

Station Usage in Seattle Bike Share

Arman Bilge^{1,2}, Yanbo Ge³, and Xiaoliu Wu²

¹Computational Biology Program, Fred Hutchinson Cancer Research Center and Departments of ²Statistics and ³Civil and Environmental Engineering, University of Washington



Background

- ▷ Pronto! has been Seattle's bike share system for 2 years
- ▷ Users must start and end their trips at various stations across the city
- ▷ Program is shutting down at the end of March due to insufficient use
- ▷ Many reasons suggested for its failure; e.g., poor placement of stations

Questions

- ▷ Which features of a station affect its daily usage?
- ▷ Can we use these to decide where to place stations?

Data

- ▷ Pronto! (October 2014–August 2016)
 - ▷ record for every trip taken (197810 trips in total)
 - ▷ daily weather report (temperature, precipitation, wind, etc.)
- ▷ Google Maps API
 - ▷ elevation
 - ▷ nearby points-of-interest (transportation, etc.)
- ▷ Puget Sound Regional Council Census
 - ▷ socio-economic data (population density, job density, income, etc.)

Challenges (and solutions)

- ▷ Stark contrast between weekday and weekend usage patterns
 - ▷ Only analyze data from weekdays
- ▷ Overall system usage depends heavily on daily weather
 - ▷ Control for temperature and precipitation in our model
- ▷ Some stations were decommissioned or installed at later dates
 - ▷ Remove stations that have been active for less than 20% of time frame
- ▷ System has complex, network dependencies between stations
 - ▷ Only count daily departures (and not arrivals) per station, called *outflow*
 - ▷ Try to find a subset of independent stations
- ▷ No data on total membership of the system
 - ▷ Conceded limitation of our analysis

Assumptions

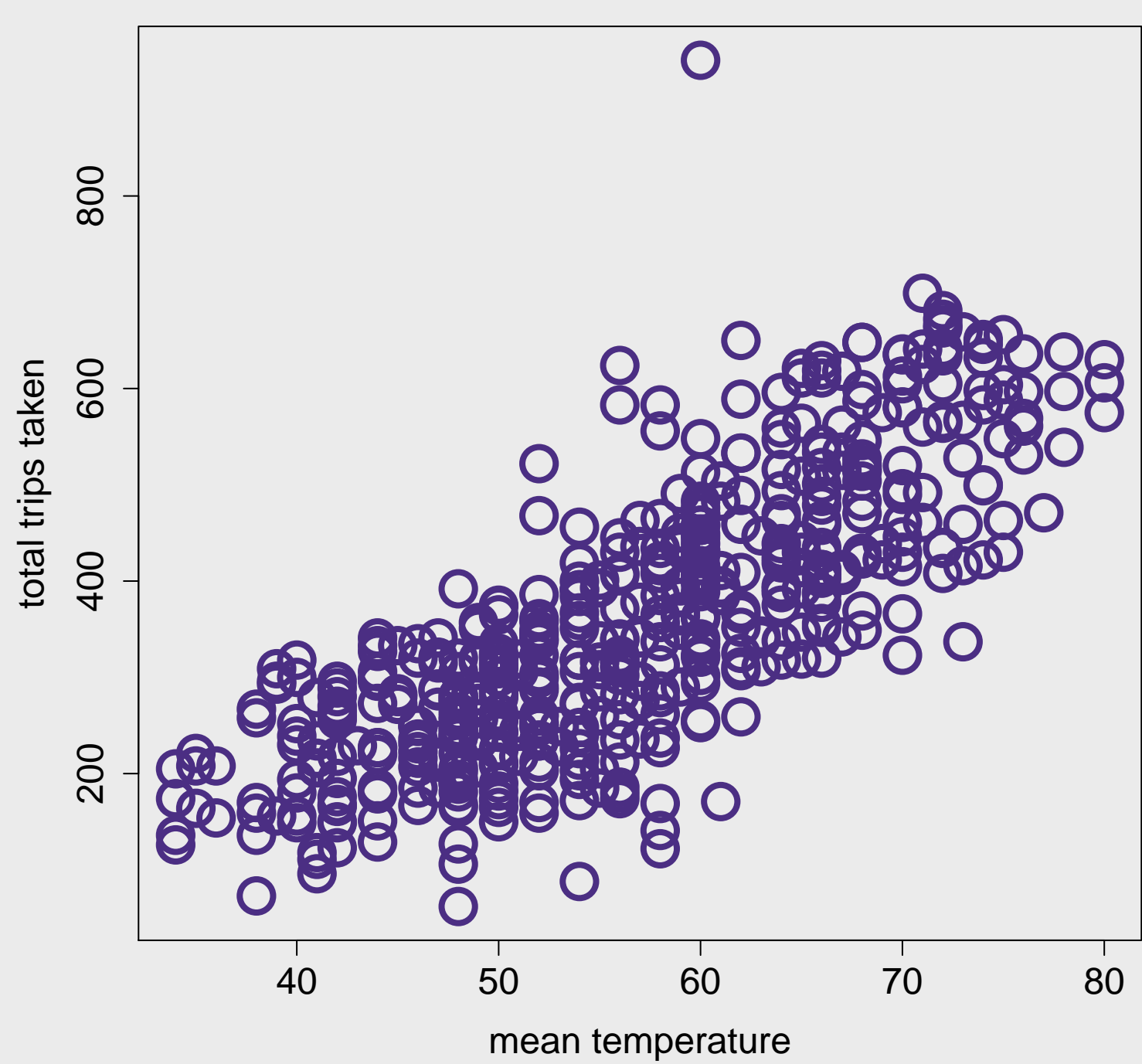
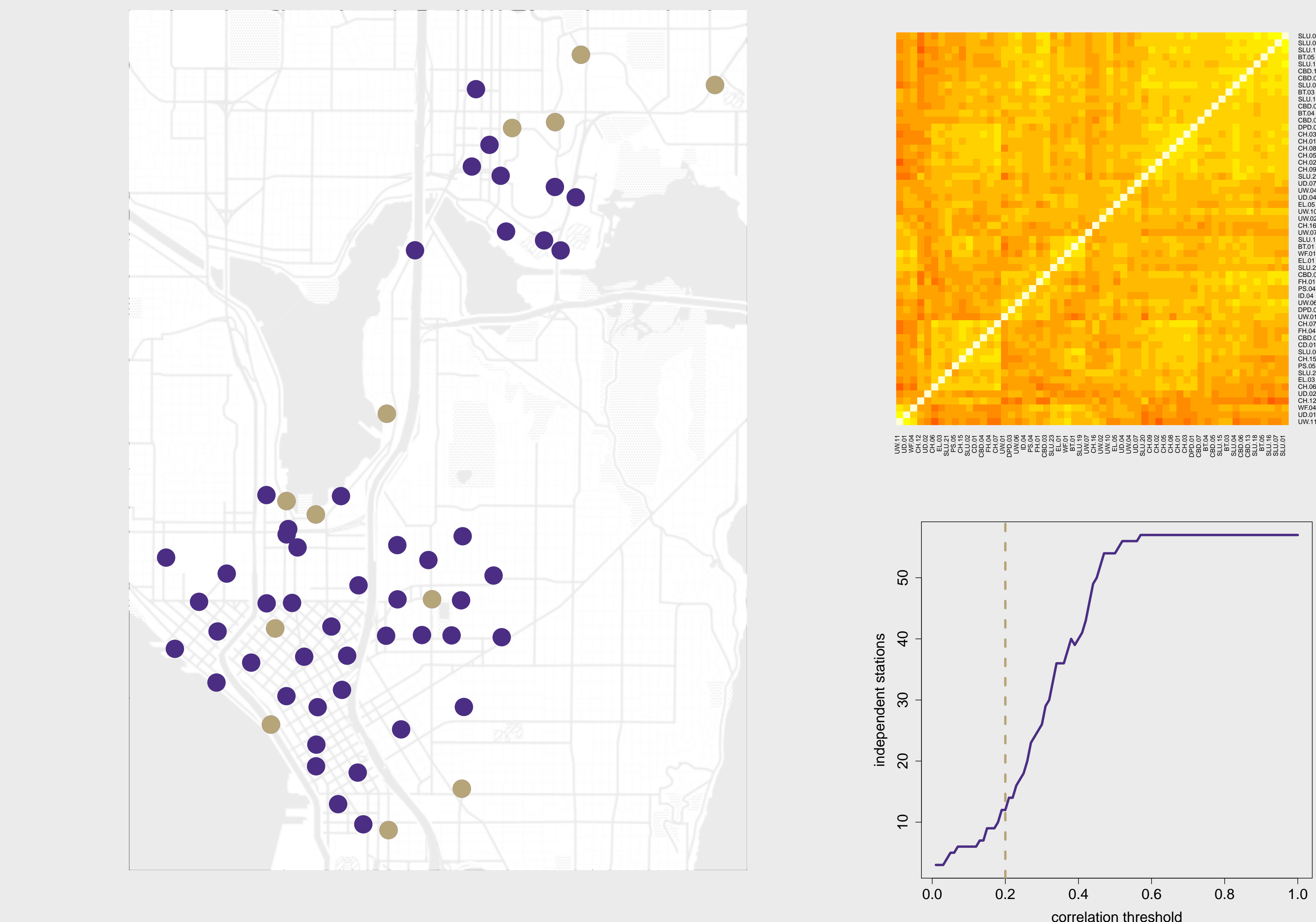
Let $y_i = [y_{ij}]^T$, where y_{ij} is the (transformed) total outflow of the j th station on the i th day. Then we assume that

$$y_i = X\beta_s + w_i^T\beta_w + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, \Sigma)$$

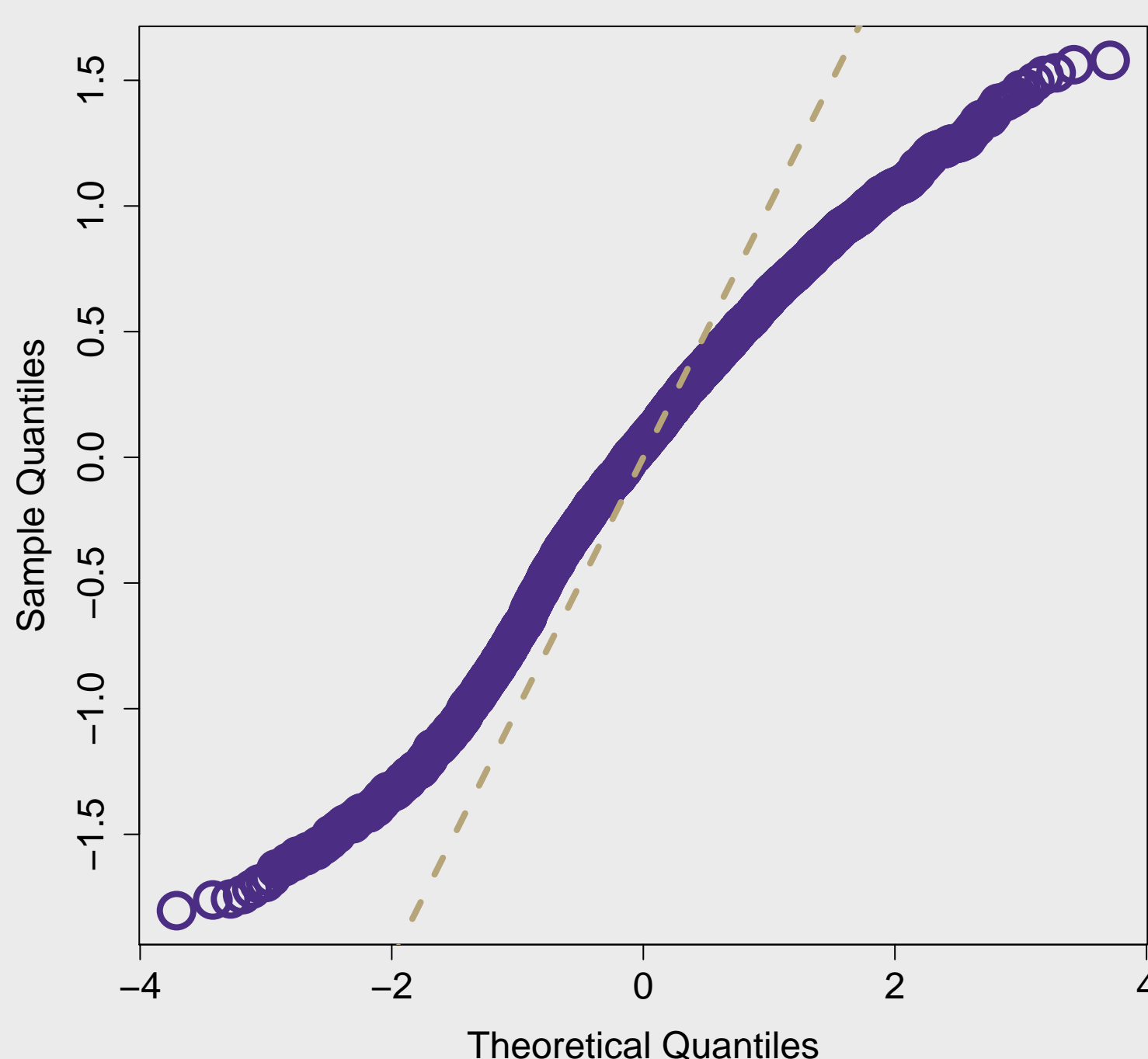
where X are the station features and w_i is the weather for the i th day. Note that we do not assume Σ is diagonal to account for inter-station correlations. Suppose we find a set of stations S_δ such that $\forall i, j \in S_\delta : \sigma_{ij} < \delta$. Then our model simplifies to

$$y_{ij} \approx x_j\beta_s + w_i^T\beta_w + \epsilon_{ij}, \epsilon_{ij} \sim \mathcal{N}(0, \sigma_{ij})$$

Results

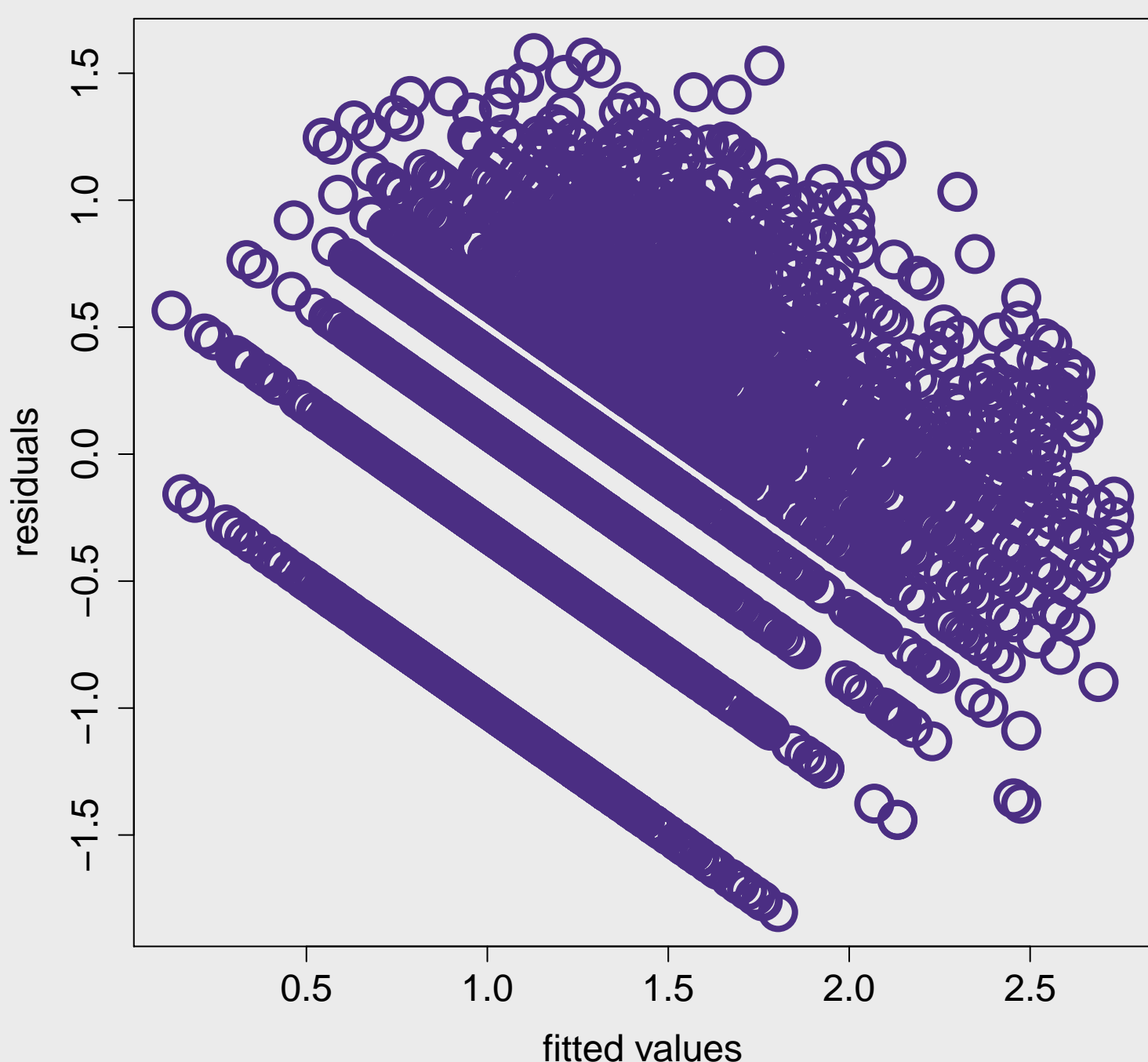


Normal Q-Q Plot



Model 1	
(Intercept)	0.924 (0.083)***
elevation	-6.239 (0.532)***
population density	-0.035 (0.003)***
job density	0.007 (0.000)***
income	-0.007 (0.004)
number metro routes	-0.010 (0.001)***
distance to metro	4.428 (0.609)***
listed on Google maps	-0.601 (0.033)***
mean temp.	0.021 (0.001)***
precip.	-0.536 (0.042)***
R ²	0.327
Adj. R ²	0.325
Num. obs.	4831
RMSE	0.628

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$



Methods

1. log-transform the per-station daily outflow to stabilize variance
2. Following technique for seemingly unrelated regression (SUR):
 - (a) Regress on weather covariates using ordinary least squares (OLS)
 - (b) Estimate $\hat{\Sigma}$ from the residuals
3. Use $\hat{\Sigma}$ to select independent subset of stations S_δ (e.g., with a connected-components algorithm)
4. Fit model with station features and weather covariates to data from selected stations using OLS

Validation

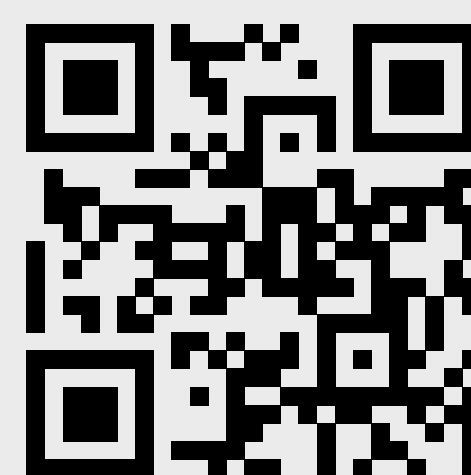
- ▷ Two independent code bases gave identical results
- ▷ Although not exactly normally distributed, residuals are symmetric
- ▷ Transformed observations demonstrate homoscedasticity (constant variance)

Concluding Remarks

- ▷ Linear regression is a challenging framework for analyzing a network; we appropriately subsampled our data to minimize violation of assumptions
- ▷ Controlling for weather, all station features considered were significant
 - ▷ Surprisingly, elevation was negative correlated with departures
 - ▷ Station proximity to transit decreased departures
 - ▷ Also surprisingly, stations listed on Google maps had decreased outflow
- ▷ Future work should analyze *inflow* (i.e., daily number of arrivals per station) and compare to our results

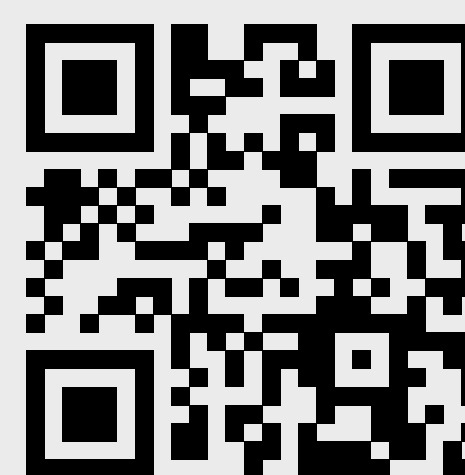
References and Acknowledgements

T Amemiya (1985). *Advanced econometrics*. Harvard University Press.
E Fishman (2016). *Transport Reviews*, 36:92–113.
K Gebhard and RB Noland (2013). *Transportation Research Board* 92.
D Kahle and H Wickham (2013). *The R Journal* 5:144–161.
DJ Kim et al. (2012). *Transportation Research Board* 91.
P Leifeld (2013). *Journal of Statistical Software*, 55:1–24.
A Small (2017). The Four Horsemen of the Bike Share Apocalypse. *CityLab*.
Thanks to Elena Erosheva and Michael Karcher for helpful discussions and our STAT 504 class for their feedback.



Download this poster

<https://git.io/vyPjX>



Fork the source code

<http://git.io/vyPjw>