

SENSING AND PREDICTING THE PULSE OF THE CITY THROUGH SHARED BICYCLING

International Workshop on Spatio-temporal Data Mining for a Better Understanding of Human Mobility:
The Bicycle Sharing System Case Study
December 5th, Paris, France



A quick story re: **traveling** to this workshop.

*This is an operations problem
a user satisfaction problem
a perception problem*



I never take Vélib'. The redistribution is completely broken.

— **Girl On Train**

Parisian Exchange Student from England

bikeshare research stakeholders & benefits



1

Operators: benefit from more accurate models of demand for load balancing

2

End-users: benefit from understanding and forecasting how system will be used for trip planning

3

Urban planners: can use bike models to improve the bikeability of the city

4

Social scientists & human geographers: can study and better understand human mobility and routine

“The social sciences can finally have access to masses of data that are of the same order of magnitude of their older sisters, the natural sciences”

Bruno Latour, 2007

French philosopher and sociologist

data

this talk

my
research

bikeshare data

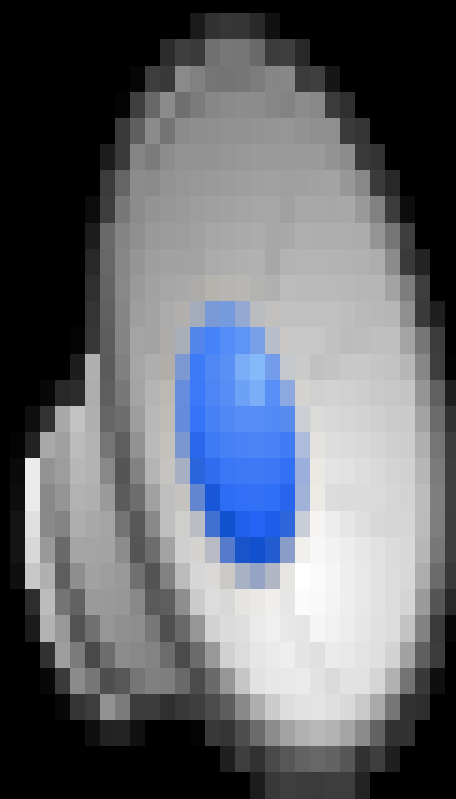
granularity



Station Capacity Data

Data collection method

- (1) Data dump from operator
- (2) Scrape bikeshare website



bikeshare data

granularity



Station Capacity Data

Data collection method

- (1) Data dump from operator
- (2) Scrape bikeshare website

bikeshare data

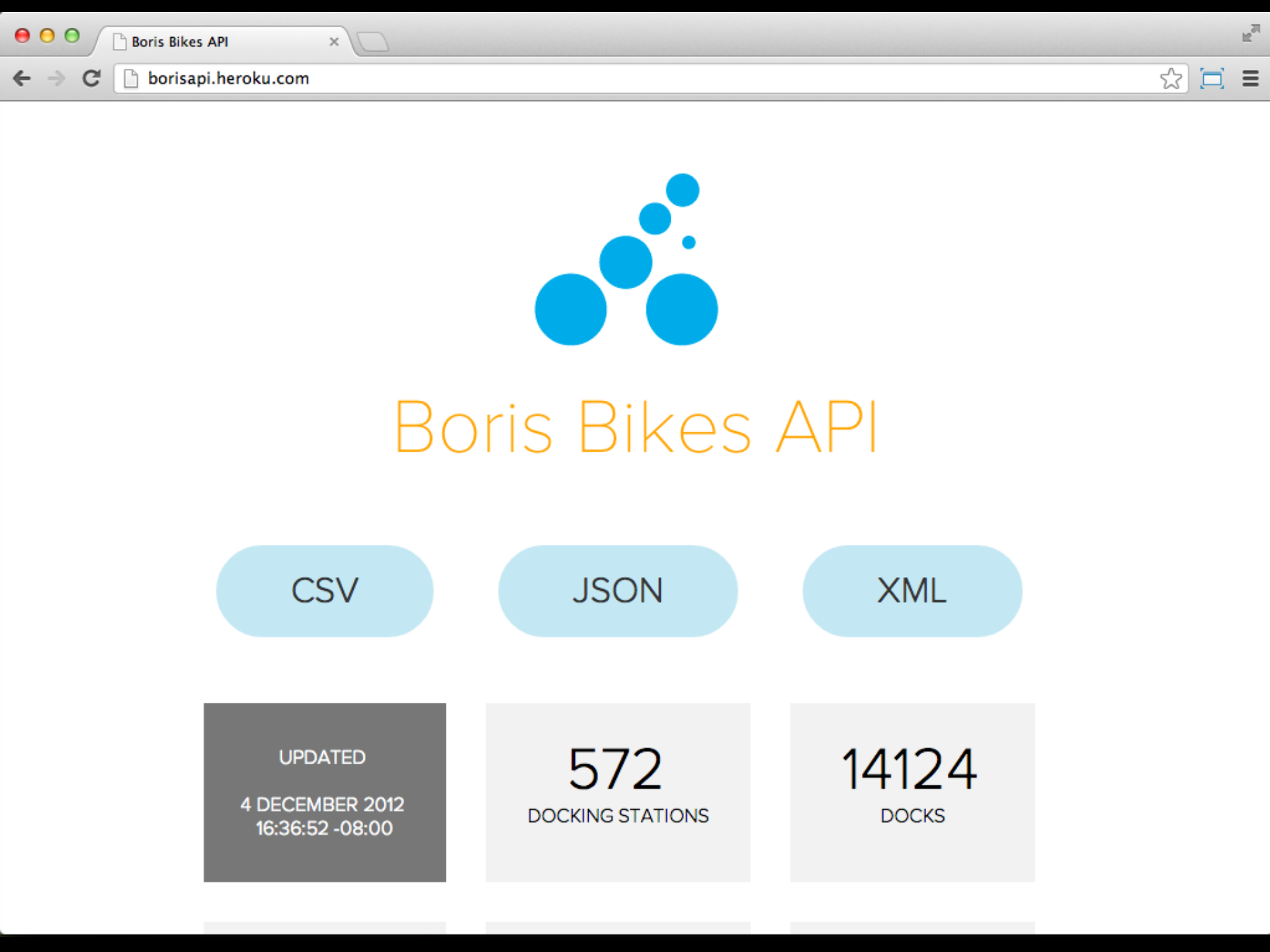
granularity



Station Capacity Data

Data collection method

- (1) Data dump from operator
- (2) Scrape bikeshare website
- (3) 3rd-party DIY APIs



Boris Bikes API

CSV

JSON

XML

UPDATED

4 DECEMBER 2012
16:36:52 -08:00

572

DOCKING STATIONS

14124

DOCKS

bikeshare data

granularity



Station Capacity Data

Data collection method

- (1) Data dump from operator
- (2) Scrape bikeshare website
- (3) 3rd-party DIY APIs



O/D Bike Data

Data collection method

- (1) Data dump from operator
- (2) ...

Greater Greater Washington

The Washington, DC area is great. But it could be **greater**.

Capital Bikeshare releases anonymous trip data

by [David Alpert](#) • January 11, 2012 3:32 pm

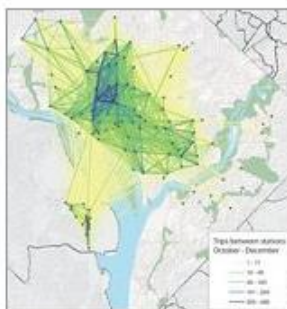
Programmers or analysts interested in studying Capital Bikeshare patterns or creating useful apps can now do a lot more. Capital Bikeshare has followed through on its promise and posted data files with individual (but anonymous) trip data.

The files, one for each quarter going back to late 2010, list individual trips, including the time each started and ended, duration, which station it started and ended at, and an identifying number for the individual bike. It doesn't say anything about the member who used the bike, except whether they are a "registered" (annual or monthly) member or a "casual" member (daily or 3- or 5-day).

Now, people can generate tables or graphics showing the most popular station pairs, or where people most often go from an individual station, or what weather patterns make usage heavier or lighter, or where the nighttime activity is, and much more.

This data has been available for some time for London, allowing people to create animations of a day's CaBi usage and diagrams of a single bike's path over several days.

The folks who built those and other tools can now even adapt their code to work for Capital Bikeshare, if they're so inclined.



Anyone can now make a map like this one of CaBi trip patterns. Image from [CommuterPageBlog](#).

Subscribe

All posts RSS • Add to Google

Comments on this post

Posts by David Alpert

Follow us on Twitter @ggwash

Become a fan on Facebook

Get a daily summary by email:

Most Active Posts

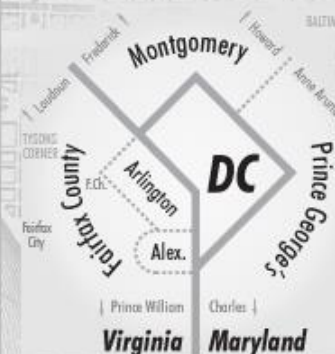
Recent Links



Breakfast links: Getting around

- Metro keeping 2 Southeast bus routes
- Need to fund transit to grow
- Transit pays off
- What's going to happen to transportation?
- More DC Bikeshare delayed

How can our region be greater?



Features



Animated history of Metrorail





Trip History Data

Overview of the Trip History Data

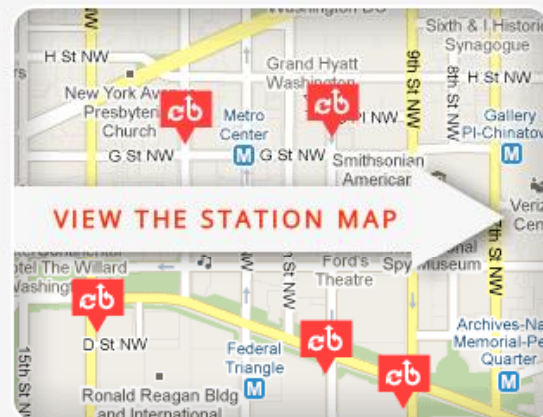
When a rental occurs within the system our software collects basic data about the trip. That data can be exported from our system and used for various types of analysis or research. By making this data available we hope to stimulate that analysis or research amongst a much wider community. Please note, all private data including member names has been removed from these files.

What is included?

Each .csv file contains data for one quarter of the year. Within each file there are 7 columns.

join

for a day, 3 days, month, or year.



Some Quick **Fun** Facts

Analysis Enabled by Anonymized O/D Bike Data

- (1) "Average" trip is downhill (by -1.94 meters)
- (2) Last mile usage: four most common trips are short and seem to cover areas that subway/bus do not
- (3) Sixth most common trip is a return trip from Smithsonian station back to the Smithsonian station
- (4) Casual vs. members usage fees: 40.7% casual riders incur fees vs. 3.3% of members incur fees



Mar 23, '12
03:50

11

Using **heuristics** to infer bike flow.

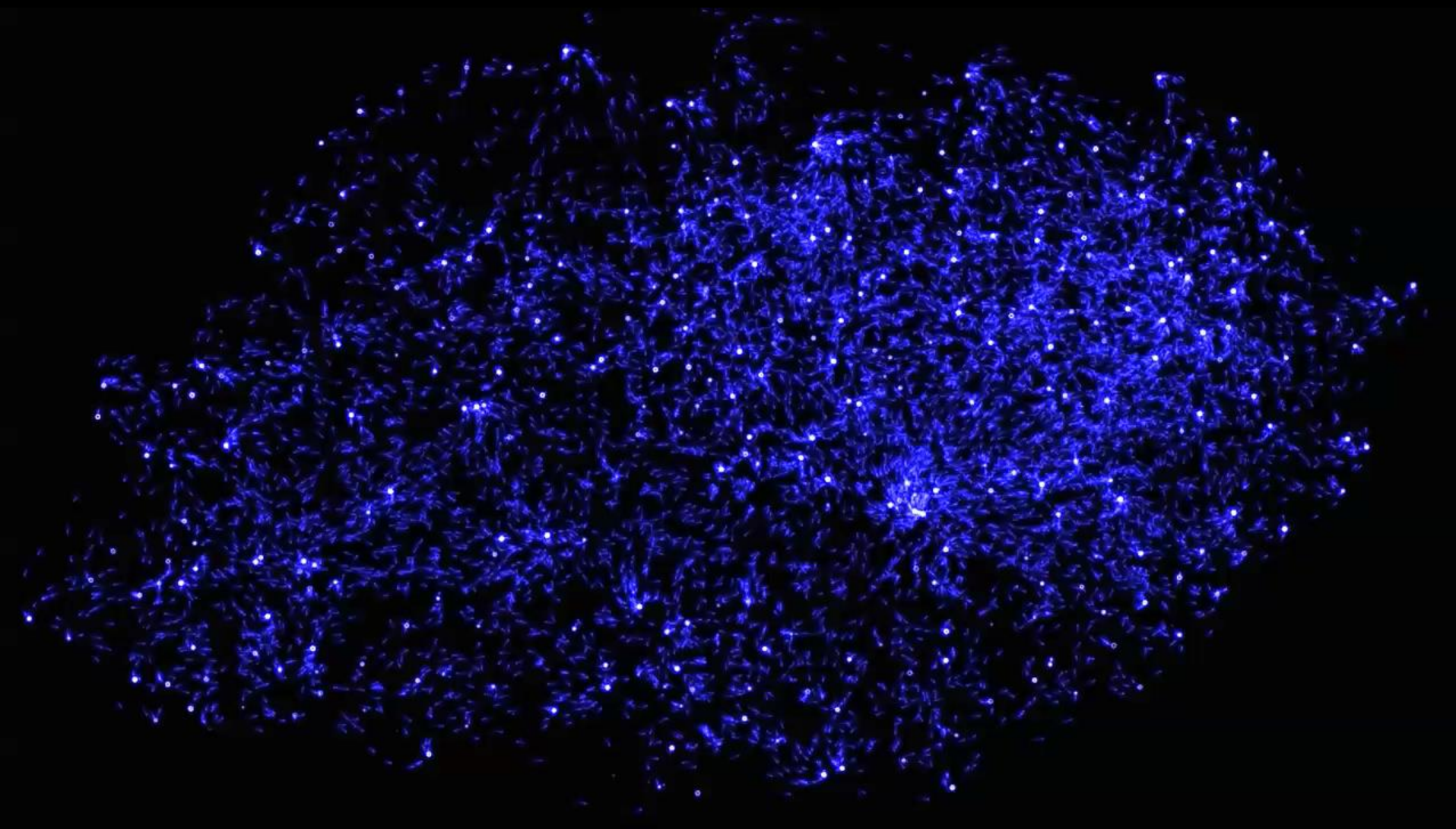


At 2010-10-04 06:01:00 there were 3 bikes in use.

The routing is done using OpenStreetMap data & Routing routing scripts optimized for bike usage (*i.e.*, **constant speed** on all road types, **obeying one-way roads** and **taking advantage of marked cycleway**). I've tweaked the desirability of road types, so that the trunk and primary roads are only slightly less desirable than quieter routes.

— **Oliver O'Brien**
UCL CASL

<http://oliverobrien.co.uk/2011/02/boris-bikes-flow-video-now-with-better-curves>



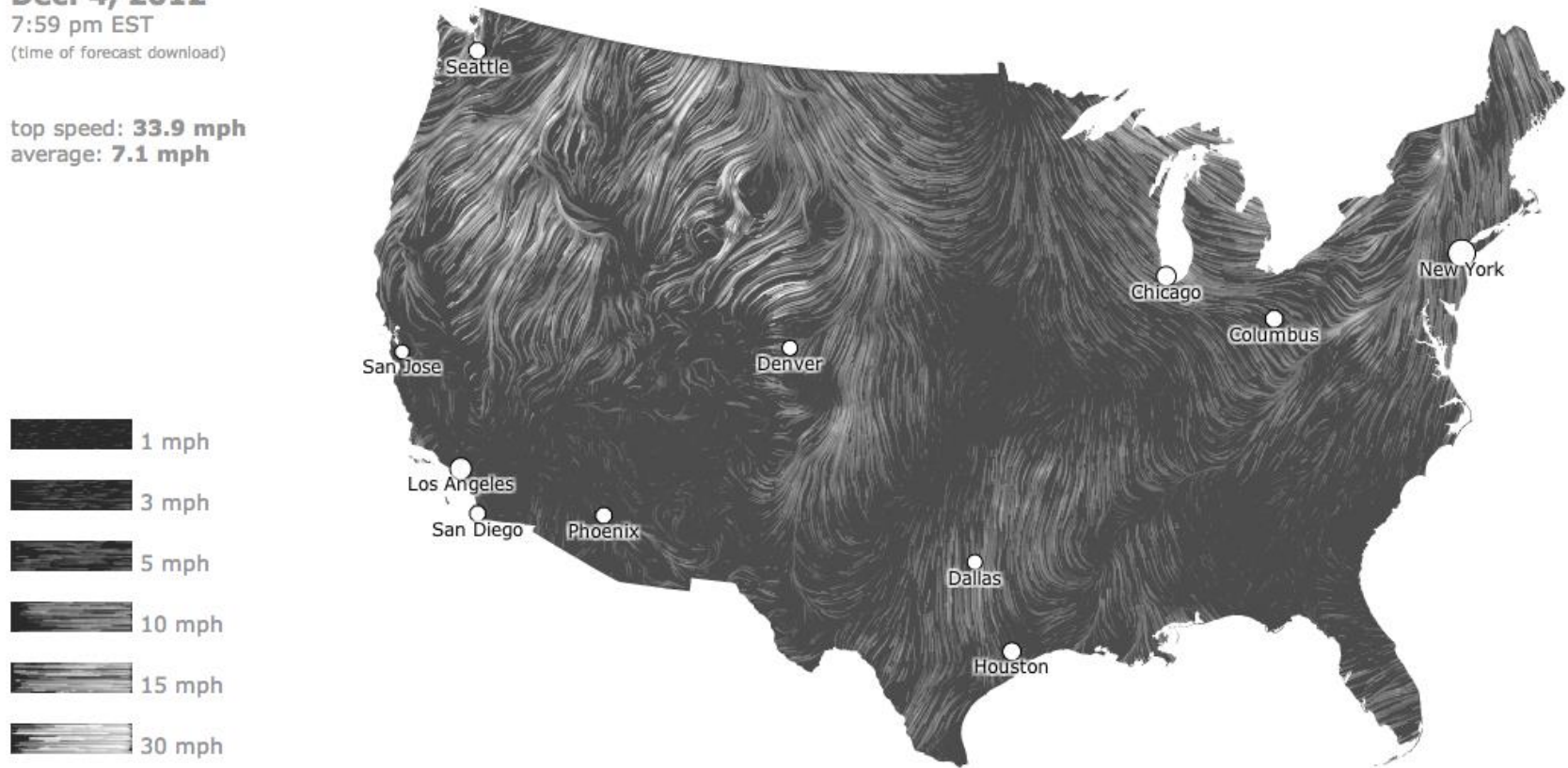
wind map

Dec. 4, 2012

7:59 pm EST

(time of forecast download)

top speed: **33.9 mph**
average: **7.1 mph**



bikeshare data

granularity



Station Capacity Data

Data collection method

- (1) Data dump from operator
- (2) Scrape bikeshare website
- (3) 3rd-party DIY APIs



O/D Bike Data

Data collection method

- (1) Data dump from operator
- (2) ...



O/D User Data

Data collection method

- (1) Data dump from operator
- (2) Scrape user account logs

Identifying and Explaining Inter-Peak Cycling Behaviours within The London Cycle Hire Scheme

R. Beecham¹, J. Wood², A. Bowerman³

¹City University London, St. John's Street, London, EC1V 0HB
Email: roger.beecham.1@city.ac.uk

²City University London, St. John's Street, London, EC1V 0HB
Email: j.d.wood@city.ac.uk

A visual analytics approach to understanding cycling behaviour

Roger Beecham*
City University London

Jo Wood, Member, IEEE †
City University London

Audrey Bowerman‡
Transport for London

ABSTRACT

Existing research into cycling behaviour

“Working collaboratively with Transport for London (TfL), customer records reporting a unique customer identifier, gender and postcode, have been made available. So too has a complete set of user journeys [for those customer ids]”

customers using the scheme to encourage inter-peak usage.

Working with colleagues at TfL, and with access to LCHS's customer databases, we attempt to identify and explain inter-peak travel context, circumstances and customer characteristics that underpin. Using techniques from information visualization and geovisualization [Roberts 2004, Wood et al. 2011], we demonstrate how this type of analysis has been achieved using a visual analytics application. After outlining our some initial findings and discuss how, in better understanding customer behaviours, we aim to provide insights that may directly inform and operational decisions around the LCHS's expansion.

2. Dataset and Analysis Techniques

2.1 The LCHS Customer and Journeys Dataset

Our analysis relies on two complementary datasets: a complete set of journey records. For every customer registered with the scheme, a unique customer identifier is generated. For every journey made, a pair representing the docking station that journey started at and the docking station that journey ended at, are recorded within a journey record. For every journey made, a pair representing the docking station that journey started at and the docking station that journey ended at, are recorded within a journey record.

Existing research into cycling behaviour has focussed on actual rather than self-reported behaviour have generally relied on GPS logs [3] or automated traffic counts at fixed sites [6]. Due to their complexity, the former have been relatively small in scale and the latter, only meaningfully identify behaviour.

Shared-bicycle schemes offer new research possibilities. In many recent schemes, data on usage are continually reported to central databases. Researchers working within data mining [5] [7], and information visualization [13] have queried these data to identify patterns of scheme use, as well as more nuanced space-time journey flows. This analysis has nevertheless been constrained by the level of detailed information made publicly available. Whilst a complete set of journey records, including journey origin-destination (OD) and start and end times was used by Wood et al. [13], these data could not be linked back to individual customers. Cyclists' journey histories, and the context framing those journeys, could not be identified. This limits the extent to which such data can be used to engage with the more complex questions around motivations and barriers to cycling [7] [13]. Working collaboratively with Trans-

*e-mail: roger.beecham.1@city.ac.uk

†e-mail: j.d.wood@city.ac.uk

‡e-mail: audreybowerman@tfl.gov.uk

¹Transport for London's cycle hire website.



Figure 1: An early visual analytics prototype, which links customer segments (top left), with a spatial (centre) and temporal (bottom) view.

port for London (TfL), customer records reporting a unique customer identifier, gender and postcode the customer registered with, have been made available. So too has a complete set of user journeys. Linking with geodemographic and other contextual information, and querying these attribute rich data within a visual analytics application, we attempt to explore and explain cycling behaviour from an individual customer perspective.

2 OBJECTIVES AND APPROACH

The research project aims to:

- Classify bike share customers according to the journeys they make.
- Validate, rewrite and add qualitative descriptions to these classifications paying attention to journey context; by querying the dataset at particular space-times, in response to changes internal and external to the scheme.
- Furnish social scientists, and strategists within TfL, with generalisable insights into the barriers, incentives and conditions that motivate cycling behaviour.

This approach sits comfortably within a visual analytics framework [11]. We will take a large and attribute rich dataset, query it to identify general and distinct space-time journey patterns, before creating modelled data and subsequently querying the models. Attempting to engage with difficult questions around cycling behaviours, it will be necessary to quickly consider many combinations of contextual variables and customer attributes. This might be best achieved through a highly flexible visual analytics application.

3 EARLY ANALYSIS AND FIRST VISUAL ANALYTICS PROTOTYPE

At the time of writing, the dataset contained 114,947 valid customer records, linked to 6,490,479 journeys. After loading data into an SQLite database, relevant derived variables were computed. Linking customers with their journeys, recency-frequency (RF) segmentation, a technique used within direct marketing [8], was performed. Matching customers' postcodes to geographic coordinates,

Rentals (72)

Capital Bikeshare is currently undergoing a system update which affects your rental history statistics from being calculated accurately. During this update, your distance, calories burned and CO2 lbs saved statistics will not be displayed. We are working to finalize the update and will provide correct statistics for your rental history again upon its completion. Your other rental history data has not been affected and is accurately displayed on this page. Thanks for your patience!

Start Station	Start Date	End Station	End Date	Duration	Cost
Lamont & Mt Pleasant NW	10-27-2012 10:40 am	New York Ave & 15th St NW	10-27-2012 11:22 am	42 minutes, 35 seconds	\$ 1.50
Connecticut Ave & Newark St NW / Cleveland Park	10-23-2012 9:09 pm	Lamont & Mt Pleasant NW	10-23-2012 9:20 pm	10 minutes, 49 seconds	\$ 0.00
Adams Mill & Columbia Rd NW	10-21-2012 10:51 am	Adams Mill & Columbia Rd NW	10-21-2012 10:51 am	9 seconds	\$ 0.00
20th & Crystal Dr	10-20-2012 11:25 am	Braddock Rd Metro	10-20-2012 11:55 am	29 minutes, 30 seconds	\$ 0.00
Lamont & Mt Pleasant NW		New York Ave & 15th St NW		32 minutes	

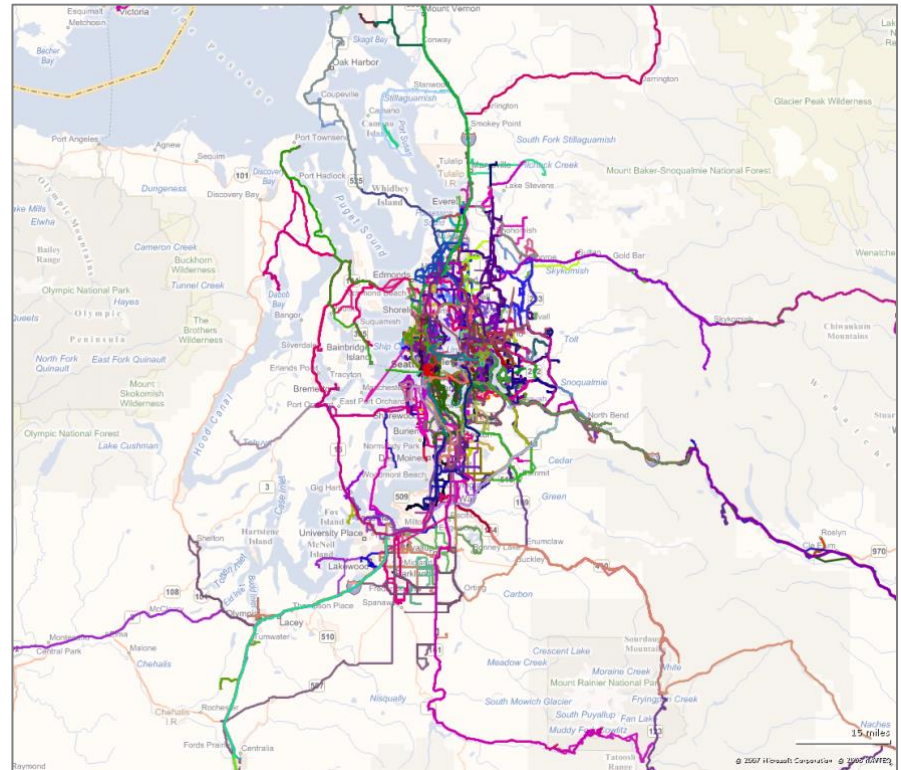
Human route repetition in **car driving**.

Route Prediction from Trip Observations

14,468 trips / 240 subjects

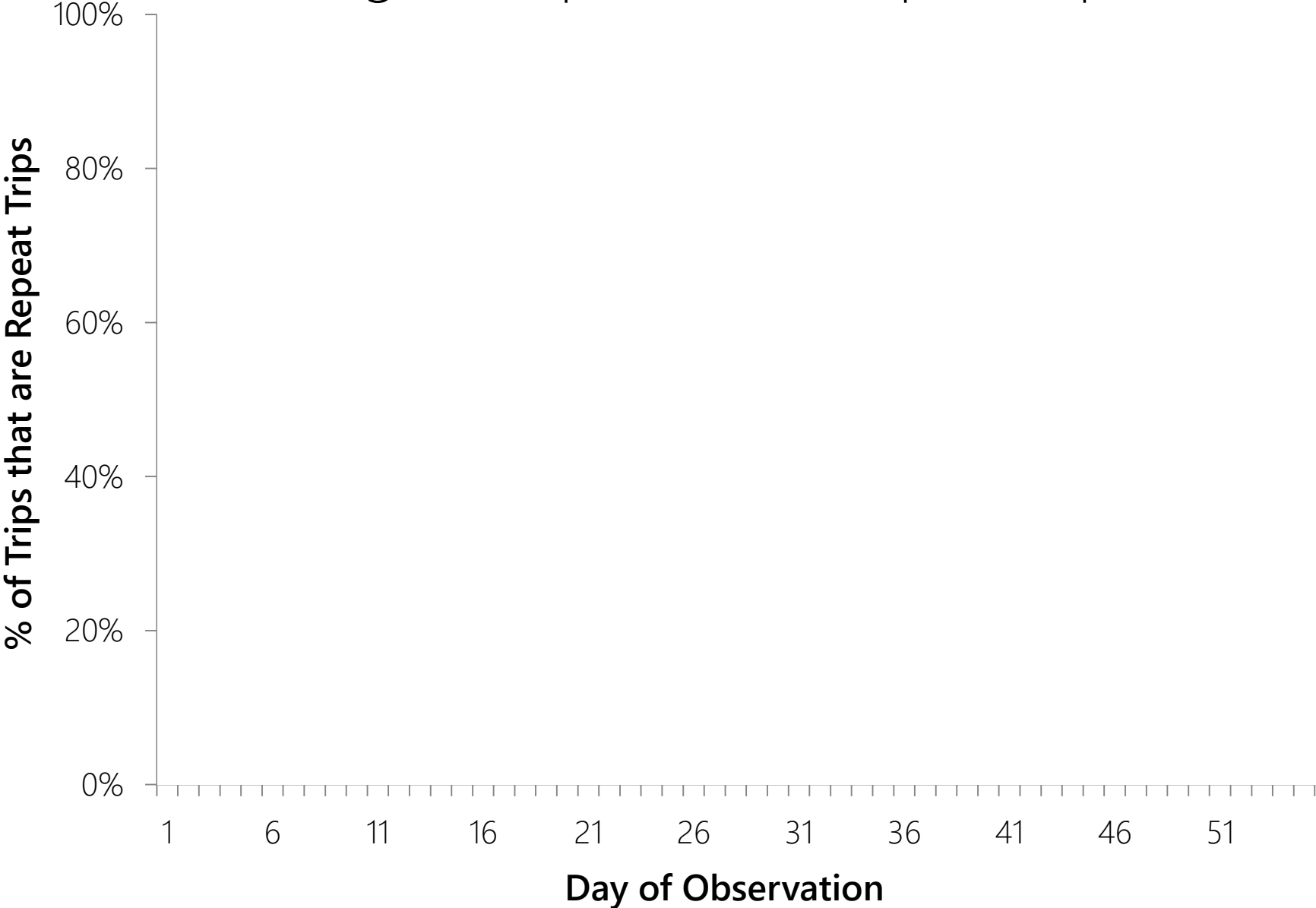
Description	Average	Median
trip distance (miles)	7.7	4.2
trip time (min)	16.3	11.5
num trips / day	4	3.9
num trips / subject	60.3	50
num days of data / subject	15.1	13

High Level Trip Stats

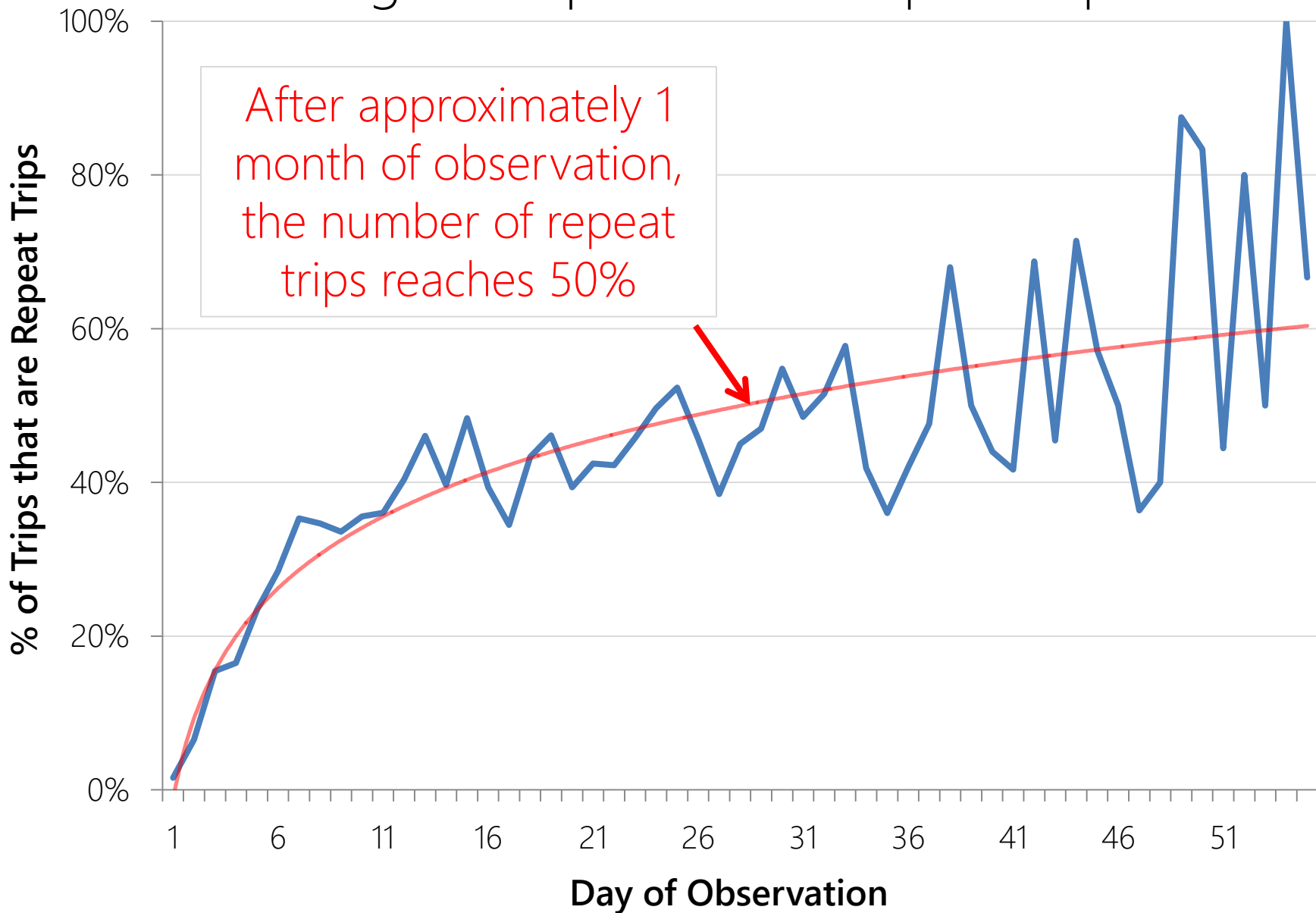


Greater Seattle Area

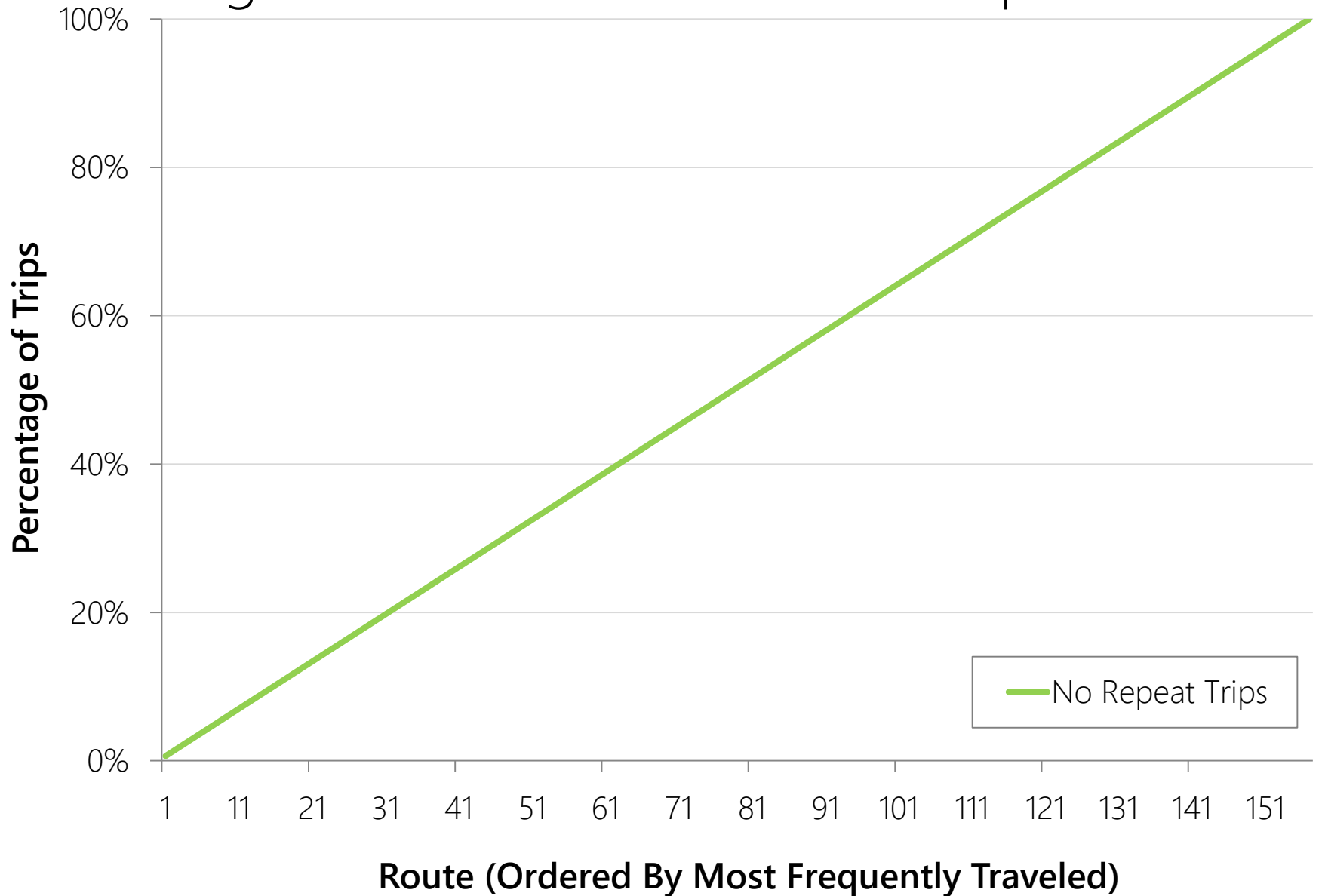
Percentage of Trips that are Repeat Trips



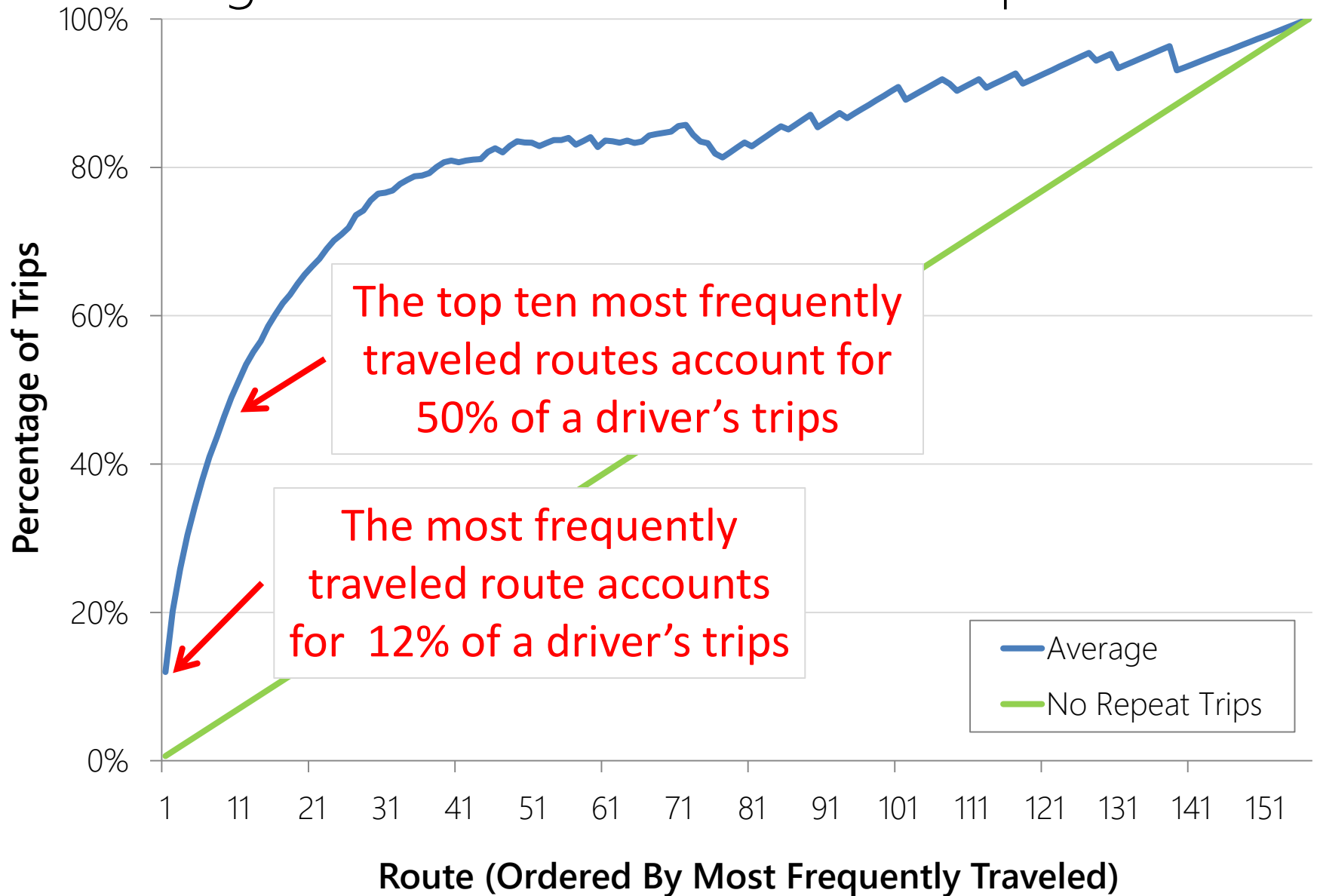
Percentage of Trips that are Repeat Trips



Avg Cumulative Distribution of Trips in Routes



Avg Cumulative Distribution of Trips in Routes



bikeshare data

Event Based

Station I/O with cloud
when bike docked/hired



Station Capacity Data



O/D Bike Data



O/D User Data

Continuous

Use sensors to track
movement



Instrumented Cities
(e.g., CCTVs)



Instrumented Bicycles
(e.g., GPS sensor)



Instrumented Users
(e.g., mobile phones)

bikeshare data

Event Based

Station I/O with cloud
when bike docked/hired



Station Capacity Data



O/D Bike Data



O/D User Data

Continuous

Use sensors to track
movement



Instrumented Cities
(e.g., CCTVs)



Instrumented Bicycles
(e.g., GPS sensor)



Instrumented Users
(e.g., mobile phones)

A person wearing a bright yellow jacket is riding a red shared bicycle. The bicycle has a white fender with the 'bicing' logo and the website 'www.bicing.com'. In the background, there is a large, historic stone building with arched windows and a courtyard filled with many other red shared bicycles. Another person in a dark jacket is standing near a bicycle on the right. The scene is set in an urban environment with cobblestone pavement.

sensing and
predicting the
movement of a
city via shared
bicycling

[Froehlich *et al.*, UrbanSense2008; IJCAI2009]

bicing

barcelona, spain

Summer 2008:

- 373 stations
- 6,000 bicycles
- 150,000 subscribers





bicing

jueves 07 de agosto de 2008

Catalán

Mapa de
estacionesInformación del
servicio

Zona de usuarios

Contacto

Noticias

Inicio
Mapa Web
Añadir a Favoritos

acceso usuarios

Usuario:

jogineumann

Contraseña:

Entrar

Regístrate [aquí](#)

Copyright © 2007 BICING · Todos

los derechos reservados [[Aviso](#)[Legal](#)]

Mapa de estaciones



A continuación se pueden visualizar en el siguiente plano las estaciones de bicicletas actualmente en funcionamiento. Así como, ver en tiempo real las disponibilidades de bicicletas en cada una de ellas.

Distrito: Código postal:
 Dirección: Estaciones Vacías: ☐ Estaciones Llenas: ☐



Estaciones con mas de una bici
 Estaciones sin bicis
 Estaciones cercanas

RESULTADOS

Total estaciones activas: 375

Ciutat Vella

Pedralbes

Sants-Montjuïc

[Mostrar todas](#)

Eixample

Poble Nou

Sarrià Sant Gervasi

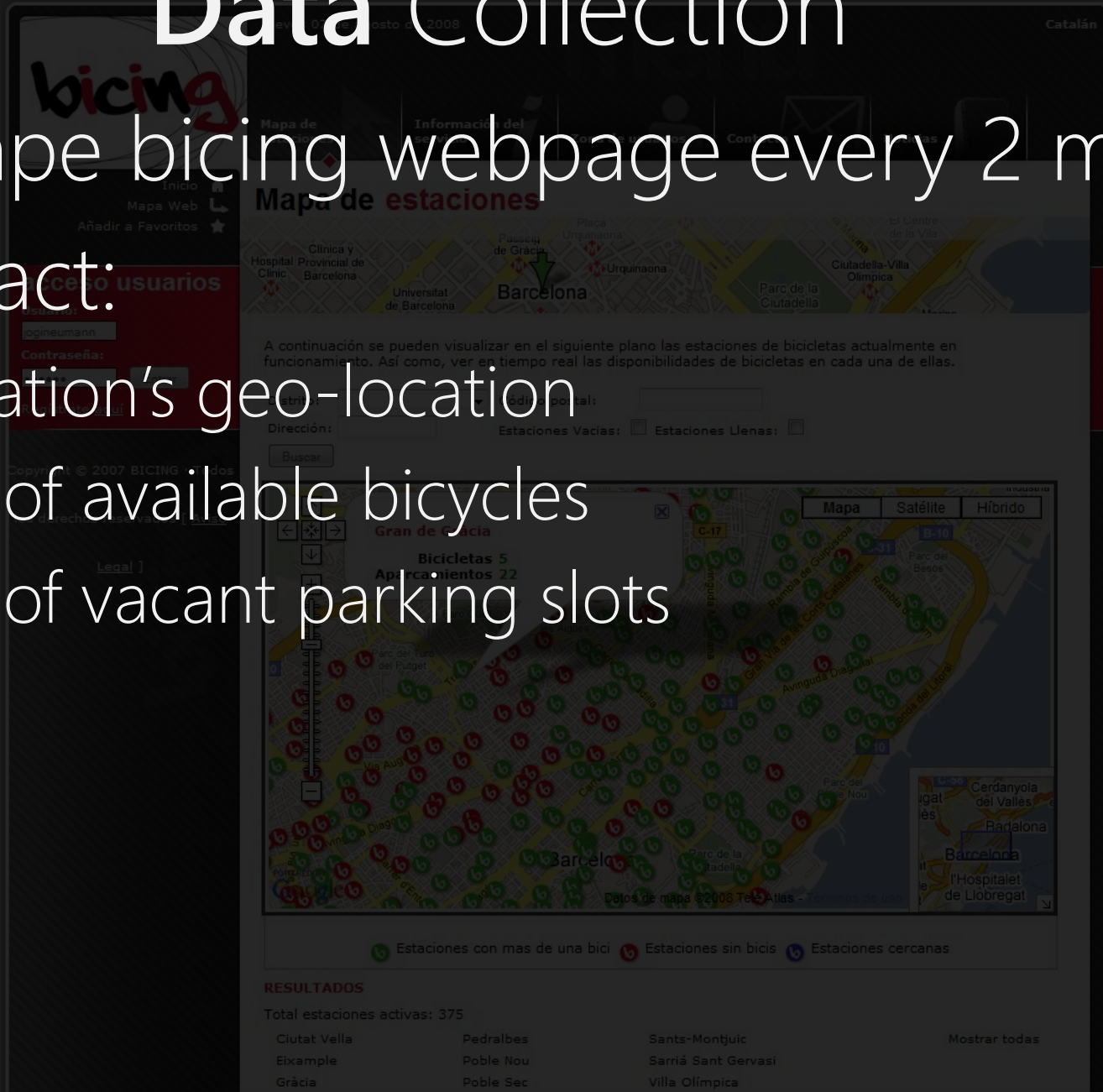
Gràcia

Poble Sec

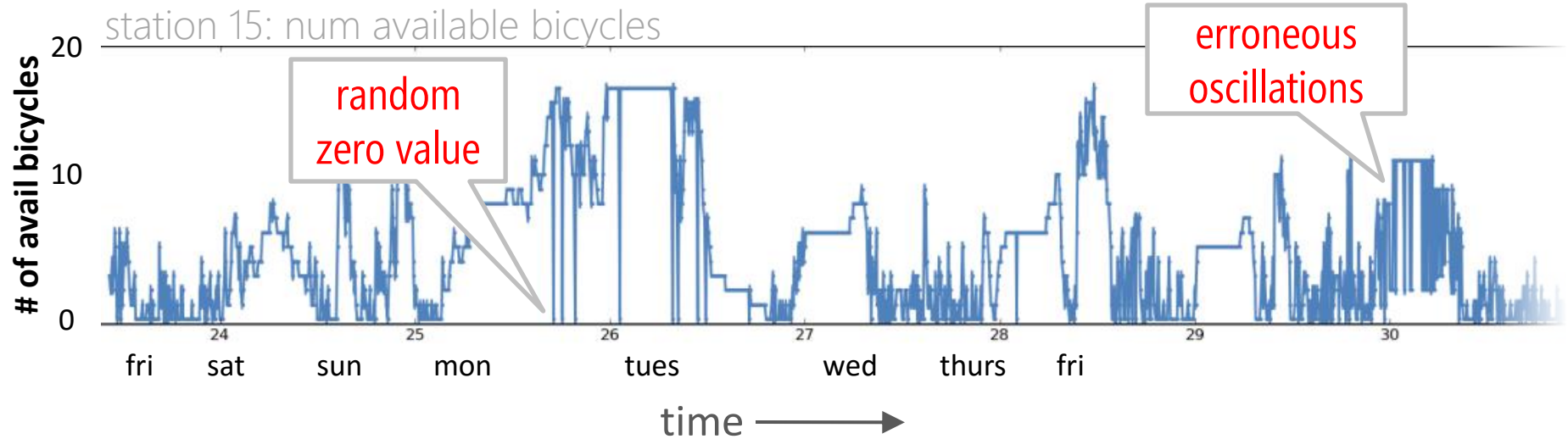
Villa Olímpica

Data Collection

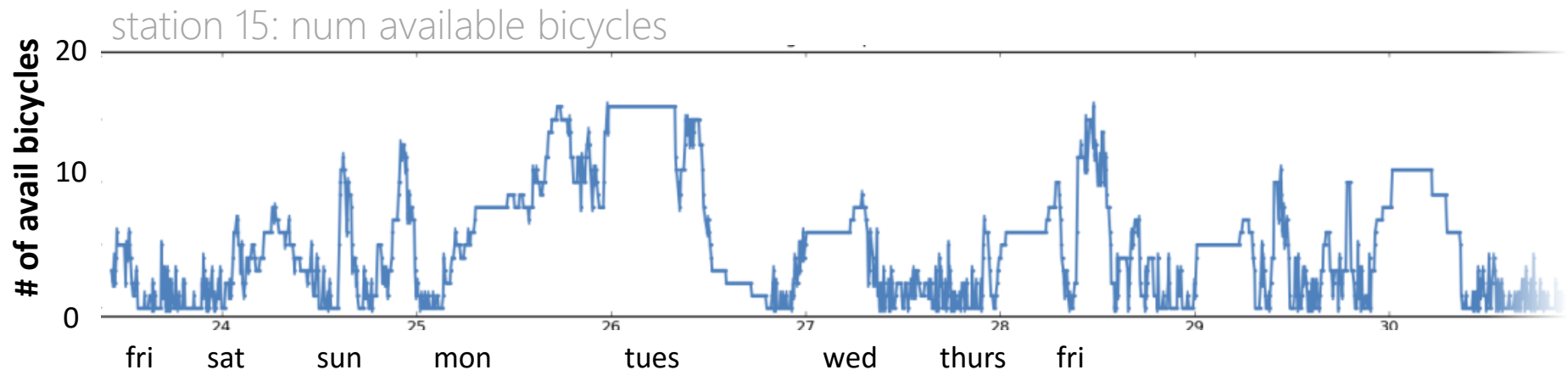
- Scrape biking webpage every 2 mins
- Extract:
 - station's geo-location
 - # of available bicycles
 - # of vacant parking slots



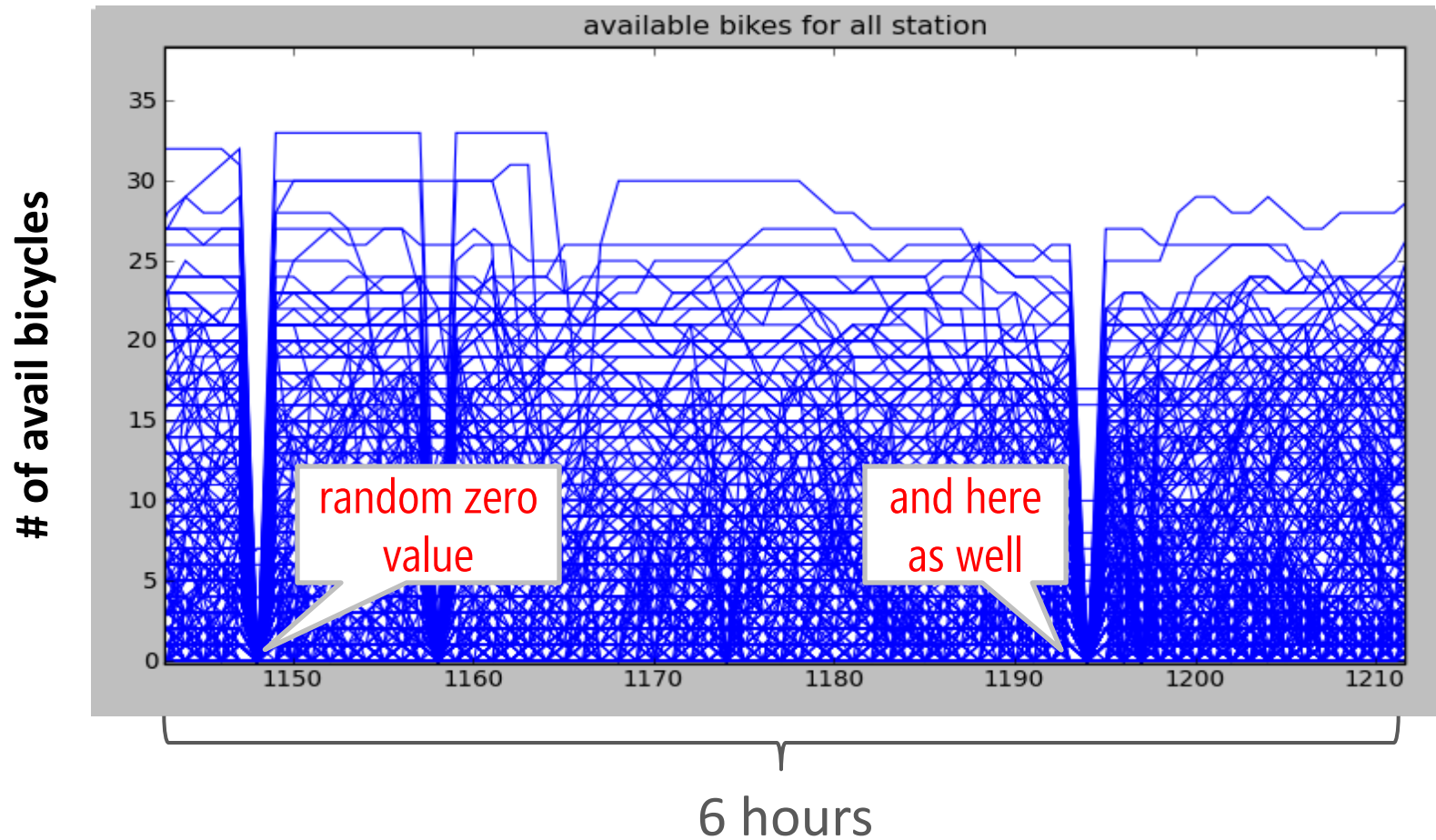
before cleansing



after cleansing



How do we know what to clean?



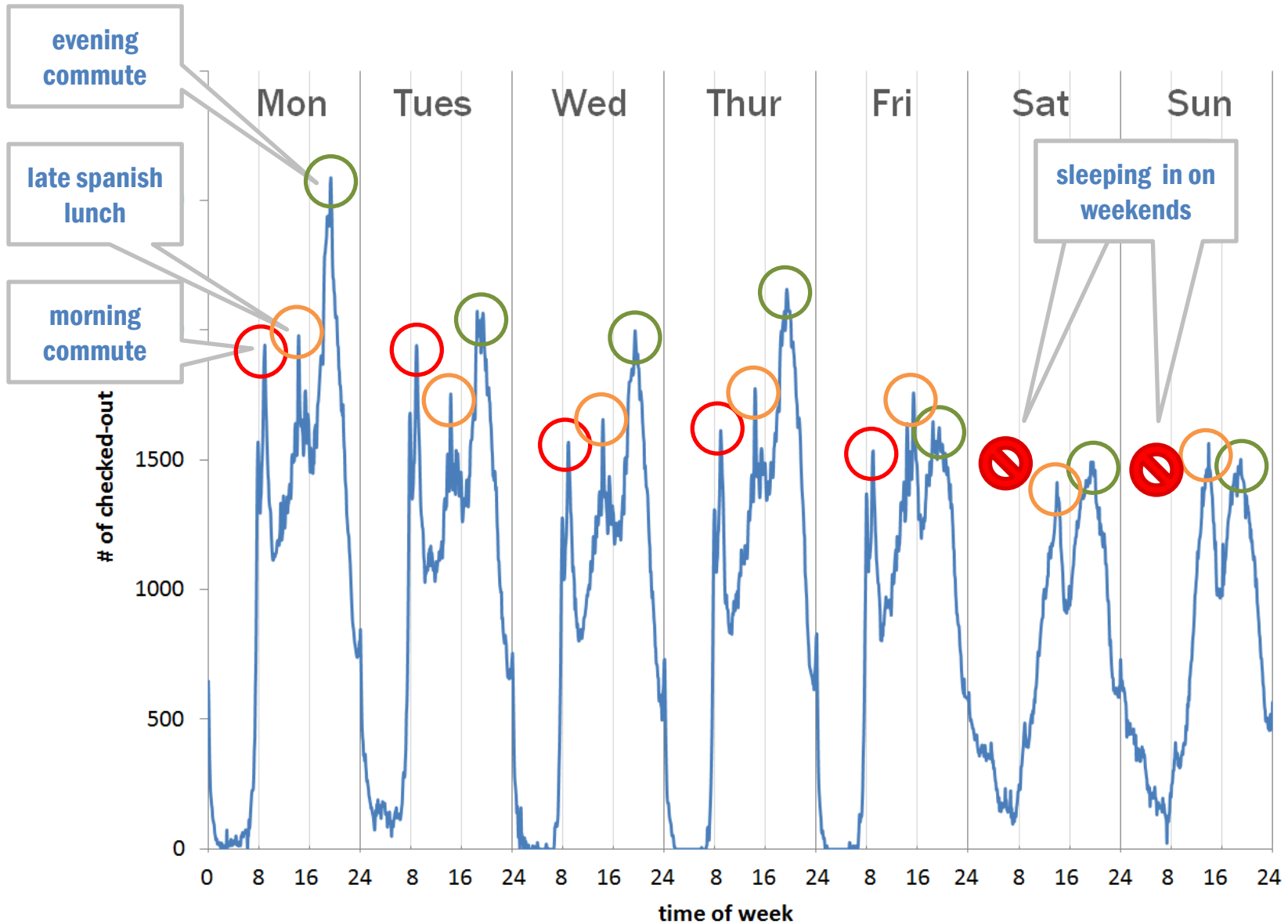
dataset

13 weeks of observations

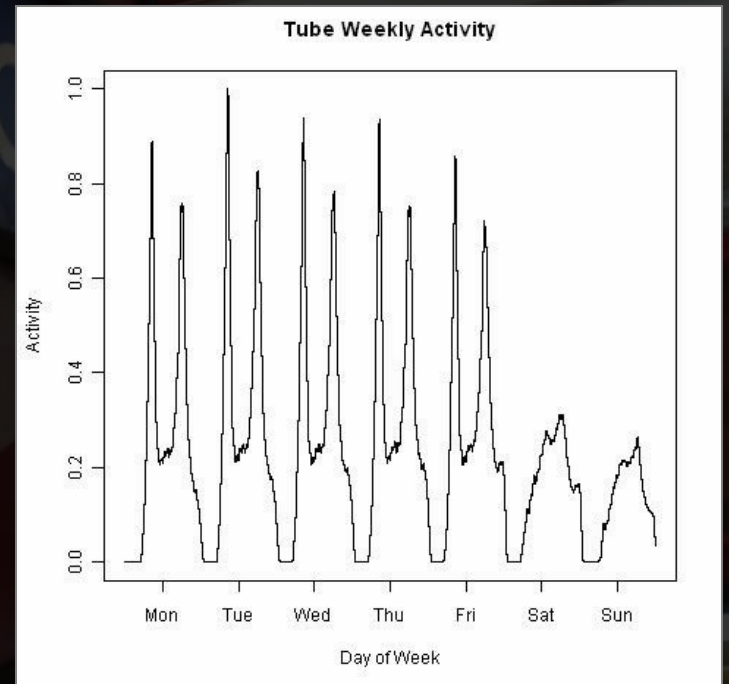
Aug 27 – Dec 1, 2008

	raw dataset	cleaned dataset
stations	390	370
days	25K	22.7K
observations	26.1M	20.2M
parking slots	9831	9315

Num checked-out bicycles across all stations



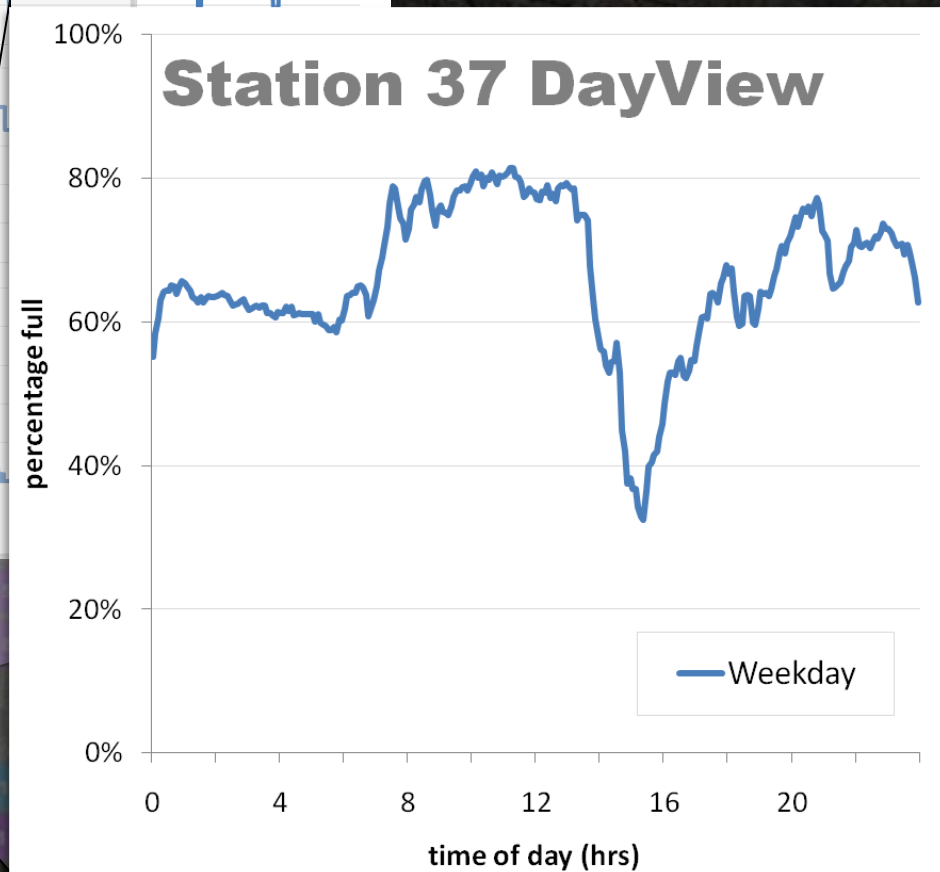
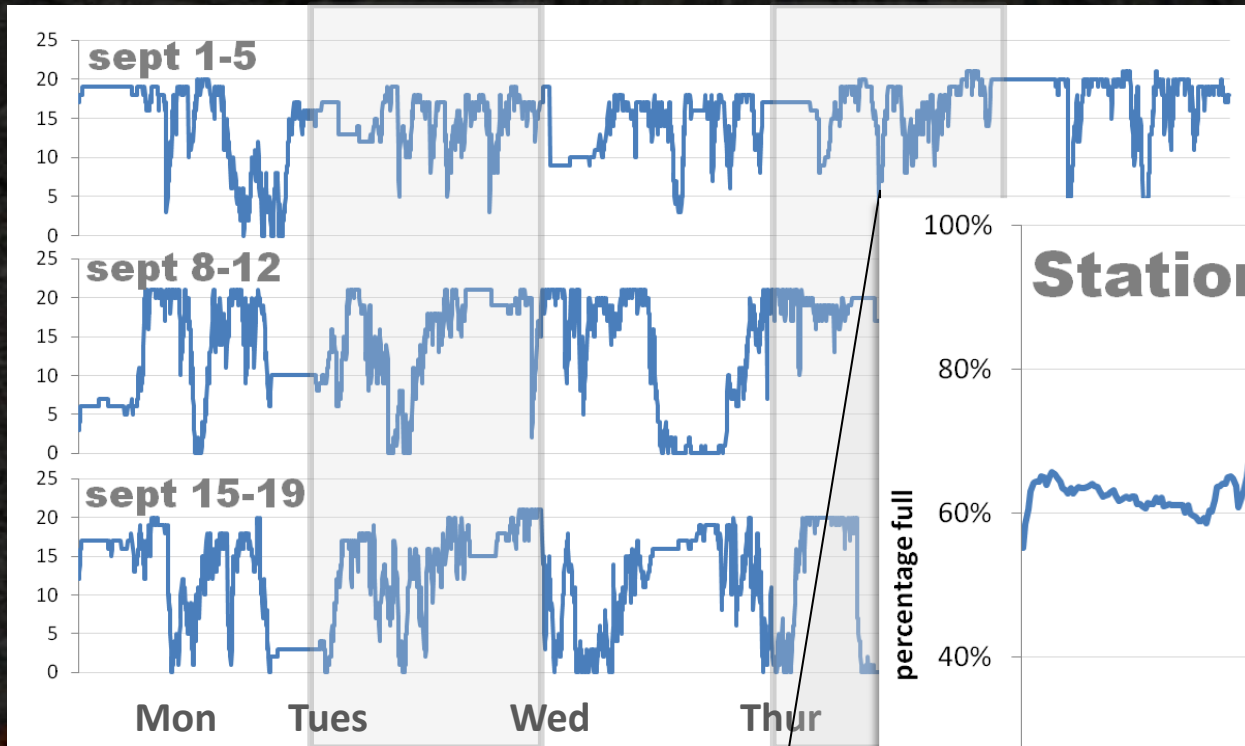
Two-spike pattern found in study of London Underground

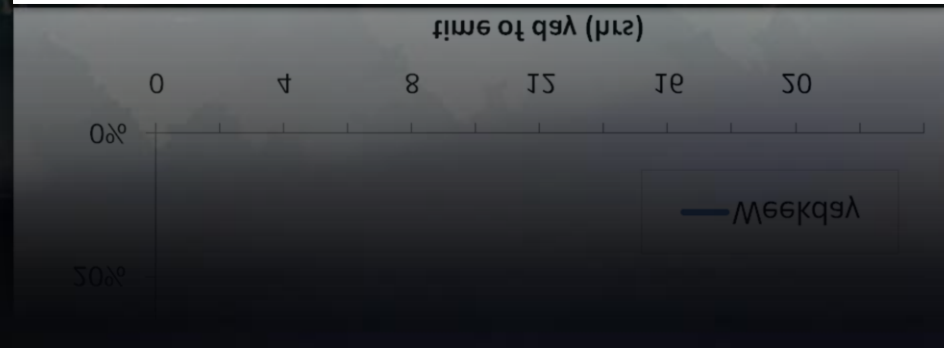
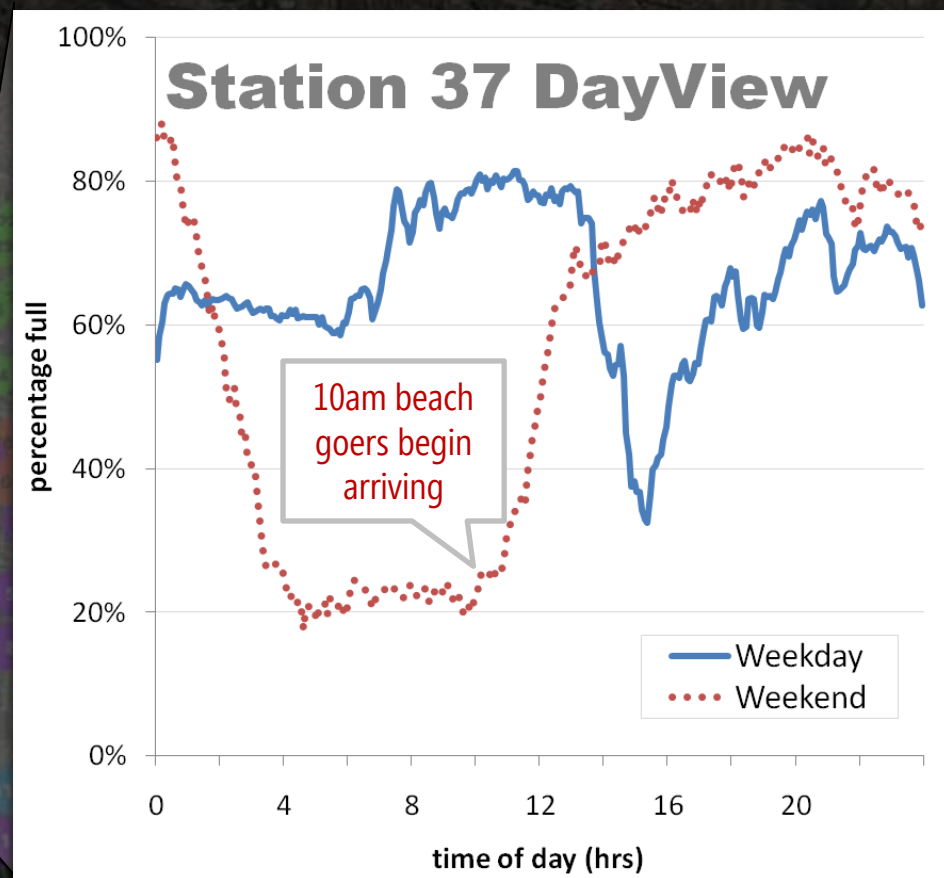


[Lathia, Froehlich & Capra, ICDM2010]

Introducing DayViews





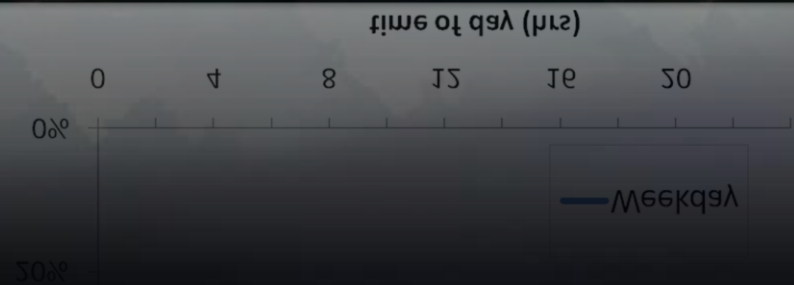
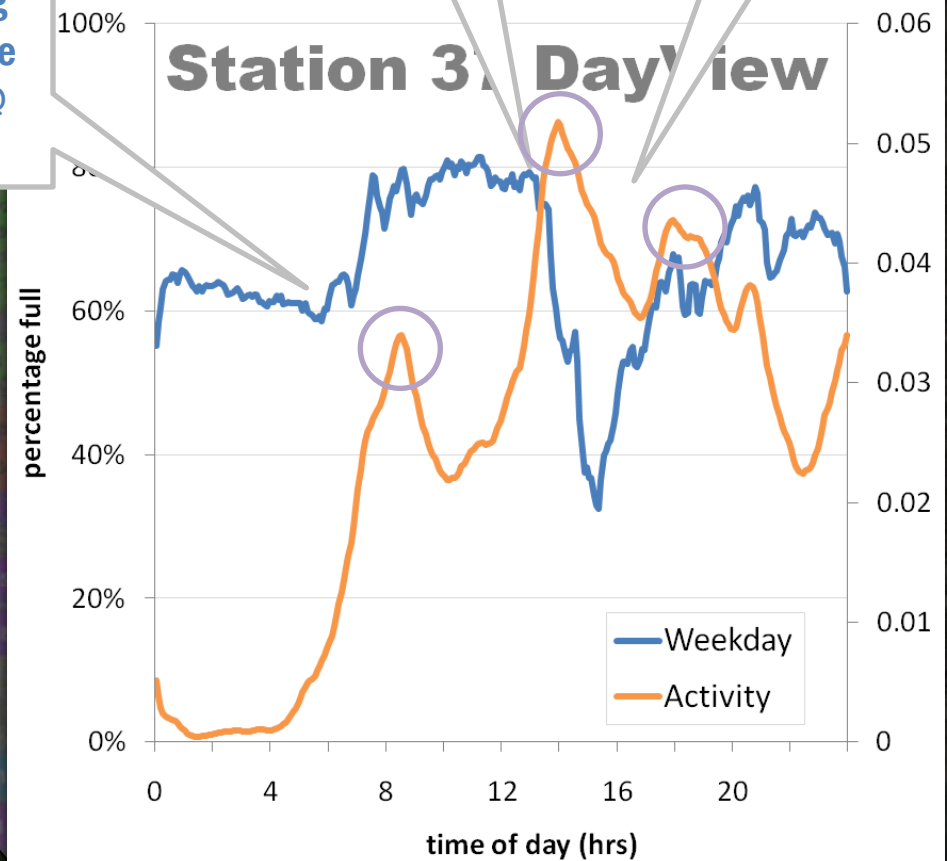


Activity Score:
 $AS(t) = |B_t - B_{t-1}|$

morning
commute
starts @
7am

lunch
rush
begins @
1pm

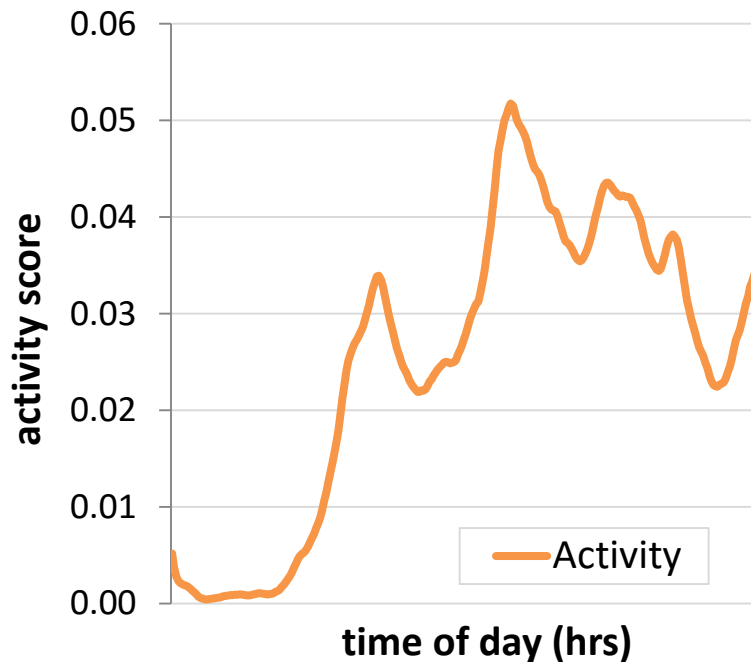
return from
lunch/work
at 3:30pm



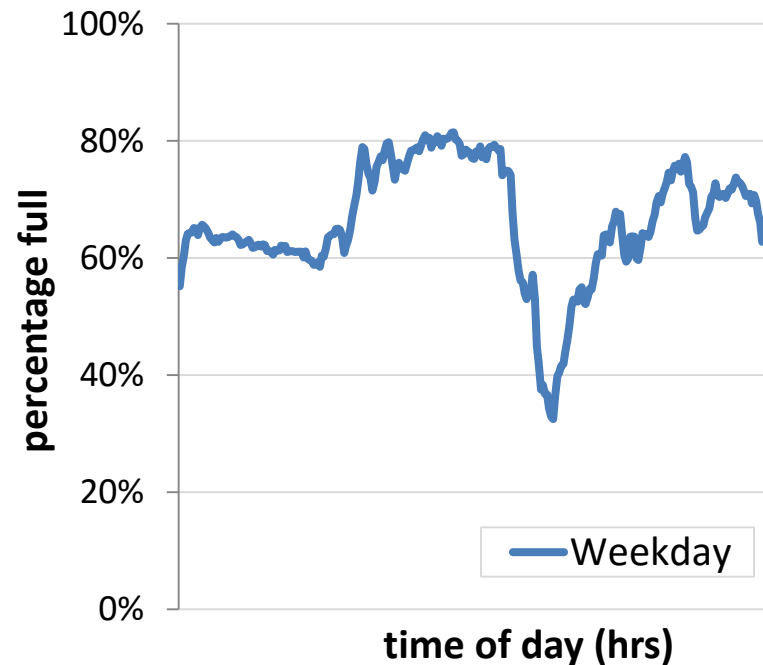
How are Bicing patterns shared across stations and distributed in the city?

Temporal Clustering

activity clusters

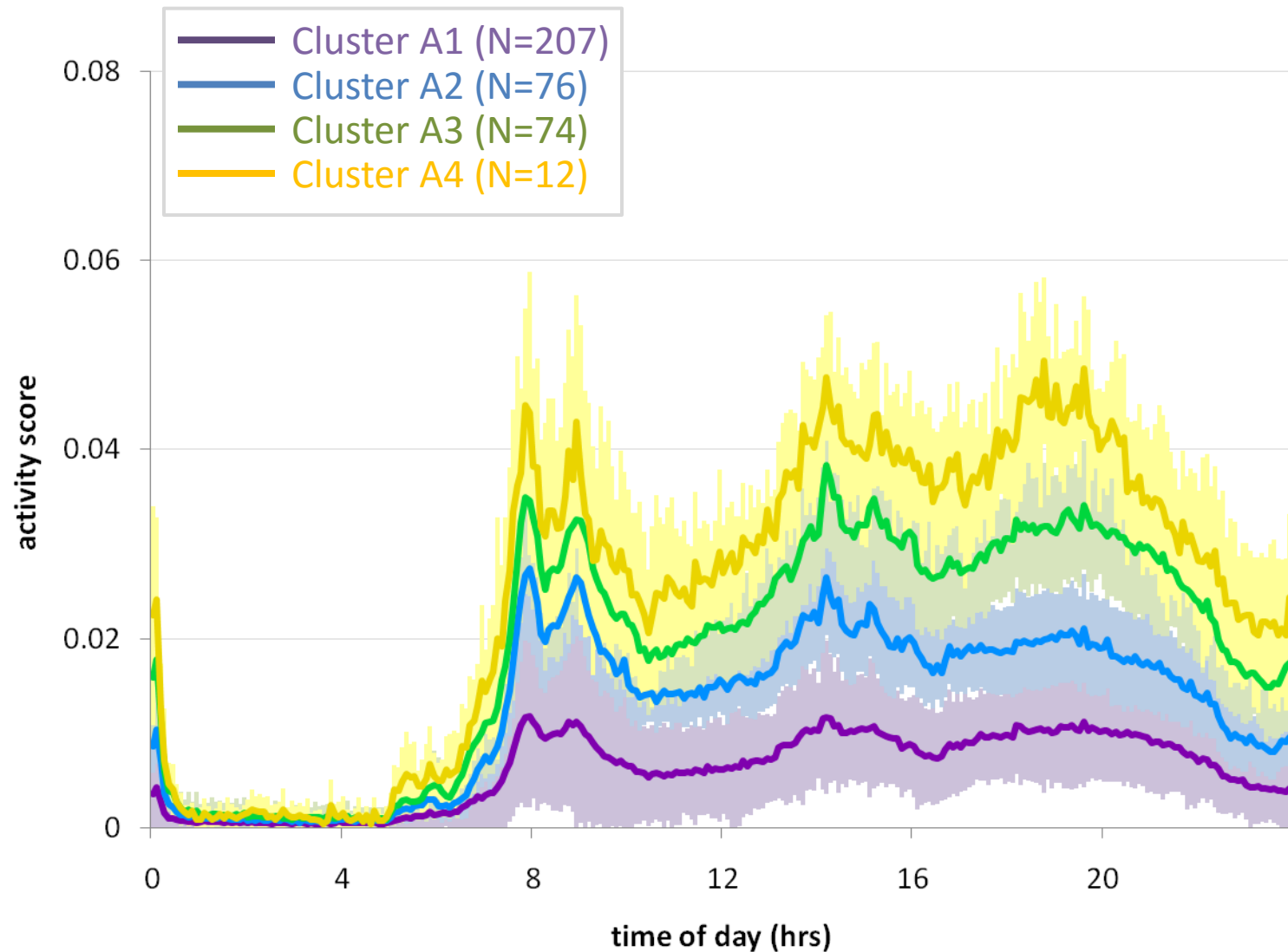


available bicycle clusters

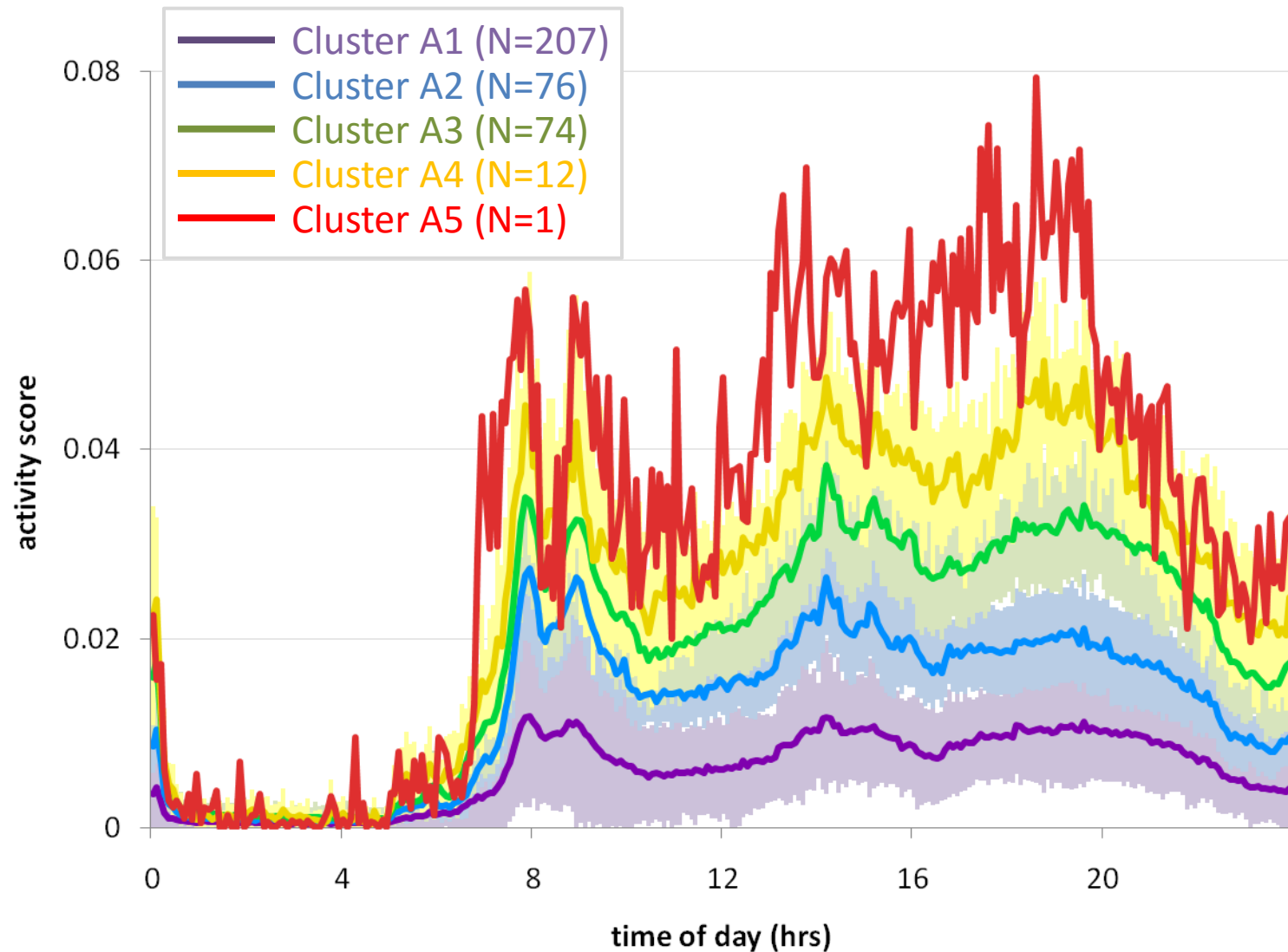


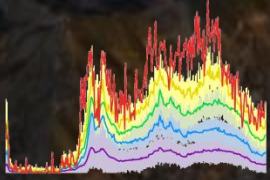
Applied dendrogram clustering with
dynamic time warping as distance metric

Activity Clusters



Activity Clusters

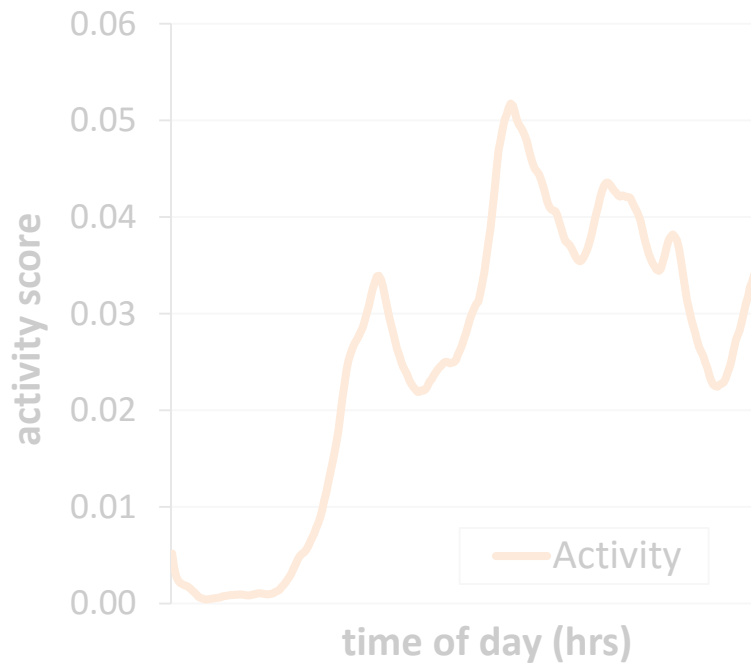




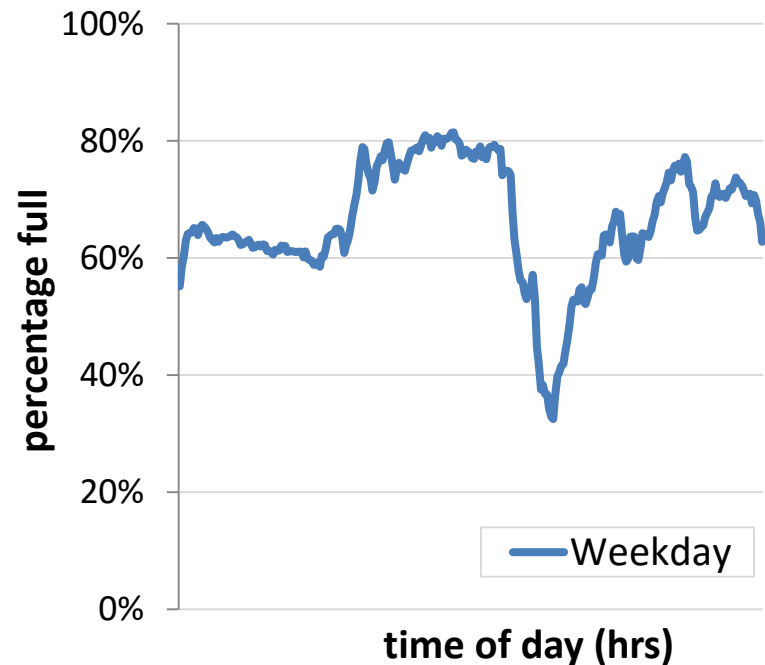
2724 m

Temporal Clustering

activity clusters



available bicycle clusters

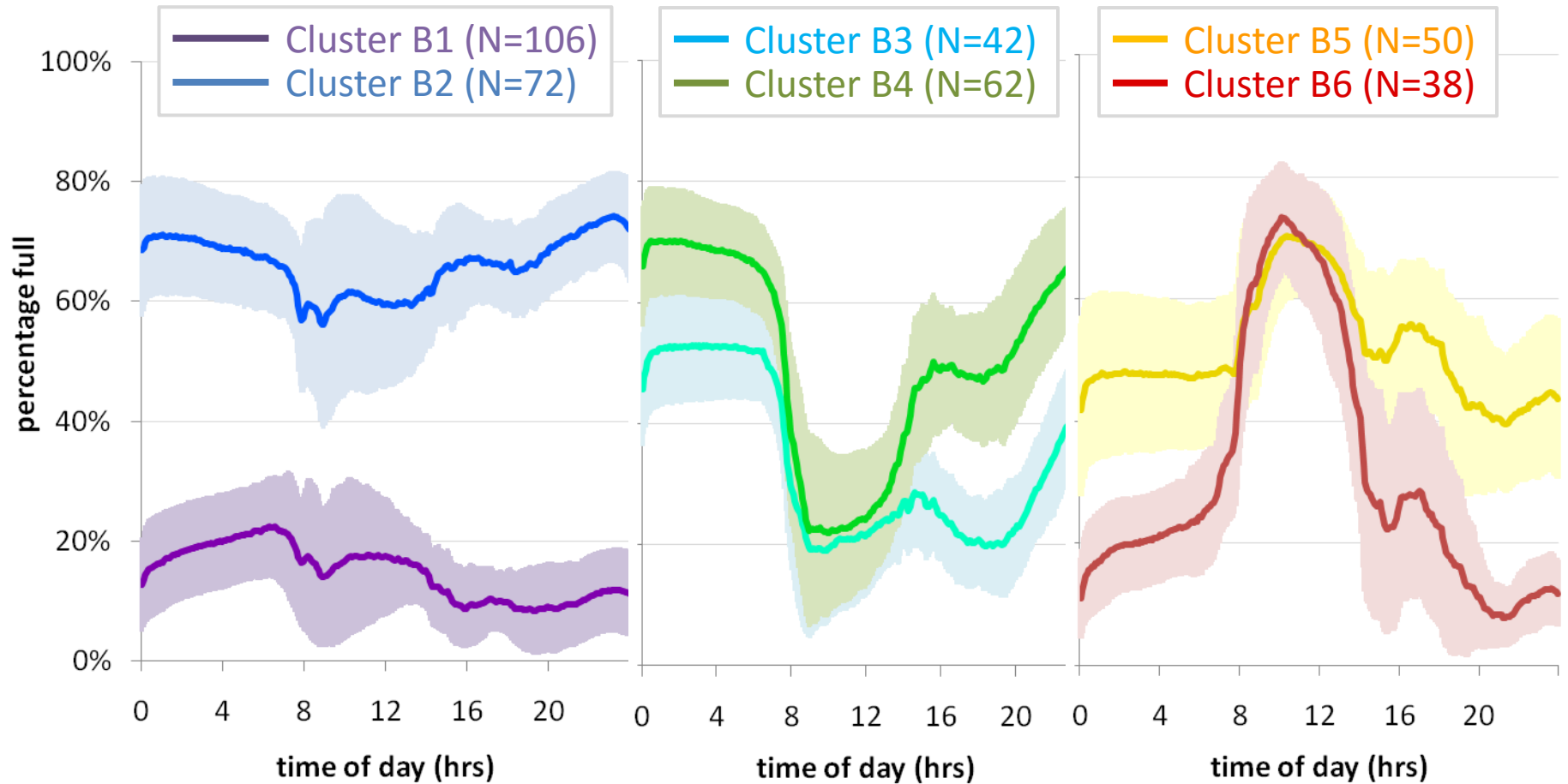


Available Bicycle Cluster

flat

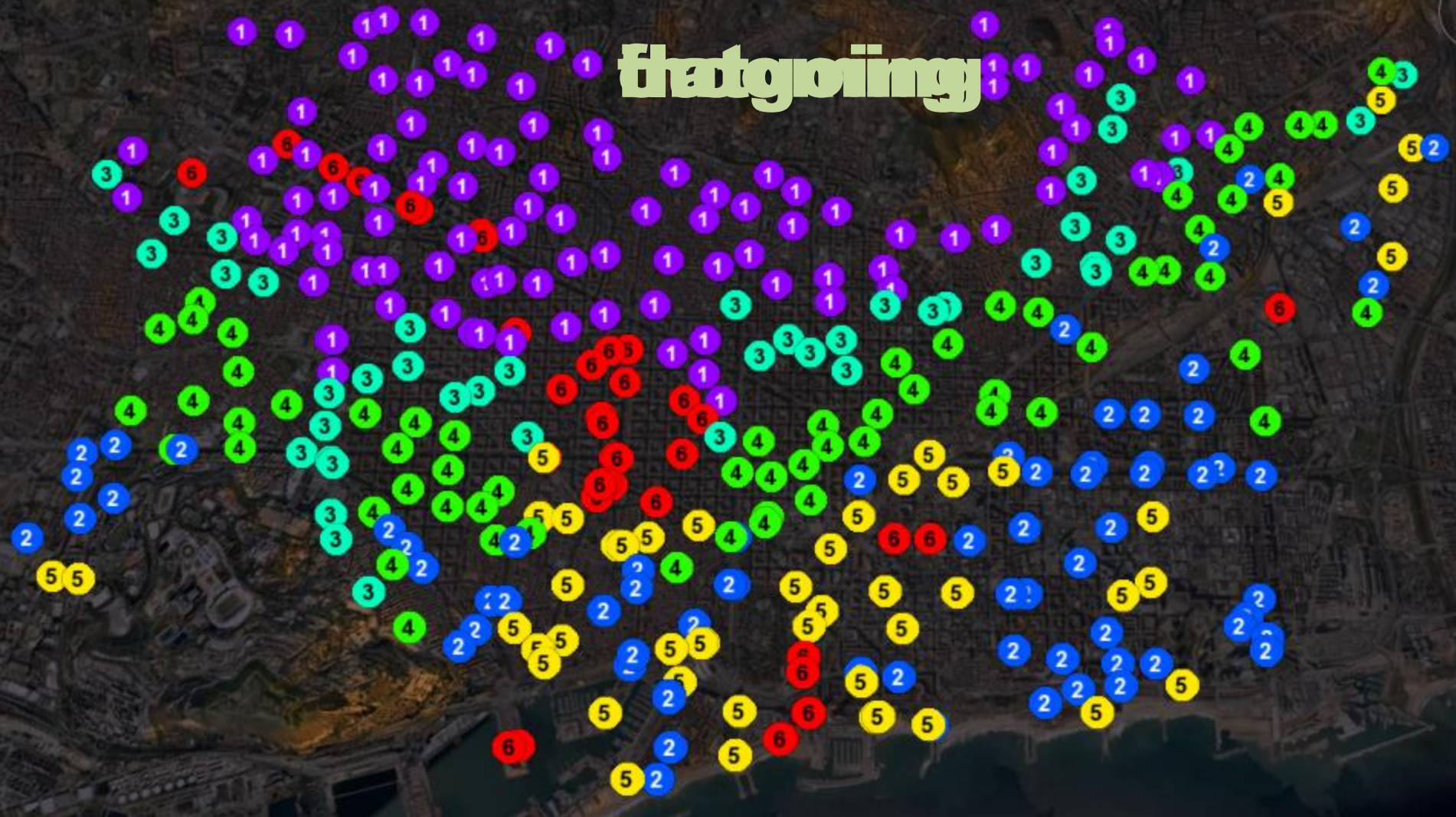
outgoing

incoming





thatology



2683 m

Can Bicing station usage be predicted?

Why Care?

- load balancing
- assist urban planners / city officials about expected activity
- provide new web/mobile services to biking users

Uphill Station (midday)



Downtown Station (night)



76% of respondents had difficulty *finding a bicycle*

Downtown Station (morning)



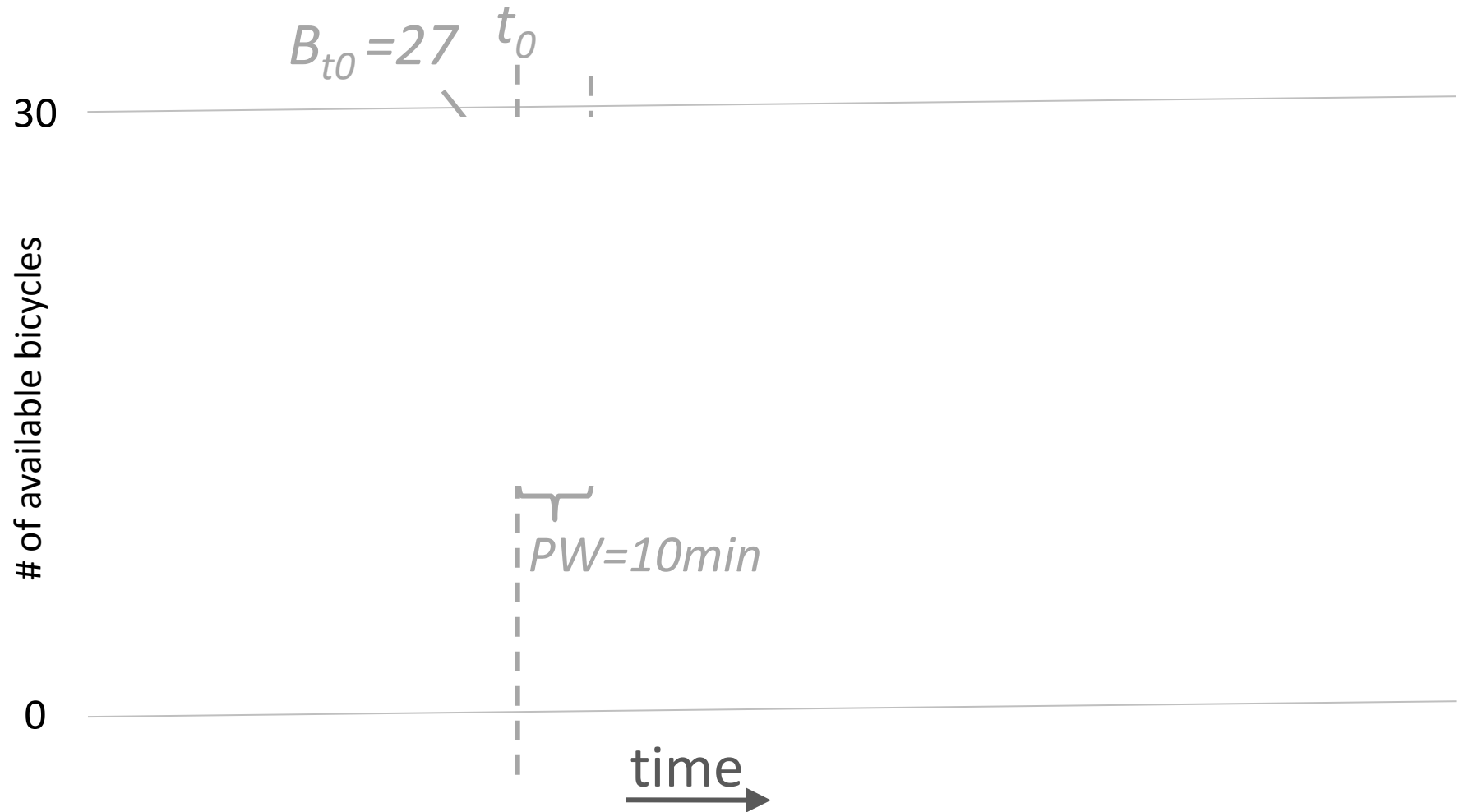
66% of respondents had difficulty *finding a parking slot*

Beach Station (evening)



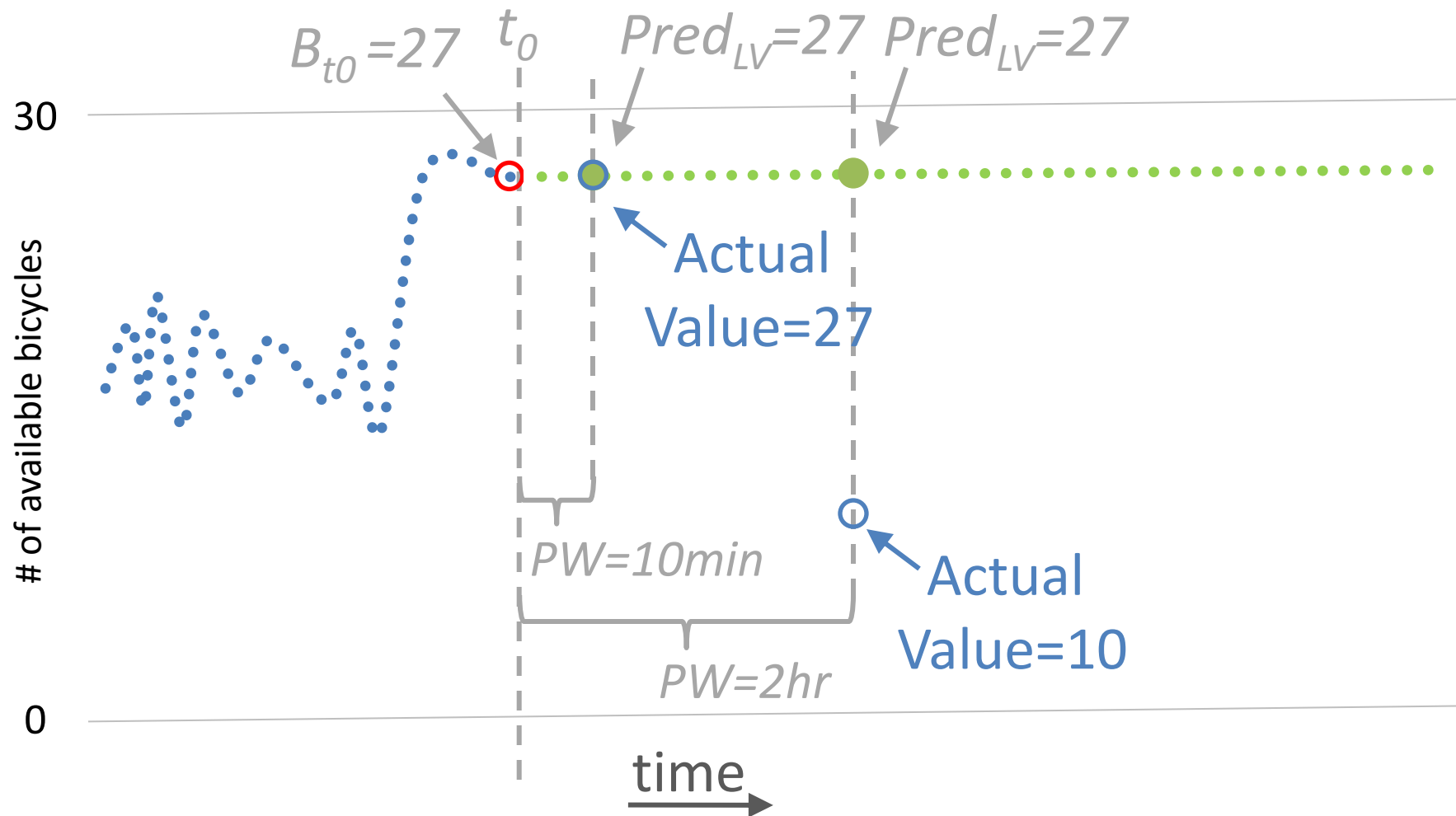
50% of respondents avoid Bicing when they are traveling to a place where they must be on time

Station Models



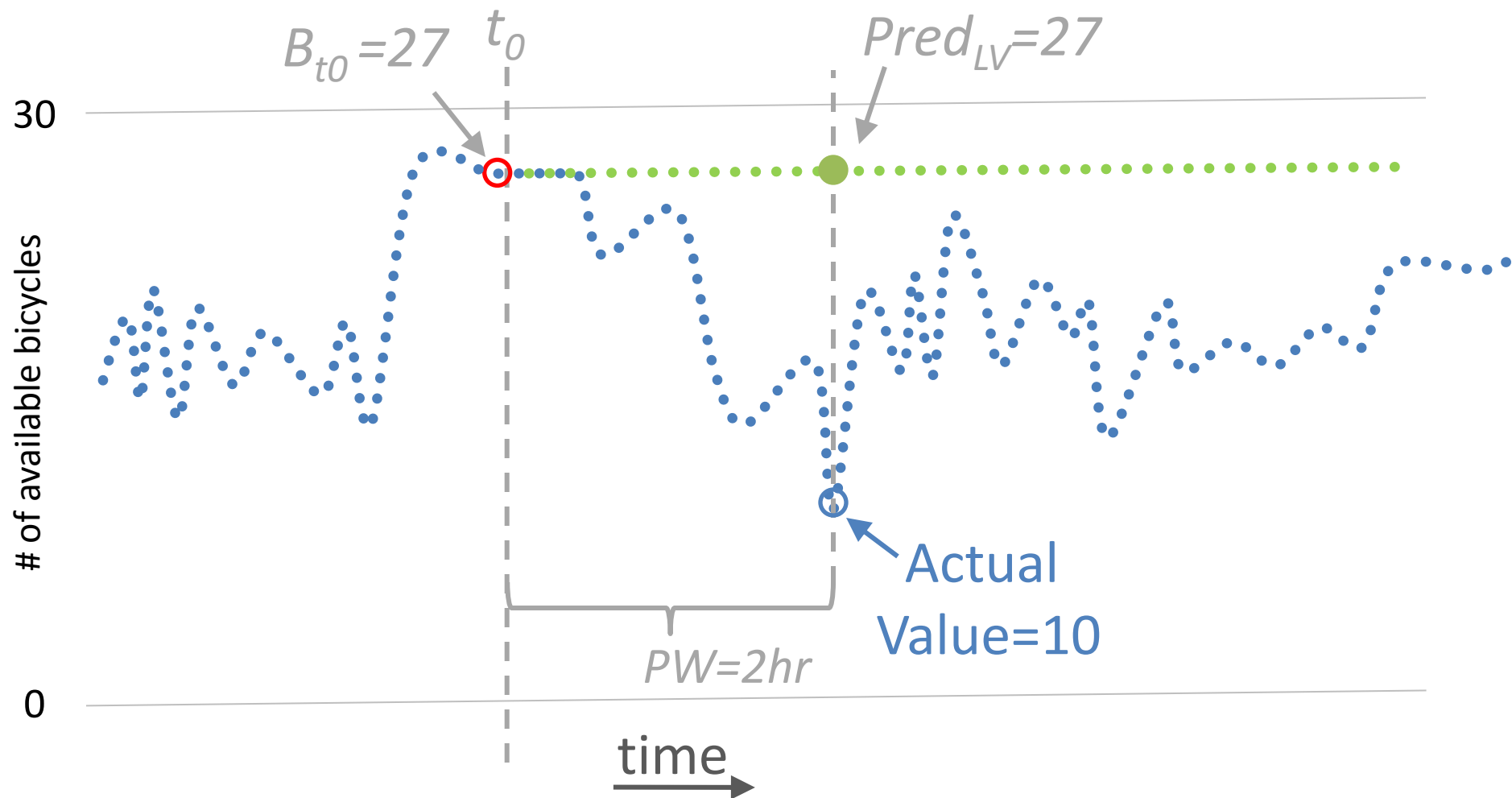
Last Value

$$Pred_{LV} = (t_0, B_{t_0}, PW) = B_{t_0}$$



Last Value

$$Pred_{LV} = (t_0, B_{t_0}, PW) = B_{t_0}$$



Historic Mean

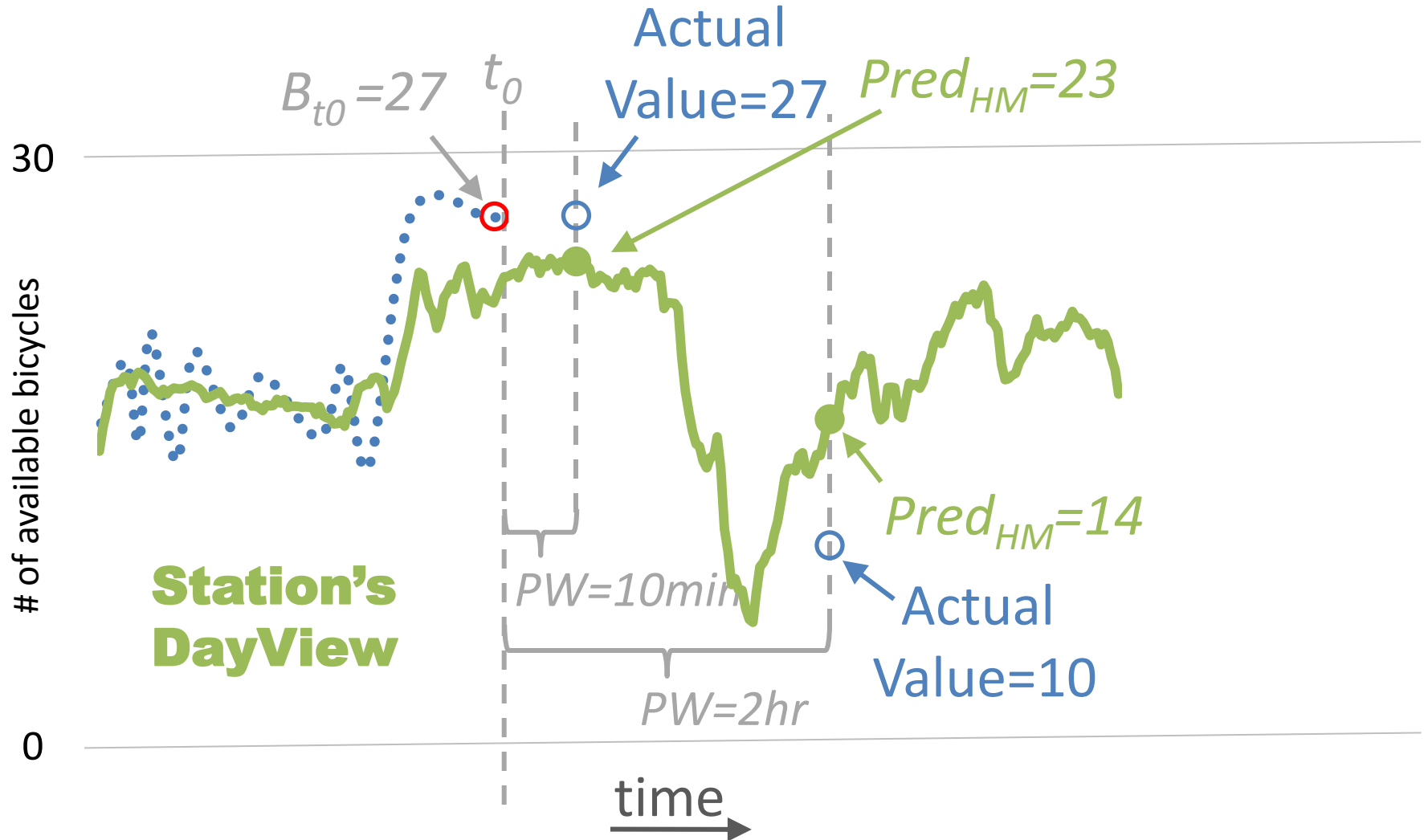
$$Pred_{HM} = (t_0, B_{t_0}, PW) = \overline{B}_{TB_{t_0} + PW}$$



**Station's
DayView**

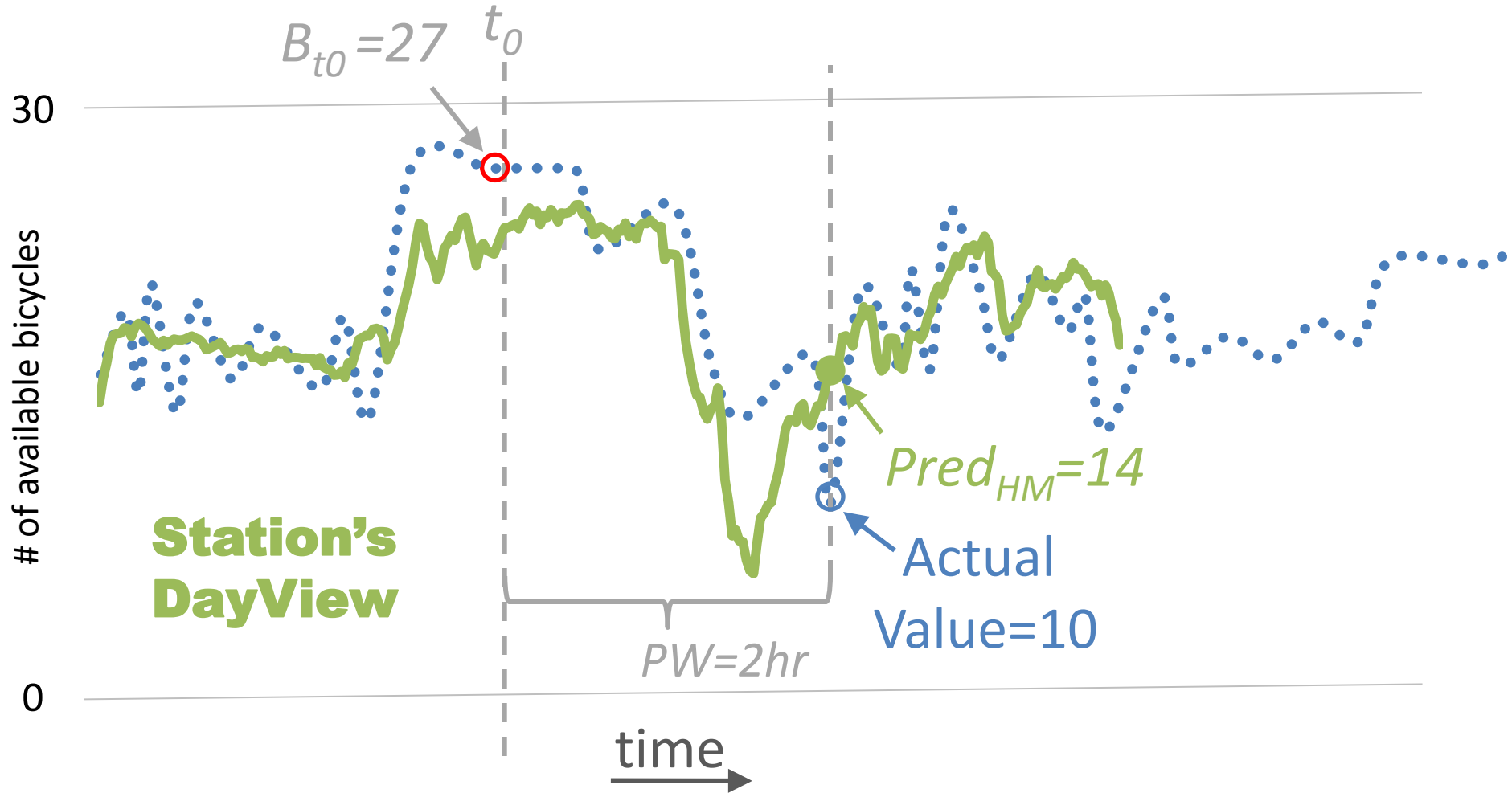
Historic Mean

$$Pred_{HM} = (t_0, B_{t_0}, PW) = \bar{B}_{TB_{t_0} + PW}$$



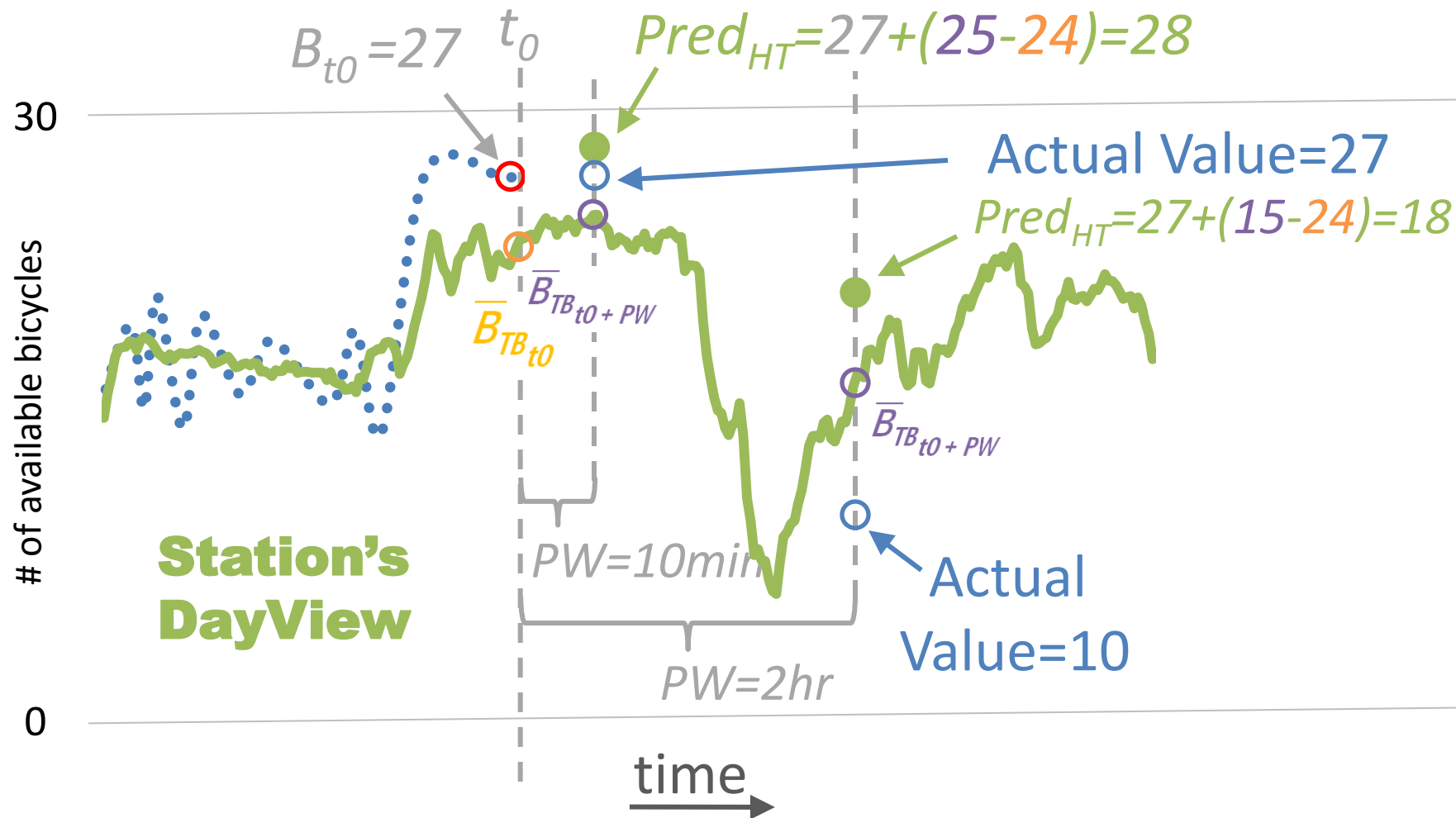
Historic Mean

$$Pred_{HM} = (t_0, B_{t_0}, PW) = \bar{B}_{TB_{t_0} + PW}$$



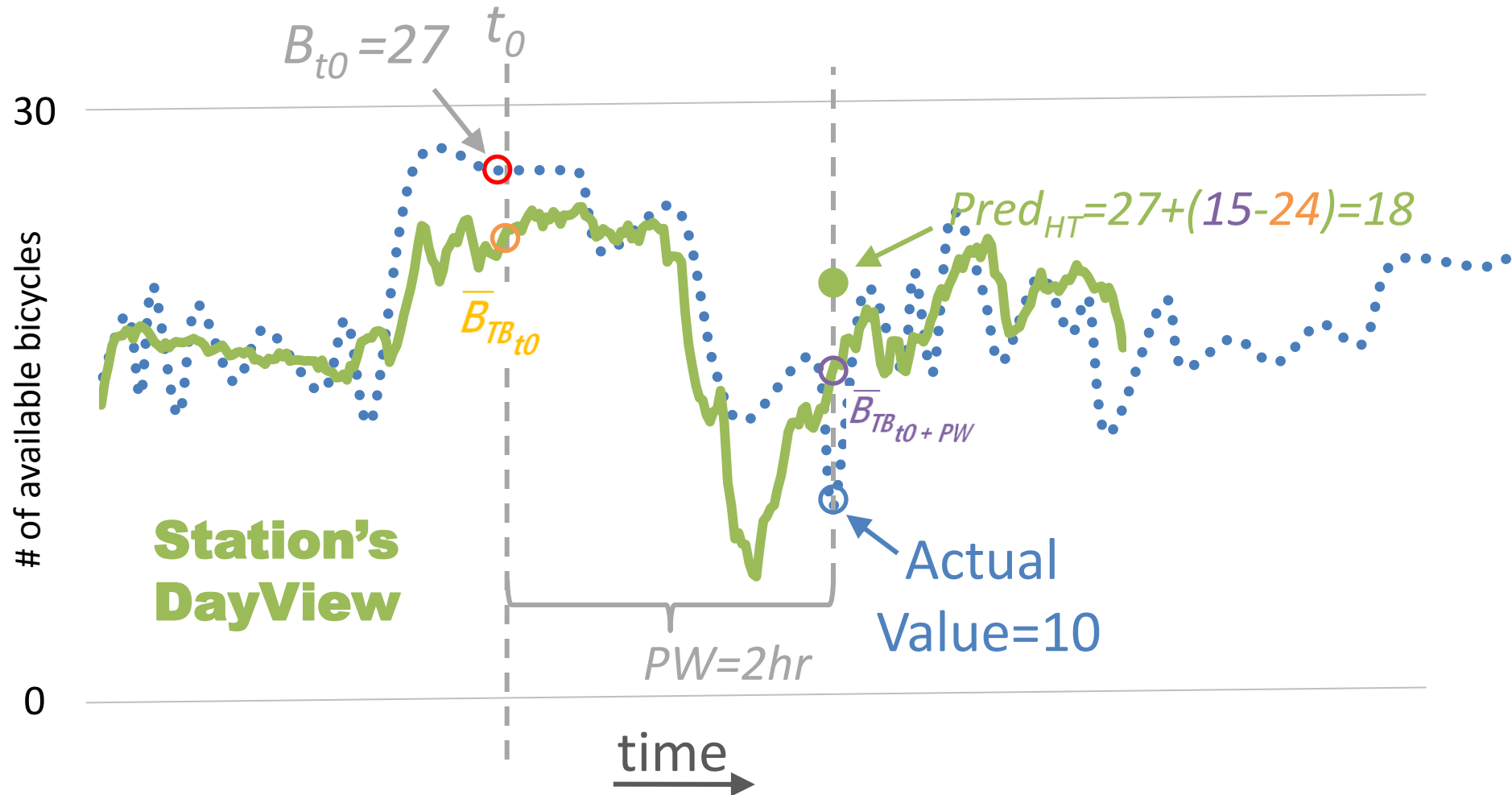
Historic Trend

$$Pred_{HT} = (t_0, B_{t_0}, PW) = \bar{B}_{t_0} + B_{TB_{t_0} + PW} - \bar{B}_{TB_{t_0}}$$



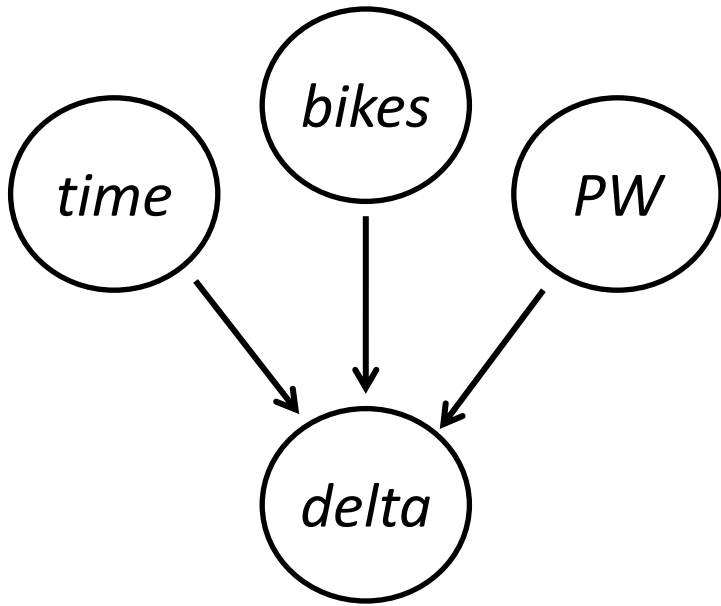
Historic Trend

$$Pred_{HT} = (t_0, B_{t_0}, PW) = \bar{B}_{t_0} + B_{TB_{t_0} + PW} - \bar{B}_{TB_{t_0}}$$



$$Pred_{BN}=(t_0, B_{t_0}, PW)=B_{t_0}+ \text{delta}$$

Bayesian Network



- *time*: discrete observed node corresponding to hours in the day
- *bikes*: the # of avail bikes at time t
- *PW*: the prediction window
- *delta*: continuous Gaussian var that represents change in number of bikes at time $t + PW$

Prediction made by adding the value of the *delta* node to the most recent observation

Prediction Evaluation

november

Sun	Mon	Tue	Wed	Thu	Fri	Sat
26	27	28	29	30	31	Nov 1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29
30	Dec 1	2	3	4	5	6

3 weeks of data to build models

1 week of test data



models were fed:

- the current time
- current # of avail bicycles
- each of the six pw values (10,20,30,60,90,120 mins)

Prediction Error Metric

- Absolute difference between the predicted number of bicycles & the ground truth observation at time $t_0 + PW$
- Error is in number of bicycles
 - normalized by the station's size

High Level Results

Model	Avg Error	Stdev of Error
Random	0.37	0.27

0.37 corresponds to
roughly 9 bicycles

*error is in normalized available bicycles (nab)

High Level Results

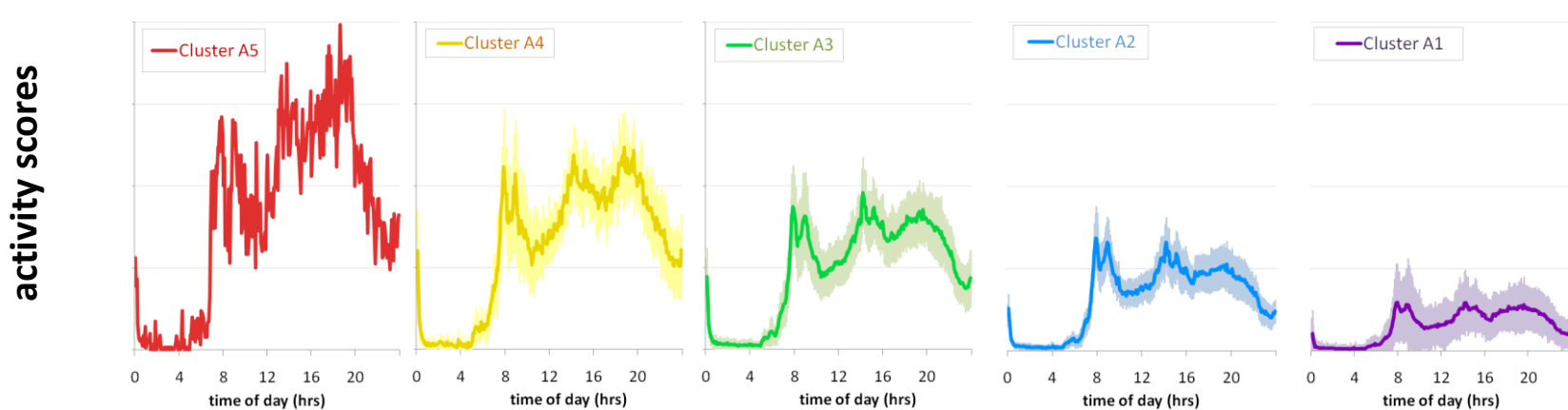
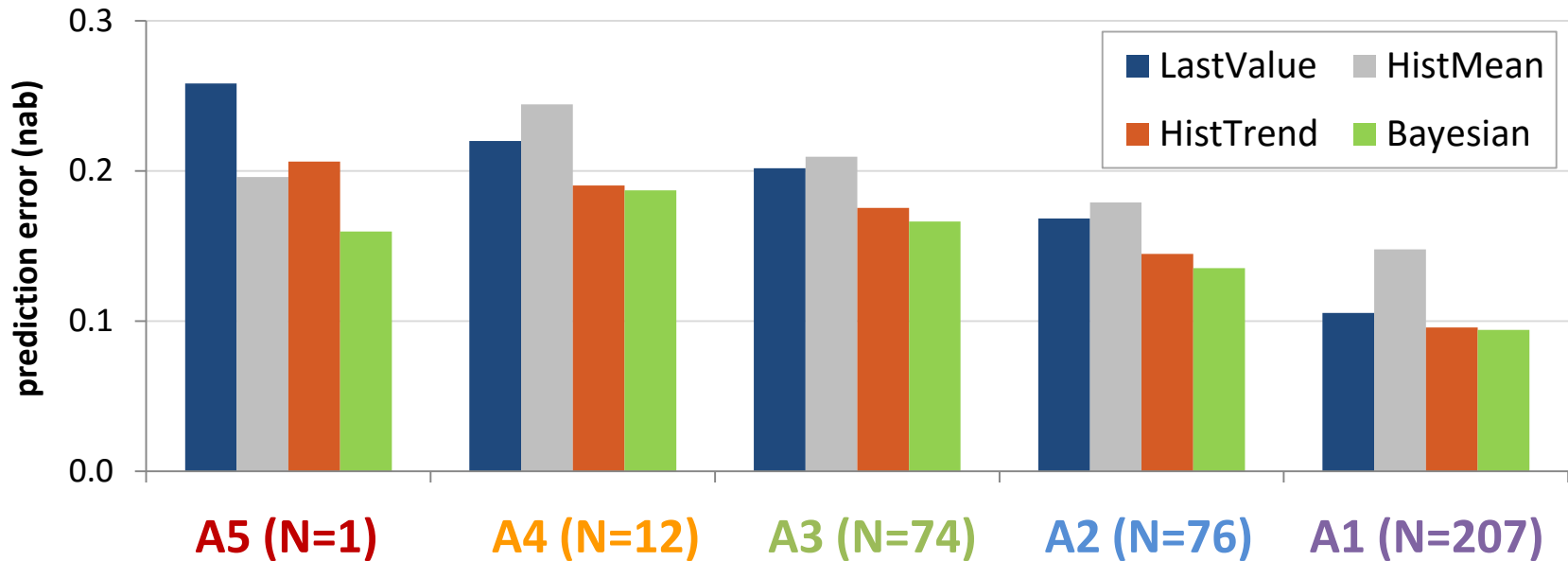
Model	Avg Error	Stdev of Error
Random	0.37	0.27
HistoricMean	0.1	0.1
LastValue	0.09	0.14
HistoricTrend	0.09	0.13
Bayesian Network	0.08	0.12

0.08 corresponds to roughly 2 bicycles

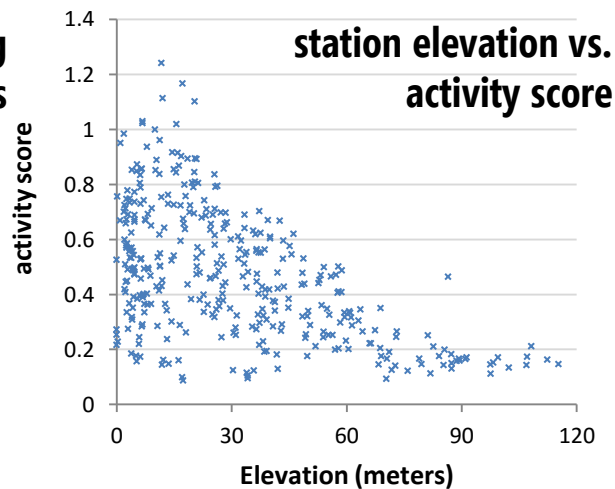
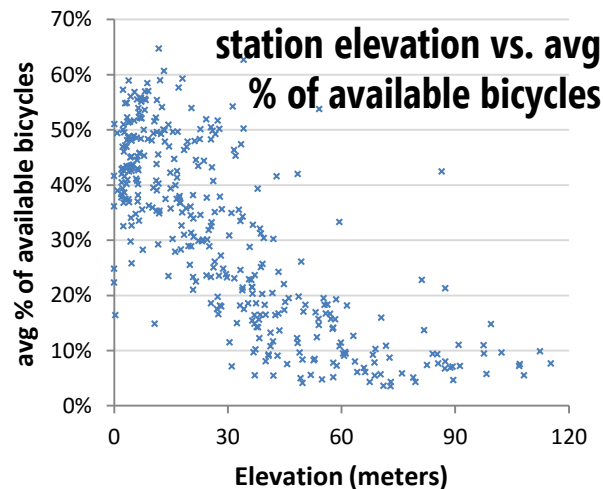
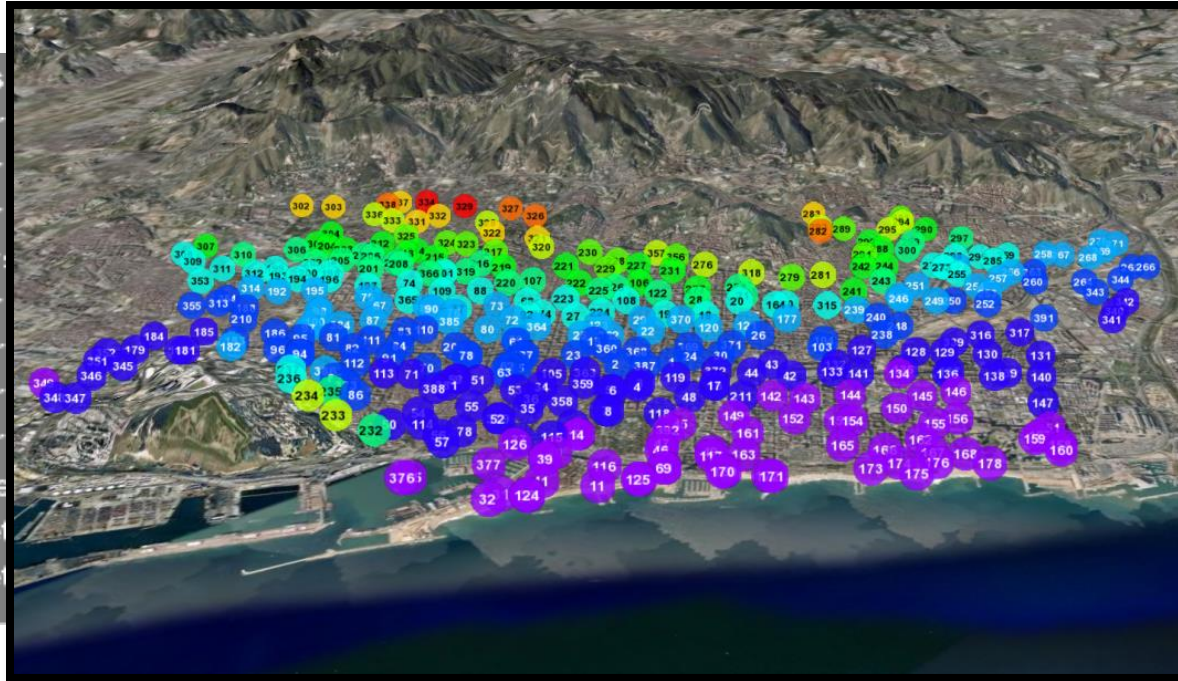
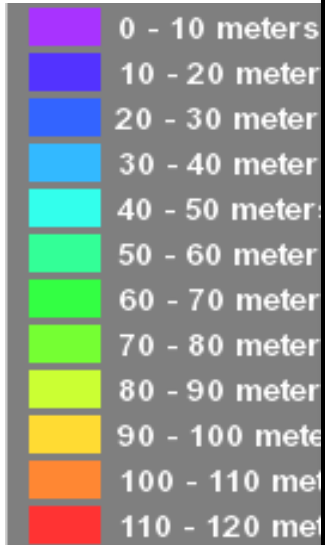
*error is in normalized available bicycles (nab)

Prediction vs. Activity Cluster

0.1 corresponds to ~2.5 bicycles at a station with 25 slots



Topographical Influences?



weather



other transit sources



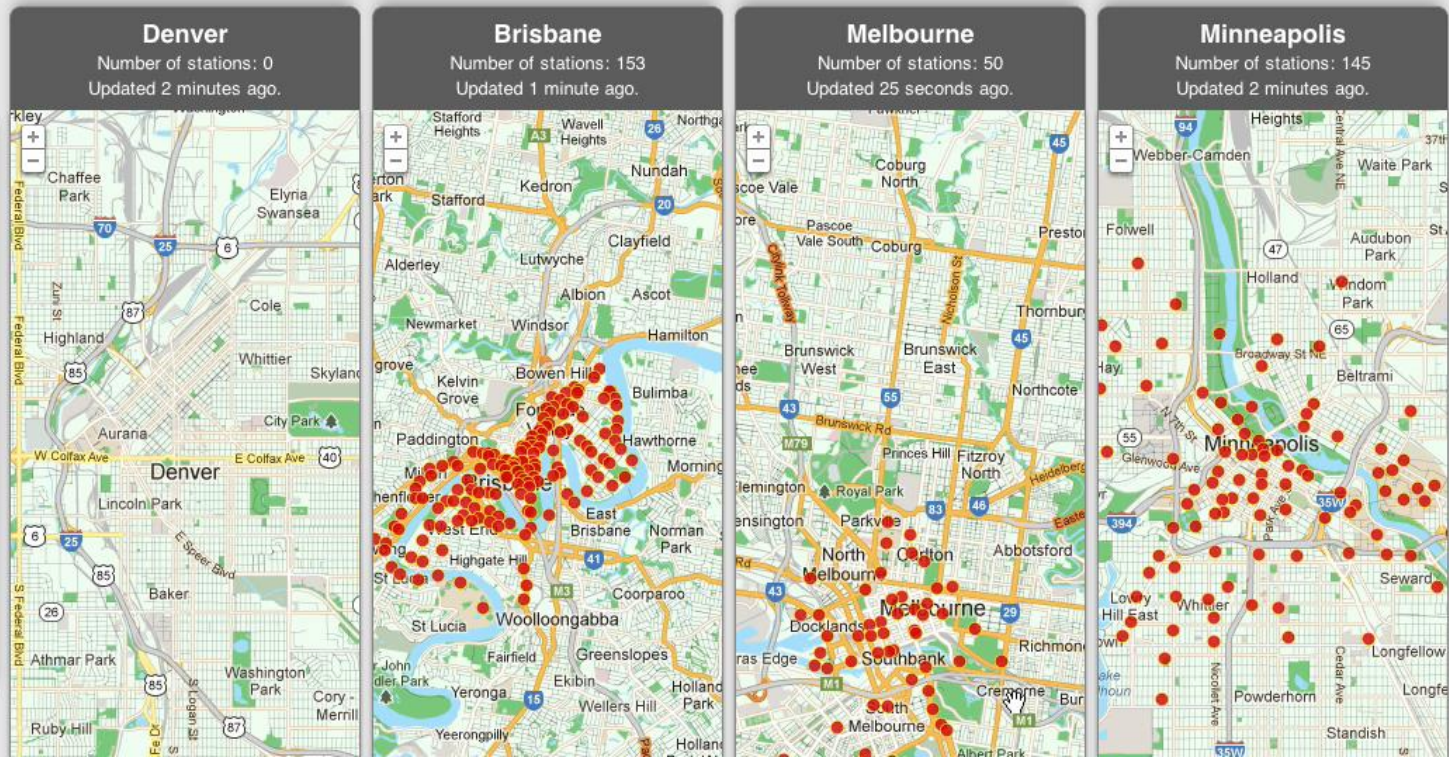
cross city examination

BIKESHARE project



Wednesday, December 5, 2012

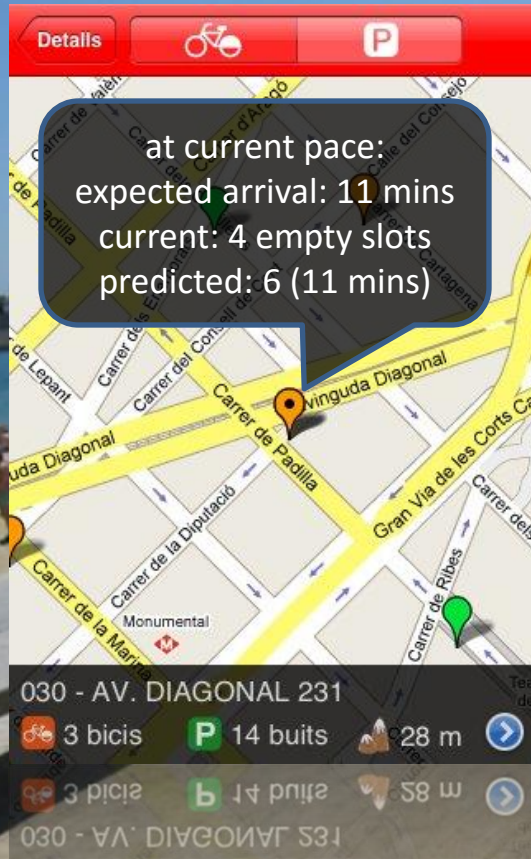
02:26:00 EDT



self-sustainable system



promote usage



Longitudinal Study

Related Collaborators



Related Publications

Individuals Among Commuters: Building Personalised Transport Information Services From Fare Collection Systems

Neal Lathia, Chris Smith, Jon Froehlich, Licia Capra, *Journal of Pervasive and Mobile Computing (PMC) 2012*

Mining Public Transport Usage for Personalised Intelligent Transport Systems

Neal Lathia, Jon Froehlich, Licia Capra, *Proceedings of ICDM2010*

Sensing and Predicting the Pulse of the City Through Shared Bicycling

Jon Froehlich, Joachim Neumann, Nuria Oliver, *Proceedings of IJCAI2009*

Measuring the Pulse of the City Through Shared Bicycle Programs

Jon Froehlich, Joachim Neumann, Nuria Oliver, *Proceedings of UrbanSense2008*

Route Prediction From Trip Observations

Jon Froehlich, John Krumm, *Proceedings of SAE2008*

Download publications here: <http://www.cs.umd.edu/~jonf/publications.html>

SENSING AND PREDICTING THE PULSE OF THE CITY THROUGH SHARED BICYCLING

International Workshop on Spatio-temporal Data Mining for a Better Understanding of Human Mobility:
The Bicycle Sharing System Case Study
Co-organized by GERI Animatic, le Labex Futurs Urbains et l'Ecole des Ponts ParisTech
December 5th, Paris, France

