

Bike Sharing in Washington D.C.

Statistical Programming - Python

MBD OCT 2018 - 017 - Group G

Context

2011

1,500 bicycles

165 stations

18,000 members



2012

1,650 bicycles

175 stations

22,200 members



Objectives

- 1. Predict the amount of users on an hourly basis**
- 2. Ensure high level of service and availability**
- 3. Optimize Logistics and Maintenance Teams**

1.

Project Structure

Project Organization

Data Preparation and Features Construction

Based on Exploratory Data Analysis and Machine Learning principles



4

Model and Predictions

Using a Linear Regression algorithm, test the impact of the features on the model score (R^2)



GitHub + GitKraken

Teamwork improved using collaborative developer tools



Machine Learning Process

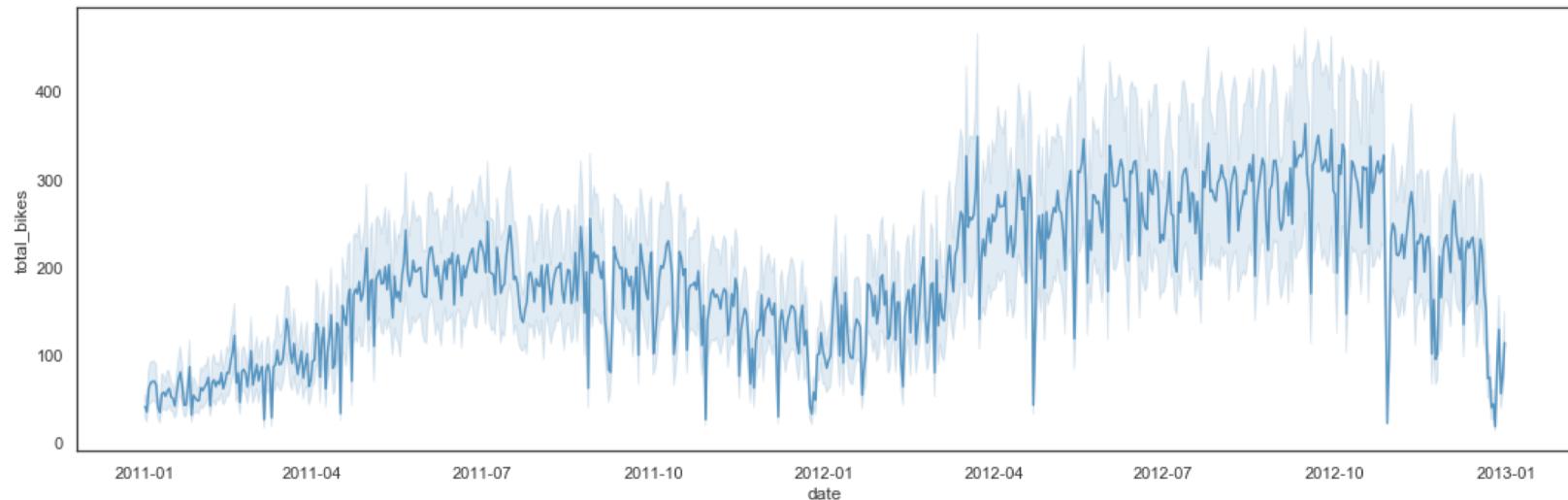
01	EDA and Data Preparation	<ul style="list-style-type: none">• Remove <i>Casual, Registered, Holiday, Feeling Temperature</i>• Scaling, Skewness, Encoding
02	Machine Learning Strategy	<ul style="list-style-type: none">• Train set: Jan 2011 - Jul 2012• Test set: Aug 2012 - Dec 2012• Time Series Cross Validation (10 folds)
03	Feature Engineering	<ul style="list-style-type: none">• Patterns on Dates and Hours• Peak Detection• Exceptional Weather Conditions• Polynomials
04	Selection and Final Metric	<ul style="list-style-type: none">• Recursive Feature Elimination• Manual Selection• Model Predictions vs Reality

2.

Data Exploration Key Insights

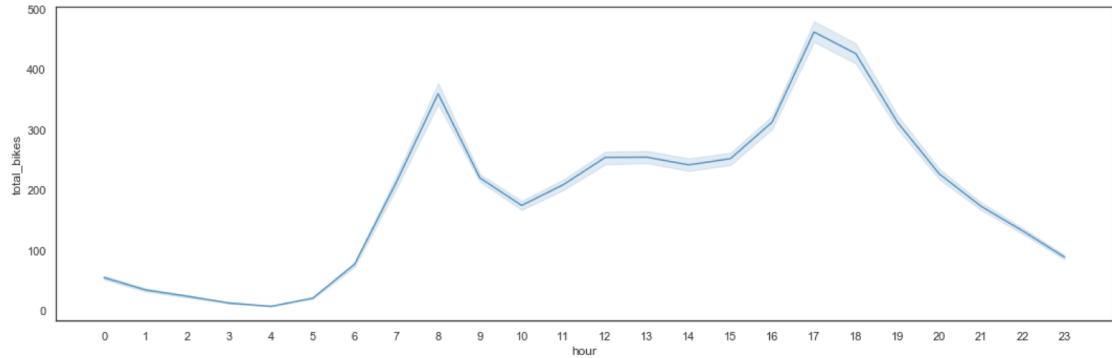
2011-2012 Utilization

Our bike sharing system gets more users every year.



Utilization by Hour

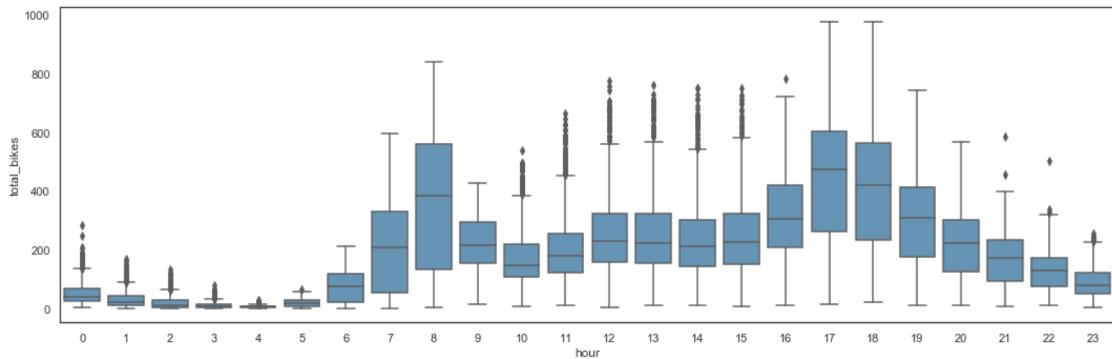
■ Day time usage



■ One peak around 8am

■ One peak between 5-6pm

■ Up to 1000 bikes within an hour

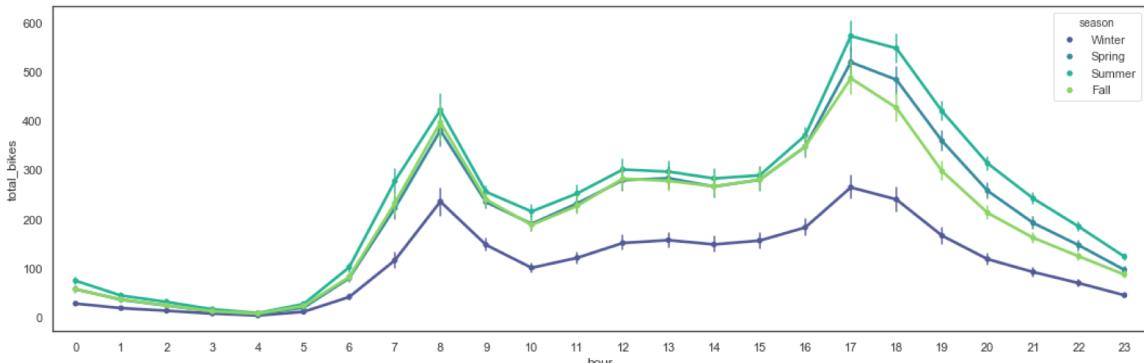
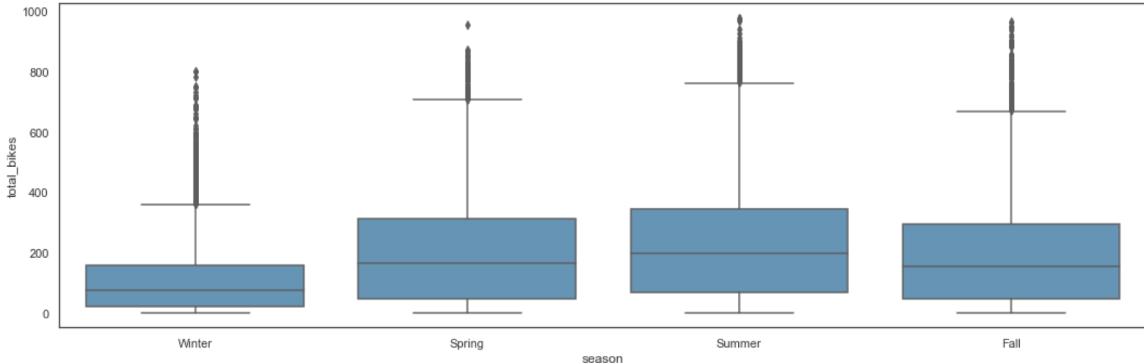


Utilization by Season

■ Summer is the high season

■ Winter is the low season

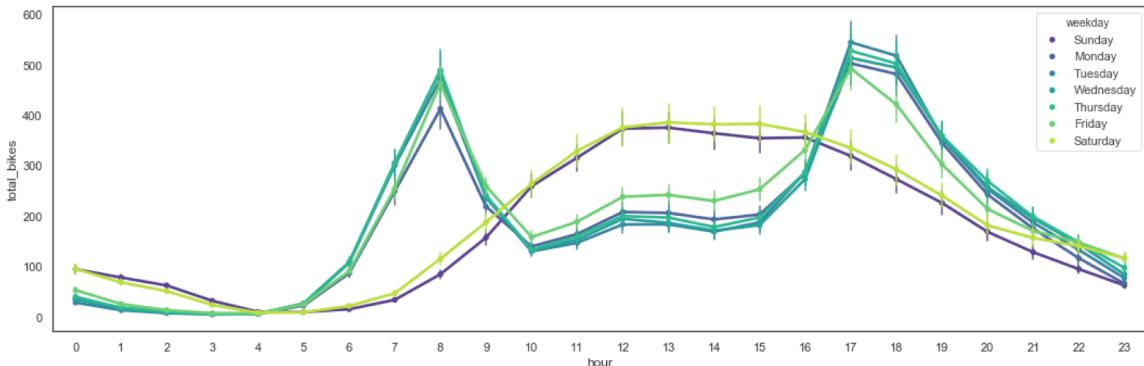
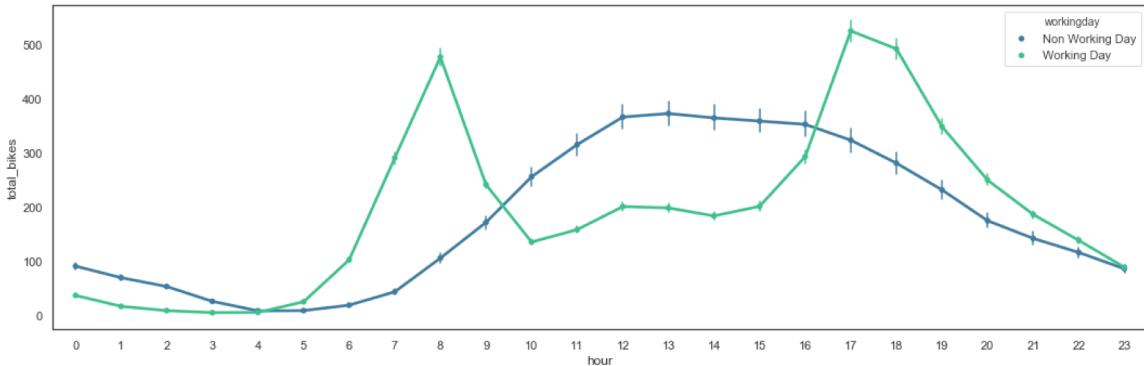
■ Spring and Fall have similar utilization shapes



Overall and Hourly Utilizations by Season

Working Days

- 2 peaks on working days during commuting hours
- No peak during non working days, but higher overall utilization in the afternoon
- Slight change of shape on Fridays, maybe because people leave work earlier on that day

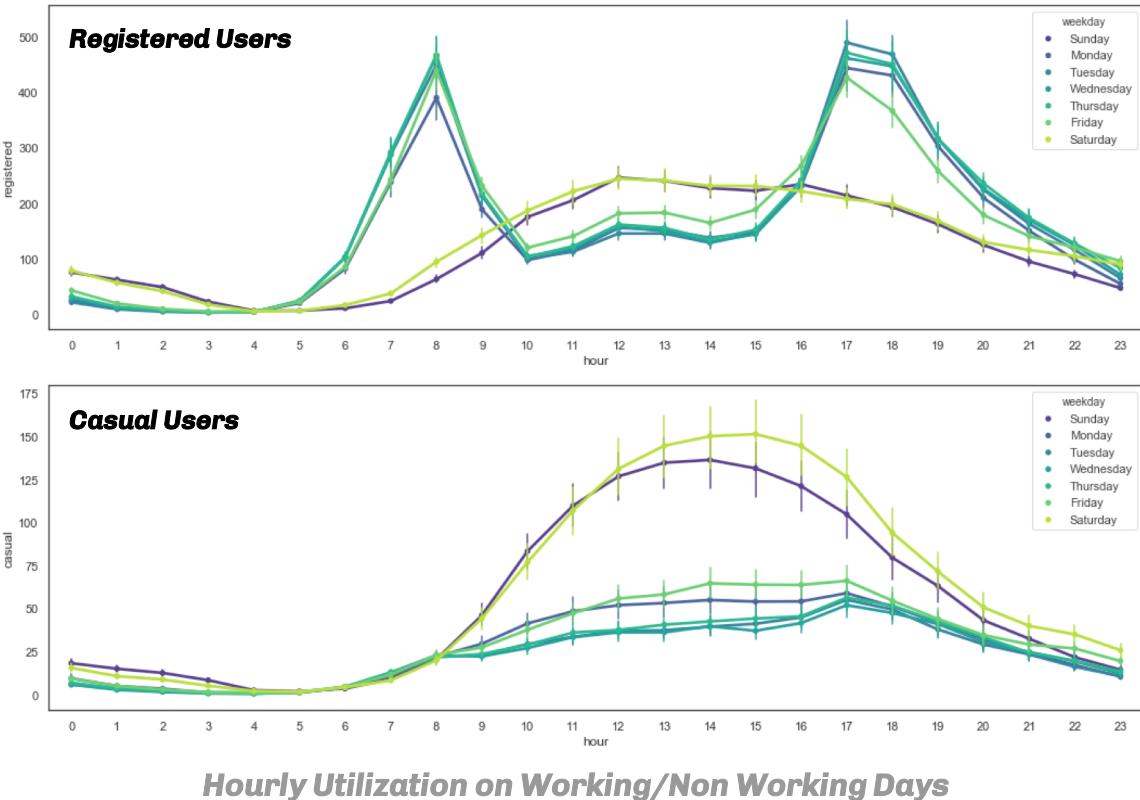


Hourly Utilization on Working/Non Working Days

Working Days

■ Clear difference in behaviours between our registered users and the casual users

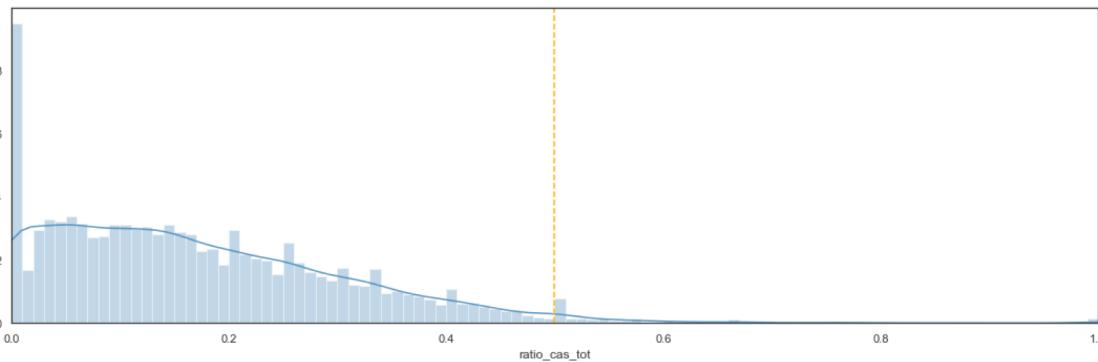
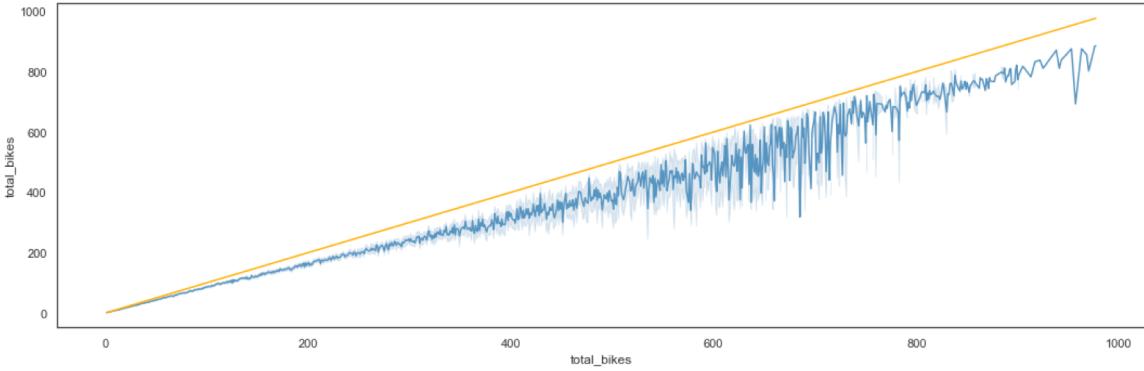
■ Commuting and Leisure effects



Utilization by User Type

- Most users are registered
- High correlation with the Total number of bikes

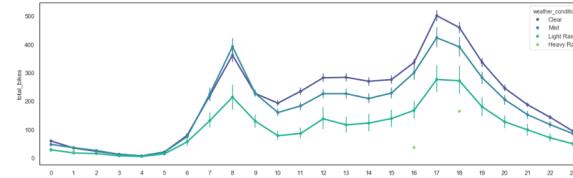
→ *Casual* and *Registered* users information removed from the dataset



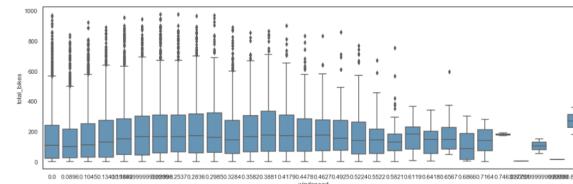
Ratio of Registered Users

Weather Conditions

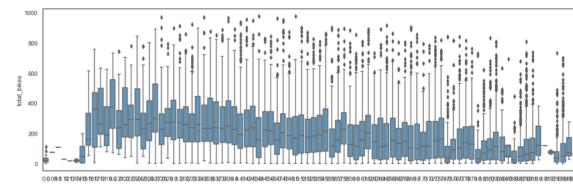
- Weather conditions have a small impact on the service utilization
- Rain has the clearest effect
- Strong Wind discourages users
- Humidity and Temperature seem to have less influence



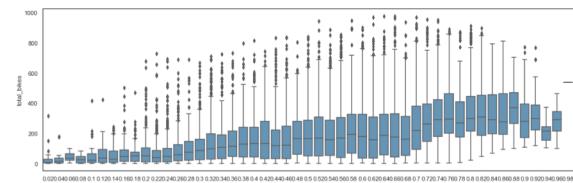
Rain



Wind



Humidity



Temperature

Utilization based on Weather Conditions

Correlations

■ Correlation between Actual and Feeling Temperatures is clear

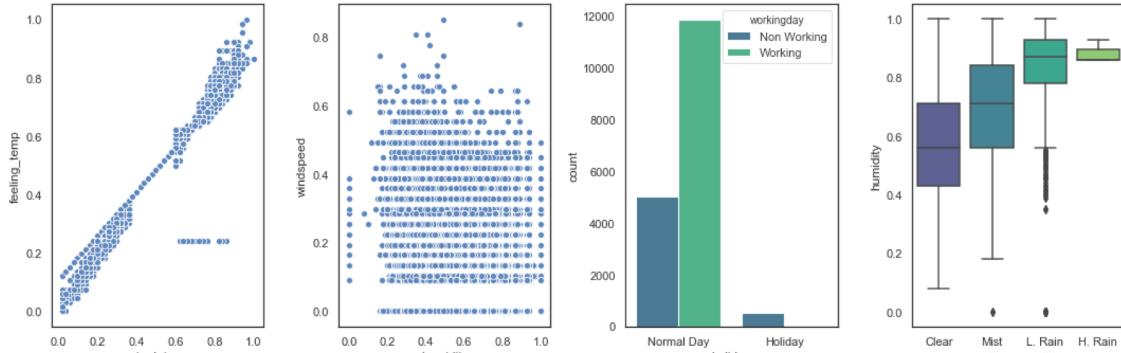
■ No other strong correlation between other variables



Correlation Matrix

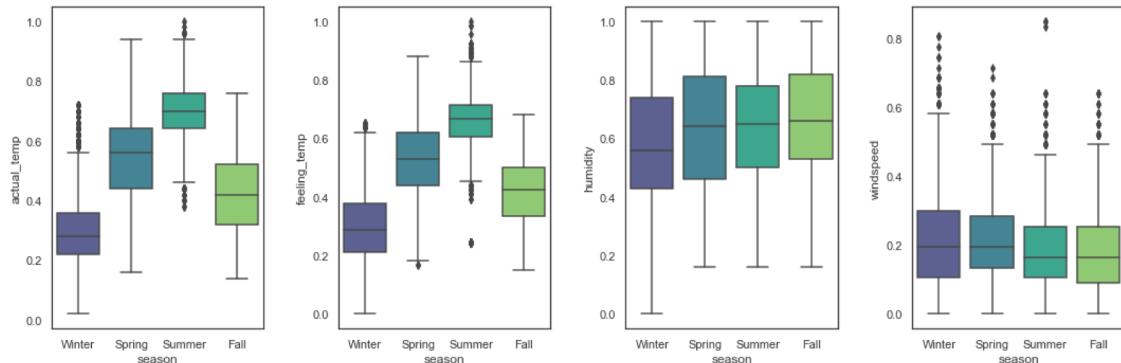
Correlations

Actual and Feeling Temperatures plot is clear



Every Holiday is a Non-Working Day

→ *Feeling Temperature and Holiday information removed from the dataset*



Pair Plots

3.

Features Construction

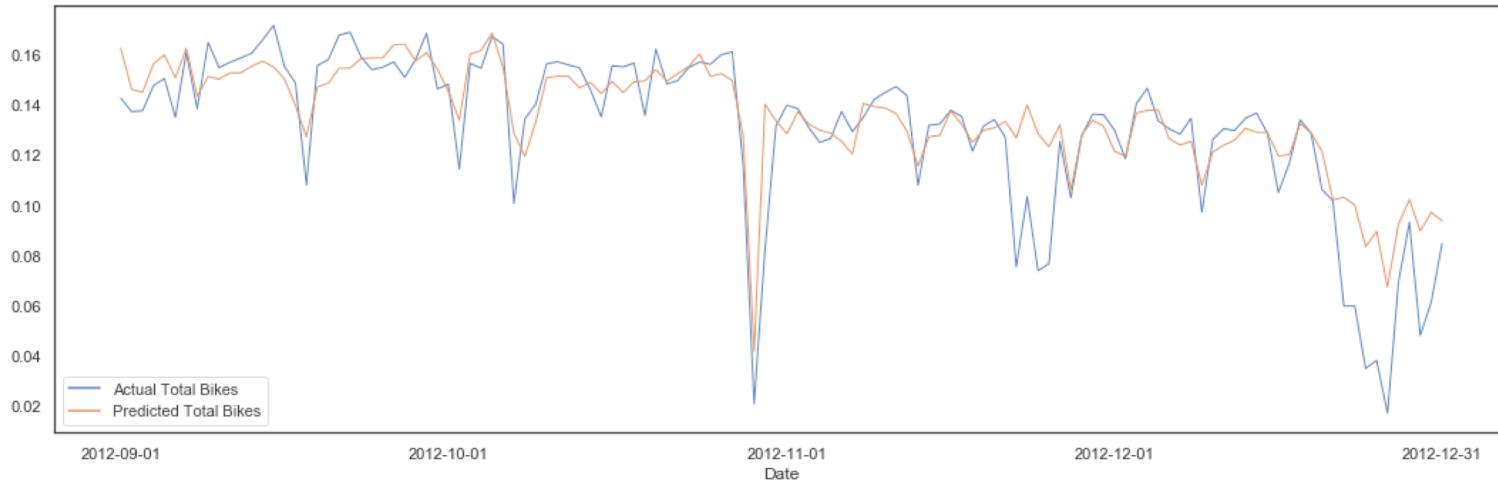
Baseline

Features:

57

R²:

0.76



Baseline Predictions vs Reality

Reminder - Features removed from dataset:

Casual, Registered, Holiday, Feeling Temperature

Calendar Features

R²



—



BASELINE

Without
Casual, Registered,
Holiday, Feeling
Temperature

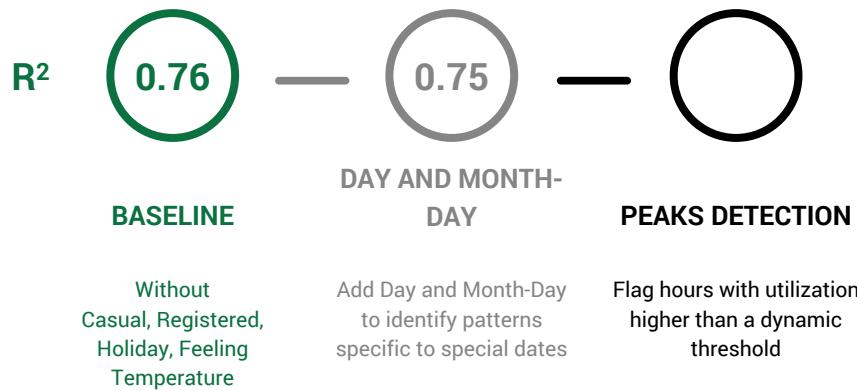
DAY AND MONTH-
DAY

Add Day and Month-Day
to identify patterns
specific to special dates

The screenshot shows a calendar interface for January 2019. At the top left, there's a "TODAY" button. Below it, the date "Tuesday, January 1st" is displayed, followed by "4 Items". A list of events follows: "8:00 AM Team Meeting", "10:00 AM Call Jane", "12:00 PM Lunch with John", and "7:00 PM Dinner with Jane". At the bottom of this section is a button labeled "+ Add new item". To the right of this is a large monthly calendar grid for January 2019, with days from Monday to Sunday. The days are numbered sequentially from 31 to 3, with 1st, 2nd, and 3rd highlighted in red. The 4th is grayed out. The 5th through 31st are in black. Navigation arrows are at the top right of the calendar grid.

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

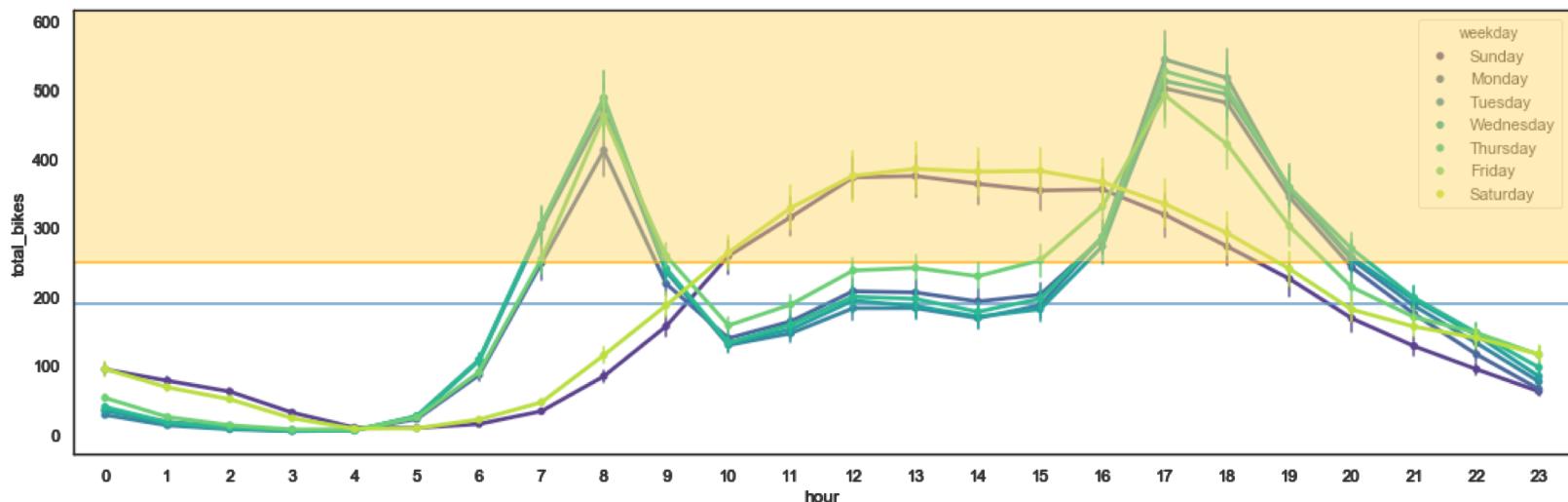
Peaks



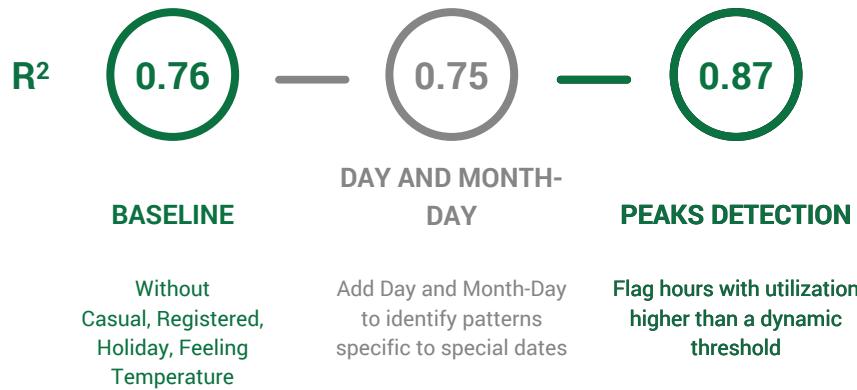
Peaks

(1+ x%)

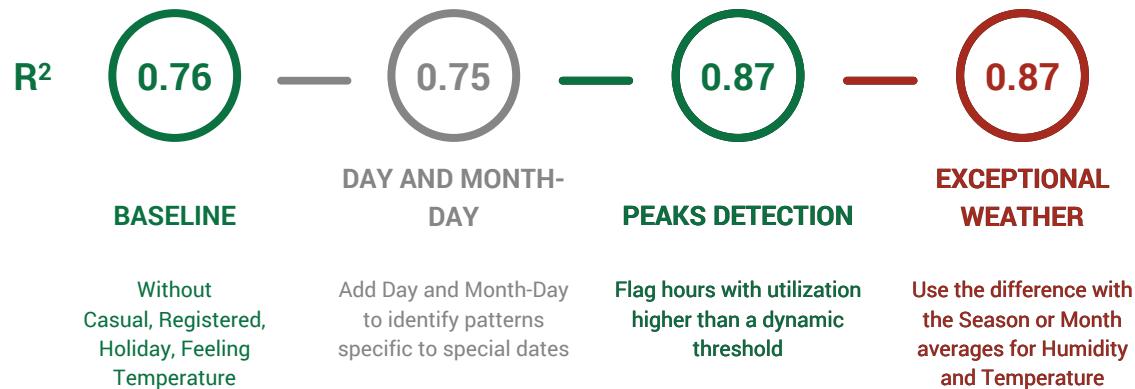
Mean Total Bikes



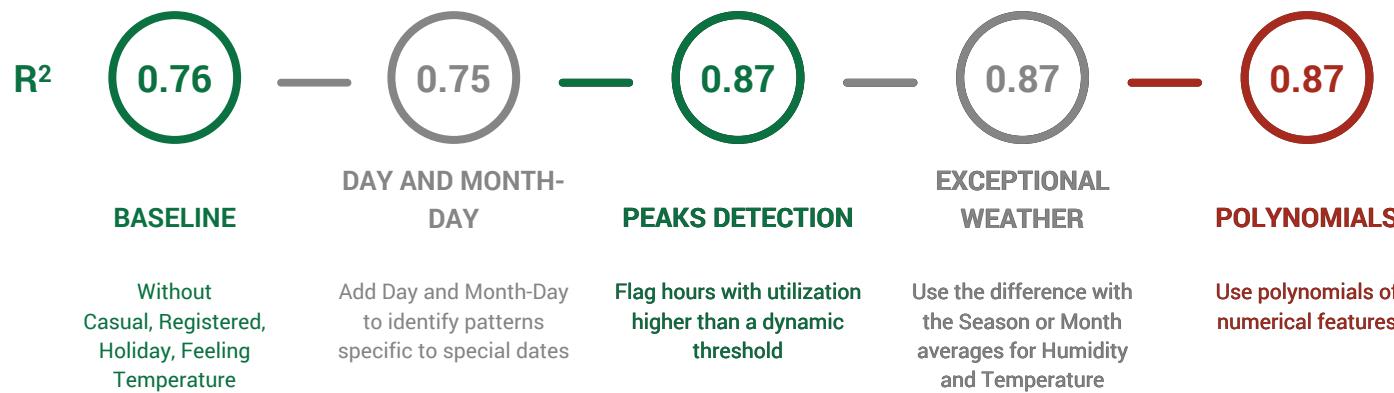
Peaks



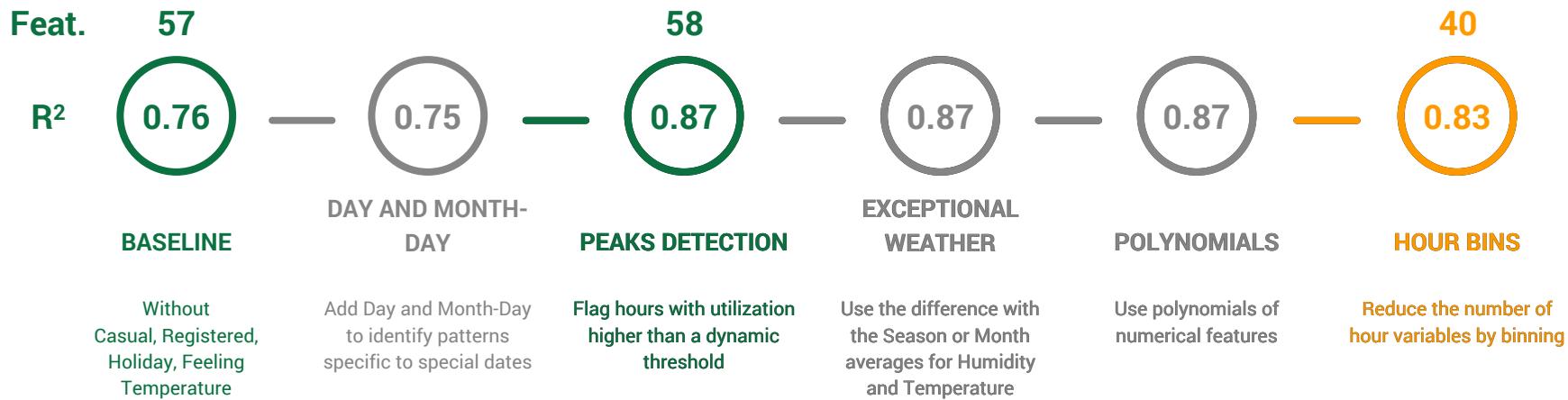
Weather



Polynomials



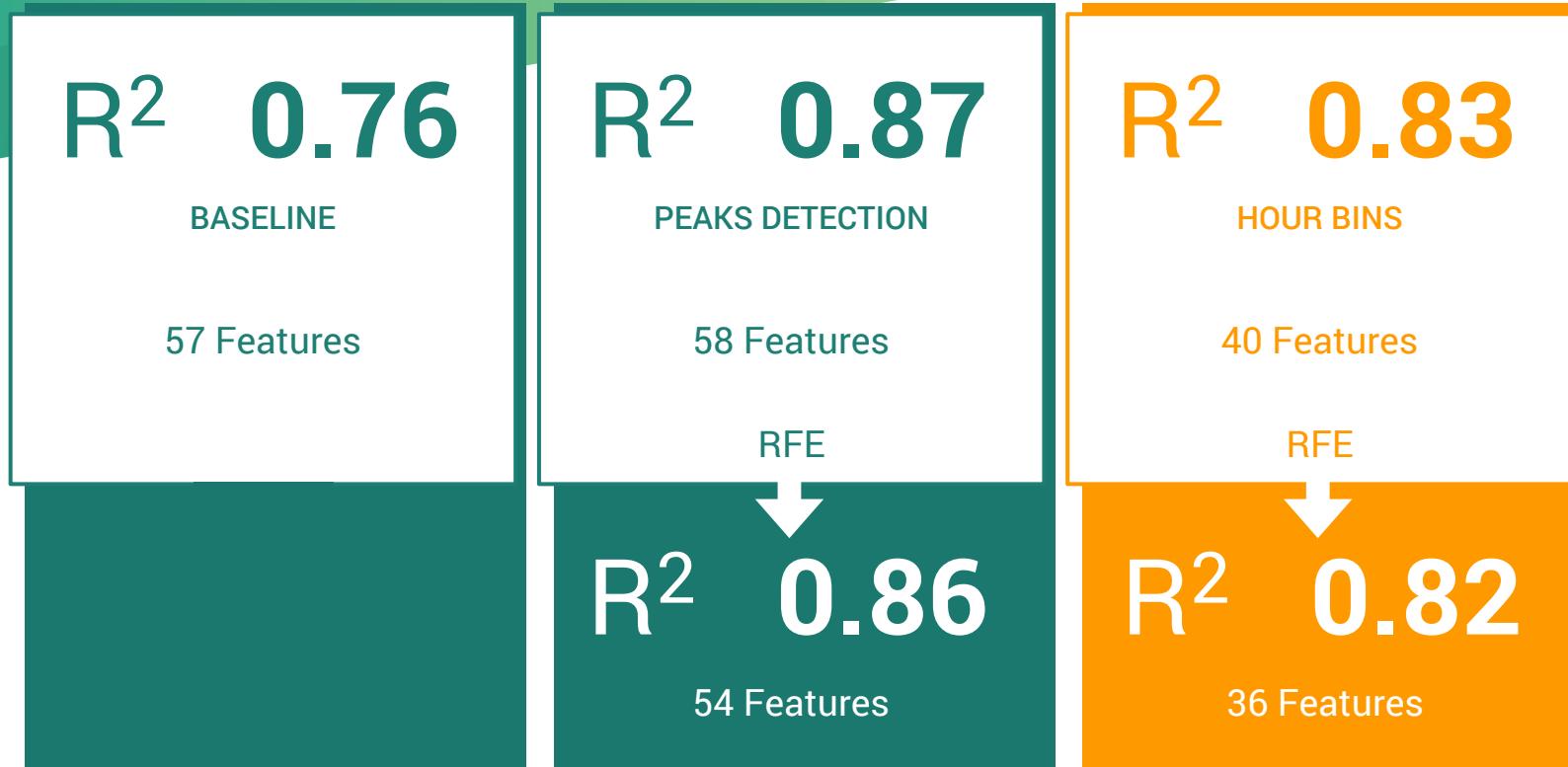
Hour Bins



4.

Model Selection

RFE



4 Features Eliminated:

Humidity | Actual Temperature | Wind Speed | Working Day

Manual Feat. Selection

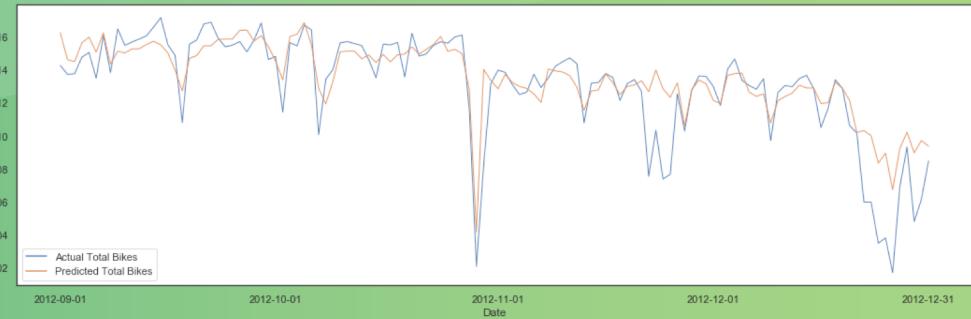
Features Kept	Features Removed
Year	Actual Temperature
Month	Humidity
Days of the Week	Windspeed
Hours	Weather Condition
Peak Detection	Working Day Flag
	Seasons

Features: 46 $R^2: 0.85$

BASELINE

Features: 57 R^2 : 0.76

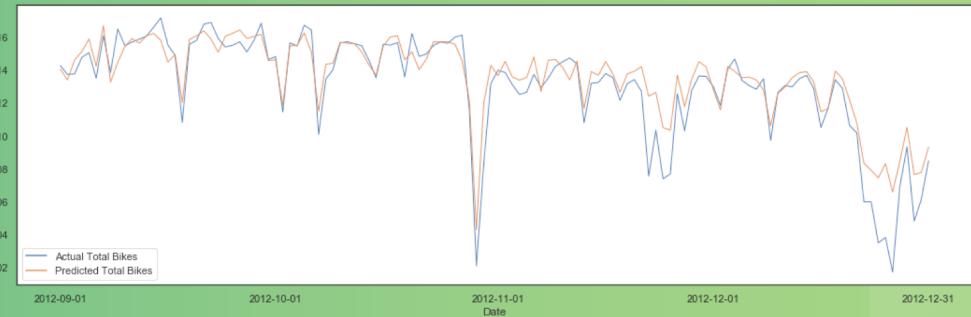
Risk of shortage during peaks



PEAKS DETECTION

Features: 54 R^2 : 0.86

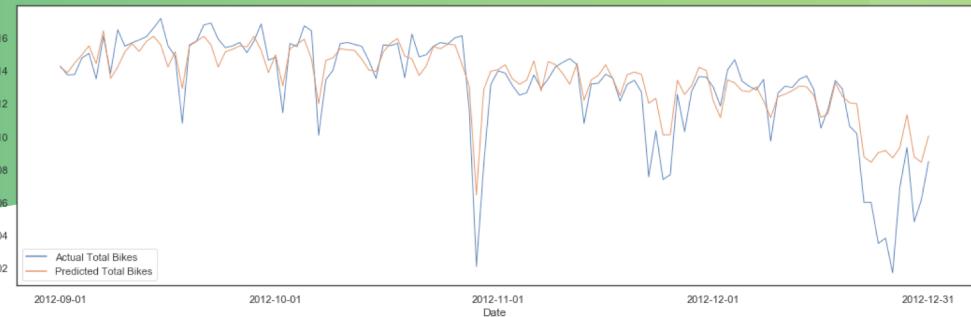
Better peaks anticipation



MANUAL SELECTION

Features: 46 R^2 : 0.85

Better general fit



5.

Business Conclusions

Optimization Using Data



Maintenance & Repair:

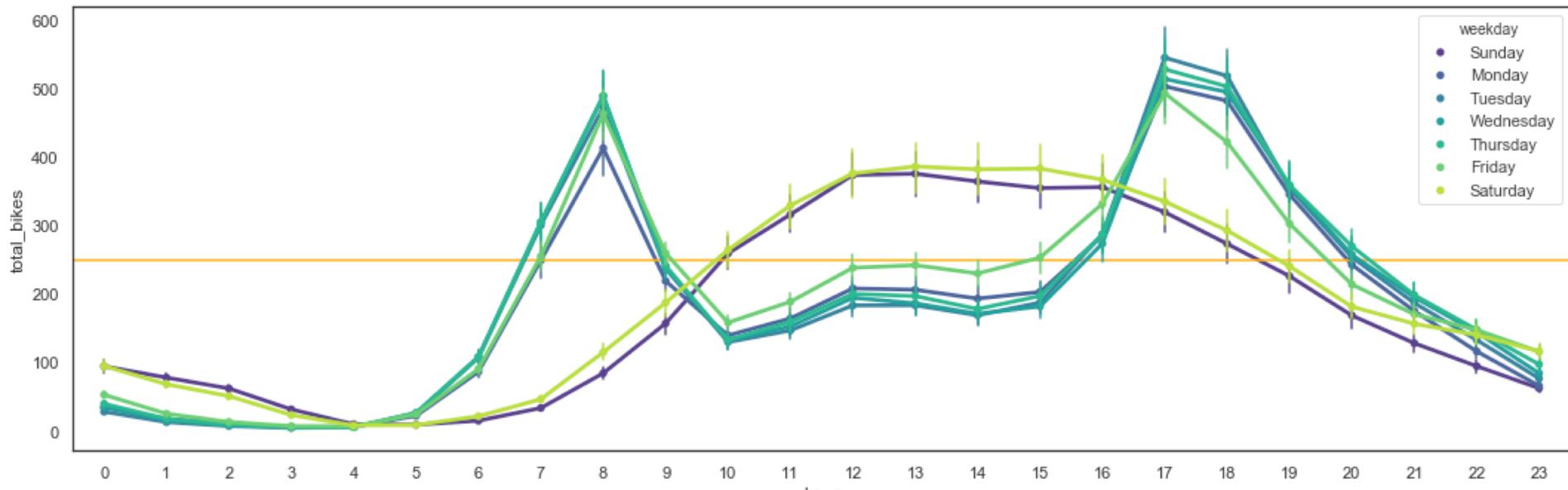
Data driven approach to optimize processes to keep bikes and docks in good repair, safe, and available.

Adapting Technologies for Future Usage:

Optimizing current operations, and the “bike valet service.”

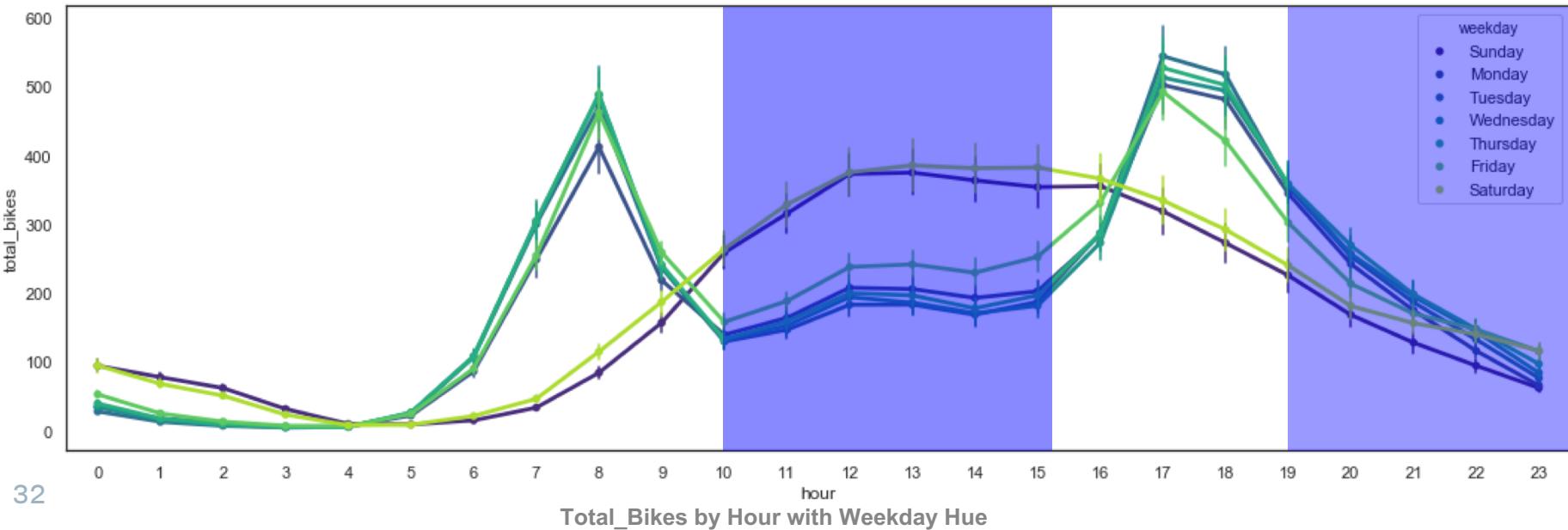
Determining Peak Times

- Peak times based on mean + 31.5%
- Process allows model flexibility
- Additional data will adapt to model



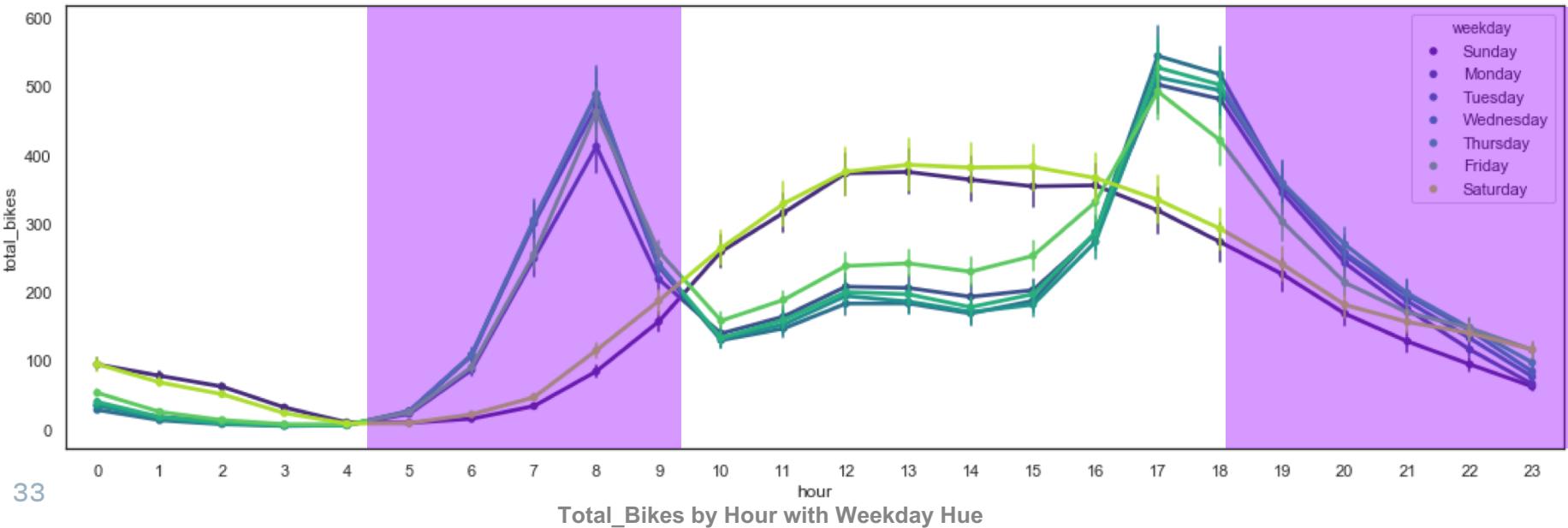
Peak Times + Maintenance Weekdays

- Weekdays/Commuting-highest usage
- Peak hours for determining maintenance time
- Goal: Least disturbance to business



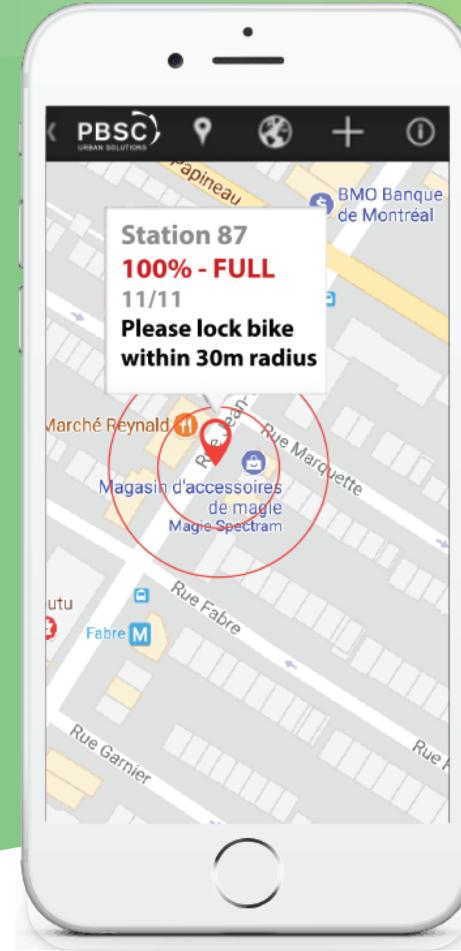
Peak Times + Maintenance Weekends

- Weekend-lower usage
- Peak hours different from weekday
- Goal: Least disturbance to business



Optimizing Operations

- Rebalancing
- Bike Valet Service
- Geofencing/Station Availability



Peak Times and Growth Optimization

- Use models in conjunction with other departments
- Avg. time increases can provide insight on inventory
- Optimize inventory based on trends





**Thank you!!
Who Has The First Question?**