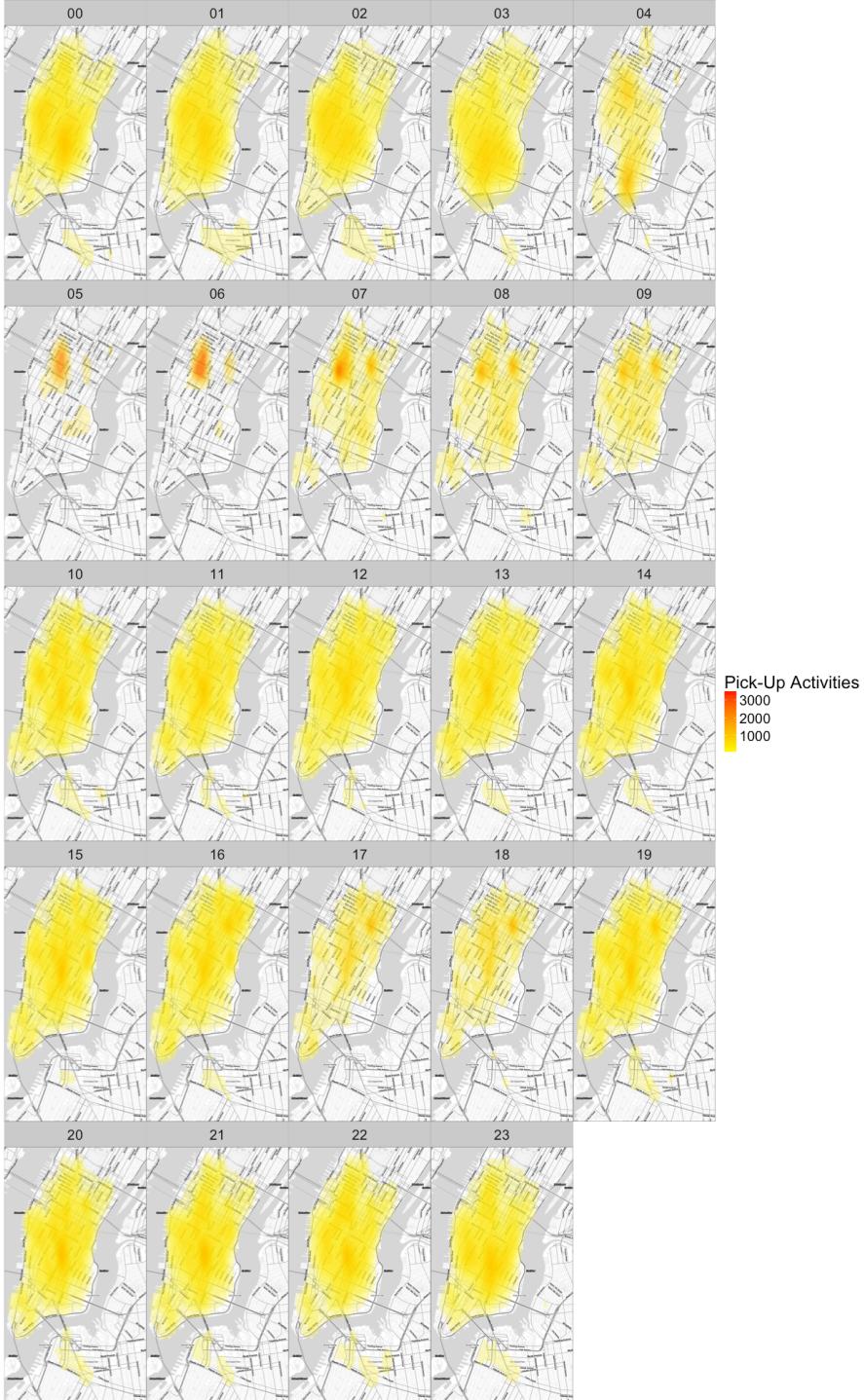
Weekly Activities

- Picking up during weekdays and weekends
- One-time customer v.s. Members



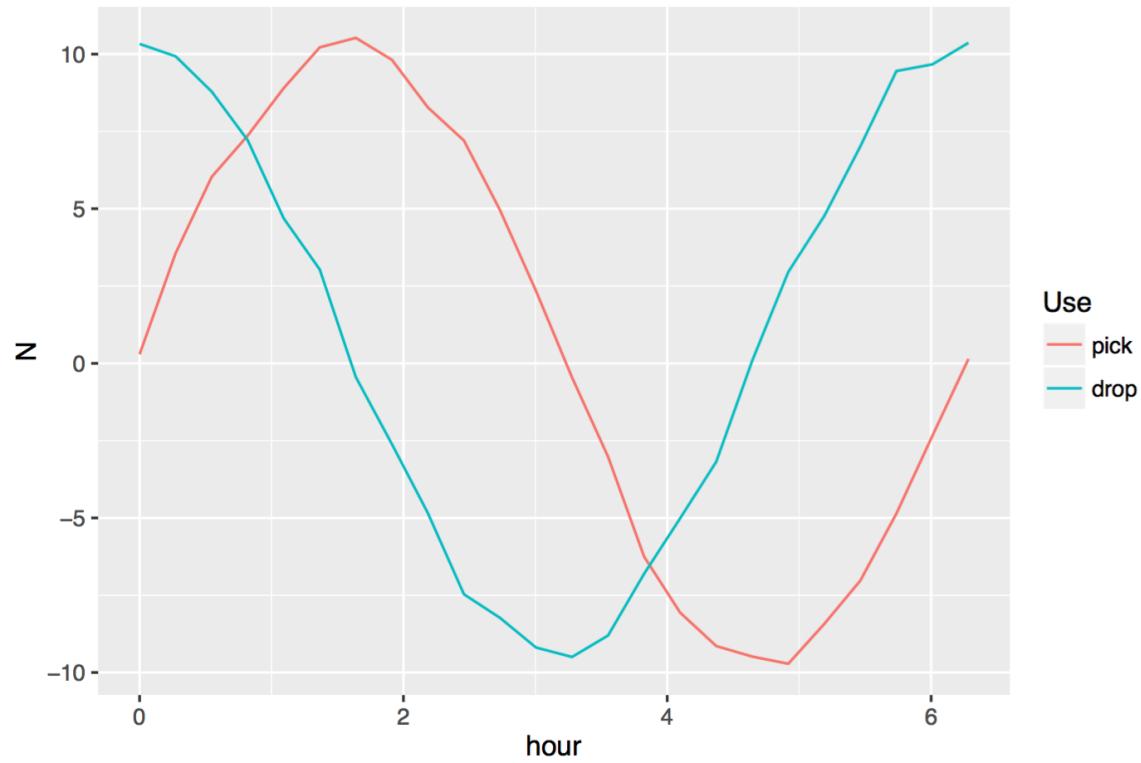
24-hour Activities

- Picking up during weekdays
- All users



Section Summary

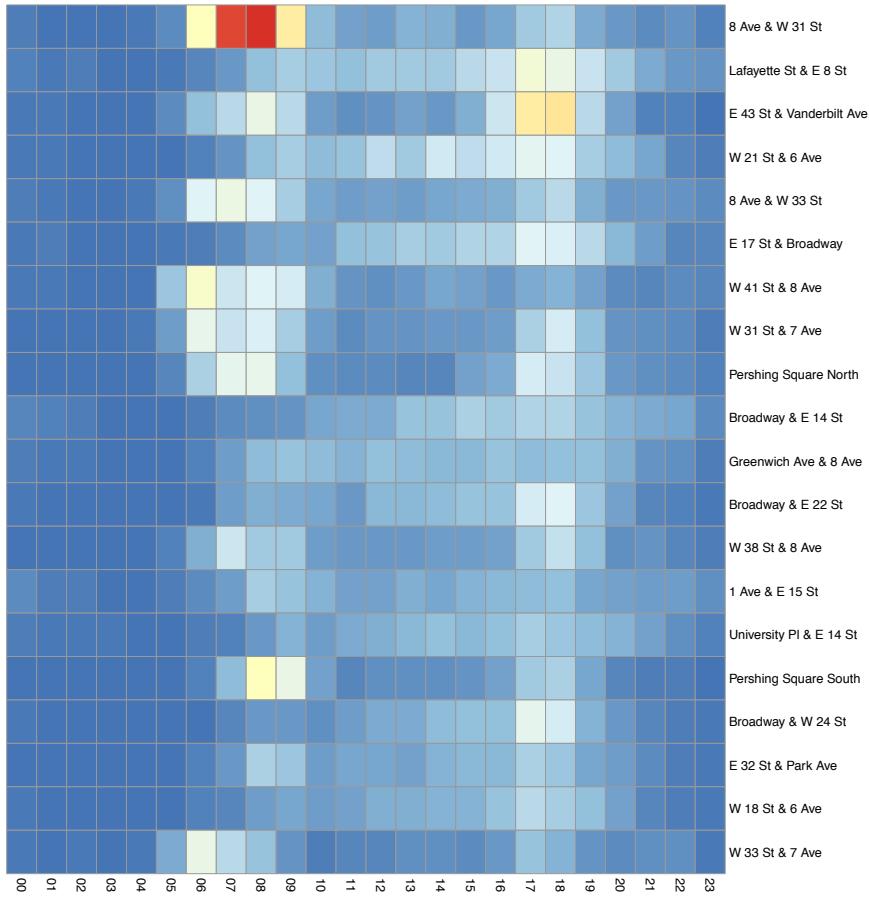
- Monitor the 24-hour activities in New York City
- Different behaviors between annual members and one-time customers
 - Member consider CitiBike for commuting
 - One-time customer ride for fun



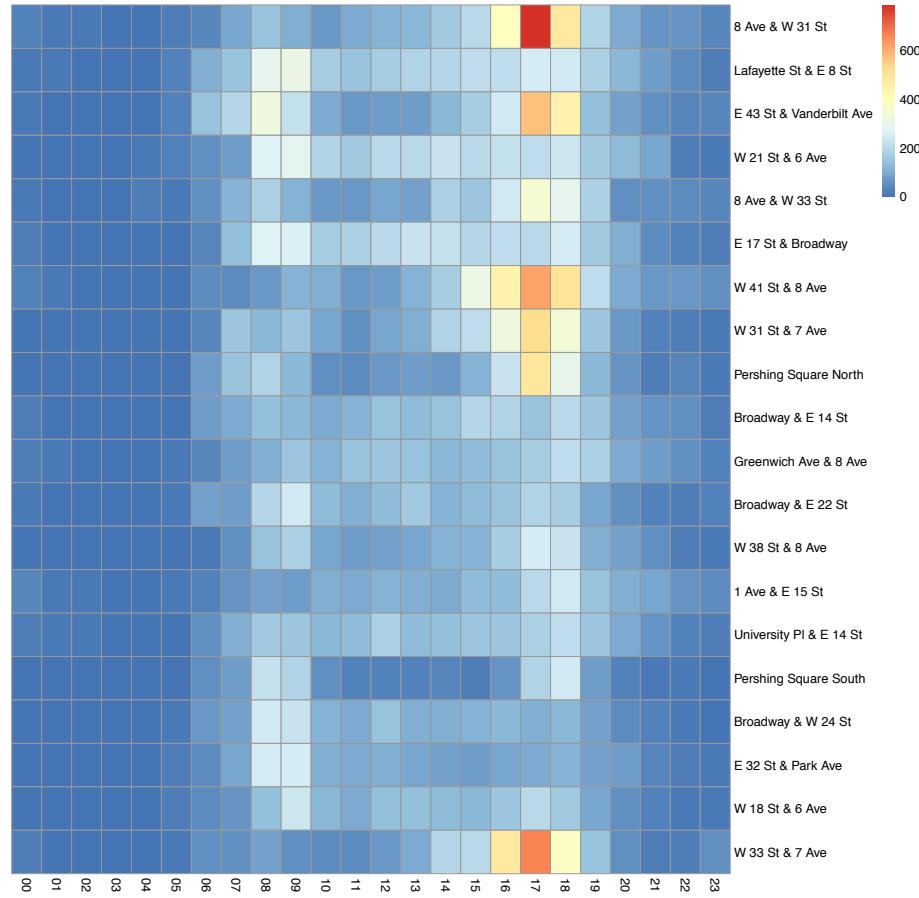
Temporal properties

Do uses have specific reasons to take bikes depending on stations?

Picking activities of Stations with Top 24-hr Picking Activities



Docking activities of Stations with Top 24-hr Picking Activities



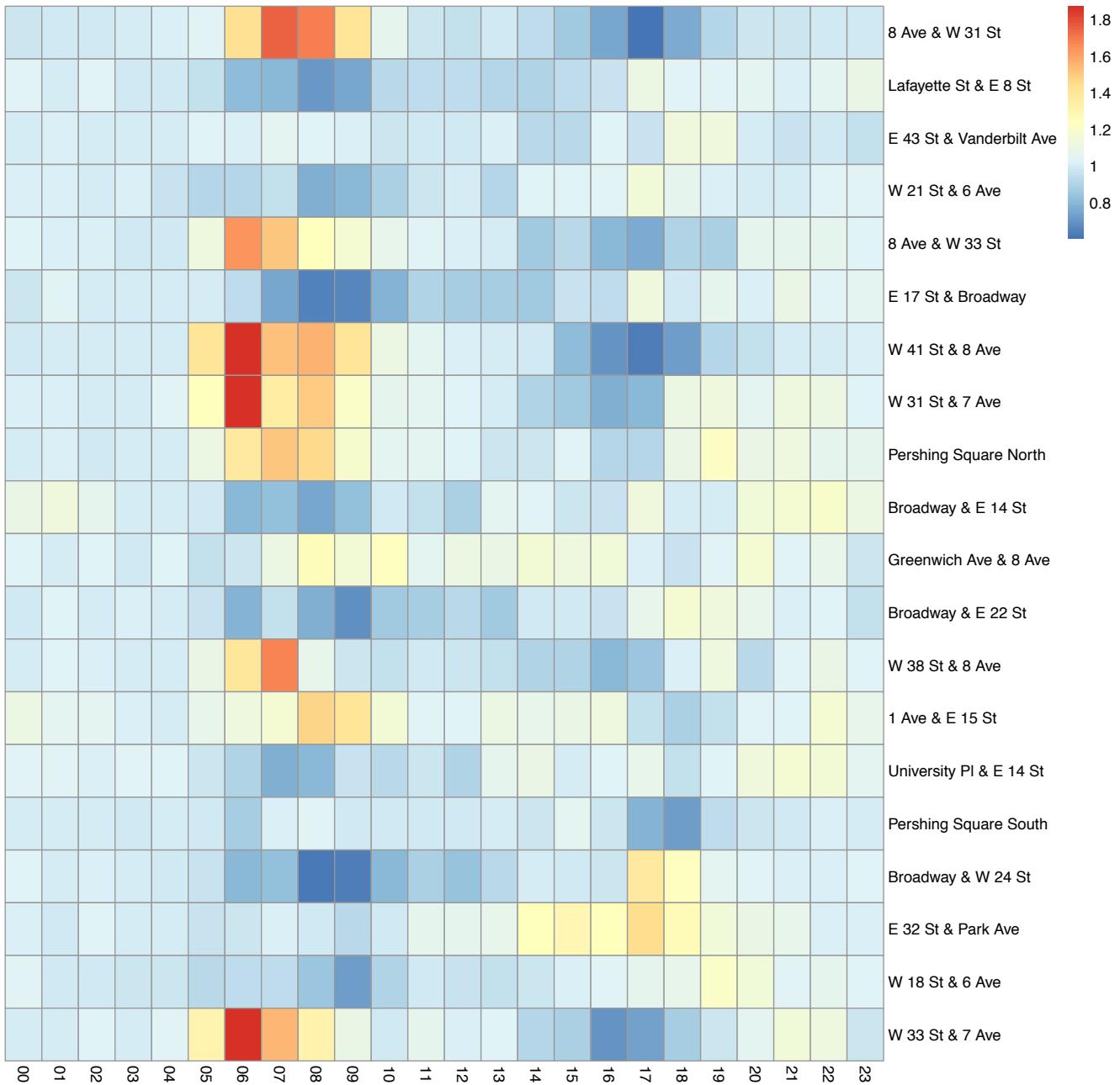
Hidden preference towards stations

- Example of “8 Ave & W 31 St”: high picking in the morning while high docking in the nightfall
- How to integrate picking and docking activities together, instead of examining two figures side-by-side?
- P/D Index = Picking / Docking
 1. Matrix with stations on rows and hour point on columns;
 2. For each row, each element is divided by its max, e.g. [0, 1, 5] -> [0, 0.2, 1.0];
 3. Perform Step1 and Step2 on picking matrix and docking matrix;
 4. P/D index = $(P + 1) / (D + 1)$.

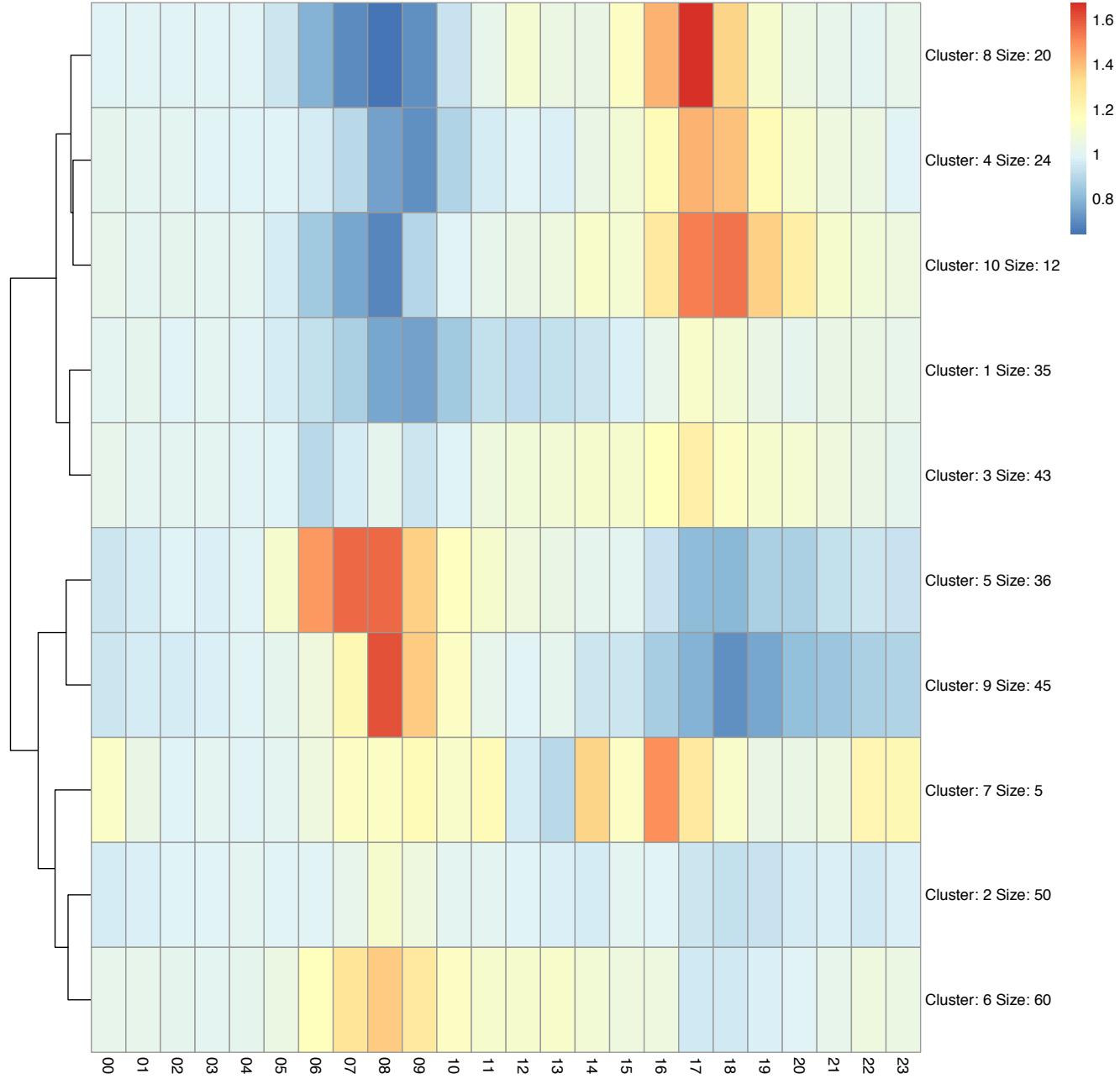
P/D Index

- Higher P/D Index, it suggests high picking or low docking, i.e. the specific station is start-point of biking trip, and vice-versa.
- Globally, the average should be (close to) 1.00.

P/D Index of Stations with Top 24-hr Picking Activities

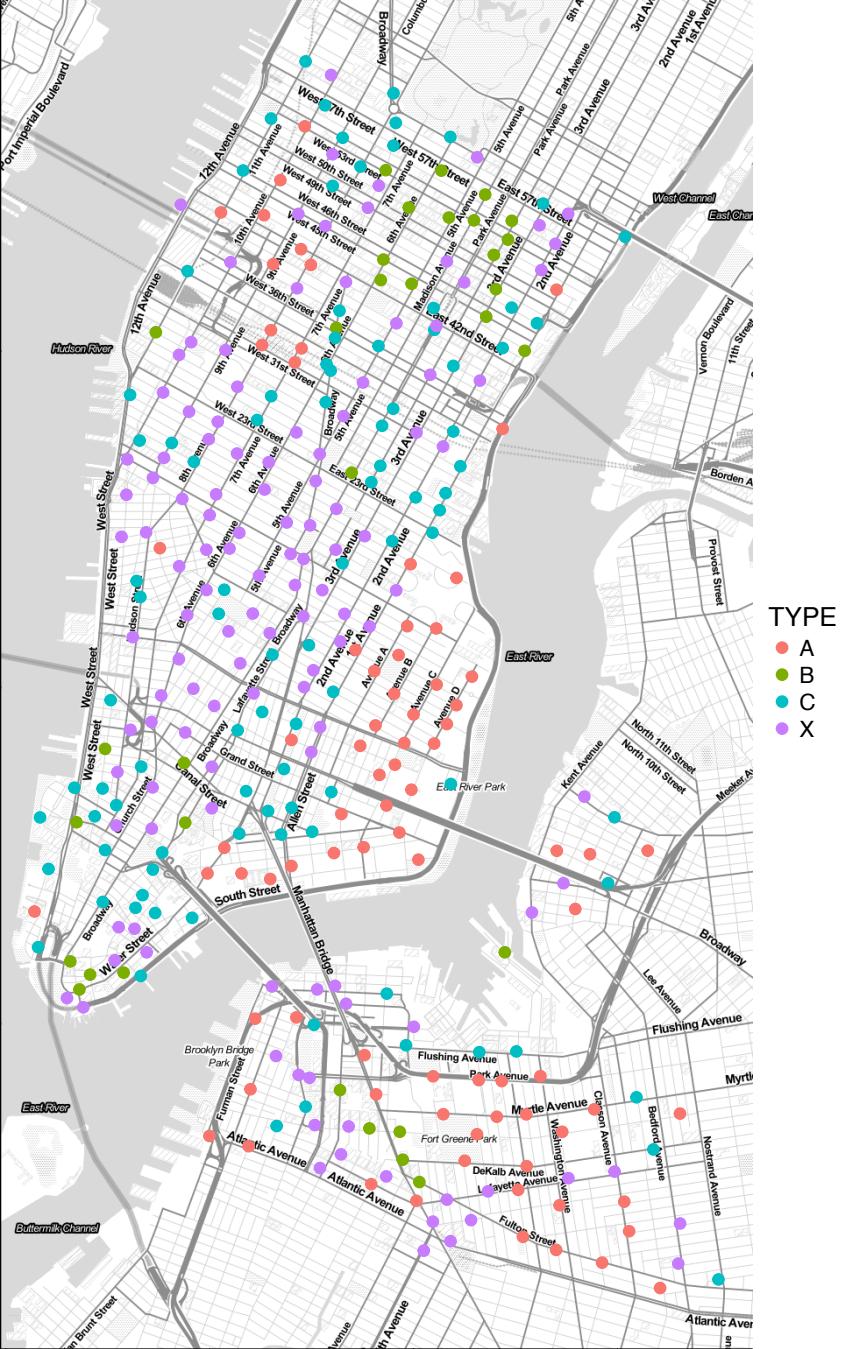


P/D Index of Stations



Type of stations

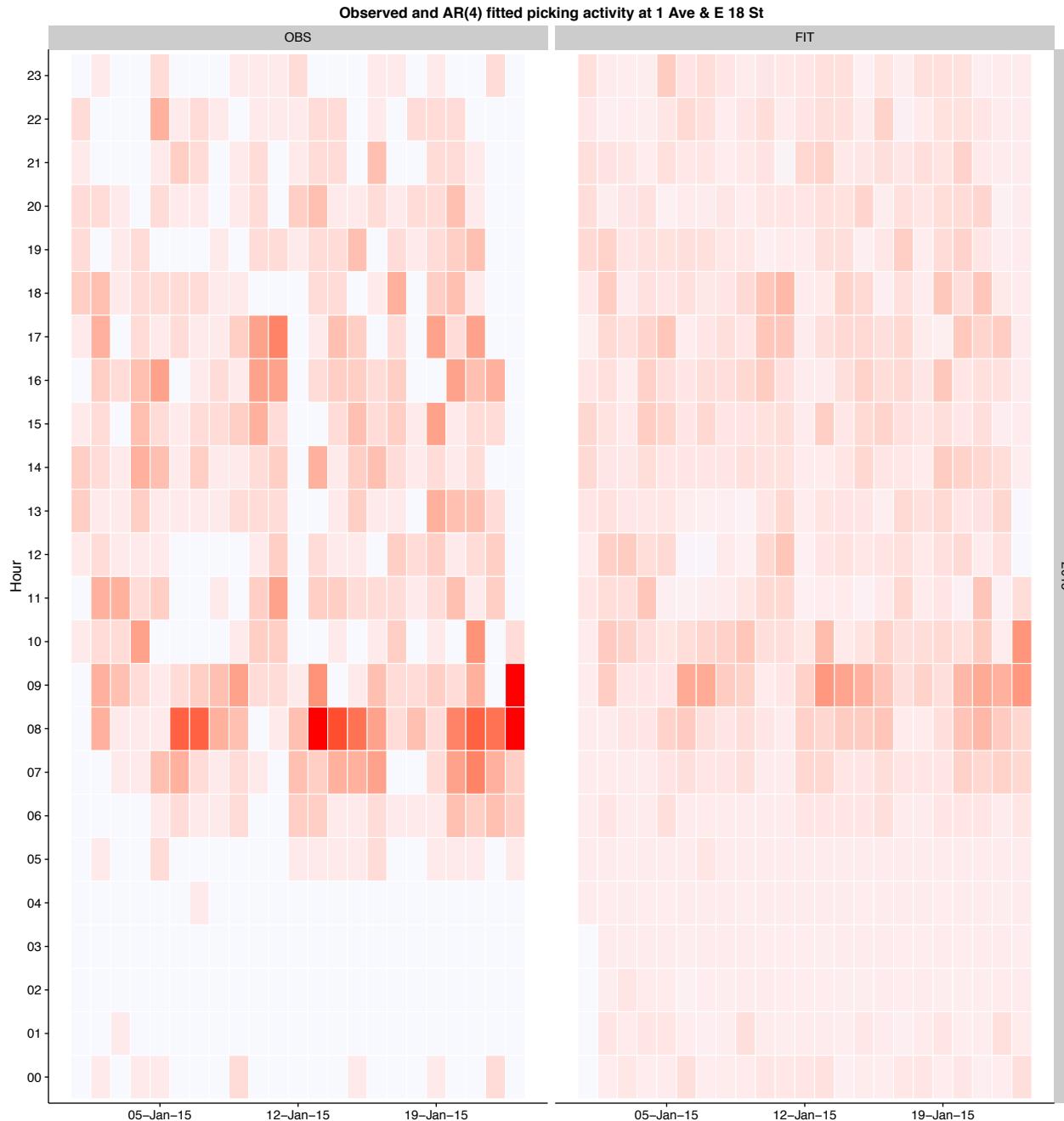
- Type A: high P/D index in the morning but low at nightfall
- Type B: low P/D index in the morning but high at nightfall
- Type C: Normal P/D index throughout the day
- Type X: the rest



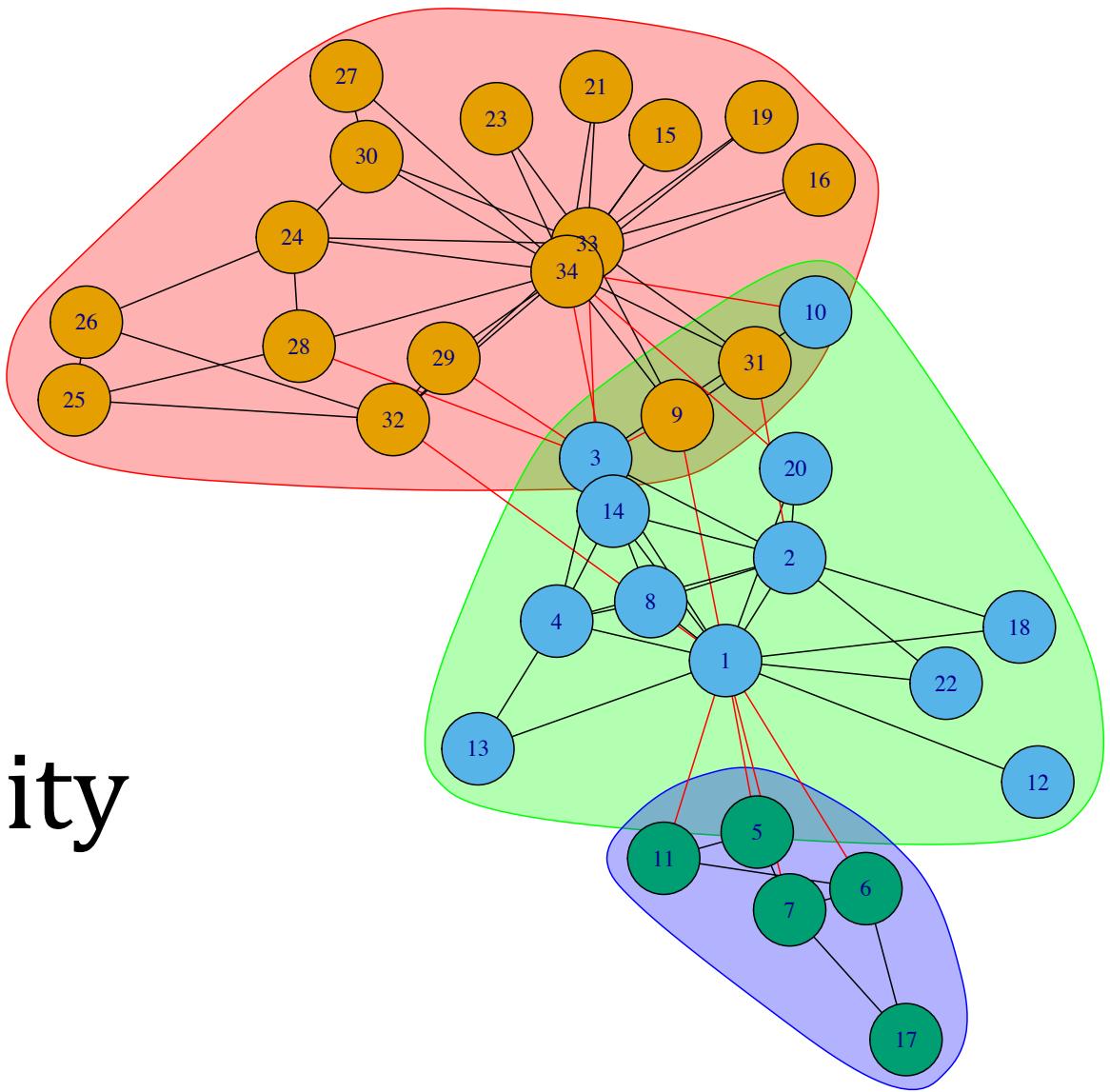
Section Summary

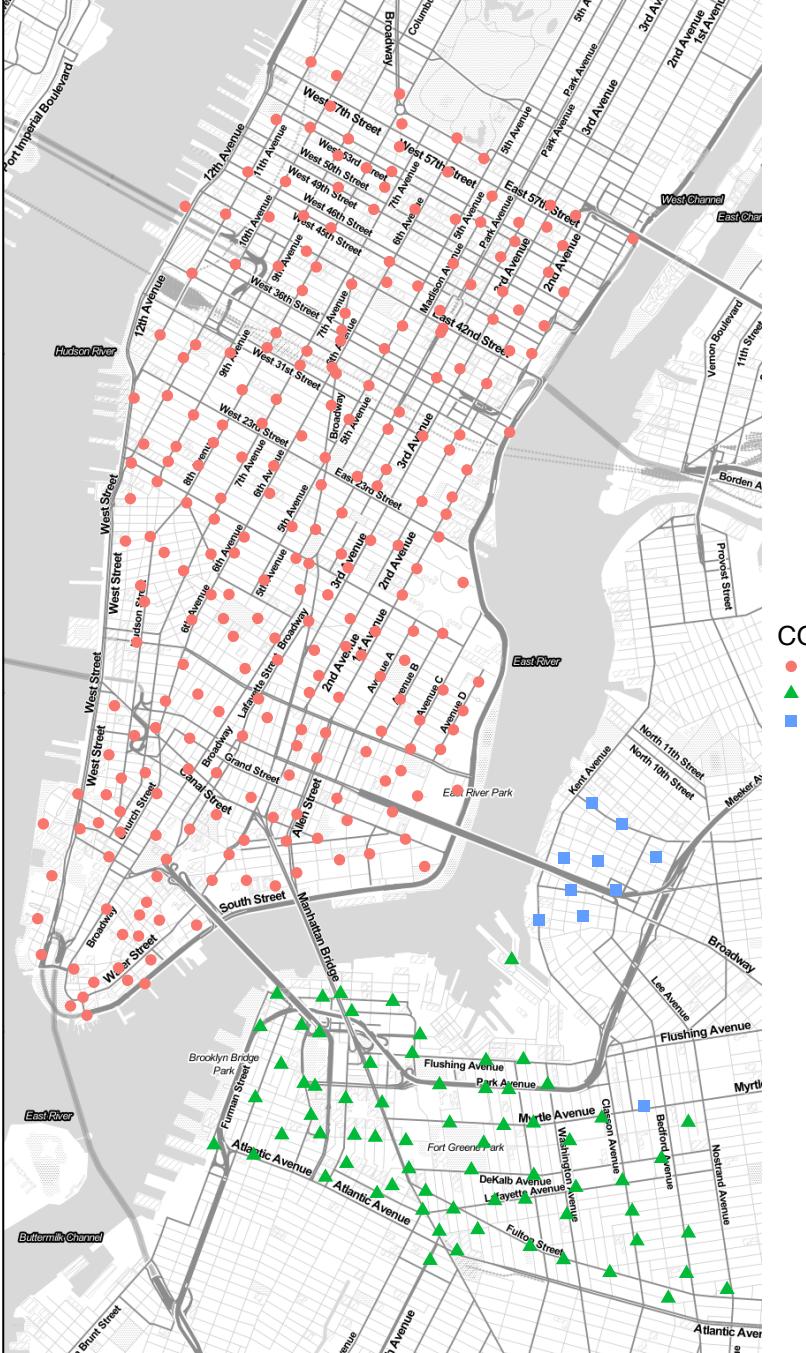
- There are hidden preferences of users towards stations
- P/D Index captured the preference differences



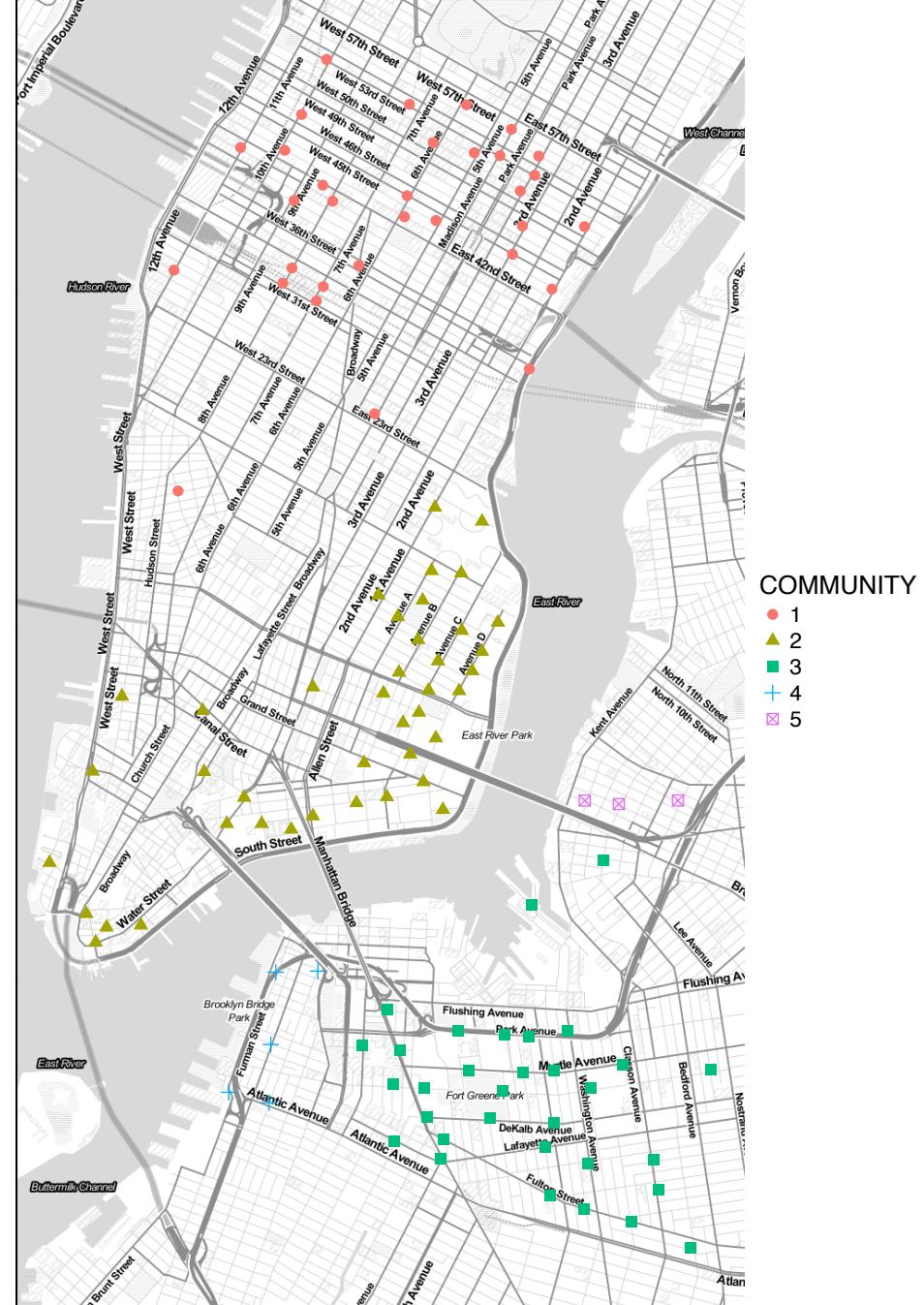
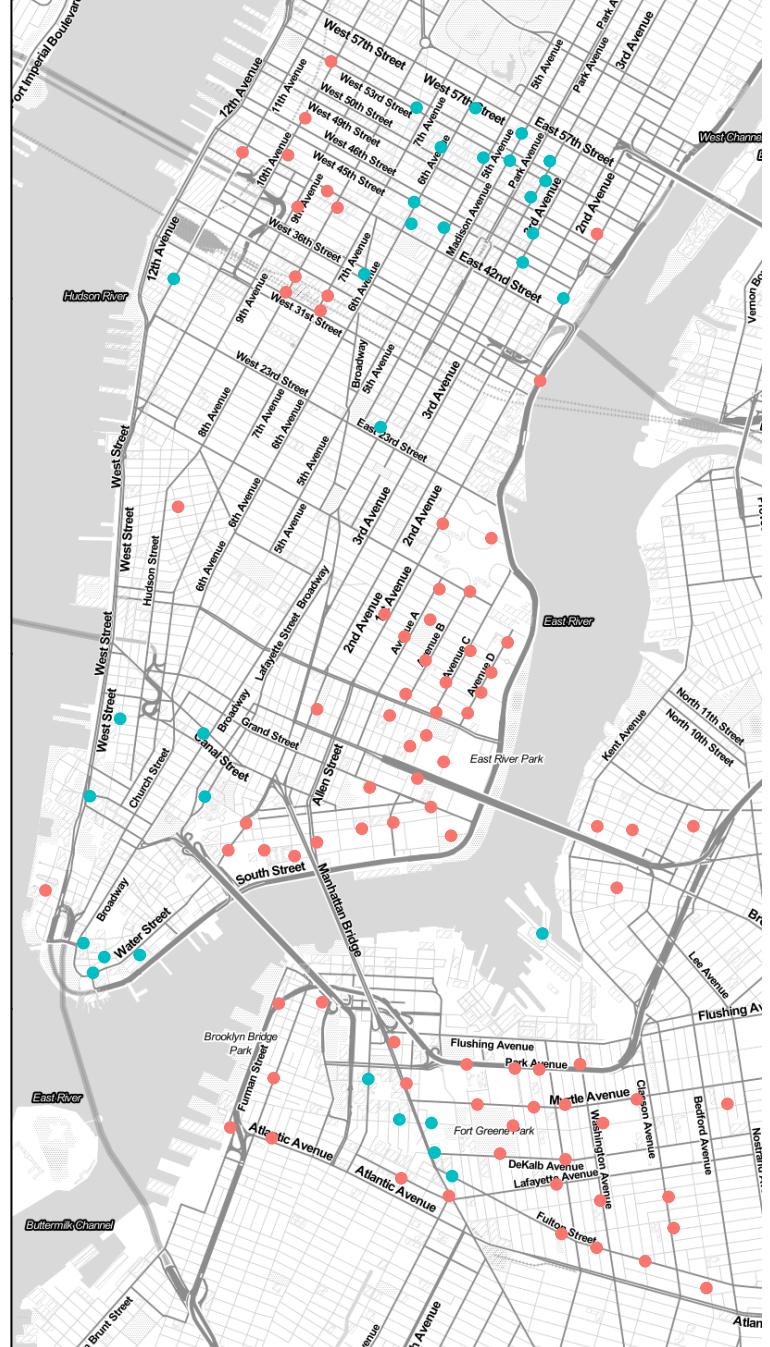


Community



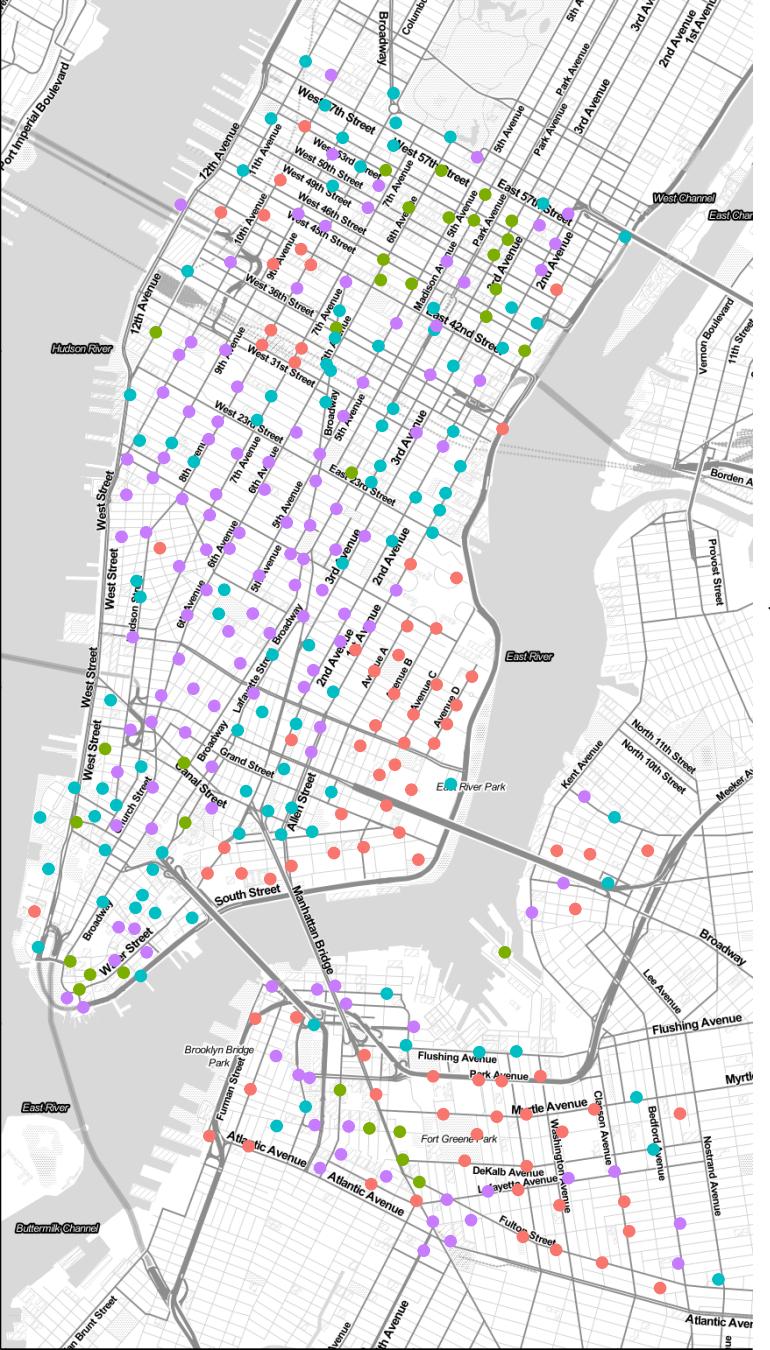


Visualization, Time-series Analysis and Community Detection on CitiBike stations



Visualization, Time-series Analysis and Community Detection on CitiBike stations





TYPE

- A
- B
- C
- X



COMMUNITY

- 1
- ▲ 2
- 3
- + 4



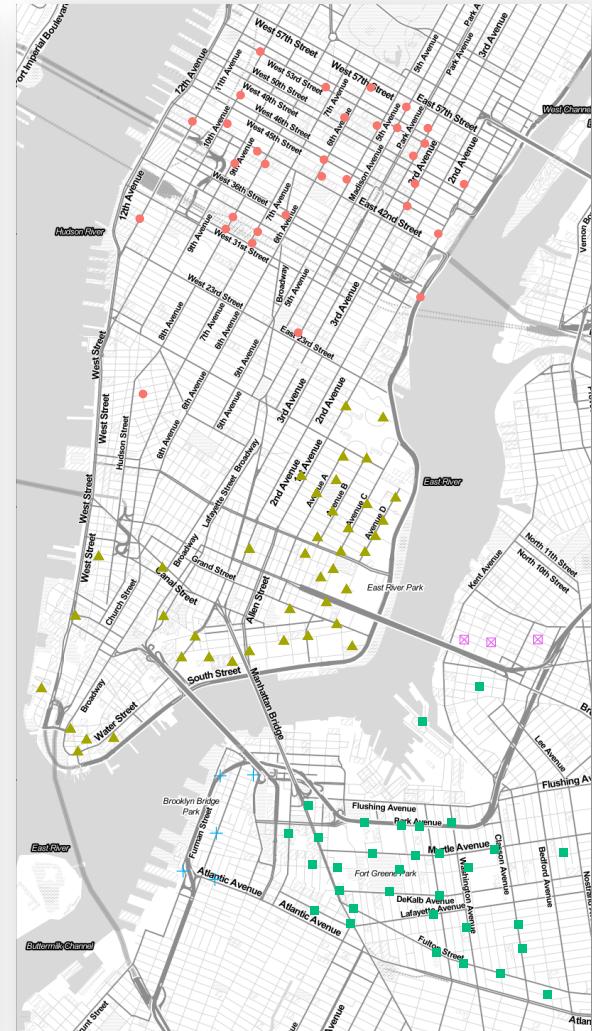
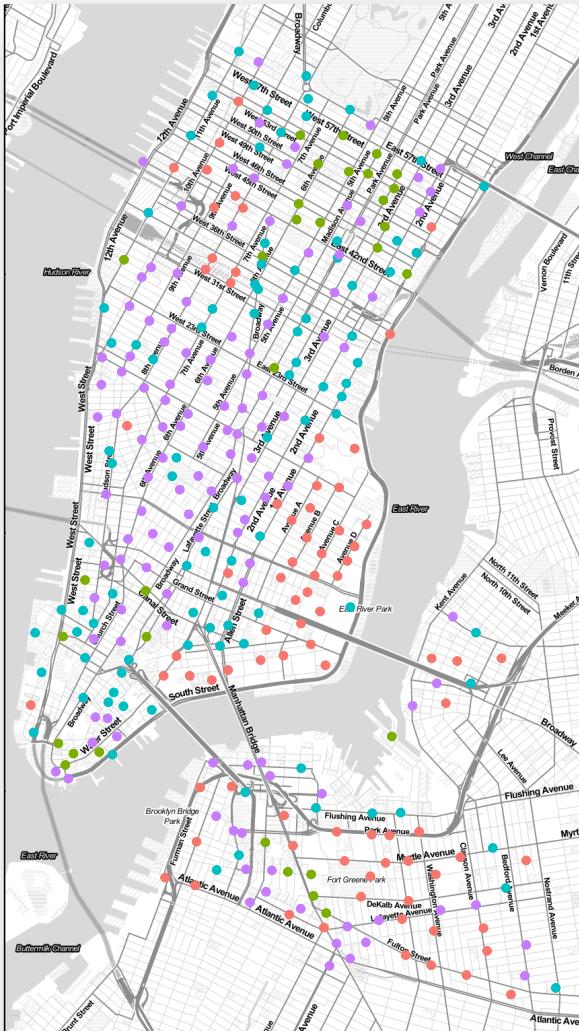
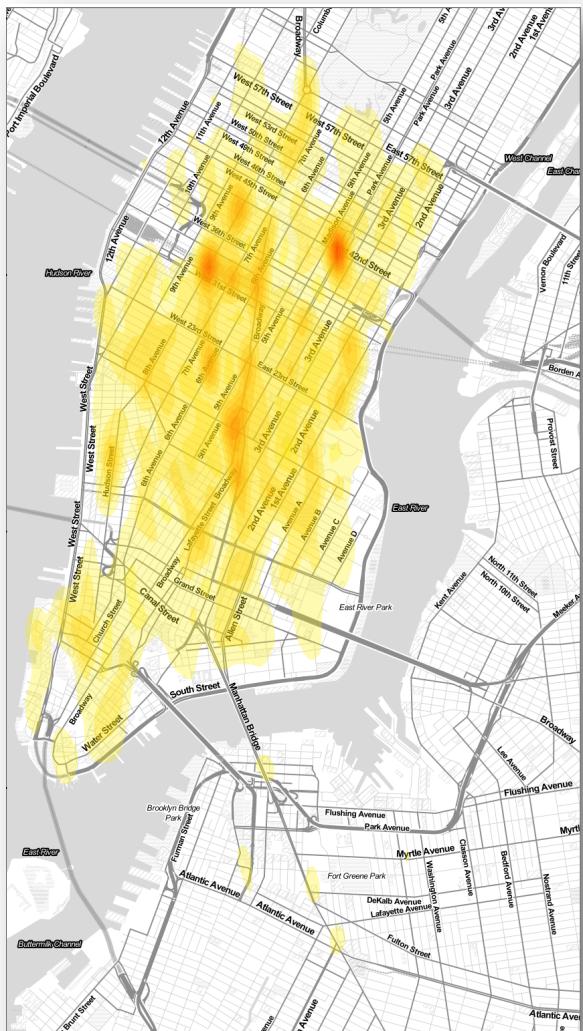
NYU

TANDON SCHOOL
OF ENGINEERING

Section Summary

- Distinct communities are identified by InfoMap algorithms of igraph package
- The general direction of flow of CitiBike

Summary



Contributions

- Y. Yan designed and finished the project;
- Y. Yan wrote the doc and slide.
- W. Zhang joined discussions.
- All are deployed at <https://github.com/Puriney/ineedabike>

Session-Info

attached base packages:

```
[1] stats      graphics   grDevices  utils  
datasets   methods    base
```

other attached packages:

```
[1] igraph_1.0.1       ggmap_2.6.1  
pheatmap_1.0.9      forecast_7.3  
timeDate_3012.100  
  
[6] zoo_1.7-13         cowplot_0.6.3  
scales_0.4.1        readr_1.0.0  
lubridate_1.6.0  
  
[11] ggplot2_2.1.0      reshape2_1.4.2  
dplyr_0.5.0
```

Questions

