

# Bay Area Bike Share



# Data Description

---

**Total Size : 2 GB**

**Station.csv** : Data about the stations where users can pickup or return bikes

**Status.csv** : Data about the number of bikes and docks available at various stations at different time stamps

**Trips.csv**: Data about individual bike trips with details such as duration, subscription type, start station name, end station name, etc.

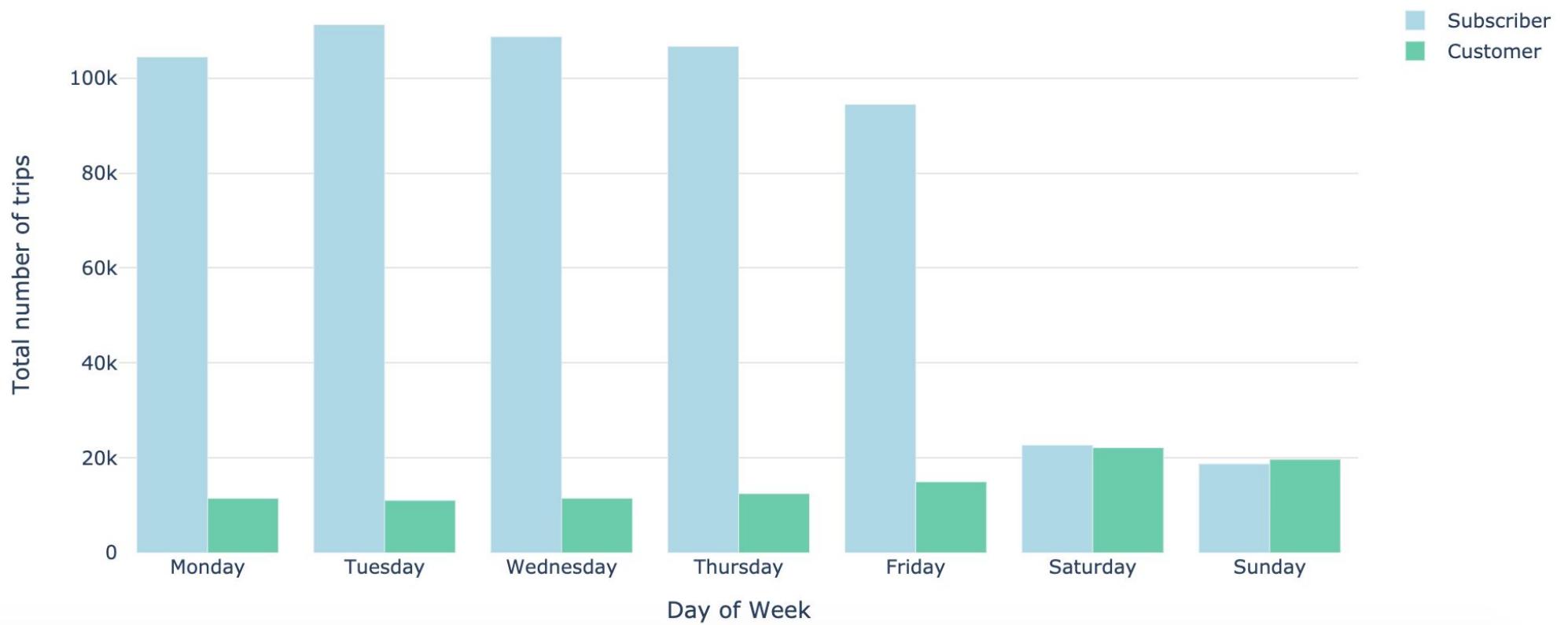
**Weather.csv** : Data about the weather on a specific day for certain zip codes. Features include temperature, humidity, visibility, wind speed, etc.

# Preprocessing & Data Visualization

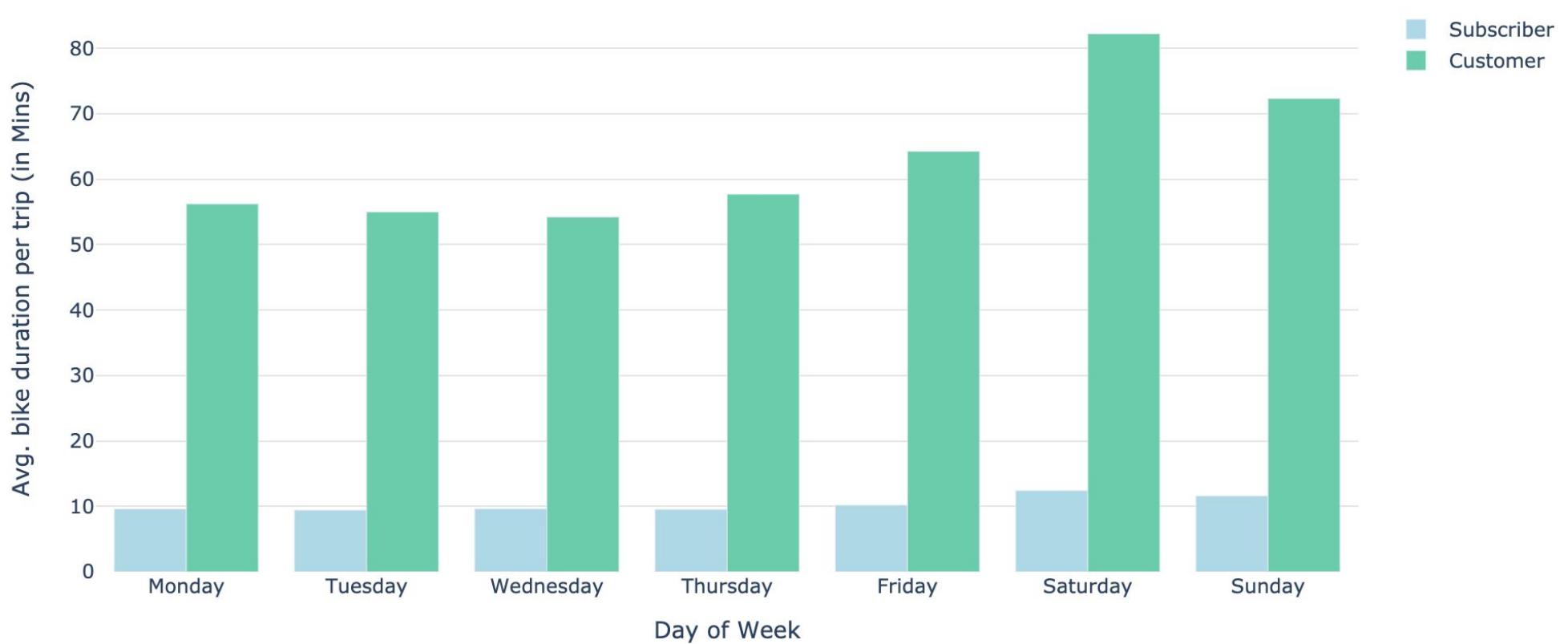
---

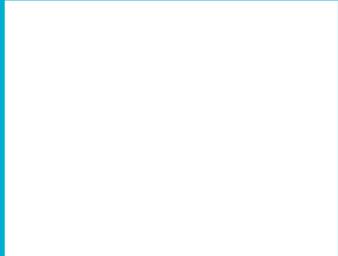
- - Segregated subscribers and customers into two different RDDs
  - Aggregated the total number of trips per day of the week for subscribers and customers
  - Computed average bike duration per trip per day of the week for both subscribers and customers

## Total Number of Bike Trips by Day of Week in Bay Area

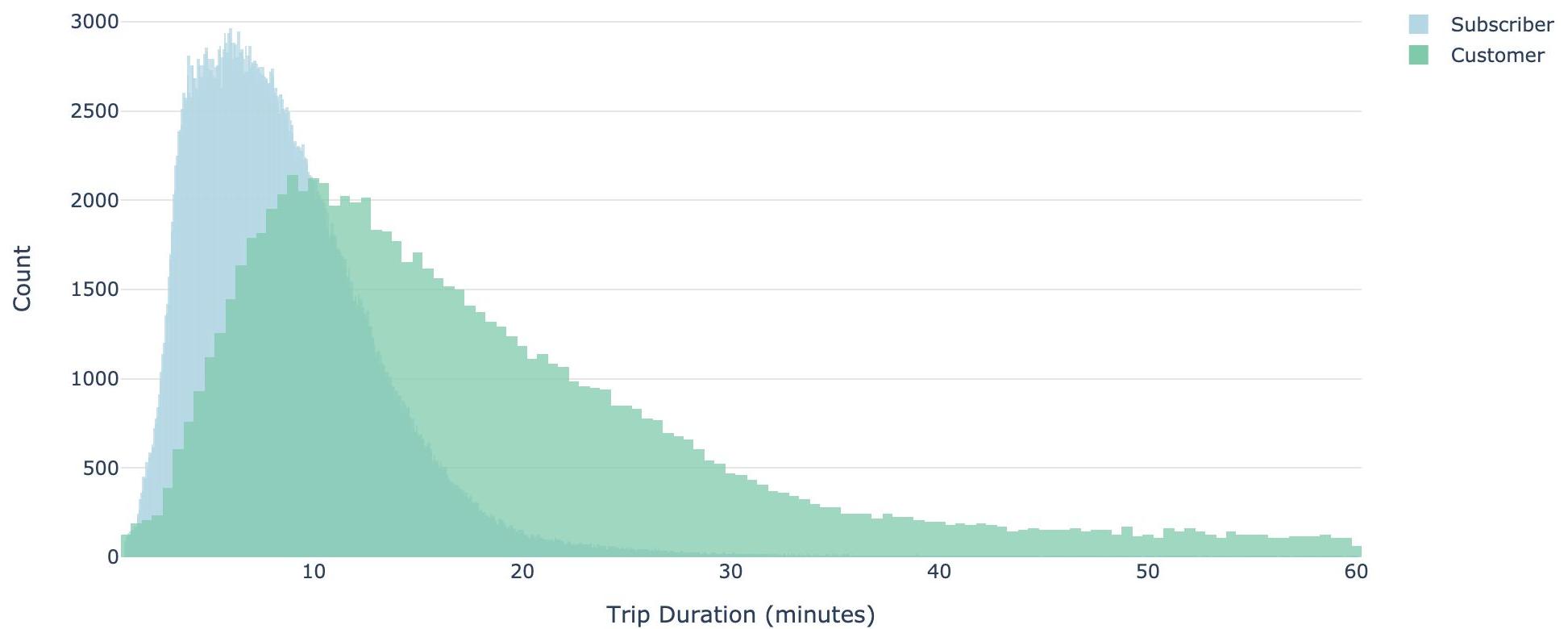


### Average Bike Duration per Trip by Day of Week in Bay Area



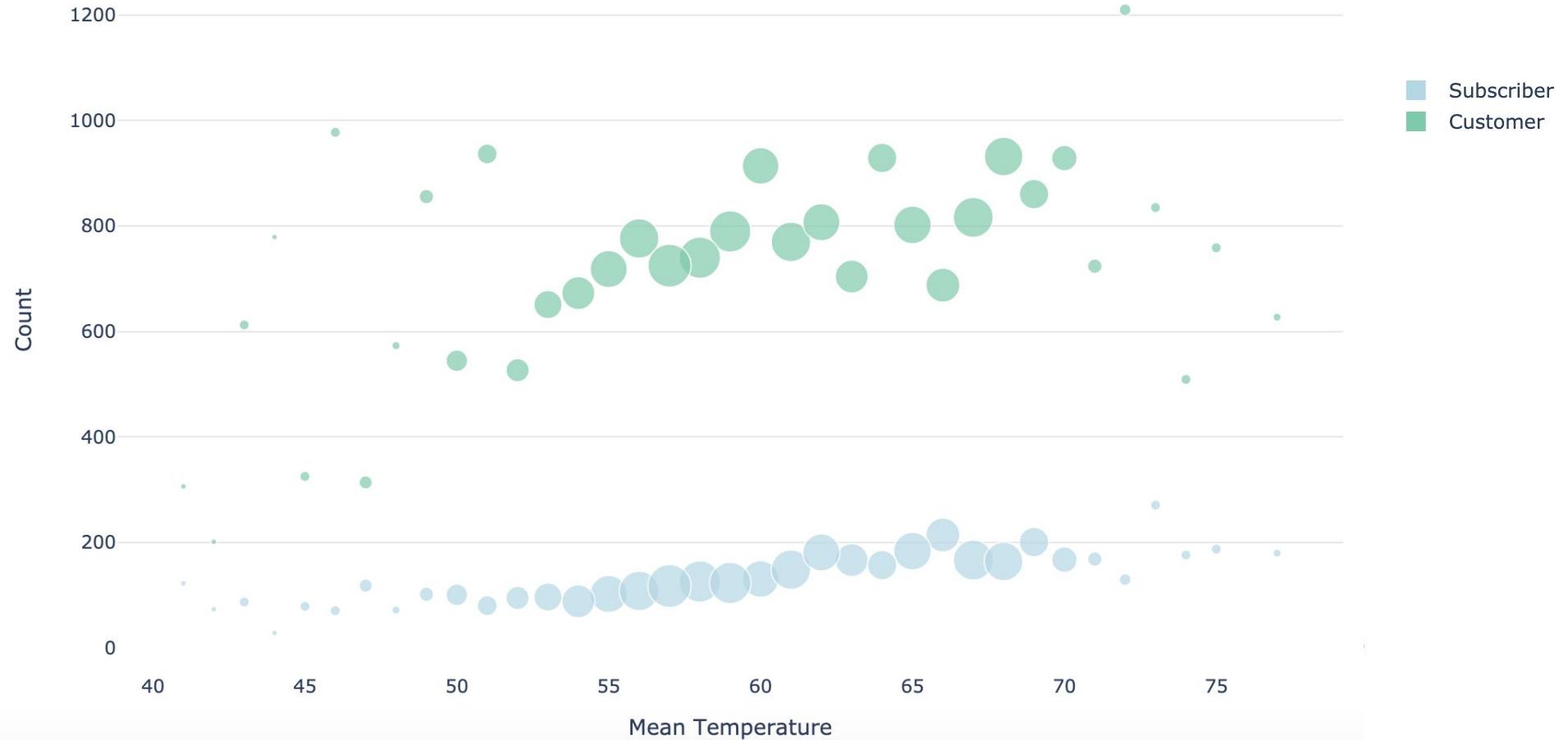
- 
- 
- For *trip.csv*, group by customer type(subscriber/customer) to count the number of trips of different durations.
  - Apply a filter to remove all the trips have a duration longer than 60 mins.

## Distribution of Trip Duration by Customer Type



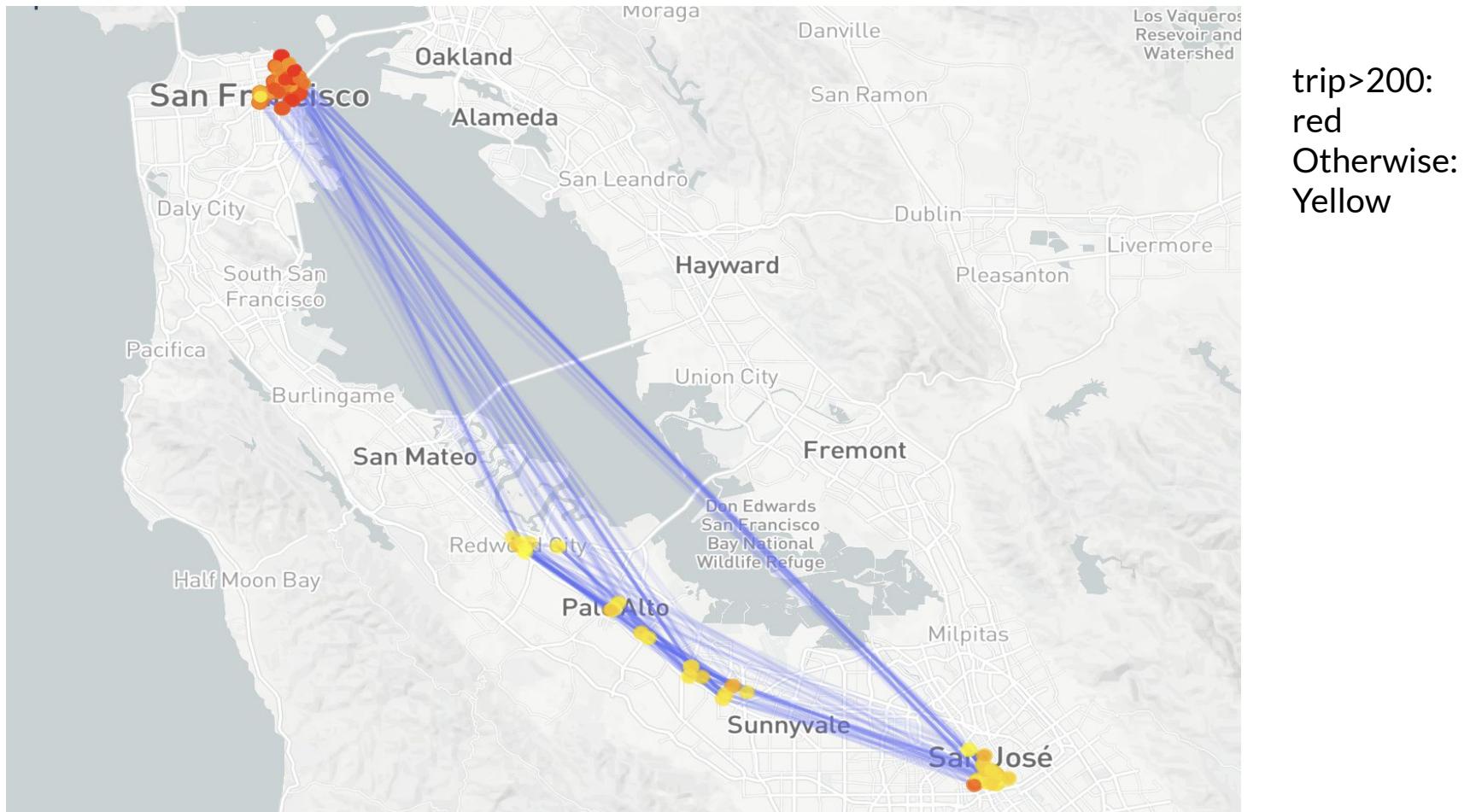
- 
- For *Trips* data, group by date to calculate total num of trips per day
  - Join *Trips* with *Weather* on the date column
  - Aggregate joined RDD to show average num of trips for different temperature

## Total Number of Bike Trips by Temperature



- Group by *start\_station\_id* and *end\_station\_id* to count the amount of trips from unique start stations and end station key pairs
- Joined *start\_station\_id* and *end\_station\_id* with longitude and latitude

## Link



## Examining average bike availability at each hour of the day on weekends and weekdays

---

- In the status dataset found the proportion of bikes available for every hour between 2014 and 2015.
- Grouped by a weekend day indicator column, hour and the name of the station
- Computed the average bike availability proportion for that hour grouped by the above variables
- [link](#)

# Cluster setting and execution time comparison

Instance Type	M5d.xlarge 16GiB memory 150 SSD GB storage	r5.8xlarge, 256 GiB memory	m4.4xlarge, 64 GiB memory	m4.4xlarge, 64 GiB memory	m4.4xlarge, 64 GiB memory
Number of Instances	1 master & 2 core nodes	1 master & 2 core nodes	1 master node & 2 core nodes	1 master node & 3 core nodes	1 master node & 4 core nodes
Execution Time	24 mins	5 mins 30s	9 mins 48s	7 mins 35s	10 mins 33 sec

# Lessons Learned

---

- Most subscribers use SF bikes for commuting to work
  - Subscribers overall take more shorted bike trips than customers.
- 
- The optimal cluster setup depends on your data processing
  - Using disk storage to shuffle data is a lot slower than memory

---



# **US Air Pollutants in the Last Decade**



# Contents

---

- Data Description
- Data Processing Goal
- Processing Outcome
- Cluster Setting and Execution Time Comparison
- Lesson Learned

# DataSet 1 : HAP

---

## Daily Hazardous Air Pollutants in United States from 1990-2017

Size: 2.3GB

Source:

[https://www.kaggle.com/epa/hazardous-air-pollutants#epa\\_hap\\_daily\\_summary.csv](https://www.kaggle.com/epa/hazardous-air-pollutants#epa_hap_daily_summary.csv)

- Arsenic:
  - highly toxic in its inorganic form
  - chronic arsenic poisoning, skin lesions and skin cancer
- Lead:
  - once taken into the body, lead distributes throughout the body in the blood and is accumulated in the bones
  - nervous system, kidney function, immune system, reproductive and developmental systems and the cardiovascular system

Pollutant name		quantity		inspection site		collection date																		
		daily_HAPS_2019																						
Parameter Code	POC	Latitude	Longitude	Datum	Parameter Name	Sample Duration	Date Local	Units of Measure	Event Type	Observation Count	Observation Percent	Arithmetic Mean	1st Max Value	1st Max Hour	AQI	Method Code	Method Name	Local Site Name	Address	State Name	County Name	City Name	CBSA Name	Date of Last Change
43502	8	33.553056	-86.815	WGS84	Formaldehyde	8 HOUR	2019-06-14	Parts per billion Carbon	None	3	100	1.466667	1.8	20		202	SILICA-ONPH-CART	North Birmingham	NO. BHAM.SOU R.R., 3009 28TH ST. NO.	Alabama	Jefferson	Birmingham	Birmingham-Hoover, AL	2019-09-24
43502	8	33.553056	-86.815	WGS84	Formaldehyde	8 HOUR	2019-06-20	Parts per billion Carbon	None	3	100	1.3	1.7	12		202	SILICA-ONPH-CART	North Birmingham	NO. BHAM.SOU R.R., 3009 28TH ST. NO.	Alabama	Jefferson	Birmingham	Birmingham-Hoover, AL	2019-09-24
43502	8	33.553056	-86.815	WGS84	Formaldehyde	8 HOUR	2019-06-23	Parts per billion Carbon	None	3	100	1.3	3.7	4		202	SILICA-ONPH-CART	North Birmingham	NO. BHAM.SOU R.R., 3009 28TH ST. NO.	Alabama	Jefferson	Birmingham	Birmingham-Hoover, AL	2019-09-24

# DataSet 2 : AQI

---

**Daily AQI Dataset in United States from 1980-2019**

Size: 915.85 MB

Source: [https://aqs.epa.gov/aqsweb/airdata/download\\_files.html#AQI](https://aqs.epa.gov/aqsweb/airdata/download_files.html#AQI)

- AQI: how polluted the air currently is, or how polluted it is forecast to become
- High AQI means high level of air pollution

AQI

State Name	county Name	State Code	County Code	Date	AQI	Category	Defining Parameter	Defining Site	Number of Sites Reporting
Alabama	Baldwin	1	3	2007-01-03	55	Moderate	PM2.5	01-003-0010	1
Alabama	Baldwin	1	3	2007-01-06	23	Good	PM2.5	01-003-0010	1
Alabama	Baldwin	1	3	2007-01-09	13	Good	PM2.5	01-003-0010	1

# DataSet 3 : Business

---

**Yearly Business Pattern in United States from 1986-2017**

Size: 84.2MB

Source:<https://www.census.gov/programs-surveys/cbp/data/datasets.html>

busPattern

Geographic area name	NAICS code	Meaning of NAICS code	Year	Number of establishments	Paid employees for pay period including March 12 (number)	county_name	state_name
Autauga County, Alabama	0	Total for all sectors	2005	872	10491	Autauga	Alabama
Autauga County, Alabama	11	Agriculture, forestry, fishing and hunting	2005	10	b	Autauga	Alabama
Autauga County, Alabama	21	Mining	2005	4	b	Autauga	Alabama

**Data Fusion:** join the above 3 datasets from 2005-2016

# Data Processing Goal

---

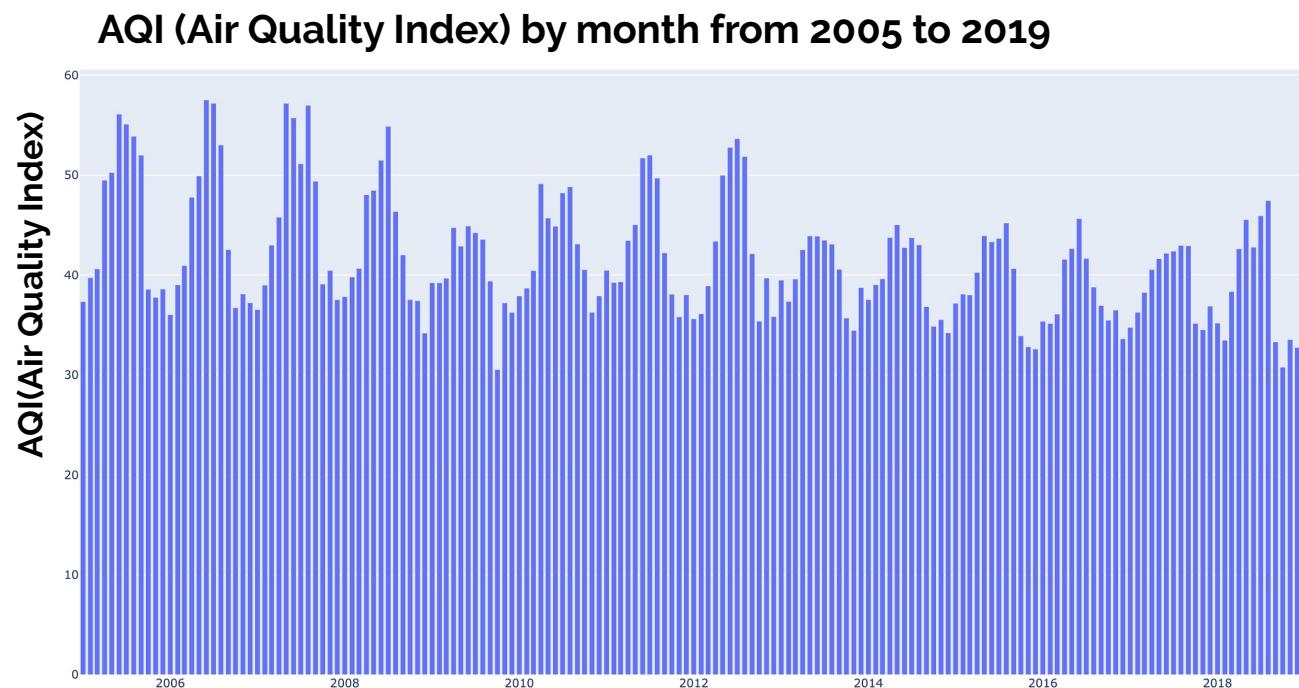
- Investigate Air Pollutants By **Time**
  - Explore air quality and hazardous air pollution over years
- Investigate Air Pollutants By **County**
  - Explore air quality and hazardous air pollution across US
- Investigate Air Pollutants By **Industry**
  - Explore the relationship of business patterns and air quality

# Processing Outcome: AQI (Air Quality Index) By Time

---

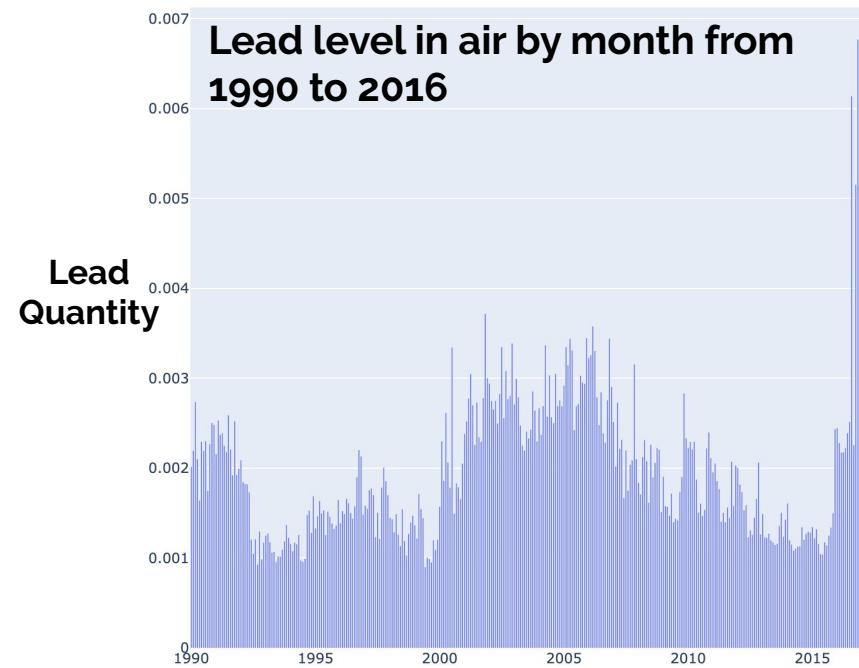
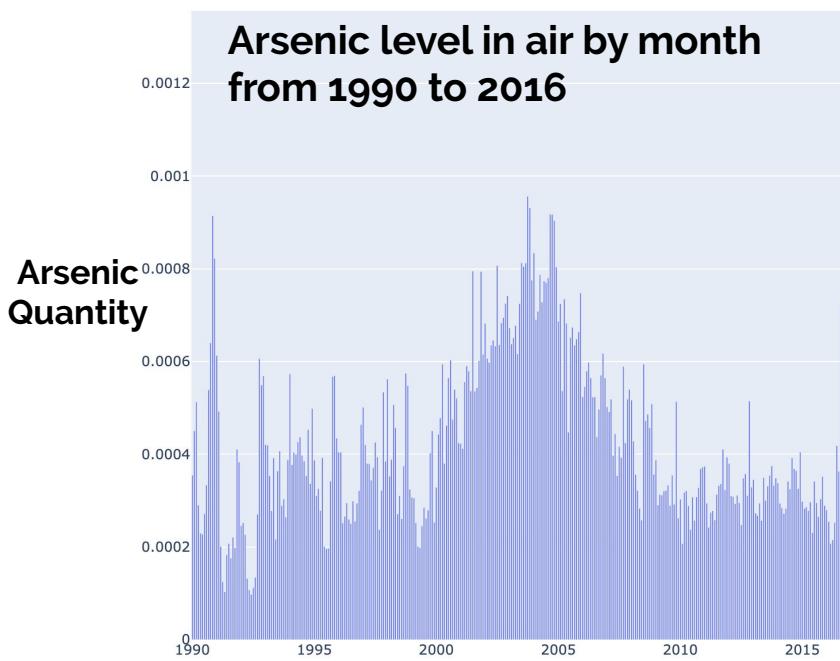
## Air Quality Index:

- obvious seasonal behavior
- yearly average AQI decreased from 2005 to 2019



# Processing Outcome: Arsenic and Lead By time

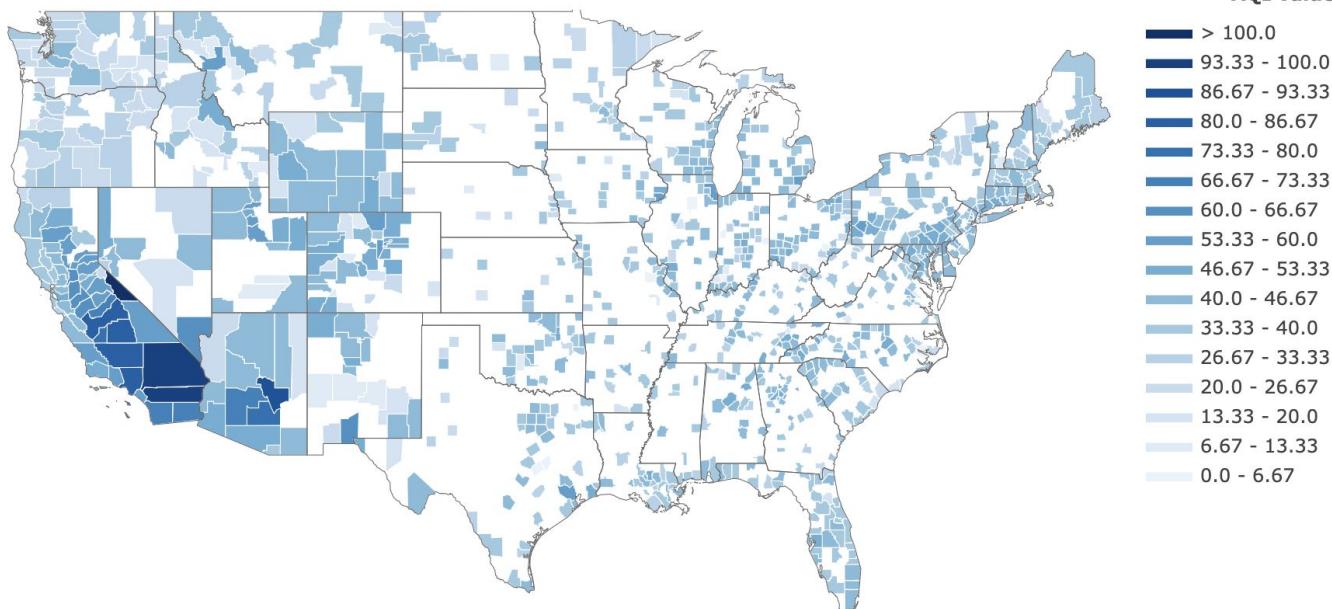
- increased from 1990 to a peak in 2005
- increased rapidly in 2016



# Processing Outcome: AQI (Air Quality Index) By County

---

AQI (Air Quality Index) by county in US in 2016



## Best 5:

- Robertson, Texas
- Titus, Texas
- Morgan, Ohio
- Beaufort, North Carolina
- Milam, Texas

## Worst 5:

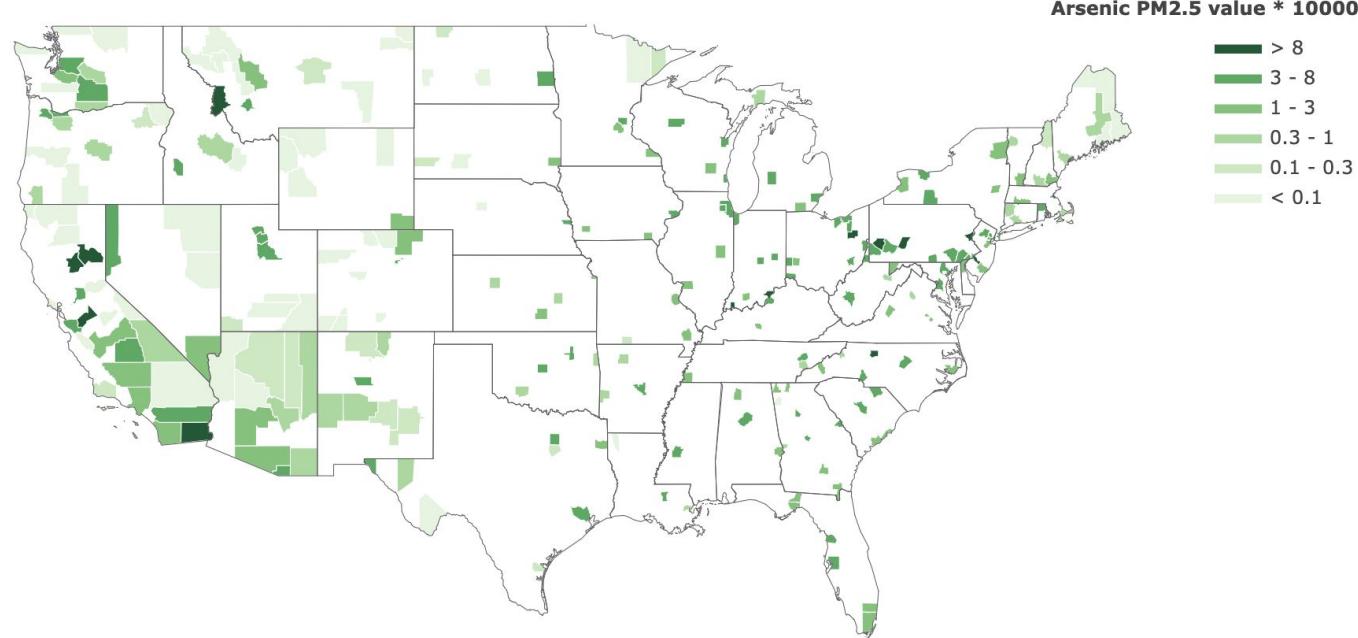
- Mono, California
- San Bernardino, California
- Riverside, California
- Gila, Arizona
- Los Angeles, California

# Processing Outcome: Arsenic PM2.5 By County

---

## Arsenic PM2.5 by county in US in 2016

Unit: Micrograms/cubic meter (LC)



### Worst 5:

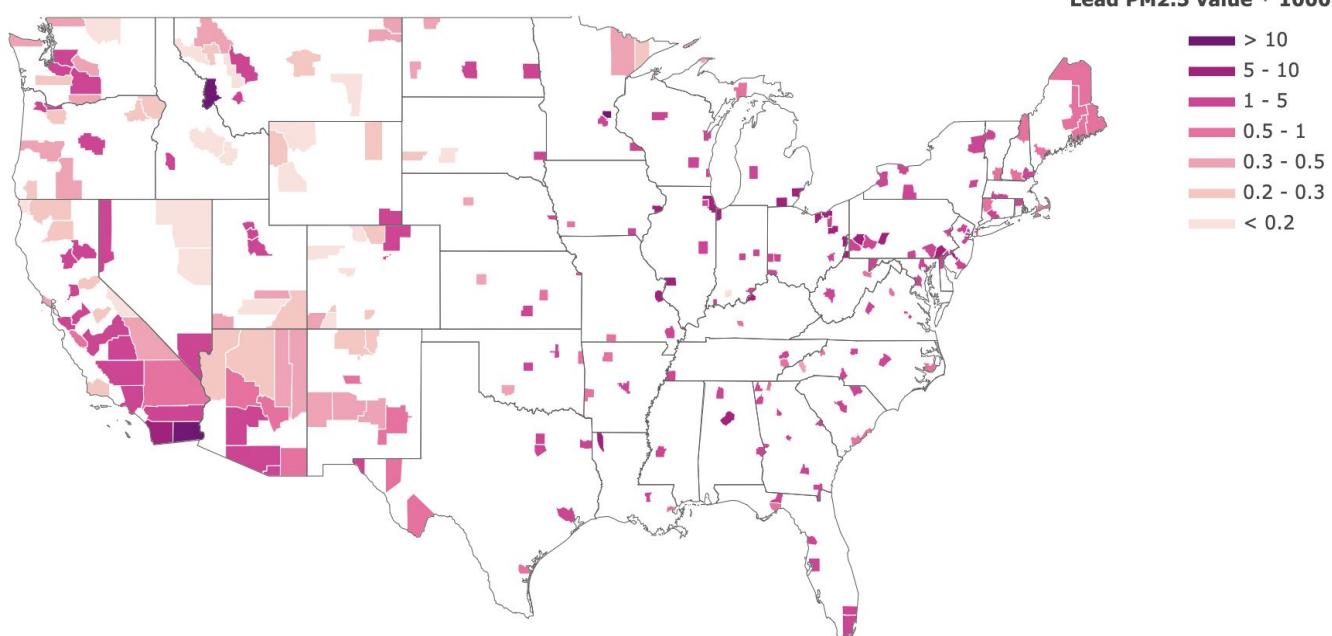
- Ravalli, Montana
- Allegheny, Pennsylvania
- Stark, Ohio
- Vanderburgh, Indiana
- Clark, Indiana

# Processing Outcome: Lead PM2.5 By County

---

## Lead PM2.5 by county in US in 2016

Unit: Micrograms/cubic meter (LC)



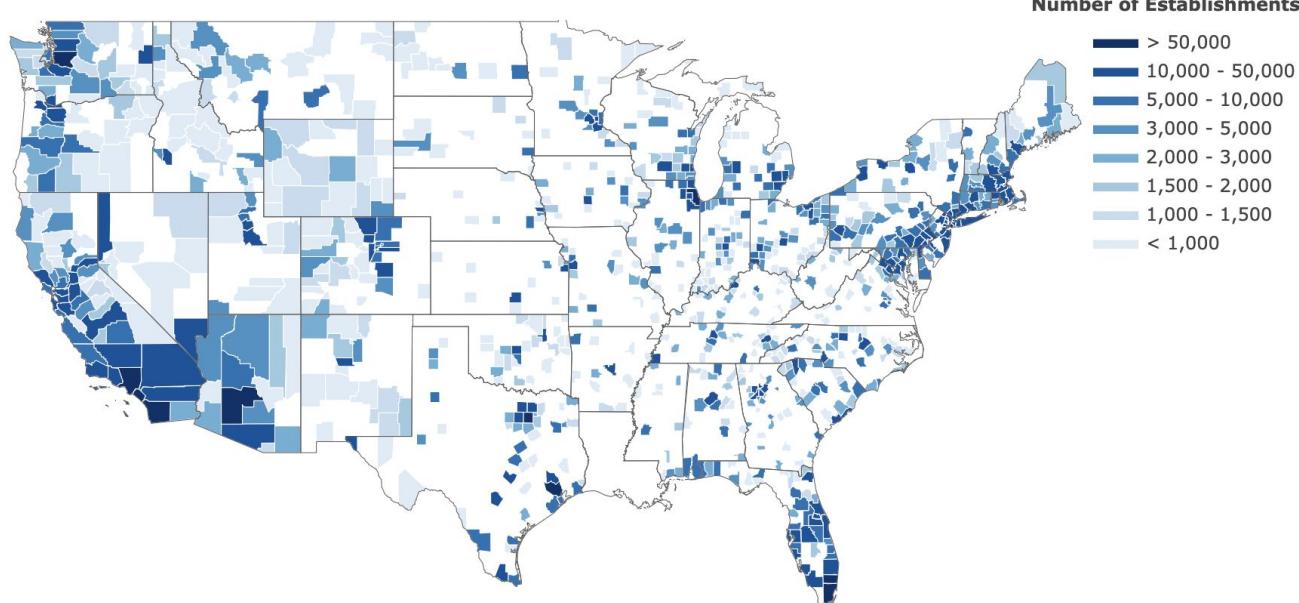
### Worst 5:

- Ravalli, Montana
- Imperial, California
- Anoka, Minnesota
- Allegheny, Pennsylvania
- Cambria, Pennsylvania

# Processing Outcome: Business Establishment By County

---

## Number of Establishment by county in US in 2016



### Largest 5:

- Los Angeles, California: 269,489
- Cook, Illinois: 133150
- New York, New York: 104691
- Harris, Texas: 100884
- Orange, California: 94703

### Smallest 5:

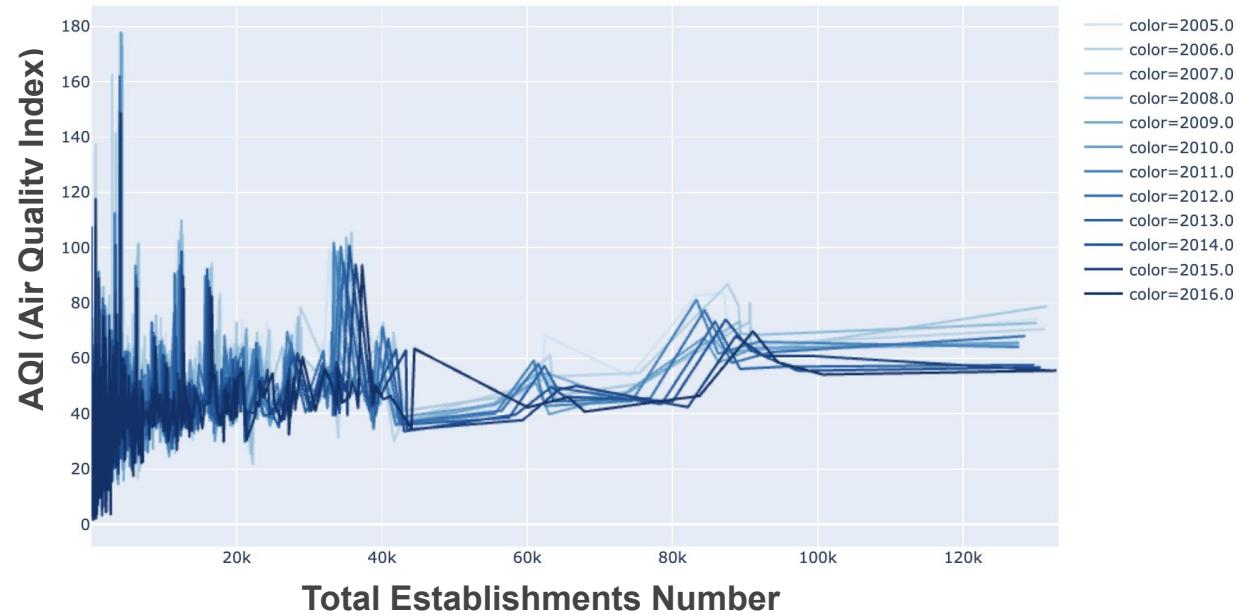
- Thomas, Nebraska: 23
- Alpine, California: 35
- Oliver, North Dakota: 42
- Jackson, South Dakota: 49
- Culberson, Texas: 54

# Processing Outcome: Business Establishments & AQI

---

- Air quality index increases with number of establishments.

AQI vs Total Establishments Number by county in US from 2005 to 2016



# Cluster Setting & Execution Time Comparison

---

Cluster Setting		Execution Time (in seconds)
machine specs	number of nodes	
m3.xlarge	3	318.74
m4.xlarge	3	272.42
m5a.xlarge	3	248.04
m5.xlarge	3	190.36
m5.xlarge	4	159.67



Execution time decreases

# Lesson Learned

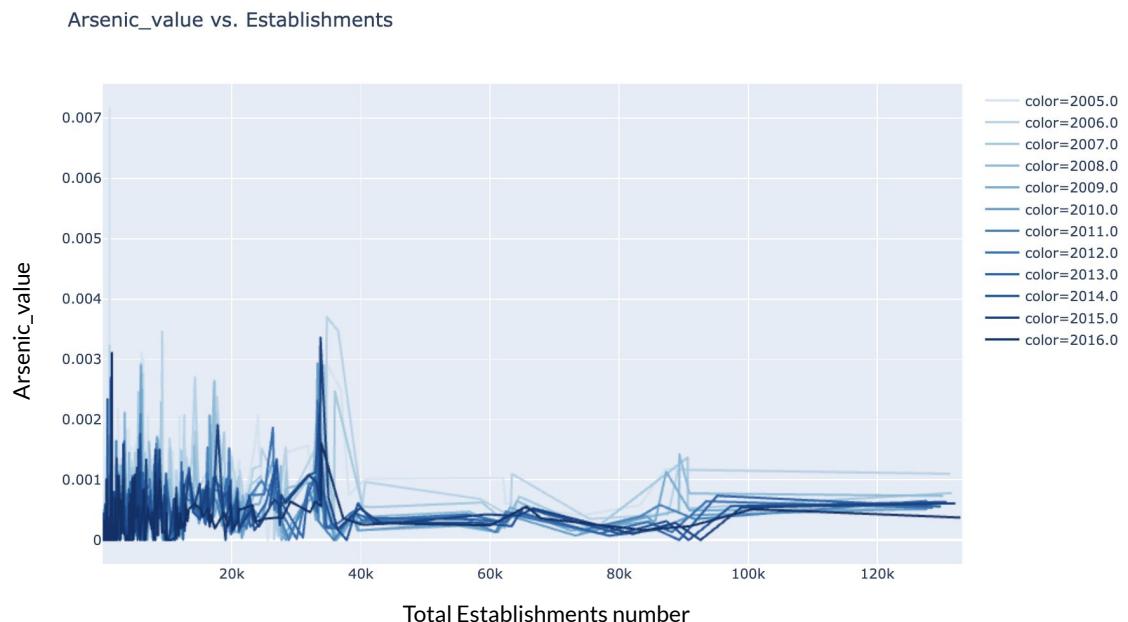
---

- Air quality is getting better over the years with seasonality.
- Air quality and pollutant level and establishment numbers are related.
- Distributed computing allows us to process large dataset.
- The larger machine specs of AWS EMR, the shorter execution time.
- The more cluster nodes of AWS EMR, the shorter execution time.

# Future Work

---

- Number of total establishments does not have strong impact on Arsenic/Lead level
- Is it possible that specific type of establishment has high impact on these hazardous air pollutants? E.g. Mining, Agricultural...
- Can we predict air pollutants level?





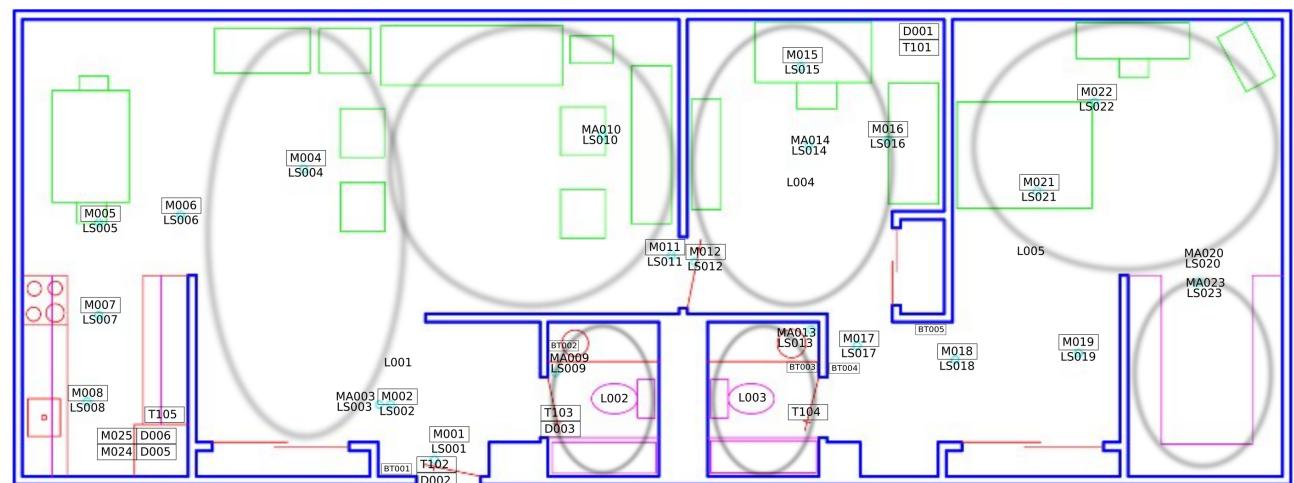
**Thank you!**



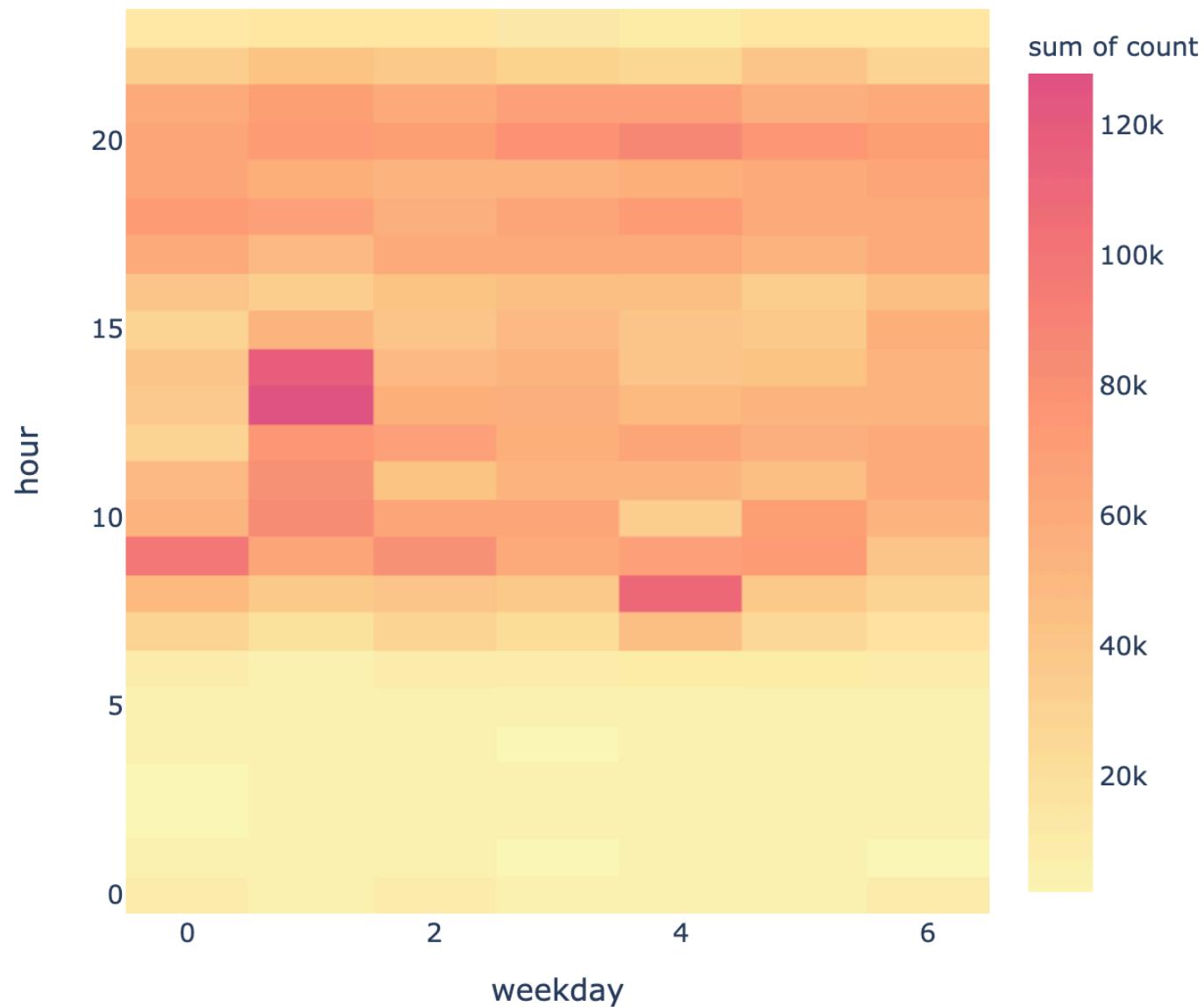
# From Sensors to Smart Homes

# Data:

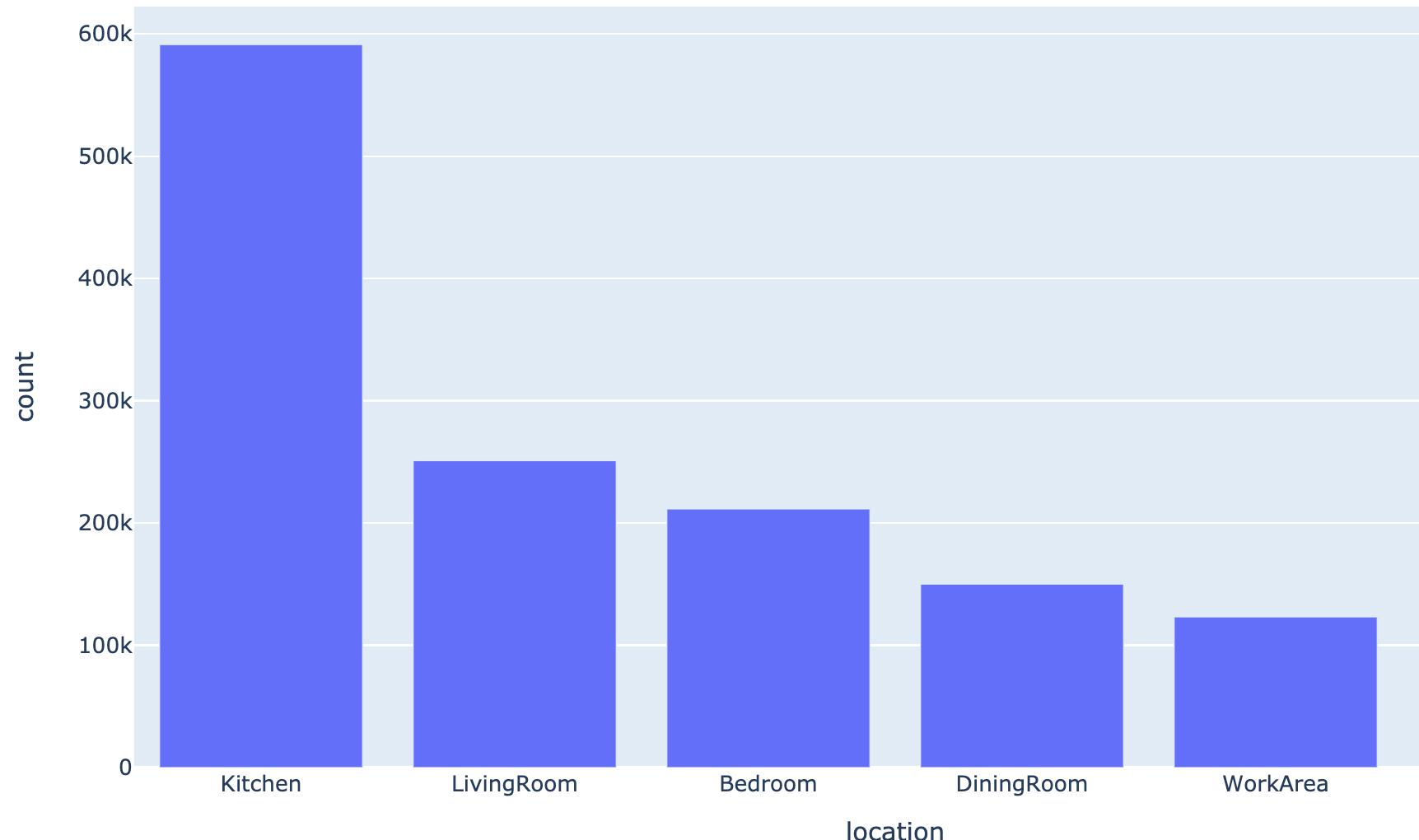
- 29 Homes
- 2 Years of Collection
- Sensors are always on and communicate to the server
- Motion sensors, Door / Temperature sensors, and Light Switch sensors
- Labels are the activity type
- Multi\_resident



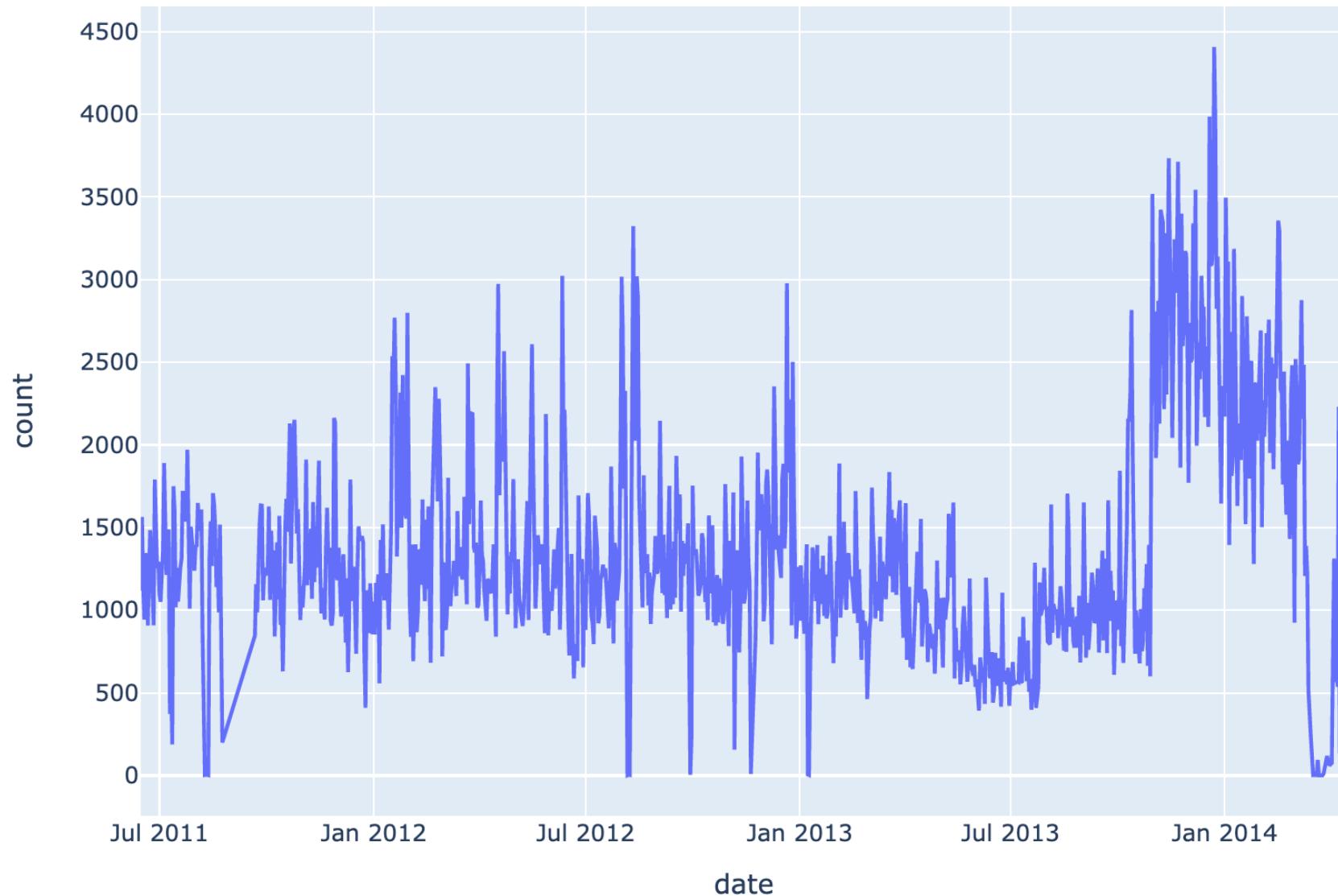
# Overall Shape of The Motion:



# Sensor Sensor on the wall...



# Tell Me When I am Gone!

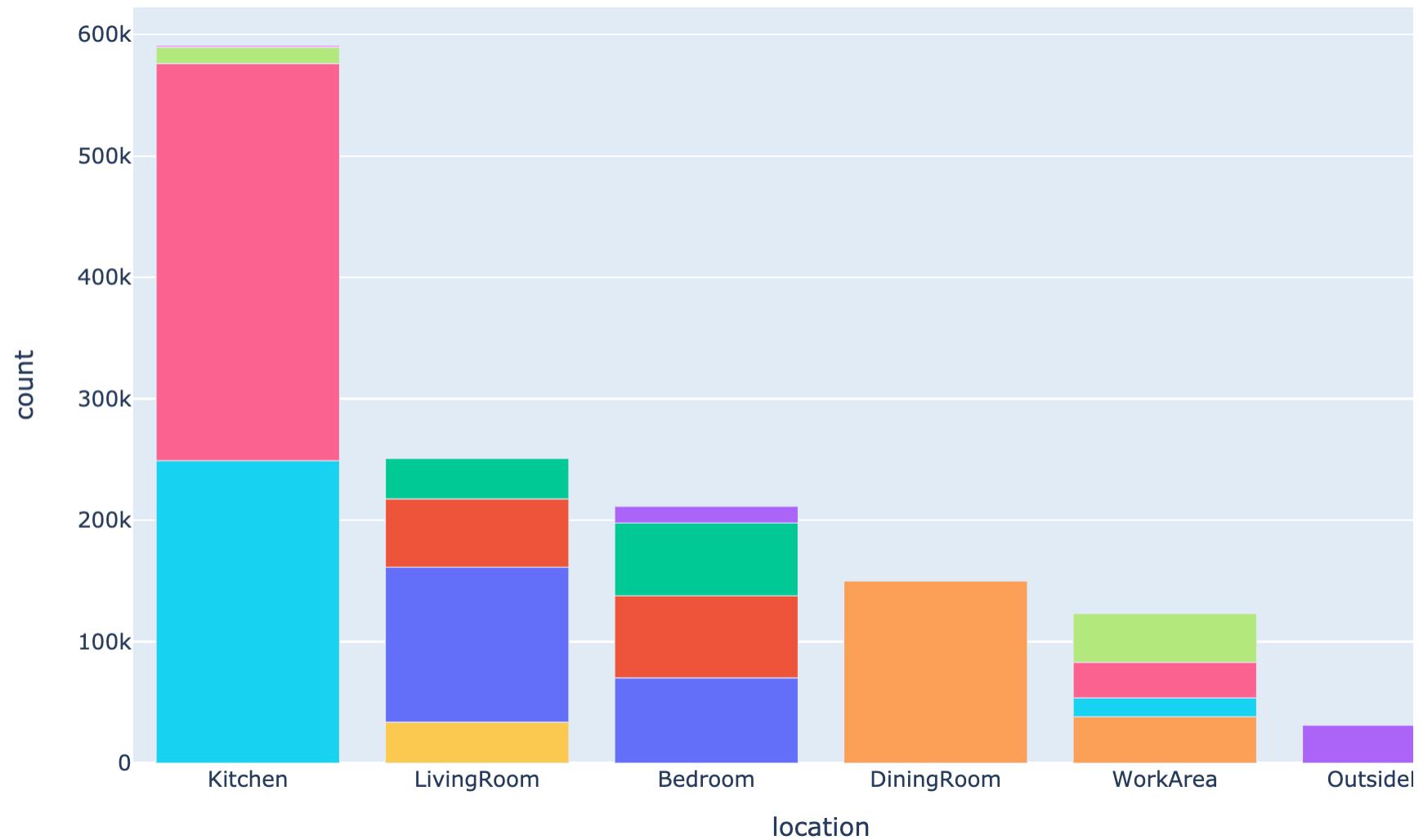


**2011-06-15**

**09:58:45.585184**

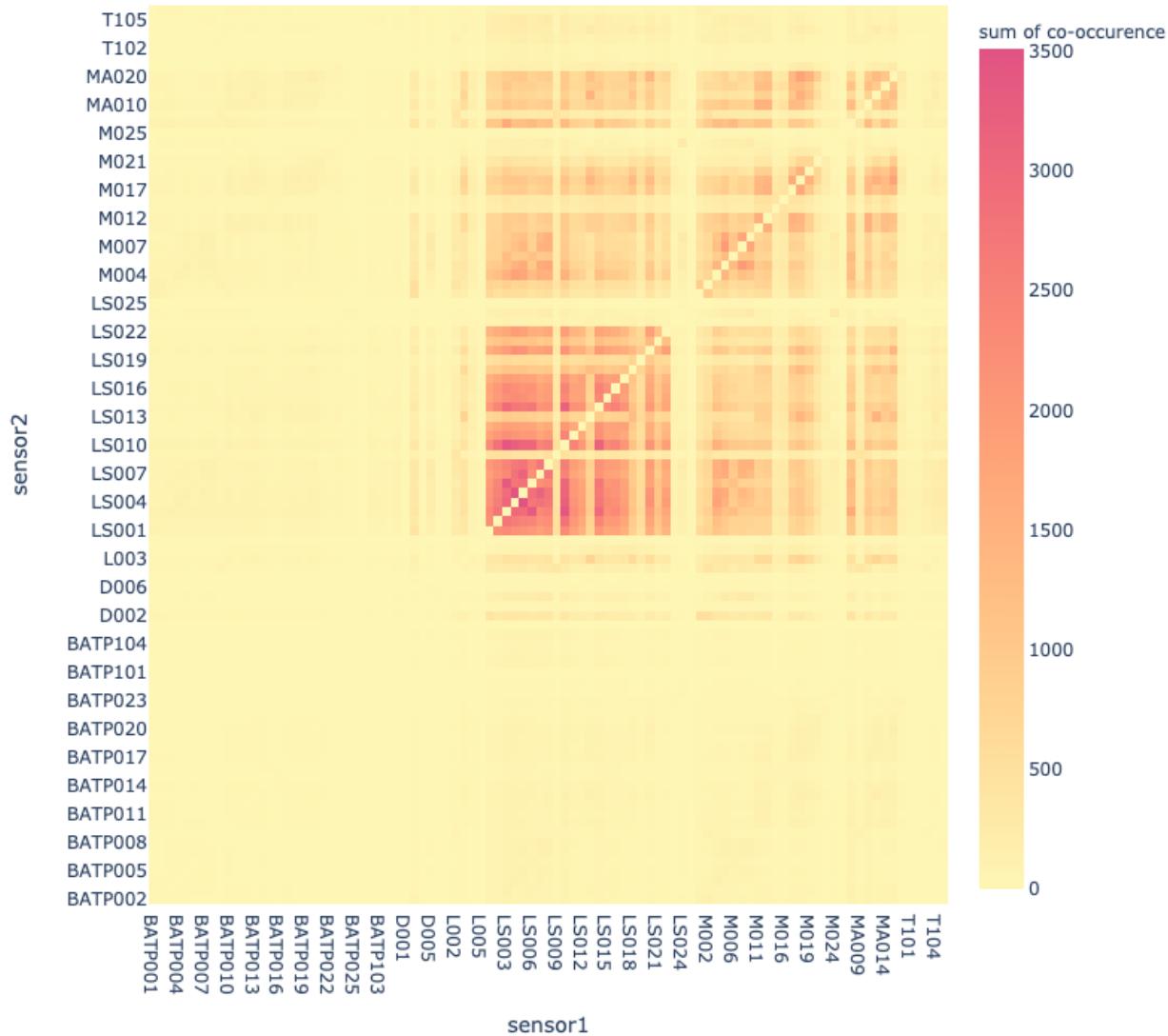
We are able to precisely detect movement, but it's no surprise as we have put too many sensors in each home.

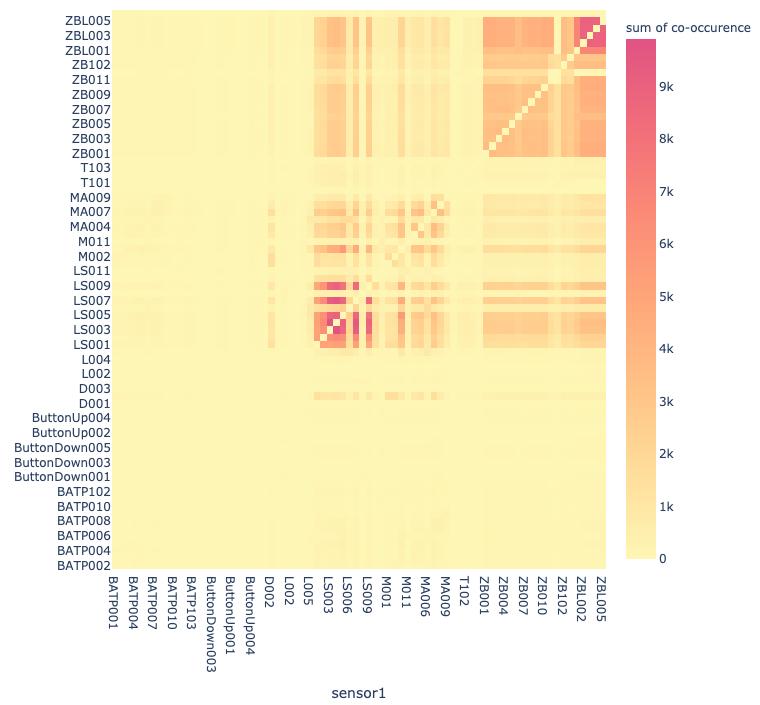
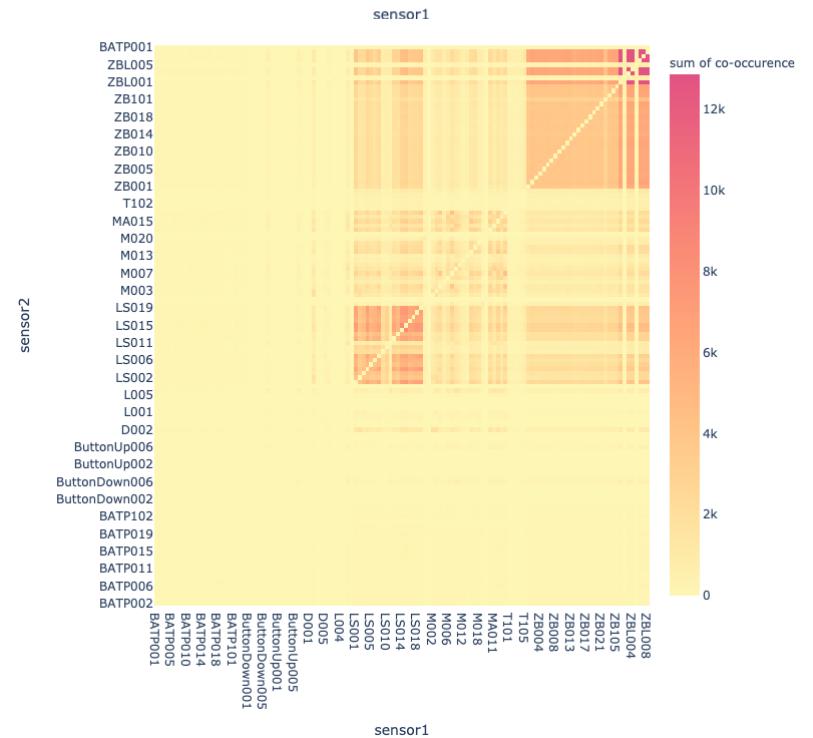
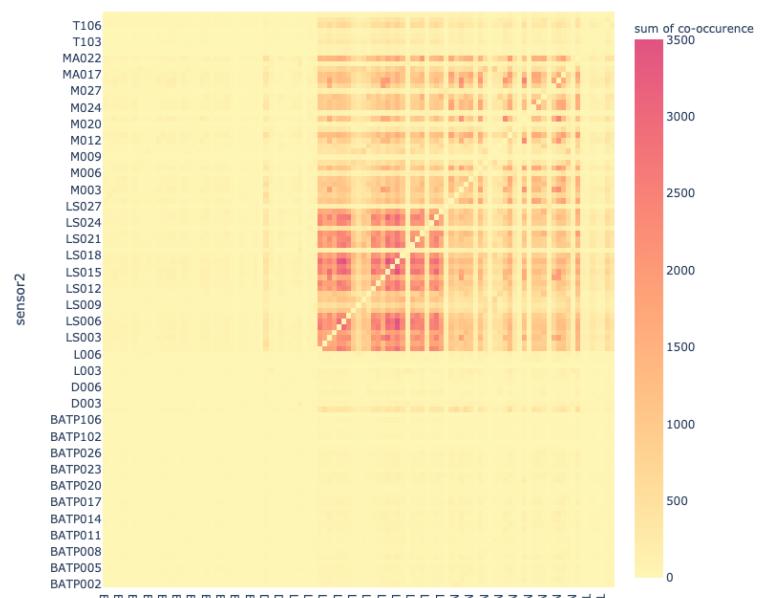
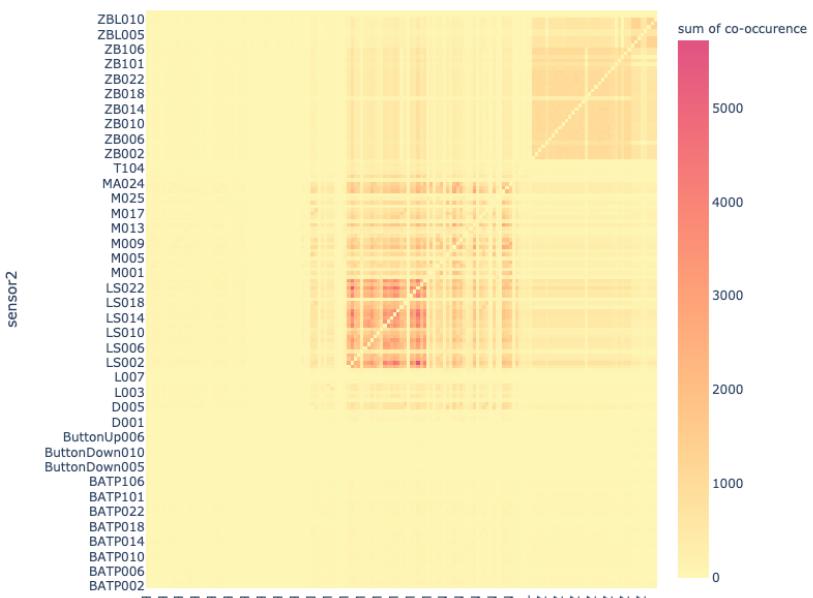
# Is This Even Smart?



Is there a correlation between sensors, and can we use less of them?

Co-occurrence Heatmap





# Running 104 jobs to process the data

## ▼ Job [104]: runJob at PythonRDD.scala:153

Progress for runJob at PythonRDD.scala:153		Job Progress: 3/3 Tasks Complete		
Stage [ID]: name at [source]:[line]	Status	Task Progress	Elapsed Time (seconds)	Failed Task Logs
Stage [250]: groupByKey at <stdin>...	SKIPPED	0/29	n/a	
Stage [251]: sortByKey at <stdin>:75	SKIPPED	0/29	n/a	
Stage [252]: runJob at Python...la:1...	COMPLETE	3/3	0.051	

	2 Workers	4 Workers
m5.xlarge	8 mins	5 mins
m5.2xlarge	5 mins	3 mins

# Conclusion:

As we can see some of these sensors are highly correlated.

We might be able to reduce the sensors, and still be able to reach the business goal, which is classifying the action type.

We can use PCA or feature importance further to see if removing some of these sensors will hurt the outcome.