

Spatial Data Analytics with Classical Data Mining and Machine Learning Algorithms

Yao-Yi Chiang

Computer Science and Engineering

University of Minnesota

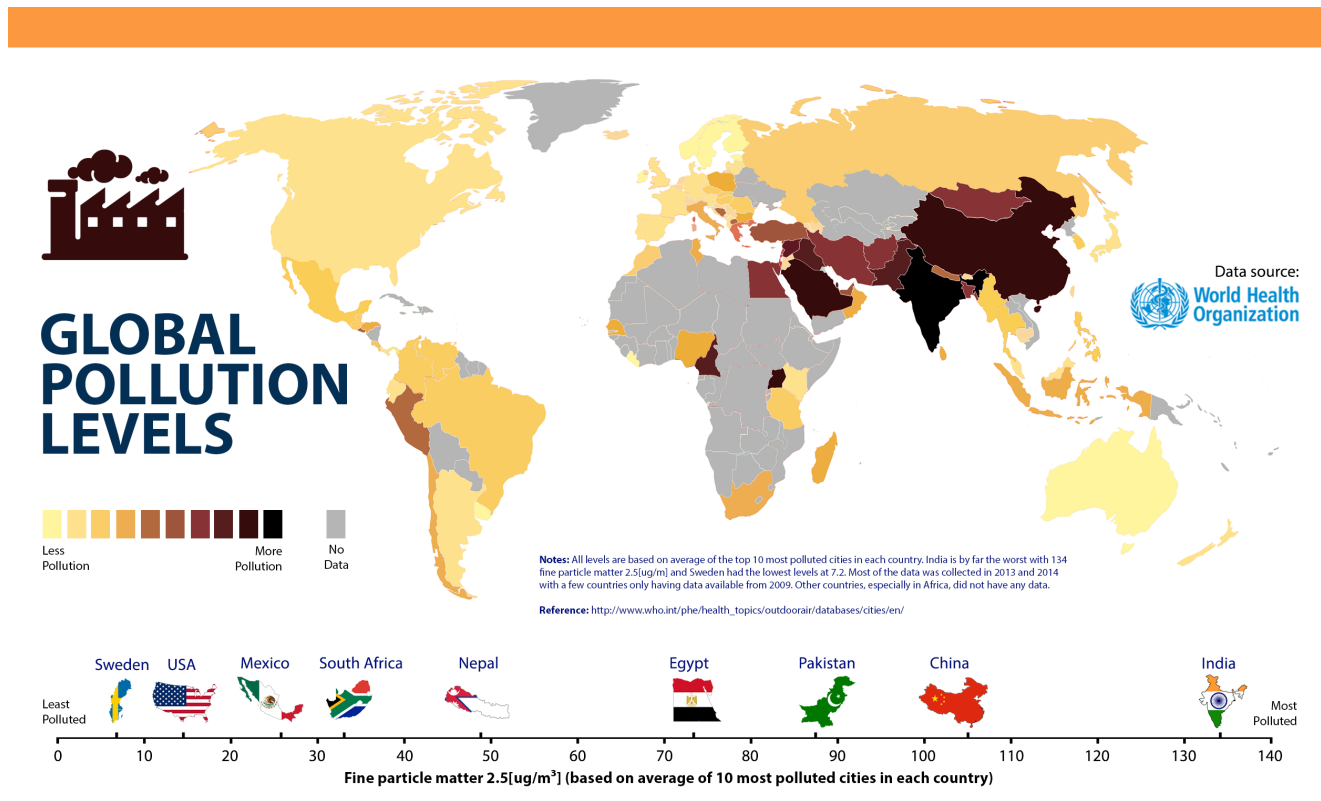
yaoyi@umn.edu



Mining Public Datasets for Modeling Intra-City PM_{2.5} Concentrations at a Fine Spatial Resolution

A motivating example

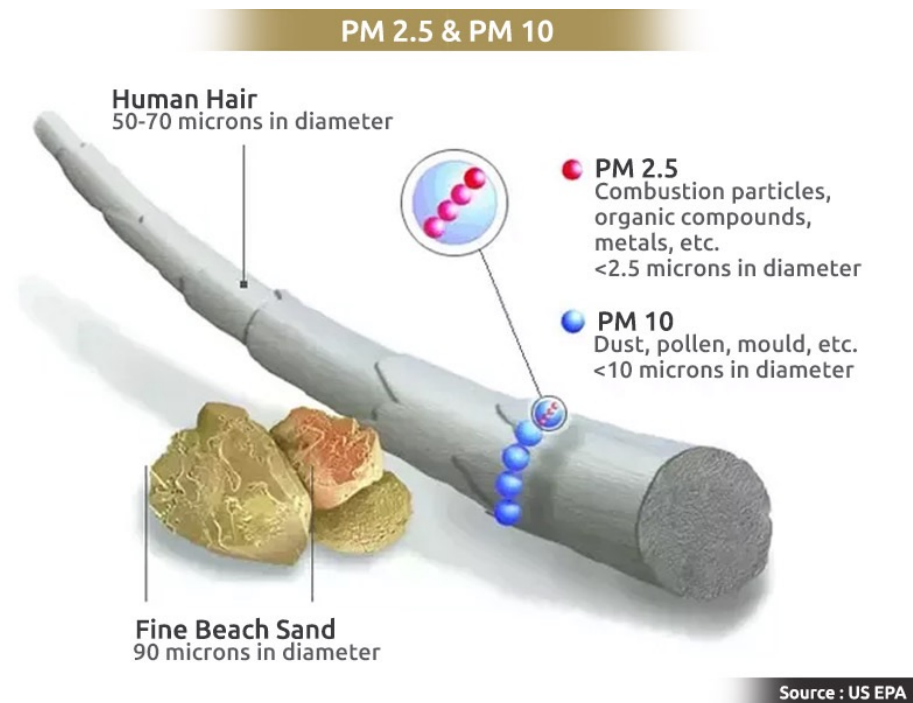
Air Pollution is a Global Problem



Air Pollutant: PM_{2.5} and PM₁₀

PM_{2.5} : fine inhalable particles, with diameters that are generally 2.5 micrometers and smaller

United States Environmental Protection Agency



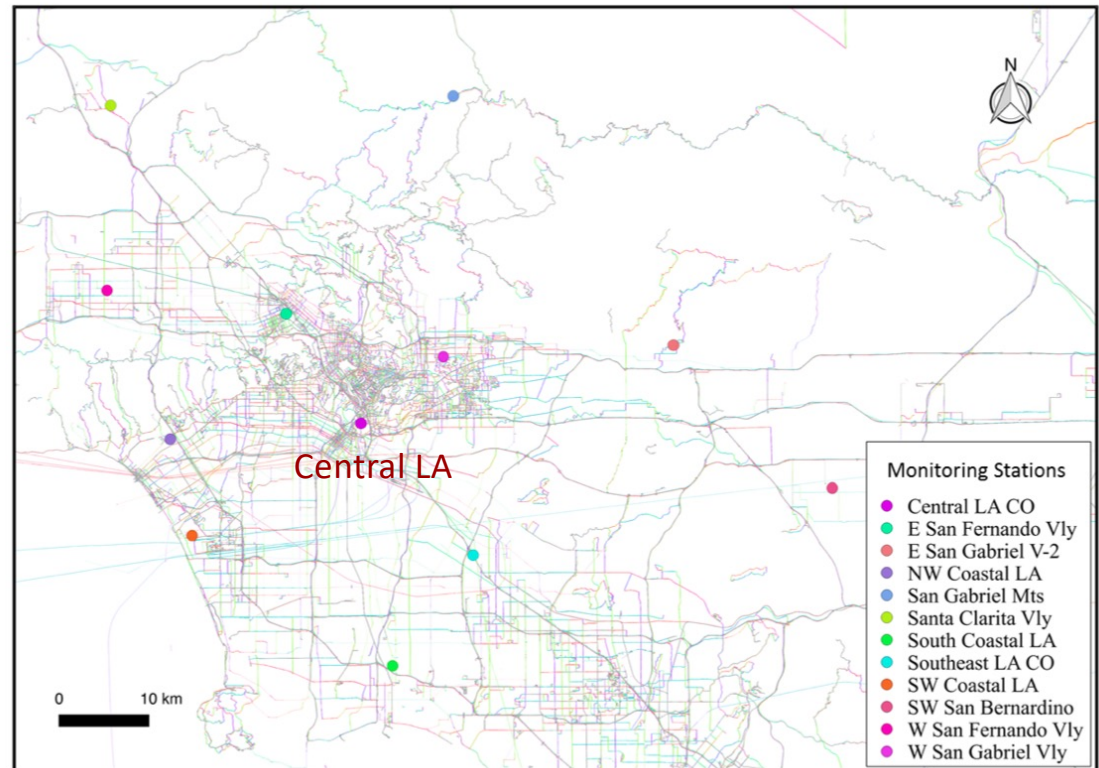
Air Quality Index

AQI: air quality index computed from a piecewise linear function of the pollutant concentration (e.g., 12.0 micrograms per cubic meter is 50 AQI for PM2.5).

Air Quality Index Levels of Health Concern	Numerical Value	Meaning
Good	0 to 50	Air quality is considered satisfactory, and air pollution poses little or no risk.
Moderate	51 to 100	Air quality is acceptable; however, for some pollutants there may be a moderate health concern for a very small number of people who are unusually sensitive to air pollution.
Unhealthy for Sensitive Groups	101 to 150	Members of sensitive groups may experience health effects. The general public is not likely to be affected.
Unhealthy	151 to 200	Everyone may begin to experience health effects; members of sensitive groups may experience more serious health effects.
Very Unhealthy	201 to 300	Health alert: everyone may experience more serious health effects.
Hazardous	301 to 500	Health warnings of emergency conditions. The entire population is more likely to be affected.

Limited Air Quality Observations

- Monitoring stations are usually sparse – 12 stations for PM_{2.5} in Los Angeles



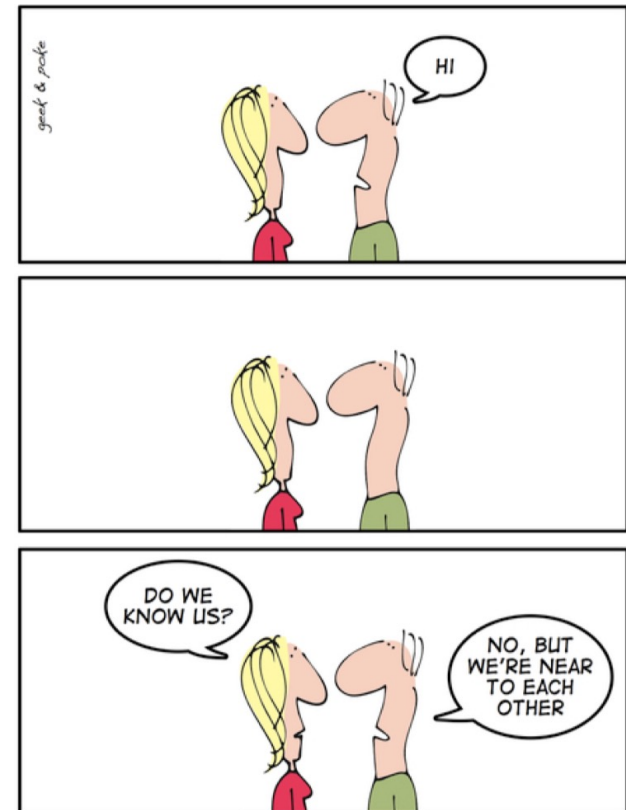
Nearby locations would have similar air quality

- Tobler's 1st law of geography:

“all things are related, but nearby things are more related than distant things”

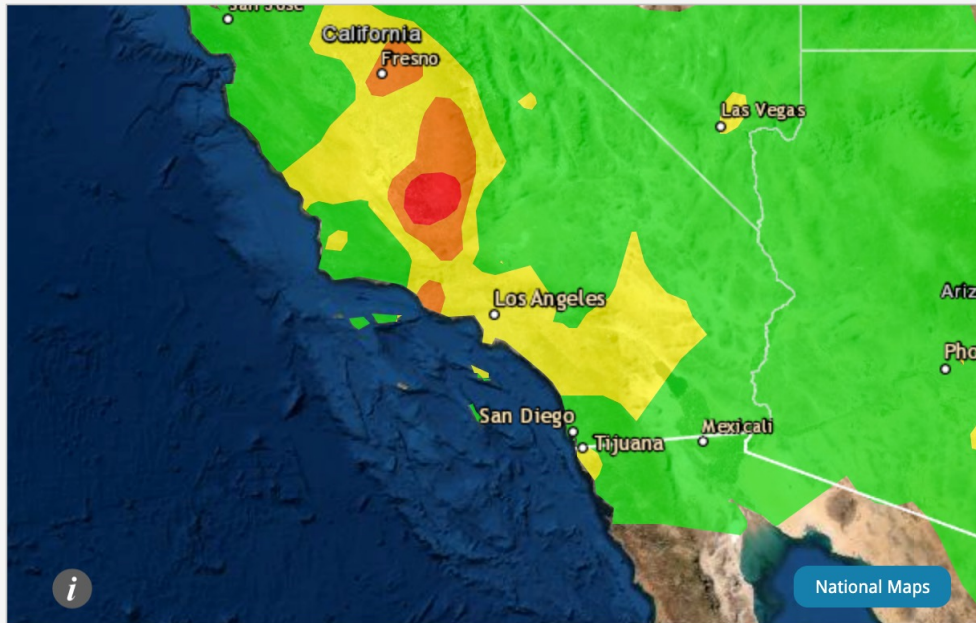
Tobler – 1970

- Spatial interpolation? IDW?



Typical Air Quality Prediction Result (AirNow)

Current Air Quality



Primary Pollutant
This pollutant currently has the highest AQI in the area.

▼ PM2.5 **59** Moderate

If you are unusually sensitive to particle pollution, consider reducing your activity level or shorten the amount of time you are active outdoors.

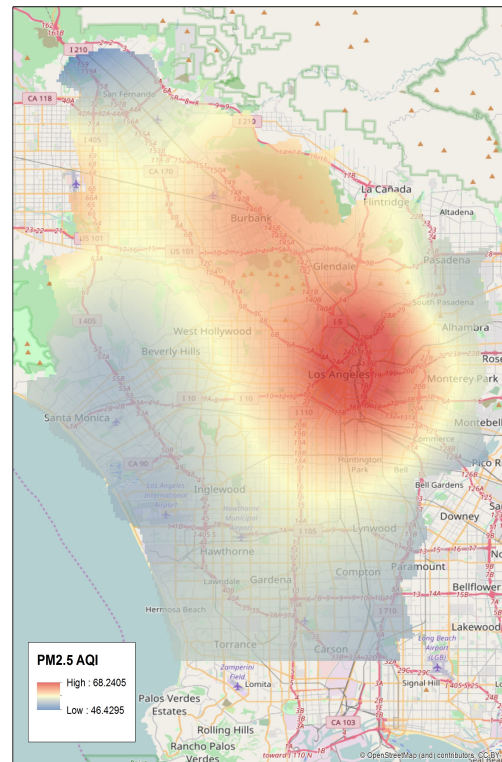
▶ PM10 **55** Moderate

▶ OZONE **31** Good

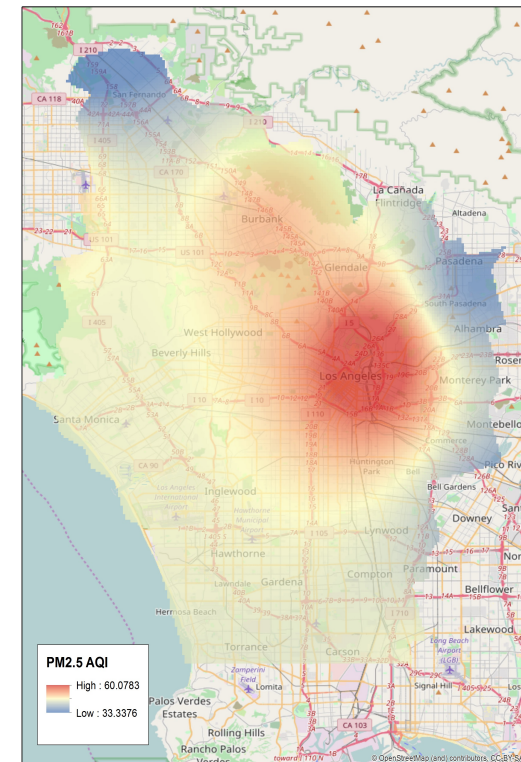
Inverse Distance Weighting

- Why are the results so smooth over space?
- Recall IDW:

$$u(\mathbf{x}) = \begin{cases} \frac{\sum_{i=1}^N w_i(\mathbf{x})u_i}{\sum_{i=1}^N w_i(\mathbf{x})} \\ u_i, \end{cases}$$



Dec 2016

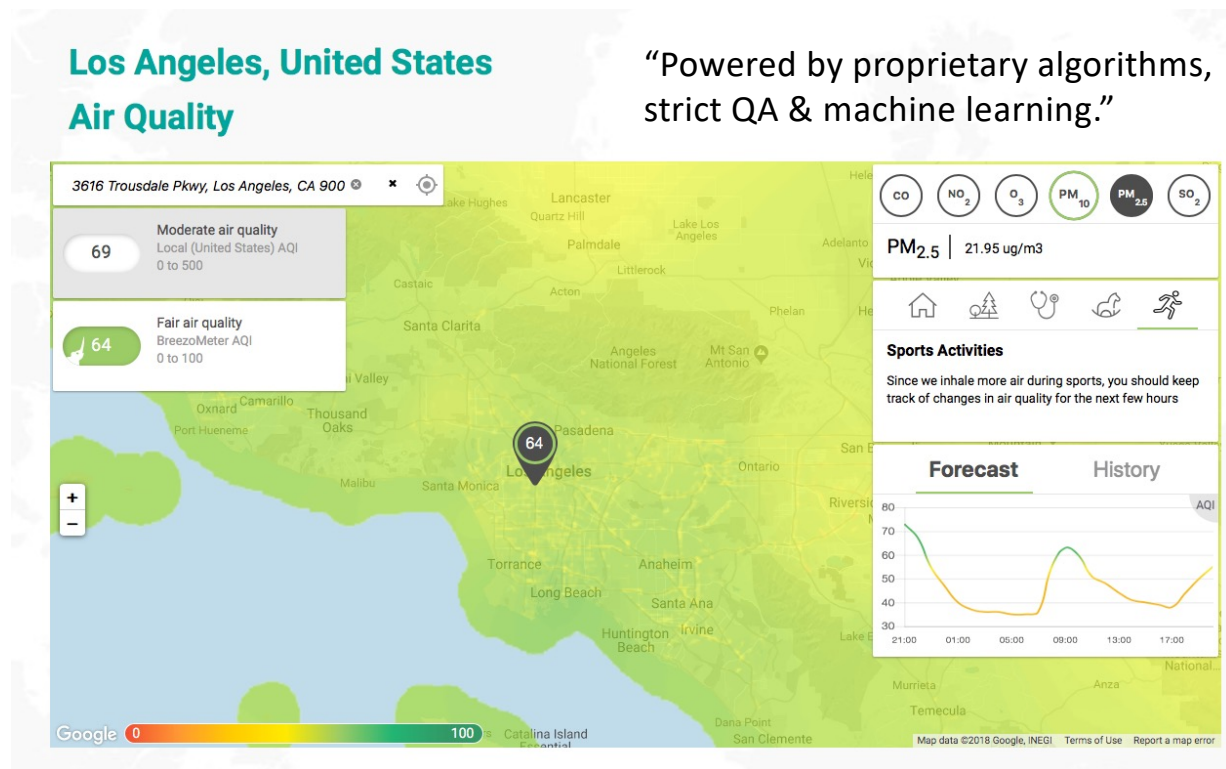


Jan 2017

Machine Learning Methods

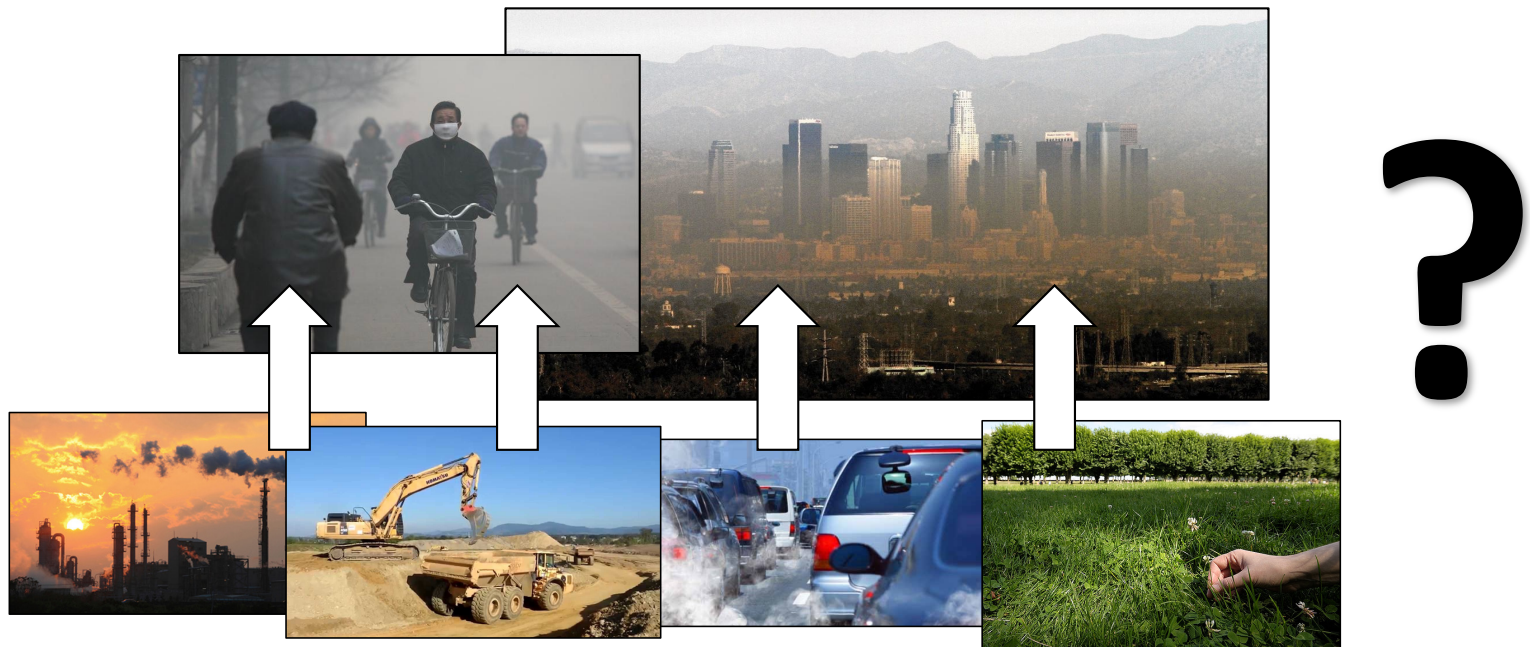
- Some prediction variations in space, e.g., the road network is obvious

“Powered by proprietary algorithms, strict QA & machine learning.”



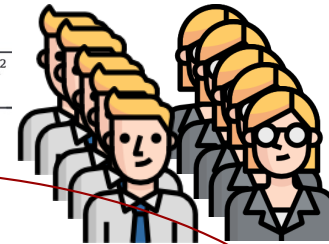
Existing Work for Air Quality Modeling

The built environment has a strong impact on air quality but *how*?



Land-use Regression Models (LUR)

Authors	Study area	Monitor counts	Dependent variables	Independent variables	Buffer size	(Adjusted) R ²
Briggs et al. (2000)	Huddersfield (UK) Sheffield (UK) Northampton (UK)	20, 28 and 35	NO ₂	Road traffic, urban land, and topography (altitudes)	300 m	0.58 to 0.76
Ross et al. (2007)	New York City (US)	28–49	PM _{2.5}	Traffic, land use, census	50, 100, 300, 500 and 1000 m	0.607 to 0.642
Su et al. (2008)	Greater Vancouver Regional District, (Canada)	116	NO/NO ₂	Road, traffic, meteorology (wind speed, wind direction and cloud cover/insolation)	3000 m	0.53 to 0.60
Mavko et al. (2008)	Portland, (US)	77	NO ₂	Traffic-related; Land use-related; Elevation; height from MSL; distance to a river; wind; direction	50, 100, 250, 300, 350, 400, 500, 750 m.	0.66 to 0.81
Rivera et al. (2012)	Girona province, (Spain)	25	Ultrafine particles (UFP)	Heavy, light and motorcy. veh in 24 h; 24 h total traffic load; length of major roads; <u>building density</u> ; distance to bus lines, highway and intersections; land cover	25, 50, 100, 150, 300, 500 and 1000 m	0.36 to 0.72
Eeftens et al. (2012)	20 European regions	20 per area	PM _{2.5} , PM ₁₀ and PMcoarse	Traffic intensity, <u>population</u> , and land-use	25, 50, 100, 300, 500, and 1000 m	0.35 to 0.94
Dons et al. (2013)	Flanders, (Belgium)	63	Traffic related air pollutant black carbon	Hourly traffic streams, daily traffic volumes, total road length; <u>population density and address density</u> ; land use variables	50, 100, 1000 m	0.44 to 0.77
Lee et al. (2014)	Taipei, (China Taiwan)	40	NO _x and NO ₂	Land use, no. of population and households, road length, altitude, distance to roads, <u>ports</u>	25, 25–50, and 50–500 m	0.63 to 0.81
Wu et al. (2015)	Beijing, (China)	35	PM _{2.5}	Traffic intensity, population, <u>bus stops, restaurants</u> , and land-use	100–3000 m	0.43 to 0.65

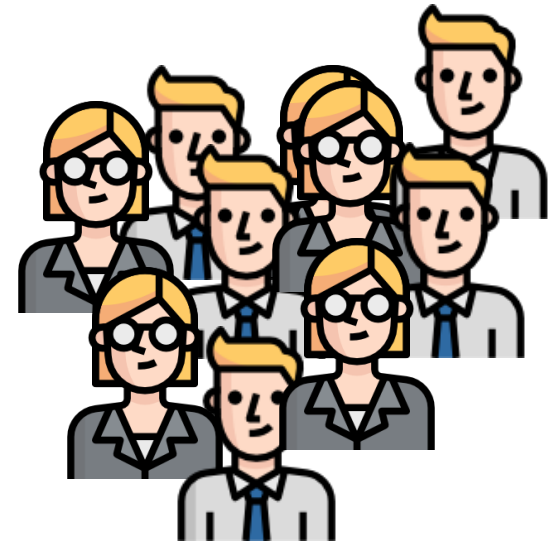


Expert-selected & Area-specific

e.g., PM_{2.5} concentrations is high near **500 meters** of **highways** in **Los Angeles**

LUR Limitations

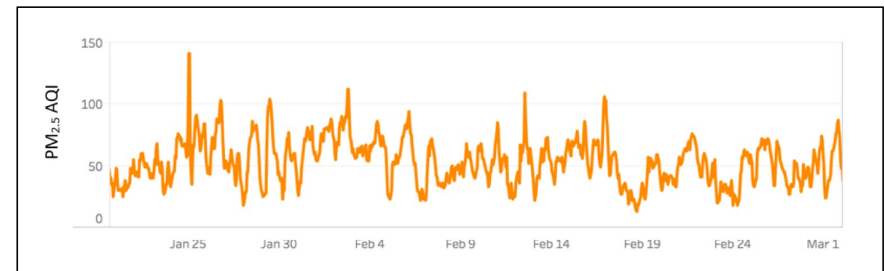
- Experts are expensive
- Do not scale well for predictions at various spatial and temporal resolutions
- Sometimes rely heavily on datasets that are not easy to obtain
 - e.g., traffic



Can we do better?

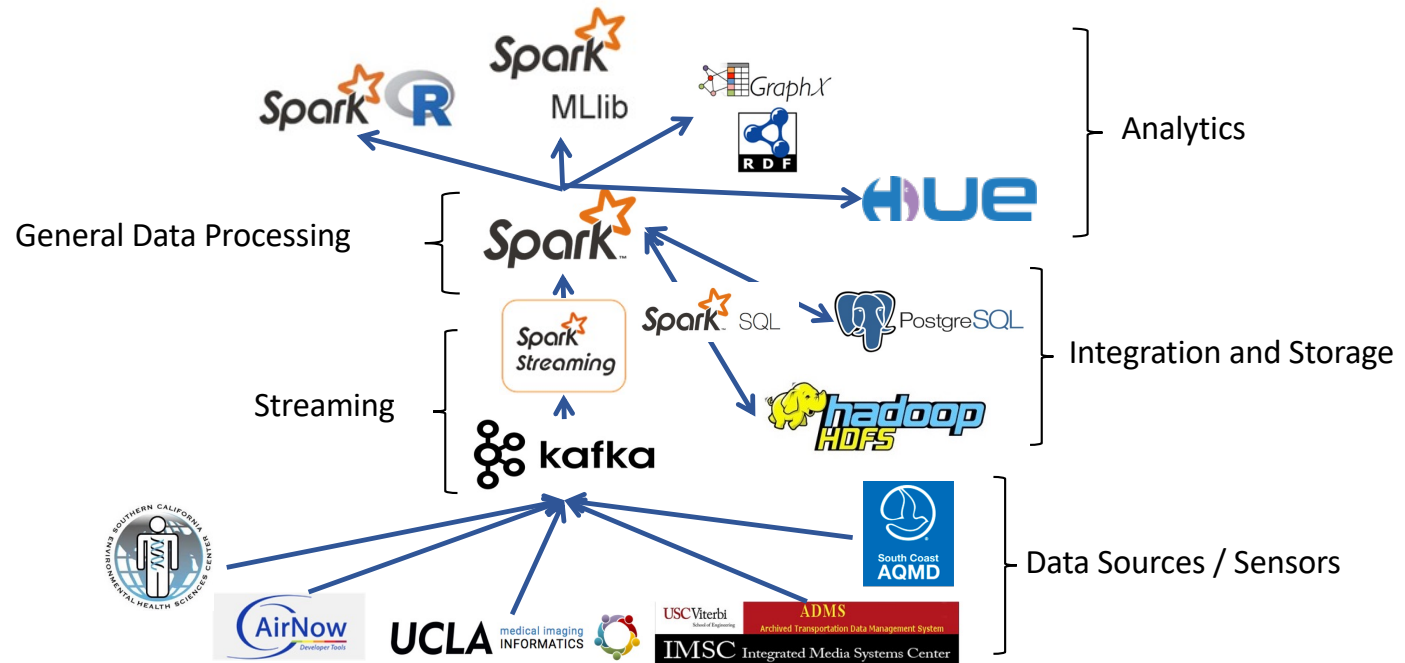
JonSnow: Data-Driven Air Quality Prediction at Fine-Spatial Scale

- Problem
 - Given **some sensors and their locations**, predicting **air quality for locations that do not have a sensor**
- Hypothesis
 - **Similar environments should have a similar air quality**



Data Collection

PRISMS-DSCIC – A scalable data integration and analysis architecture



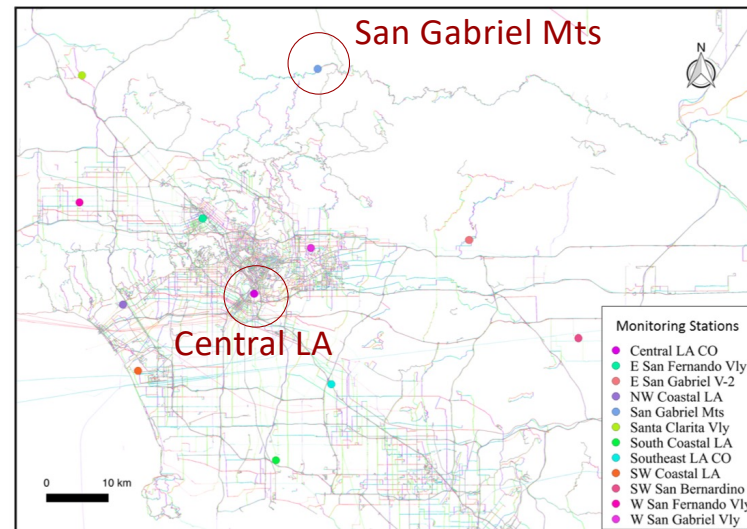
Data Sources – I

AQS (Air Quality System) Data

- Hourly PM_{2.5} AQI from **12 monitoring stations** in the Los Angeles Area from 2016-10-30 00:00:00 to 2017-08-31 23:00:00

Monitoring Station	Timestamp	PM _{2.5} AQI
San Gabriel Mts	2017-03-04 12:00:00	44
San Gabriel Mts	2017-03-04 13:00:00	54
Central LA	2017-03-04 12:00:00	60
Central LA	2017-03-04 13:00:00	68

Sample data



Data Sources – II

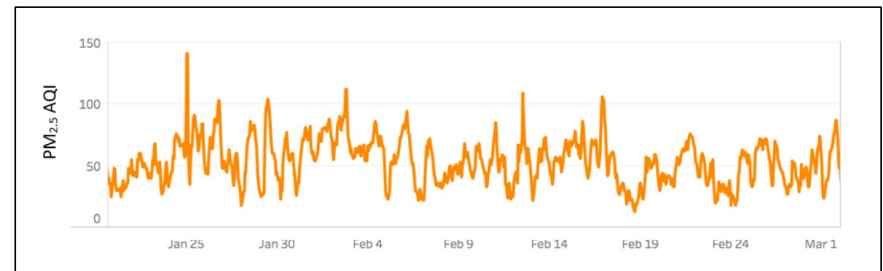
Geographic Features - OpenStreetMap (OSM)

- Land uses (67,972 polygons), Roads (544,142 lines), Water areas (11,207 polygons), Buildings (2,971,349 points), Aeroways (962 lines), etc.
- Each geographic category contains various feature subtypes
 - e.g., subtypes for “Buildings”: commercial, apartment, house, industrial, school, etc.



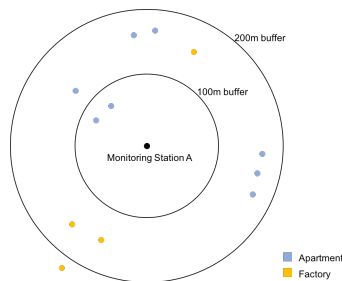
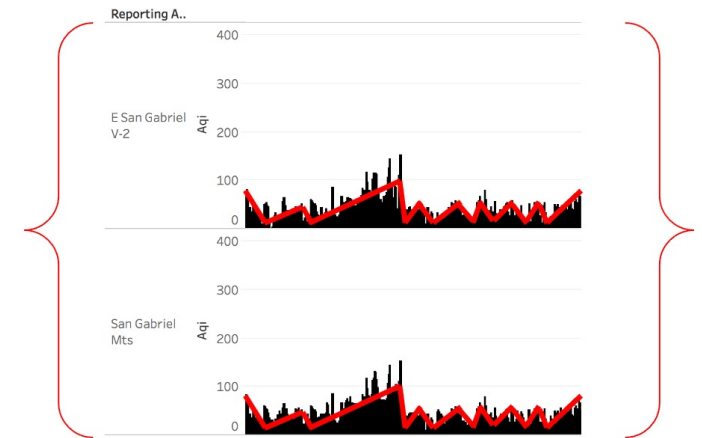
Recall: JonSnow: Data-Driven Air Quality Prediction at Fine-Spatial Scale

- Problem
 - Given **some sensors and their locations**, predicting **air quality for locations that do not have a sensor**
- Hypothesis
 - **Similar environments should have a similar air quality**



Approach Overview

- **Similar environments** should have a **similar air quality**
 - How to quantify “similar air quality”
 - Clustering of air quality measurements
 - K-Means, hierarchical clustering, dimension reduction
 - How to quantify “similar environments”
 - Train an **interpretable machine learning model** using **geographical context** to predict whether two locations would have “similar air quality”
 - Random Forest



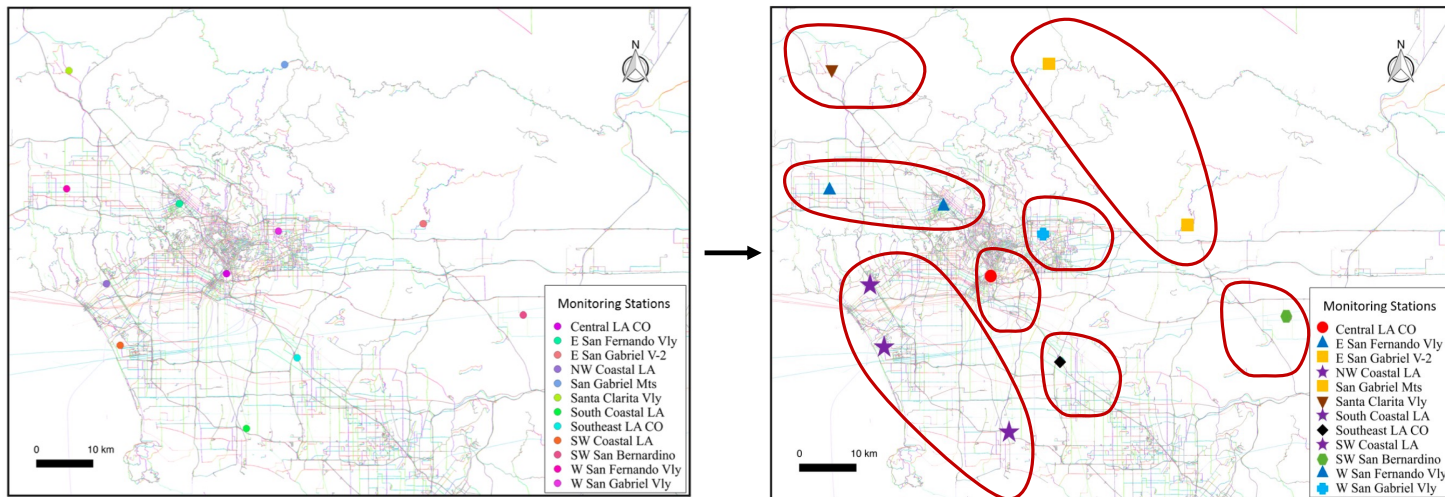
	Pedestrian 100m	Motorway 100m	... Apartment 200m	Factory 200m
Monitoring Station 0	0.1	5.28	2	0.3
Monitoring Station 1	0.0	3.27	... Apartment 200m	Factory 200m 0.432

Required Technologies

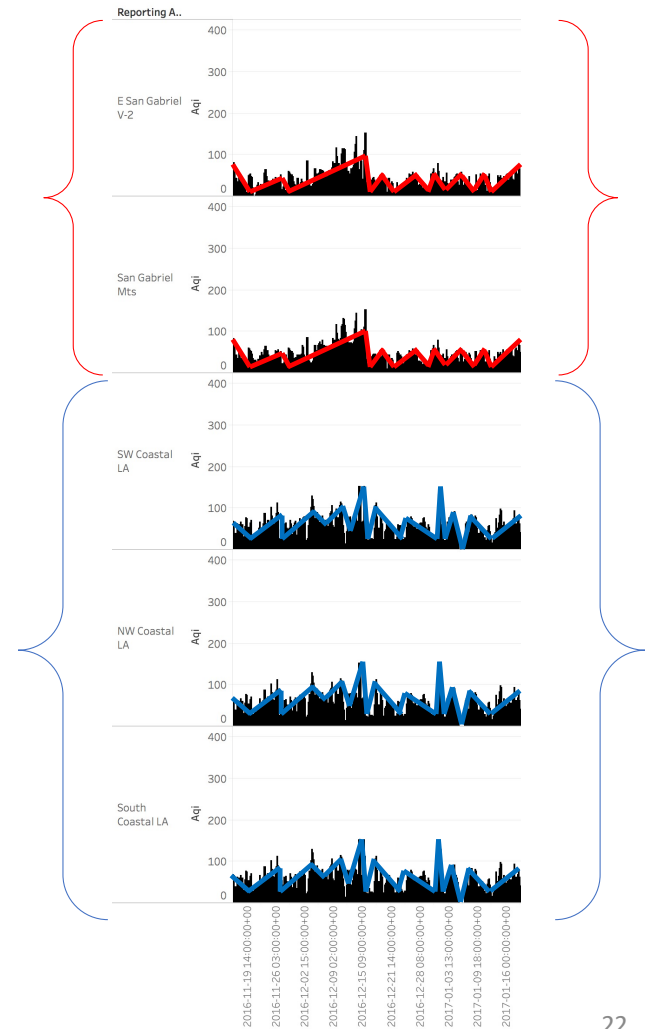
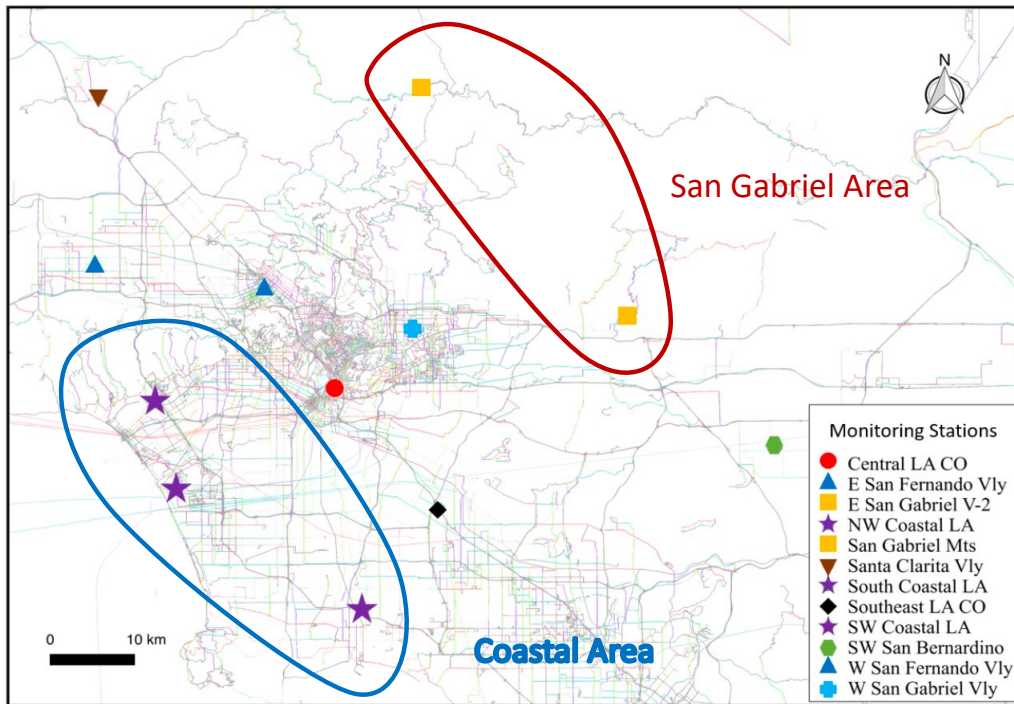
- Clustering
 - K-Means, Hierarchical Clustering
- Dimension Reduction
 - SVD (Singular Value Decomposition)
- Interpretable Machine Learning Method
 - Random Forest

Step 1. Grouping Stations based on their $PM_{2.5}$ AQIs

- To identify the monitoring stations that have **similar temporal pattern** on $PM_{2.5}$ AQIs
- These monitoring stations should have a similar environment.

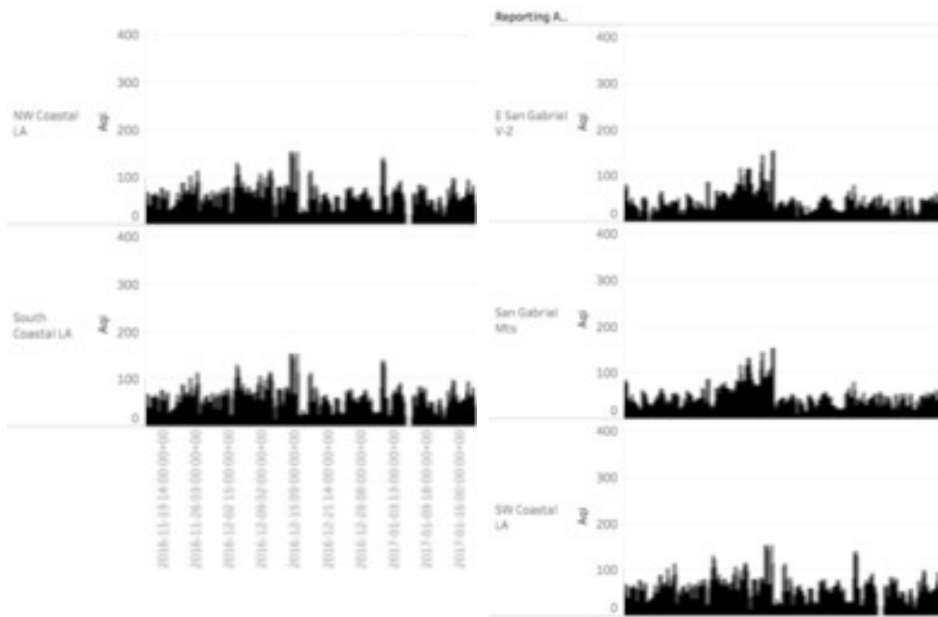


Similar Temporal Pattern

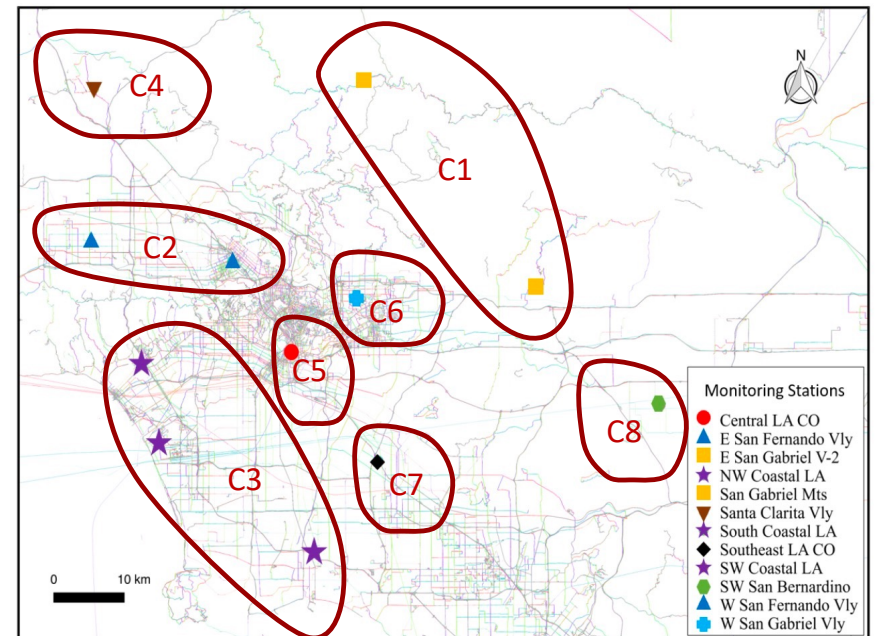


K-means Clustering

- Input: time-series observations at each station



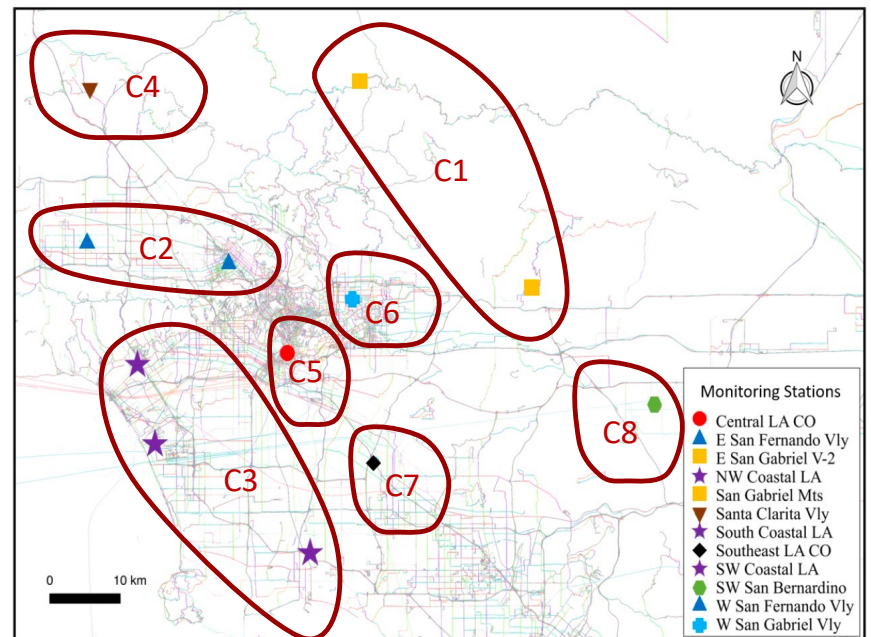
- Output: clusters of stations having a similar temporal pattern



K-means Clustering

- Recall: Hypothesis
 - Similar environments should have a similar air quality
- Stations in the same cluster have a similar temporal pattern
- How to quantify “similar environment”
 - what specific geographic feature types (e.g., primary roads, industrial areas, parks)
 - from what distance have the most impact on the clustering result?

- Output: clusters of stations having a similar temporal pattern



Step 2. Generating Geographic Abstraction

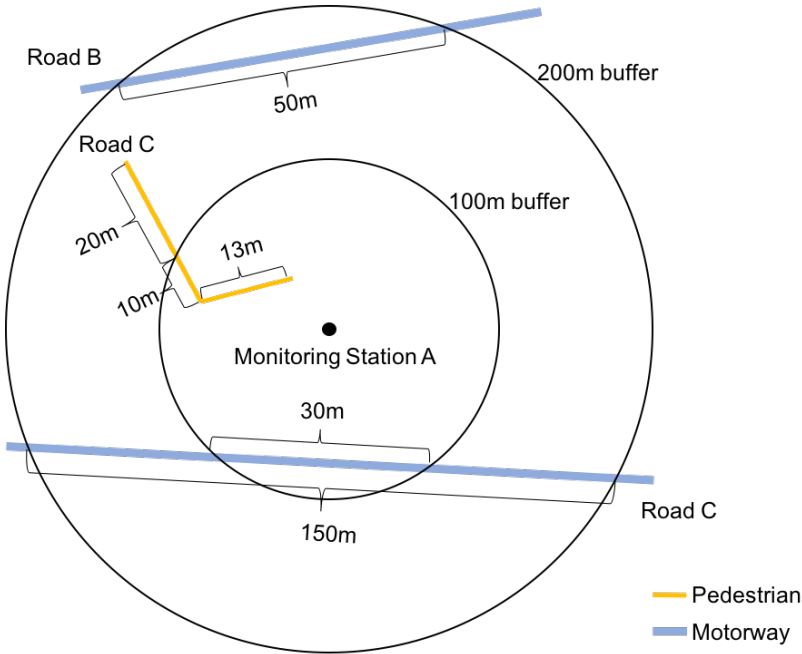
Length of line features

- e.g., Roads, Aeroways

Example Roads

	100m	200m
Pedestrian	23	43
Motorway	30	200

↓
[23, 30, 43, 200]



Step 2. Generating Geographic Abstraction

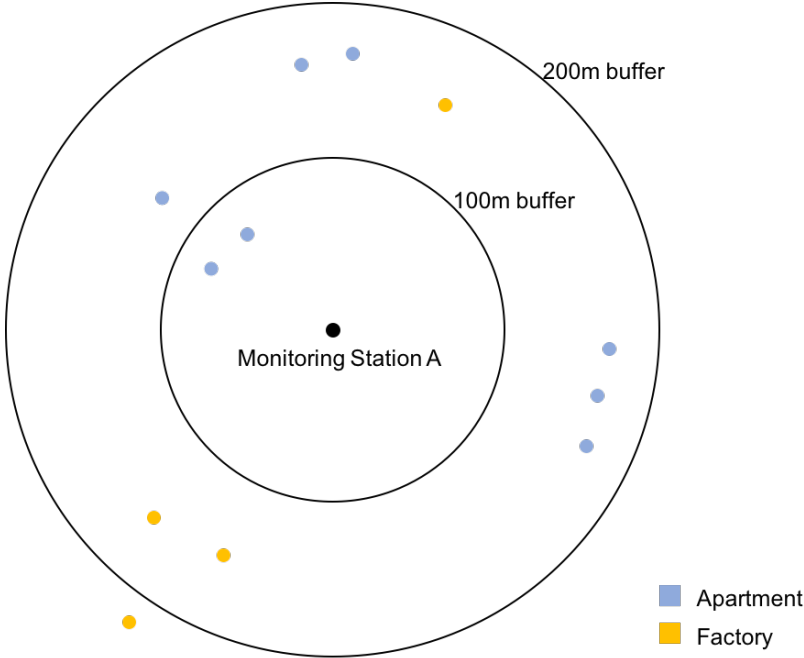
Count of point features

- e.g., Buildings

Example Buildings

	100m	200m
Apartment	2	8
Factory	0	3

↓
[2, 0, 8, 3]



Step 2. Generating Geographic Abstraction

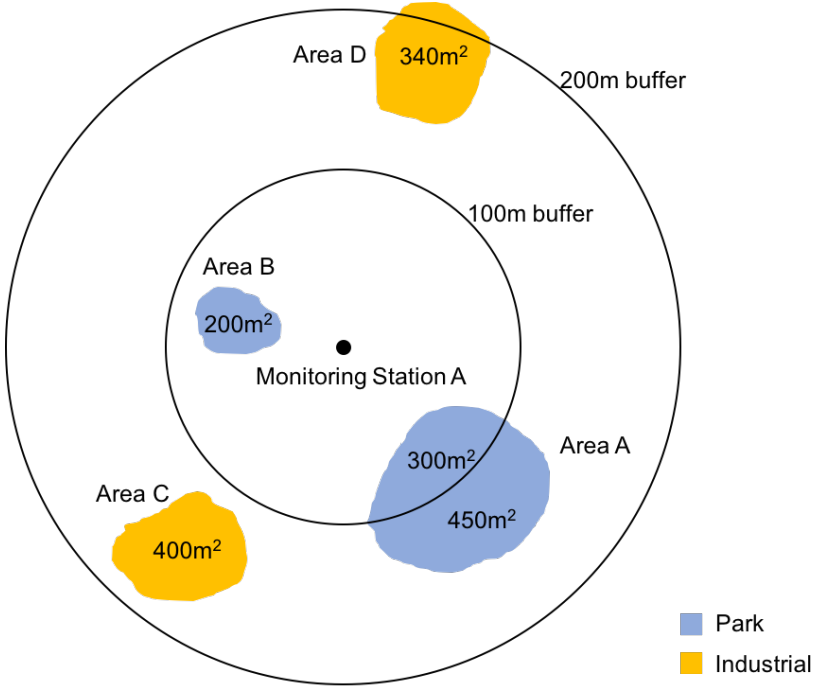
Area of polygon features

- e.g., Land uses, Water areas

Example Land uses

	100m	200m
Park	500	950
Industrial	0	740

↓
[500, 0, 950, 740]



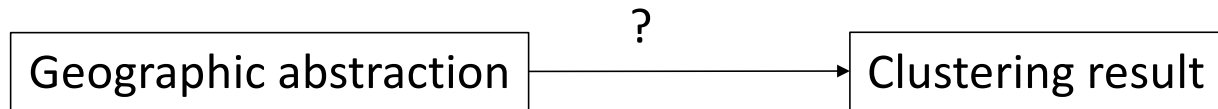
Step 2. Generating Geographic Abstraction

- Generating a large vector for each monitoring station

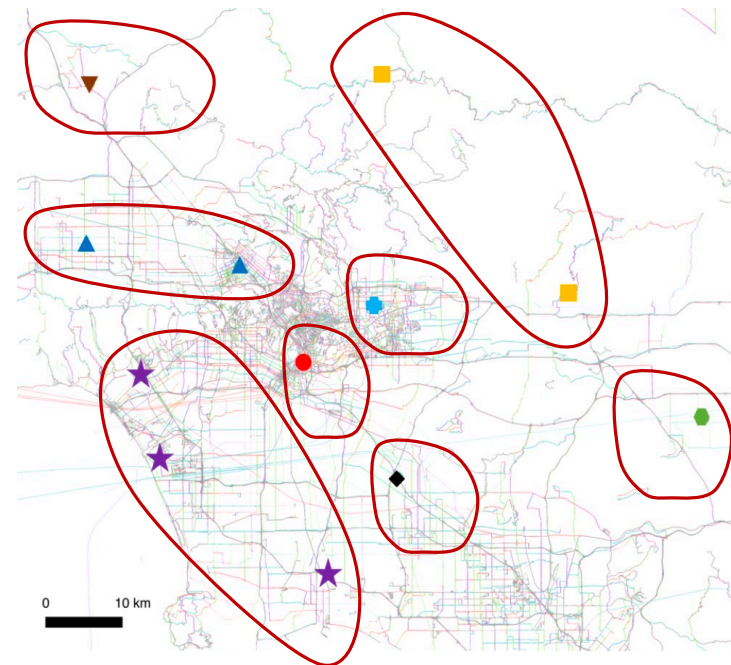
		Pedestrian 100m	Motorway 100m	Pedestrian 200m	Motorway 200m	Park 100m	Industrial 100m	
Monitoring Station X	[23	30	43	200	500	0	
		Park 200m	Industrial 200m	Apartment 100m	Factory 100m	Apartment 200m	Factory 200m	Distance to Ocean
		950	740	2	0	8	3	4000]

- In practice, we create buffers from 100 meters to 3,000 meters with an interval of 100 meters
 - 3,500+ components in a vector

How to quantify “similar environment”

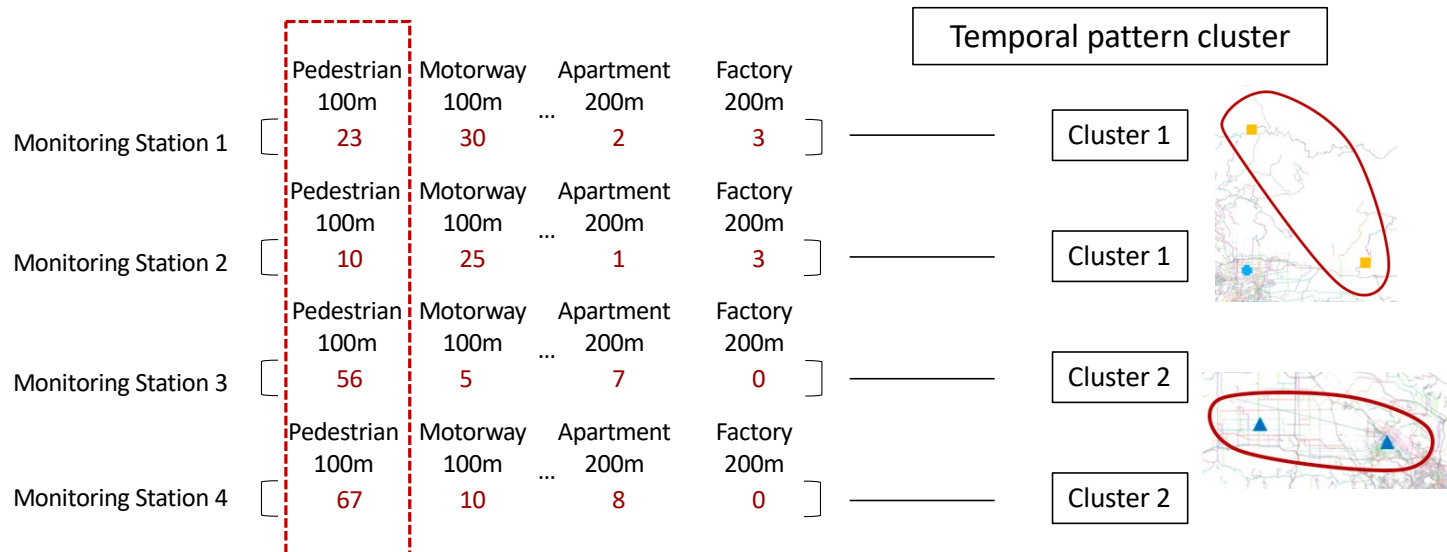


	Pedestrian 100m	Motorway 100m	Apartment ... 200m	Factory 200m
Monitoring Station 1	23	30	2	3
Monitoring Station 2	10	25	1	3
Monitoring Station 3	56	100	8	0
.....				
Monitoring Station 12	67	10	8	0



Step 3. Computing Feature Importance

- Training a **random forest model** to
 - predict cluster label using the geographic context
 - each feature component represents a **geographic feature type** within **certain distance**
 - quantify the **impact** of **each feature component**



Step 3. Generating Geo-context

- Multiplying each geographic abstraction value by its feature importance to generate geo-context

$$\text{Geographic Abstraction Vector } \mathbf{A} = [a_1, a_2, \dots, a_n]$$

$$\text{Importance Vector } \mathbf{I} = [i_1, i_2, \dots, i_n]$$

$$\text{Geo-Context Vector } \mathbf{C} = \mathbf{A} * \mathbf{I}$$

Monitoring Station 1 (Geographic Abstraction)	[Pedestrian 100m	Motorway 100m	... Apartment 200m	Factory 200m]
		23	30	2	3	
↓						
Monitoring Station 1 (Geo-context)	[Pedestrian 100m	Motorway 100m	... Apartment 200m	Factory 200m]
		0.0	3.27	0.041	0.432	

Example of Importance

Geo-feature	Importance
Pedestrian 100m	0.000
Motorway 100m	0.109
...	...
Apartment 200m	0.041
Factory 200m	0.144
...	...
Total	1.0

Step 3. Geo-context

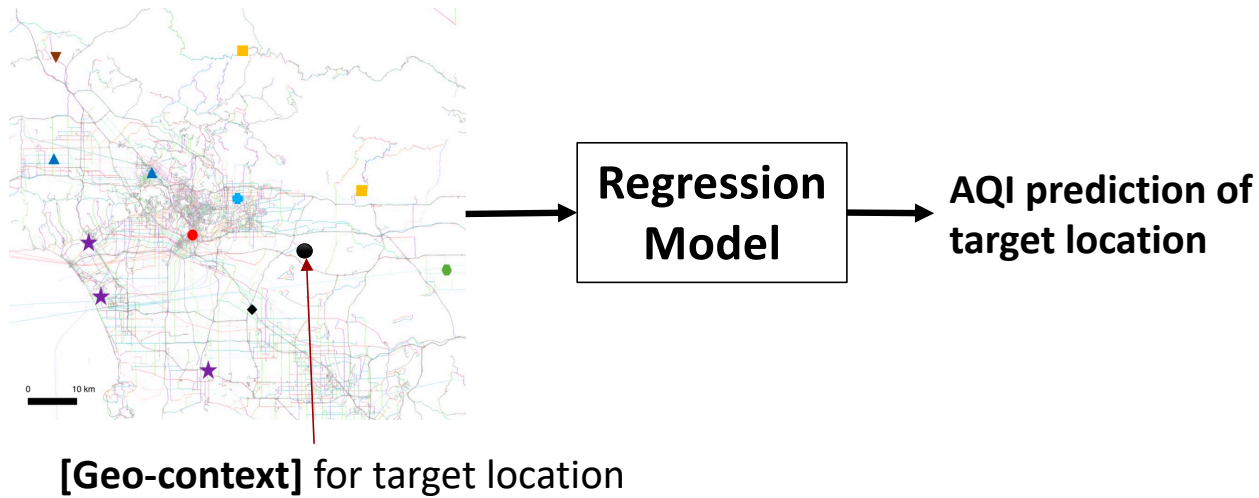
- Geo-context is an updated vector from geo-abstract for describing
 - how **each feature type** within **a certain distance** (a feature component) in **Geographic Abstraction** affects the **Temporal Pattern** (PM_{2.5} AQI)
- Reward important (relevant) features and penalize others

Monitoring Station 1 (Geographic Abstraction)	Pedestrian 100m	Motorway 100m	... Apartment 200m	Factory 200m
	23	30	2	3
Monitoring Station 1 (Geo-context)	Pedestrian 100m	Motorway 100m	... Apartment 200m	Factory 200m
	0.0	3.27	0.041	0.432

Step 4. Predicting $PM_{2.5}$ AQI

Train a regression model to predict $PM_{2.5}$ AQI for a target location at time T

[Geo-context, AQI] for each monitoring station at time T



Experiments

Leave-one-out cross-validation method

- Predict PM_{2.5} AQI for the removed station by using other 11 stations
- Compare our approach with baseline methods

Predicting **at a fine scale**

- Predict PM_{2.5} AQI of each point on an 1-mile-apart fishnet covering most of the Los Angeles area (604 points)
- Visualize the fine-scale prediction results

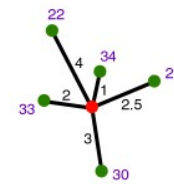
Experiment & Result – I

Leave-one-out cross-validation method

- Tested with **three methods** on **three temporal scales**
 - Geo-context, Geo-abstraction, IDW (Inverse distance weighting)
 - **Monthly** (7 months), **daily** (233 days), and **hourly** (168 hours)
 - **RMSE** - root-mean-square error; **MAE** - mean absolute error

	<i>Geo – context</i>	<i>Geo – Abstraction</i>	<i>IDW</i>
<i>RMSE (Monthly)</i>	<u>2.53984</u>	2.62391	2.88263
<i>MAE (Monthly)</i>	<u>1.86657</u>	1.93673	2.18675
<i>RMSE (Daily)</i>	4.33786	4.35857	<u>4.10172</u>
<i>MAE (Daily)</i>	3.26140	3.28176	<u>3.10185</u>
<i>RMSE (Hourly)</i>	7.38823	7.59260	<u>6.66106</u>
<i>MAE (Hourly)</i>	5.06559	5.12406	<u>4.54779</u>

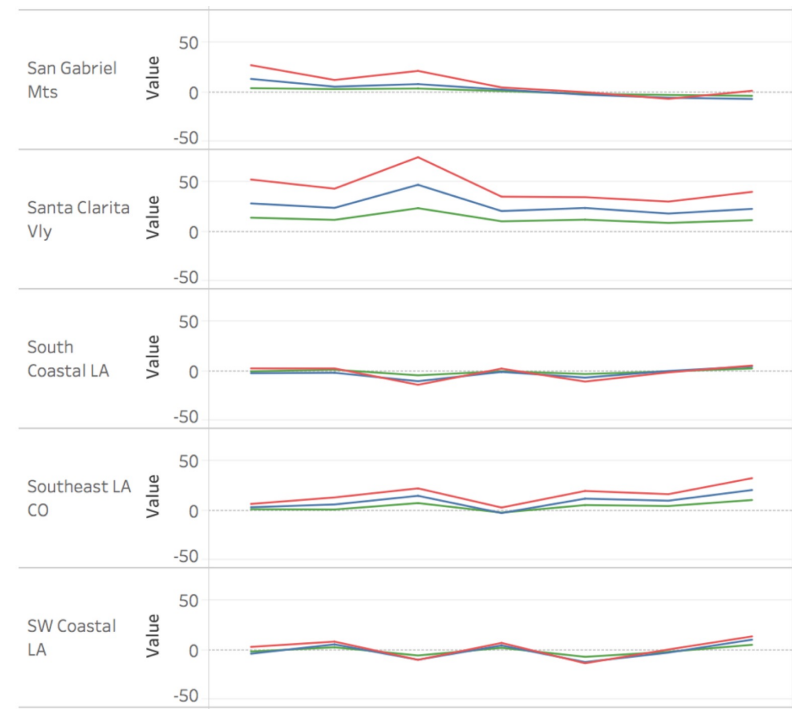
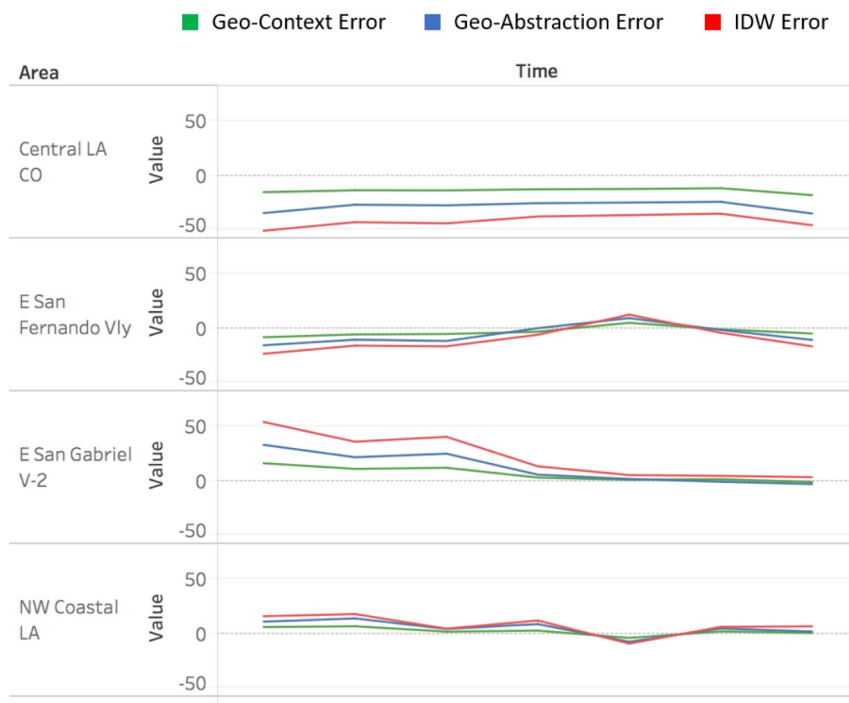
IDW method



$$Z(x) = \frac{\sum w_i z_i}{\sum w_i} = \frac{\frac{34}{1^2} + \frac{33}{2^2} + \frac{27}{2.5^2} + \frac{30}{3^2} + \frac{22}{4^2}}{\frac{1}{1^2} + \frac{1}{2^2} + \frac{1}{2.5^2} + \frac{1}{3^2} + \frac{1}{4^2}} = 32.38$$

All within 10% error margin; Significant different with 95% confidence (paired t-test)

Experiment & Result – I (Cont'd)



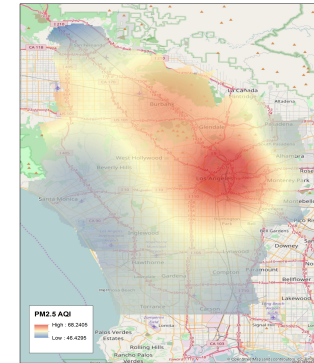
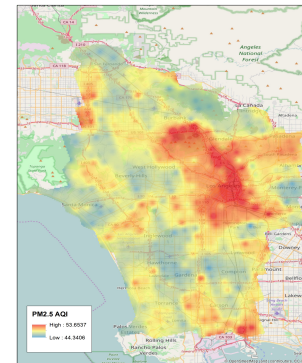
Experiment & Result – II

Predicting PM_{2.5} AQIs at a fine scale

<i>Geo Name</i>	<i>Buffer Size (meter)</i>	<i>Geo type</i>	<i>Importance (%)</i>
<i>land use</i>	1100	<i>wetland</i>	0.0051177
<i>land use</i>	1300	<i>university</i>	0.004450
<i>road</i>	600	<i>rail</i>	0.0044327
<i>land use</i>	1200	<i>village_green</i>	0.0037241
<i>road</i>	700	<i>primary</i>	0.0035520
<i>land use</i>	1900	<i>farmland</i>	0.0031458
<i>land use</i>	2700	<i>village_green</i>	0.0030063
<i>road</i>	800	<i>residential</i>	0.0028980
<i>building</i>	2000	<i>retail</i>	0.0027980
<i>building</i>	900	<i>industrial</i>	0.0027576
<i>road</i>	500	<i>tertiary</i>	0.0027357
<i>land use</i>	900	<i>pitch</i>	0.0026613
<i>building</i>	2900	<i>school</i>	0.0025681
<i>building</i>	1700	<i>garages</i>	0.0025361
<i>road</i>	1300	<i>motorway</i>	0.0023724

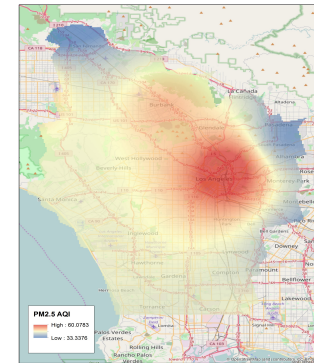
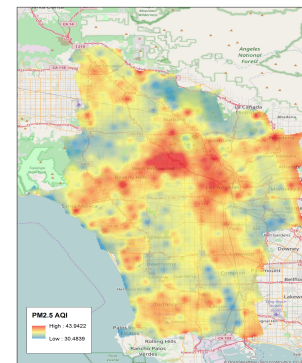
Geo-context

IDW



Dec 2016

Dec 2016



Jan 2017

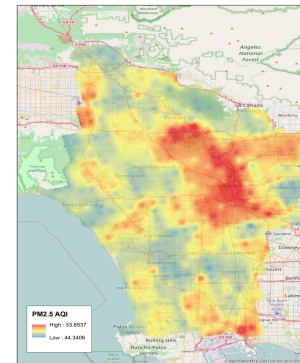
Jan 2017

Experiment & Result – II

Predicting PM_{2.5} AQIs at a fine scale

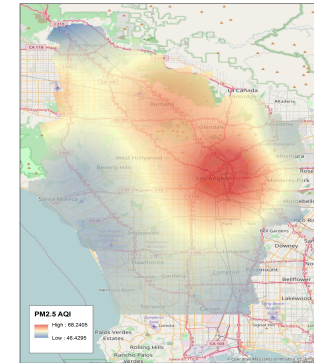
<i>Geo Name</i>	<i>Buffer Size (meter)</i>	<i>Geo type</i>	<i>Importance (%)</i>
<i>land use</i>	1100	<i>wetland</i>	0.0051177
<i>land use</i>	1300	<i>university</i>	0.004450
<i>road</i>	600	<i>rail</i>	0.0044327
<i>land use</i>	1200	<i>village_green</i>	0.0037241
<i>road</i>	700	<i>primary</i>	0.0035520
<i>land use</i>	1900	<i>farmland</i>	0.0031458
<i>land use</i>	2700	<i>village_green</i>	0.0030063
<i>road</i>	800	<i>residential</i>	0.0028980
<i>building</i>	2000	<i>retail</i>	0.0027980
<i>building</i>	900	<i>industrial</i>	0.0027576
<i>road</i>	500	<i>tertiary</i>	0.0027357
<i>land use</i>	900	<i>pitch</i>	0.0026613
<i>building</i>	2900	<i>school</i>	0.0025681
<i>building</i>	1700	<i>garages</i>	0.0025361
<i>road</i>	1300	<i>motorway</i>	0.0023724

Geo-context

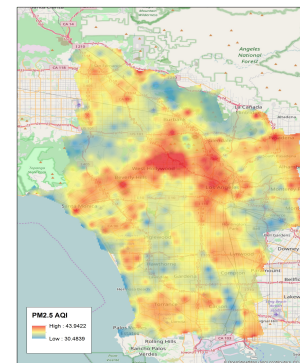


Dec 2016

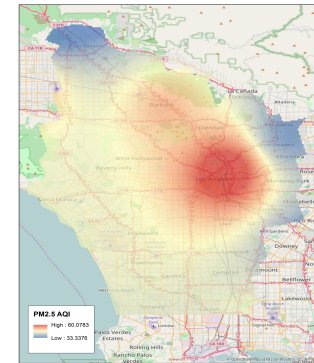
IDW



Dec 2016



Jan 2017



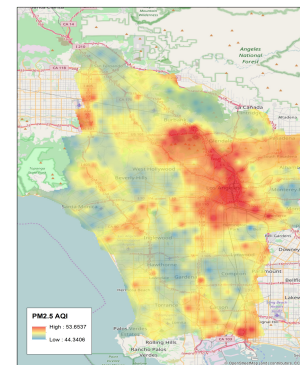
Jan 2017

Experiment & Result – II

Predicting PM_{2.5} AQIs at a fine scale

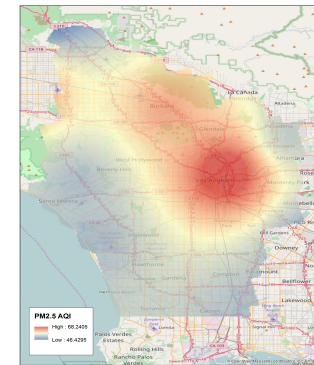
<i>Geo Name</i>	<i>Buffer Size (meter)</i>	<i>Geo type</i>	<i>Importance (%)</i>
<i>land use</i>	1100	<i>wetland</i>	0.0051177
<i>land use</i>	1300	<i>university</i>	0.004450
<i>road</i>	600	<i>rail</i>	0.0044327
<i>land use</i>	1200	<i>village_green</i>	0.0037241
<i>road</i>	700	<i>primary</i>	0.0035520
<i>land use</i>	1900	<i>farmland</i>	0.0031458
<i>land use</i>	2700	<i>village_green</i>	0.0030063
<i>road</i>	800	<i>residential</i>	0.0028980
<i>building</i>	2000	<i>retail</i>	0.0027980
<i>building</i>	900	<i>industrial</i>	0.0027576
<i>road</i>	500	<i>tertiary</i>	0.0027357
<i>land use</i>	900	<i>pitch</i>	0.0026613
<i>building</i>	2900	<i>school</i>	0.0025681
<i>building</i>	1700	<i>garages</i>	0.0025361
<i>road</i>	1300	<i>motorway</i>	0.0023724

Geo-context

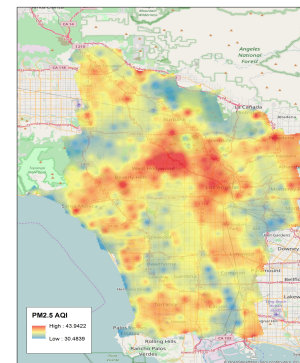


Dec 2016

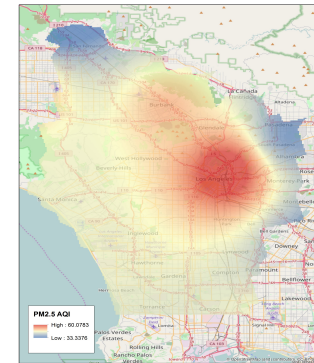
IDW



Dec 2016



Jan 2017



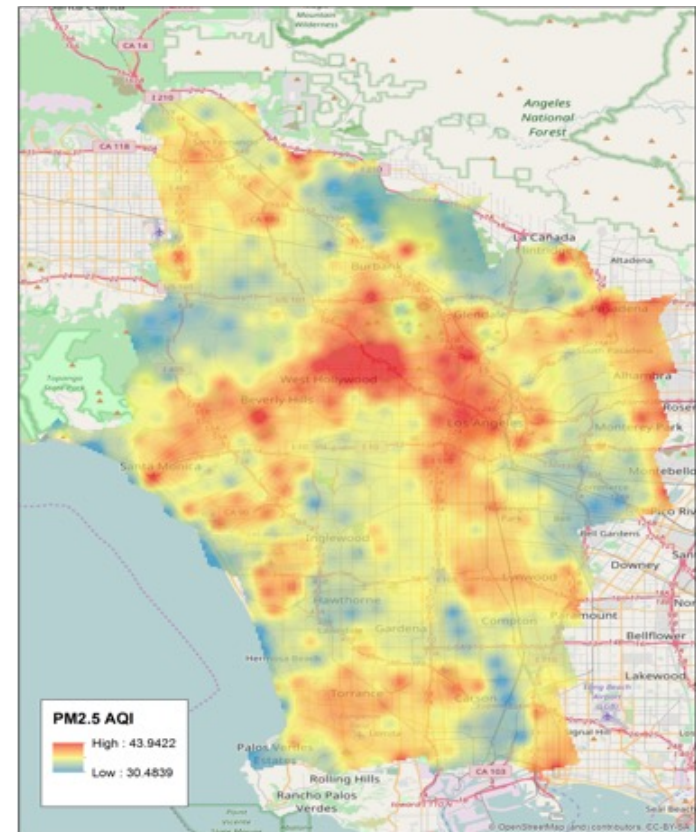
Jan 2017

Related Work

	Limitations	Advantages of our method
Spatial interpolation methods, e.g., IDW and Kriging	Not considering neighborhood characteristic	With neighboring geographic features
	Cannot generate a fine scale result with sparse monitoring stations	Can generate accurate result in a fine scale
Dispersion models	Require detailed data (e.g., building heights and distance between neighboring buildings)	Use easily accessible datasets (OpenStreetMap)
Land-use regression (LUR) methods (e.g., Hoek (2008))	Rely on expert-selected predictors, including types and spatial radii	Expert-free feature selection

Summary

- A spatial data mining approach to build an accurate model to predict $PM_{2.5}$ concentrations at a fine scale by
- **Automated selection of important geographic features** without using expert knowledge.



Additionally, Air Quality Forecasting

Goal

Build a general approach for **location-dependent time-series data** forecasting

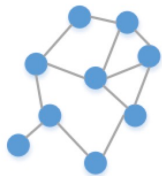
Challenges:

Existing approaches do not handle spatial correlation well
e.g., Auto-Regression Integrated Moving Average (ARIMA), Kalman filtering,
Artificial Neural Network (ANN)

Our approach

- We are building a Diffusion Convolutional Recurrent Neural Network for forecasting location-dependent time series data.
- Continuously forecasting air quality index (AQI) in next 24 hours at a fine scale using data on the PRISMS-DSCIC

DCRNN – Diffusion Convolutional Recurrent Neural Network



Graph Construction

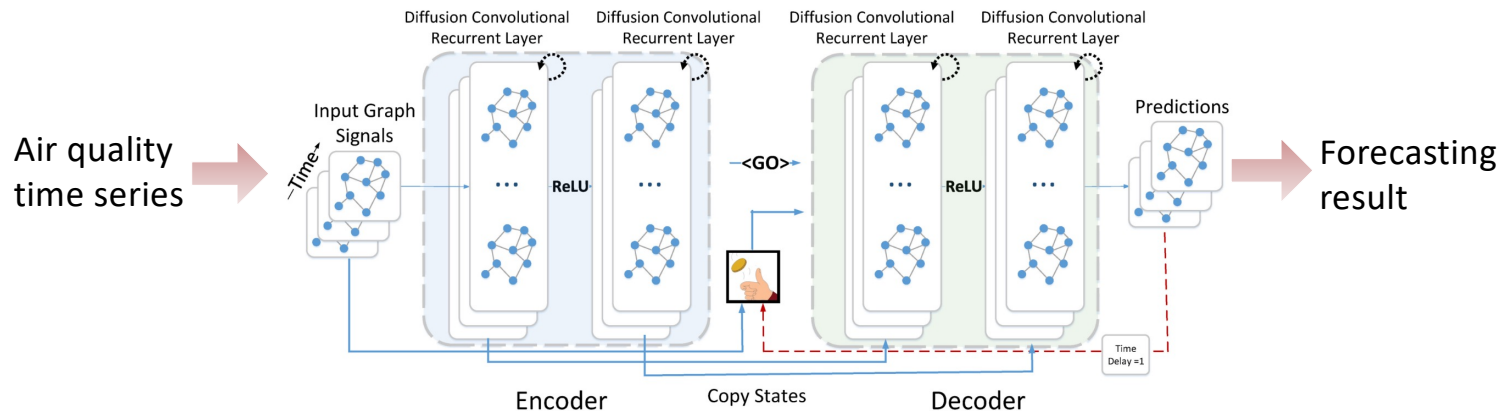
- Each point in the graph represents the time series at the station
- The link between points would be the proximity between stations (e.g., distance, geographic similarity)

Spatial Dependency Modeling

- Use diffusion convolution to learn a function that maps historical graph signal to future graph signal

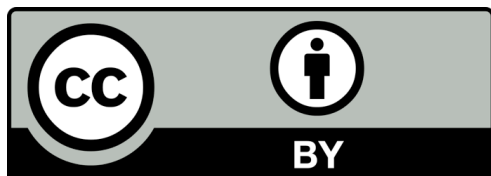
Temporal Dependency Modeling

- Use Recurrent Neural Networks



Acknowledgements

- Gil, Yolanda (Ed.) Introduction to Computational Thinking and Data Science. Available from <http://www.datascience4all.org>



<https://creativecommons.org/licenses/by/2.0/>

These materials are released under a CC-BY License

You are free to:

Share — copy and redistribute the material in any medium or format

Adapt — remix, transform, and build upon the material for any purpose, even commercially.

The licensor cannot revoke these freedoms as long as you follow the license terms.

Under the following terms:

Attribution — You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.

*Artwork taken from other sources is acknowledged where it appears.
Artwork that is not acknowledged is by the author.*

Please credit as: Chiang, Yao-Yi Introduction to Spatial Artificial Intelligence. Available from <https://yaoyichi.github.io/spatial-ai.html>

If you use an individual slide, please place the following at the bottom: “Credit: <https://yaoyichi.github.io/spatial-ai.html>”

We welcome your feedback and contributions.