# Weather & Transportation
## Streaming the Data, Finding Correlations

Provide capability to Data for Democracy democratizing_weather_data

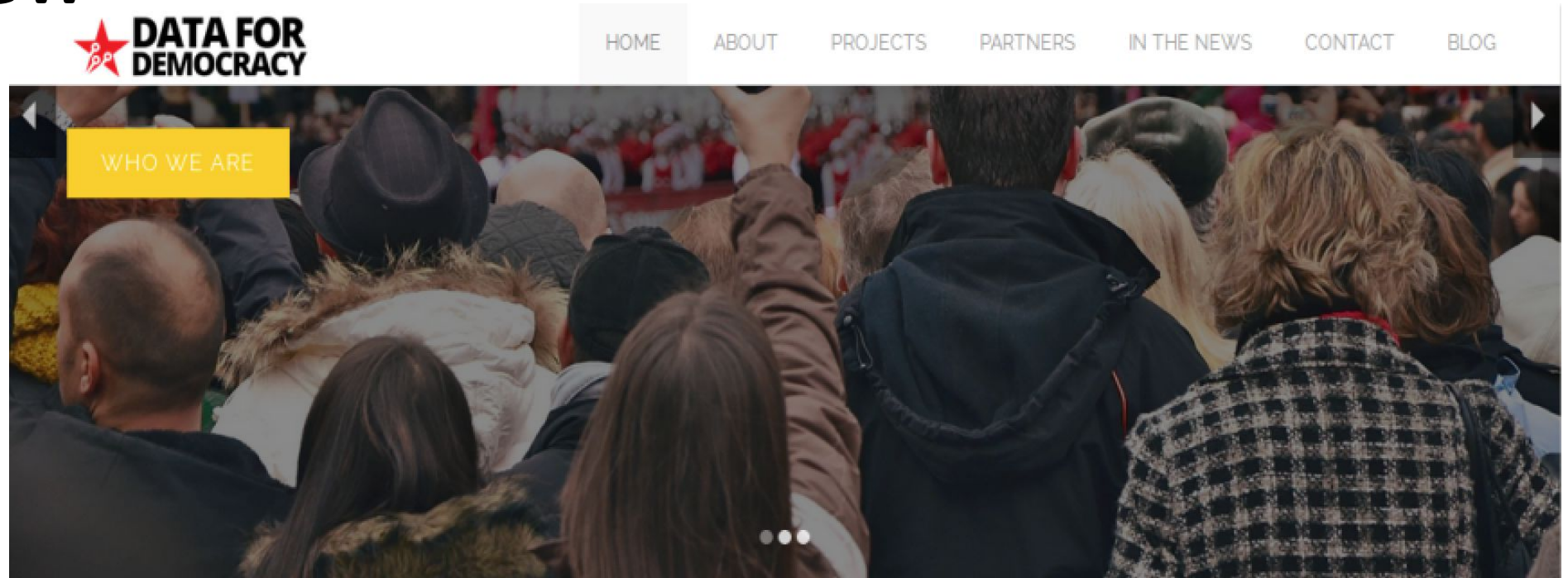University of Washington Professional & Continuing Education

BIG DATA 230 B Su 17: Emerging Technologies In Big Data

**Team D-Hawks**

John Bever, Karunakar Kotha, Leo Salemann, Shiva Vuppala, Wenfan Xu

# Overview

**Our "Client"**



**Their Mission**

To be an inclusive community for data scientists and technologists to volunteer and collaborate on projects that make a positive impact on society.
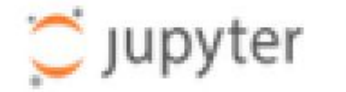
**Our Mission**

- Provide a streaming capability to extract weather and traffic data from multiple Web API's, and produce a clean merged dataframe suitable for Machine Learning and other Data Science analysis.
- Deliver code to D4D's Github Repository
- Use vendor-neutral, opensource solutions, implemented in python and Jupyter notebooks

**Learn More**   www.datafordemocracy.org  https://github.com/Data4Democracy  democratizing_weather_data/streaming
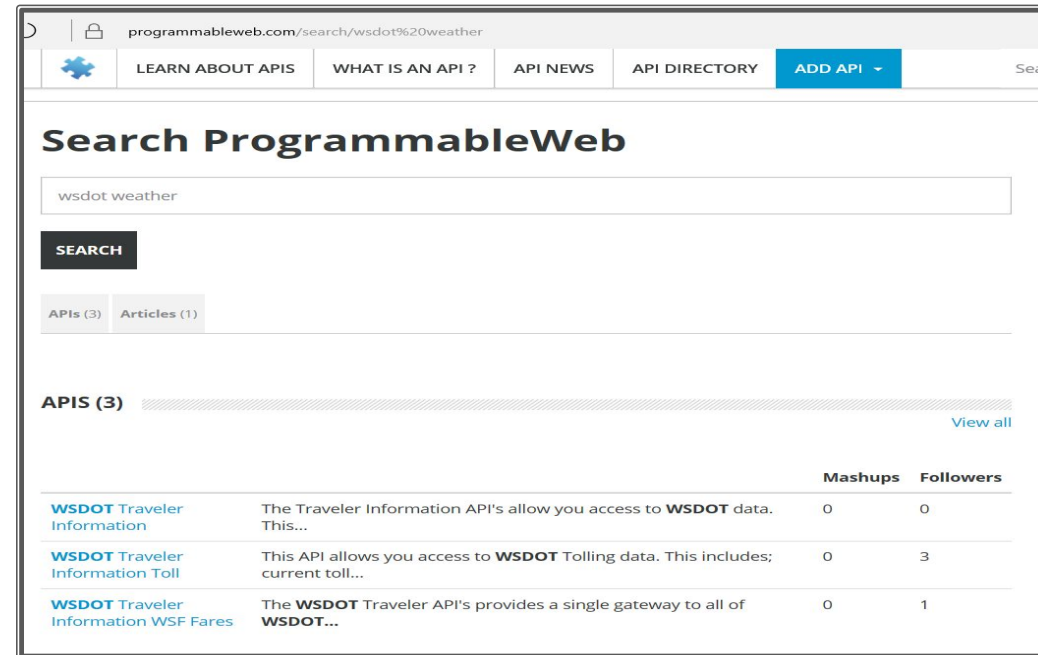
# Pipeline



- Kafka transport mechanism (vendor-neutral, open source)

- Message value is an entire JSON document

- One topic per source API, guarantees consistent schema

- Multiple json documents (sharing same schema) combined into a single dataframe

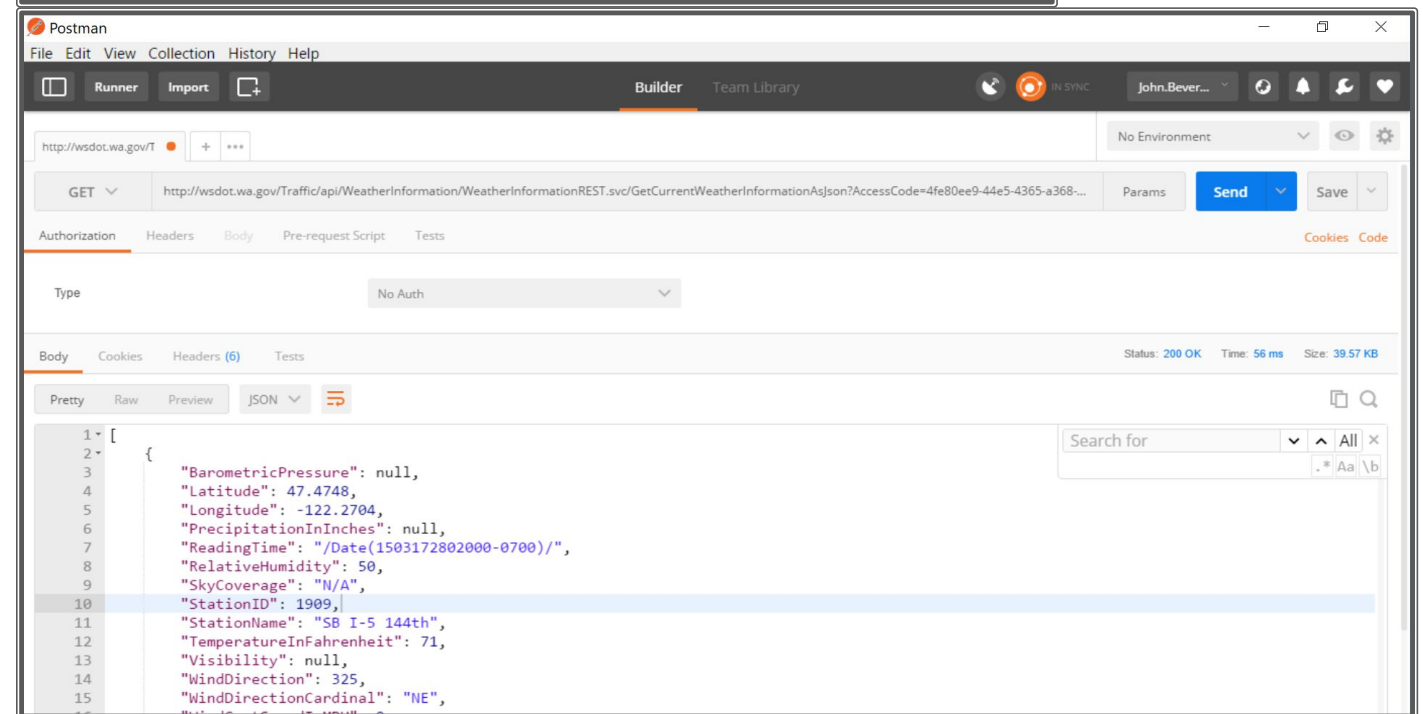- Dataframe records joined based on space and time

# Web APIs



## ProgrammableWeb.com

- A massive searchable directory of over 15,500 web APIs that are updated daily

- Includes sample source code for APIs

## Postman

- Great tool for interacting with potential APIs.

- Friendly GUI for constructing requests and reading responses.

- **Provided JSON files before pipeline was completed. Allowed analysis of data in parallel**

# Producers

```
1   import sys
2   from kafka import KafkaClient, SimpleProducer
3   import json,requests
4   from apscheduler.schedulers.blocking import BlockingScheduler
5   import logging
6
7   logging.basicConfig()
8
9
10  def pullData():
11      topic = sys.argv[1]
12      kafka = KafkaClient('localhost:9092')
13
14      producer = SimpleProducer(kafka)
15
16      #url= 'http://countdown.api.tfl.gov.uk/interfaces/ura/instant_V1'
17      url = sys.argv[2]
18      r = requests.get(url,stream=True)
19
20      for line in r.iter_lines():
21          producer.send_messages(topic,line)
22          print(line)
23
24      kafka.close()
25
26  sched = BlockingScheduler()
27  sched.add_job(pullData, 'interval', minutes=1)
28  sched.start()
```

Arguments
- Topic
- URL + Access Key

Message.Value
- JSON document

# Consumers

```python
import sys
import logging
import multiprocessing
import json
import time
from datetime import datetime
from kafka import KafkaConsumer

class Consumer (multiprocessing.Process):
    def __init__(self, topic_name):
        self.topic_name = topic_name


    daemon = True


    def run(self):
        consumer = KafkaConsumer(bootstrap_servers = 'localhost:9092',
                                 auto_offset_reset = 'latest')
        consumer.subscribe(self.topic_name)
        for message in consumer:
            print (message.value.decode('utf-8'))
            with open(datetime.now().strftime("%Y-%m-%d-%H-%M-%S"), 'w') as outfile:
                outfile.write(message.value.decode('utf-8'))

def main():
    topic_name = sys.argv[1:]
    consumer = Consumer(topic_name)
    consumer.run()
    time.sleep(10)

if __name__ == "__main__":
    logging.basicConfig(
        format = '%(asctime)s.%(msecs)s:%(name)s:%(thread)d:%(levelname)s:%(process)d:%(message)s',
        level = logging.INFO
    )
    main()
```

- One complete JSON file on disk per message

- Filename includes timestamp

- "utf-8" decoded text file

# Analysis

**7 days** of data (includes eclipse!)    **30 minutes** between readings

Load Json file, normalize, save as dataframe.
Repeat for next json file, append to prior.

**54** Weather Json Files from Yahoo    (54 rows x 31 columns)

**394** Weather Json Files from WSDOT  (40,931 rows x 16 columns)

**395** Traffic Json Files from WSDOT    (70,998 rows x 20 columns)

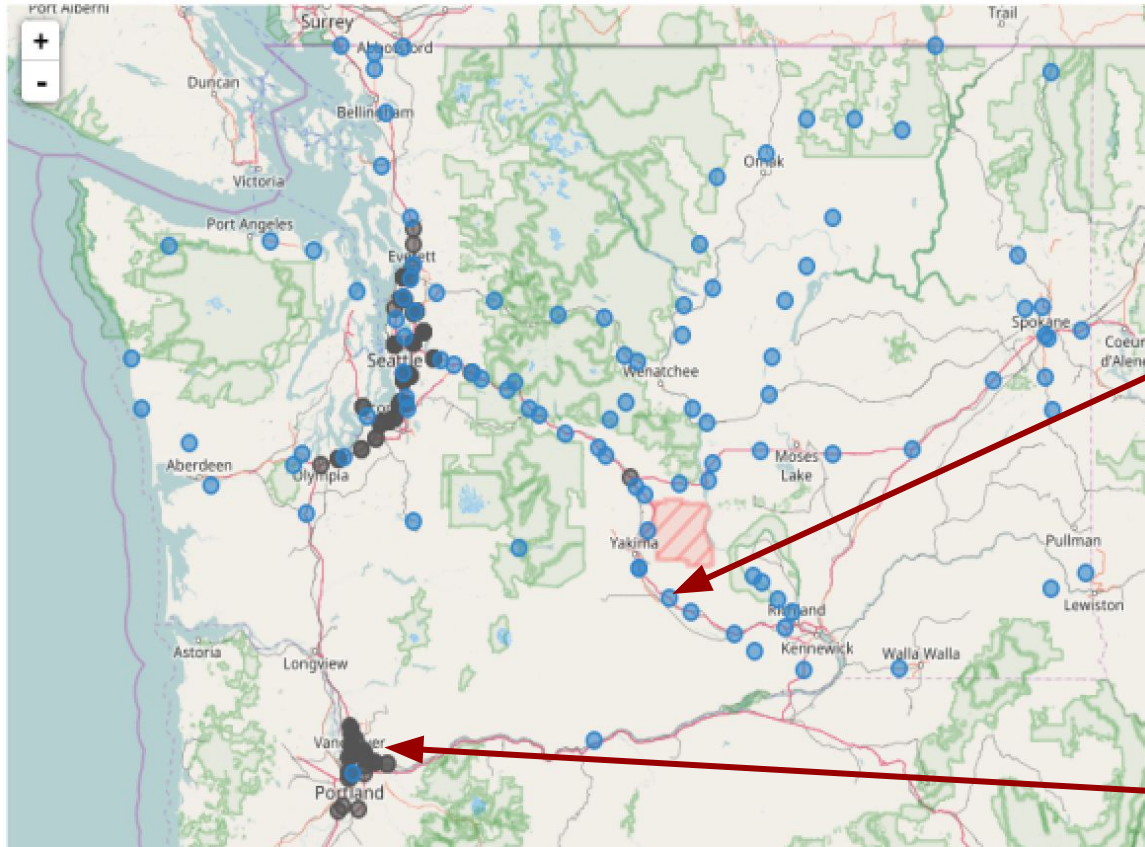Merge WSDOT & Yahoo  Weather Dataframes (use columns common to both)
Merge Traffic/Weather Dataframes.  Each Row has:
- Traffic data from a specific Traffic dataframe row
- Weather data from a weather station within *20 miles* and *30 minutes* of traffic reading.

**1** Merged Traffic/Weather Table    (52,975 rows x 30 columns)
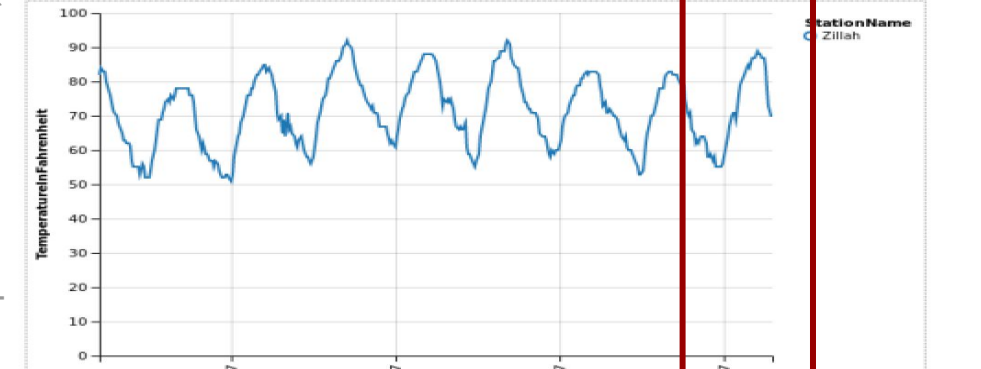
# Visualization

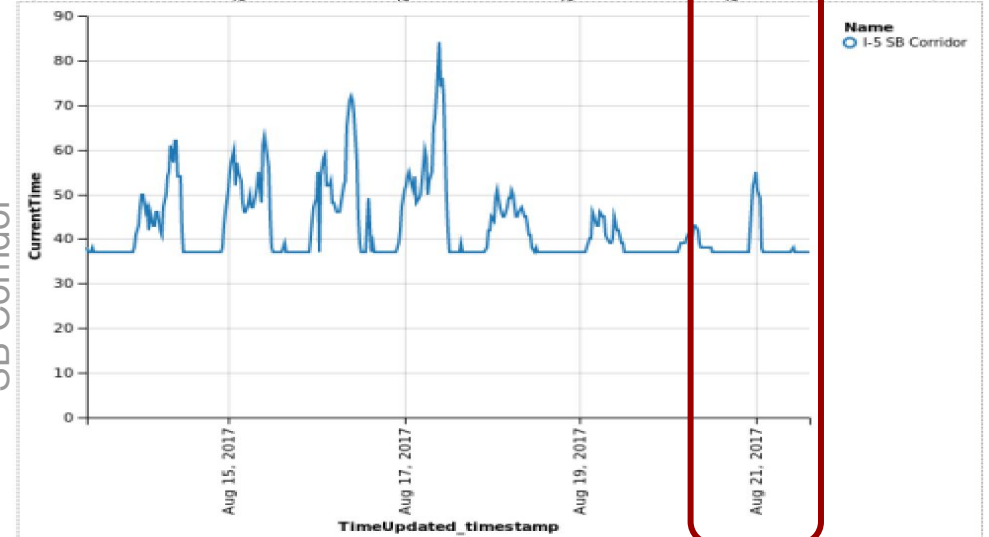Mapping with Folium (traffic in black; weather in blue)
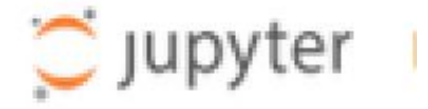
Charting with Altair



Temperature for Zillah, WA

Current Travel Time for I-5 SB Corridor

Eclipse

# Analyzing the Merged/Traffic Weather Dataset



Scatterplot Matrix with Seaborn (10% random sample)

# Wrapping Up …

## Key Takeaways
- Choose your python libraries carefully (2 lines of code for a fully-labeled lineplot vs. dozens)
- Spatial plots first, data-joins later (I-5 traffic data vs. statewide weather, also Portland)
- The fastest way to count records in a dataframe is df.shape[0]

## Conclusion
- Data for Democracy has a repeatable way to extract weather and transportation data from WSDOT and Yahoo
- Jupyter Notebook provides a teaching/coding environment
- Bitnami provides low-cost simple Kafka infrastructure

## Further Work
- Upload csv and zipped json's to data.world
- Better parameters for Producer scripts (ex. Longitude, Latitude, Date, Time)
- Config files for access keys
- More matrix plots, Data Science, Machine Learning
  - Gather data for longer time frames (fewer readings per day?)
  - Isolate matrix plots to specific locations and/or time.

THANK YOU!