



# Combining GIS and Statistics: Data Visualization to Communicate Findings

Kristen Hocutt

Lu Zhang

March 7, 2019

Washington Statistical Society

WASHINGTON, DC

**Lakers Playoff Game**  
@Staples Center

**Cyber attack**  
@Washington DC

**Power Outage**  
@Eugene area locations

**HR Incident**  
@Scranton, PA

**Conference**  
@Phoenix City Center

**Everything  
happens somewhere**

**Armed Protester**  
@Seattle Distribution  
Center

**New Starbucks**  
@Downtown Phoenix

**Amazon Prime**

**Machine down**  
@San Francisco P&DC

**Route Disruption**  
@Reno & Irvine Locations

**Winter Weather**  
@Washington DC

**Nearest Food truck**  
Landing @ Salt Lake International Airport

What are

Spatial

Statistics?

**Spatial Statistics are a set of exploratory techniques for describing and modeling spatial distributions, patterns, processes, and relationships.**

coincidence

connectivity

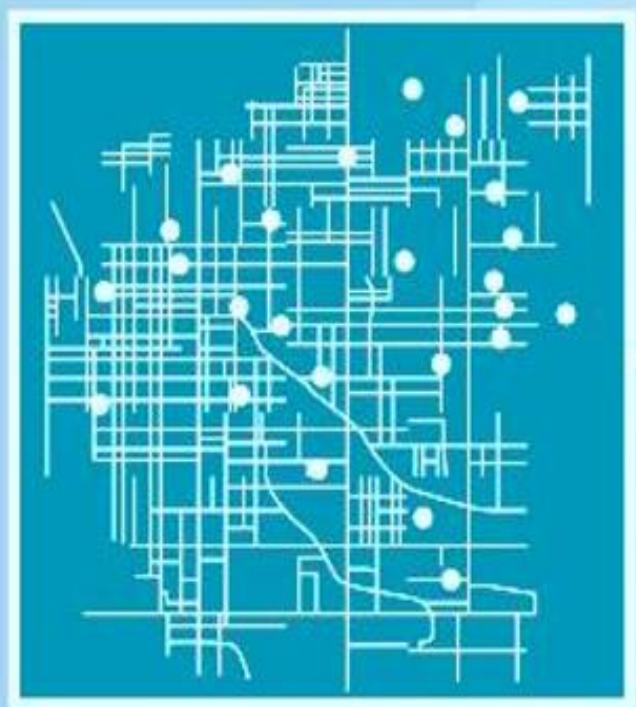
area

proximity

orientation

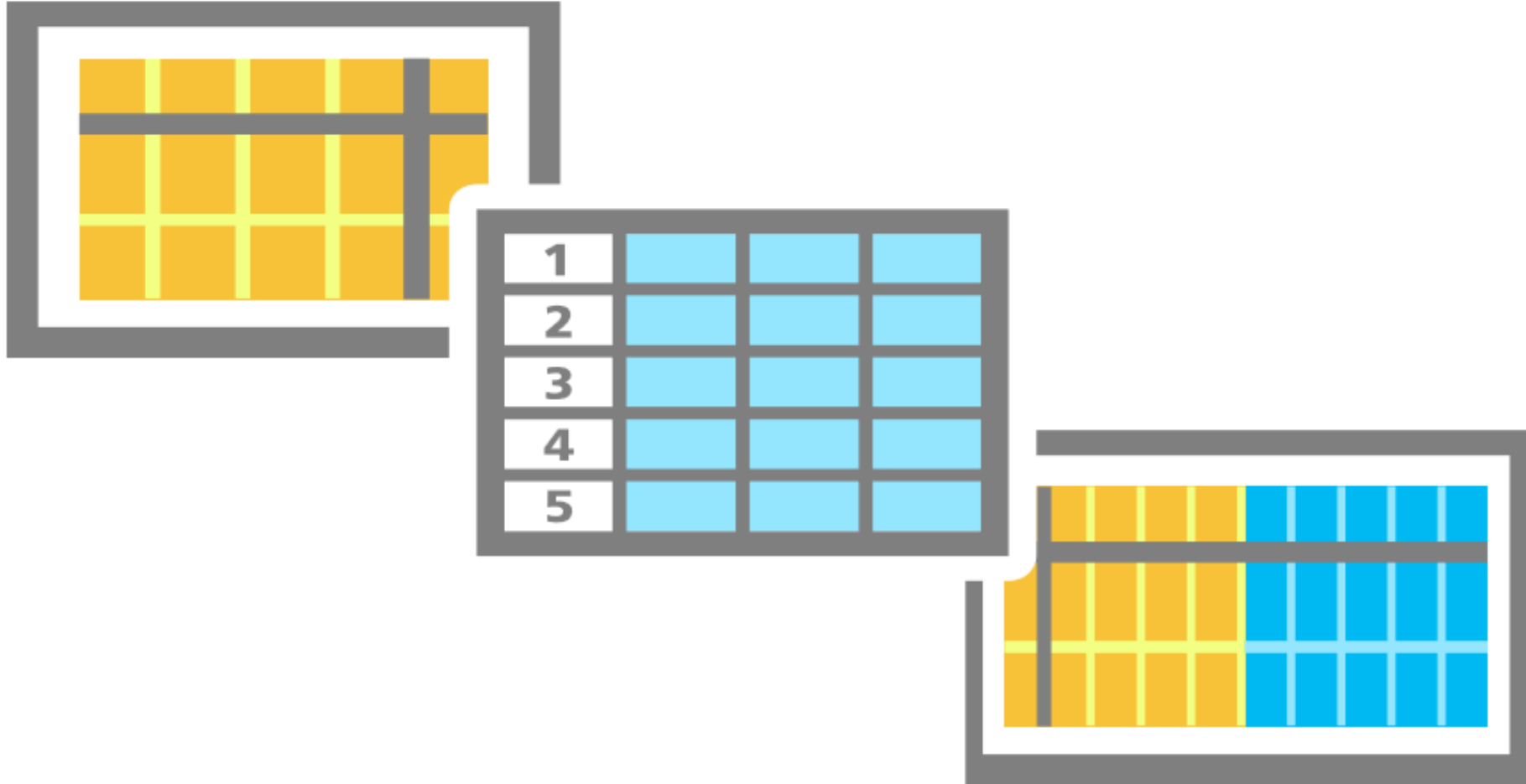
length

direction



# Spreadsheets

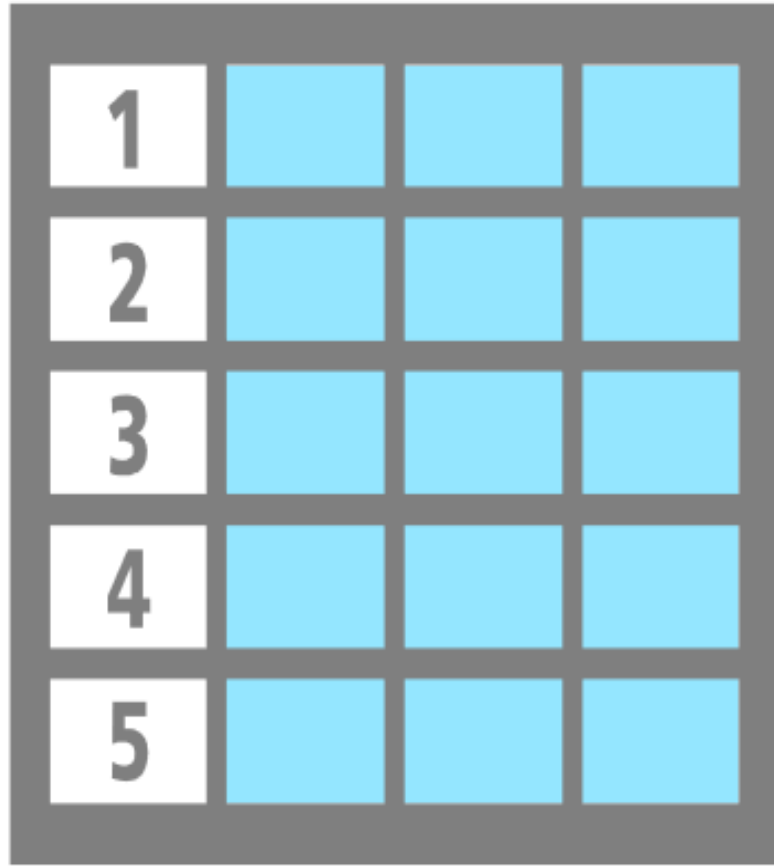
## Data or Information?



A	B	C	D	E	F	G	H
4131 021 030916	4131	21	1567 CARS	3/9/2016	Coldstream Homestead Mont	Notheaste	14 (39.32227
4139 026 030916	4139	26	2738 HARF	3/9/2016	Coldstream Homestead Mont	Notheaste	14 (39.32126
0244 034 031418	244	34	1708 W PF	3/14/2018	Union Square	Southern	9 (39.28546
1125 027 031418	1125	27	1526 N CA	3/14/2018	Oliver	Eastern	12 (39.30801
1127 051 031018	1127	51	1516 N BR	3/10/2018	Oliver	Eastern	12 (39.30782
1138 114 031418	1138	114	1424 N BE	3/14/2018	Oliver	Eastern	12 (39.30694
1145 032 031418	1145	32	1336 AISQ	3/14/2018	Oliver	Eastern	12 (39.30581
1149B010 031418	1149B	10	1308 N CA	3/14/2018	Oliver	Eastern	12 (39.30545
7056A018 030816	7056A	18	3438 6TH S	3/8/2016	BROOKLYN	SOUTHERN	10 (39.23955
1488 010 031018	1488	10	2718 E OLI	3/10/2018	Berea	Eastern	13 (39.30825
1499 030 031418	1499	30	2058 E HO	3/14/2018	Broadway East	Eastern	13 (39.30685
1504 014 031518	1504	14	2407 E OLI	3/15/2018	Broadway East	Eastern	13 (39.30771
1505 055 031018	1505	55	1417 N MI	3/10/2018	Berea	Eastern	13 (39.30754
1521 037 031018		37	2525 E HO	3/10/2018	Berea	Eastern	13 (39.30672
1542 076 031218	1542	76	1224 N PO	3/12/2018	Berea	Eastern	13 (39.30563
1544 028 031018	1544	28	3113 E PRI	3/10/2018	Berea	Eastern	13 (39.30593
1661 014 031318	1661	14	3027 MCE	3/13/2018	Ellwood Park/Monument	Southeast	13 (39.29786
0007 053 122915	7	53	1726 APPL	12/29/2015	EASTERWOOD	WESTERN	7 (39.30841
1667 014 031318	1667	14	429 N WA	3/13/2018	CARE	Eastern	13 (39.29616
1655 083 031318	1655	83	515 N POR	3/13/2018	McElderry Park	Southeast	13 (39.29725
1656 063 031318	1656	63	512 N ROS	3/13/2018	McElderry Park	Southeast	13 (39.29726
1657 062 031418	1657	62	514 N GLO	3/14/2018	McElderry Park	Southeast	13 (39.29736
1658 074 031318	1658	74	531 N BEL	3/13/2018	McElderry Park	Southeast	13 (39.29776
1659 002 031318	1659	2	503 N KEN	3/13/2018	McElderry Park	Southeast	13 (39.29704
4581 010F 030918	4581	010F	3223 SPAL	3/9/2018	Central Park Heights	Northwest	6 (39.34726
4609 055 031418	4609	55	3701 MAN	3/14/2018	Central Park Heights	Northwest	6 (39.34307
4616 005 031318	4616	5	3008 WOC	3/13/2018	Central Park Heights	Northwest	6 (39.34655
5902A004 031518	5902A	4	3307 LAKE	3/15/2018	Belair-Edison	Notheaste	13 (39.32087
5907 031 031518	5907	31	3044 CHES	3/15/2018	Belair-Edison	Notheaste	13 (39.32395
5929 061 031518	5929	61	4340 SHAM	3/15/2018	Belair-Edison	Notheaste	2 (39.32625
0006 032 011509	6	32	2002 PRES	1/15/2009	EASTERWOOD	WESTERN	7 (39.30800
0031 057 032316	31	57	1834 LAUF	3/23/2016	SANDTOWN-WINCHESTER	WESTERN	7 (39.30327
0055C063 032316	0055C	63	1104 N CA	3/23/2016	SANDTOWN-WINCHESTER	WESTERN	9 (39.30135
0073 024 031616	73	24	930 N MO	3/16/2016	SANDTOWN-WINCHESTER	WESTERN	9 (39.29980
0075 016 032416	75	16	1515 MOS	3/24/2016	SANDTOWN-WINCHESTER	WESTERN	9 (39.29977
0082 058 032416	82	58	805 N BRIC	3/24/2016	MIDTOWN-EDMONDSON	WESTERN	9 (39.29785
0096 052 031516	96	52	707 N MO	3/15/2016	HARLEM PARK	WESTERN	9 (39.29682
1109 040 031018	1109	40	1704 N BR	3/10/2018	Oliver	Eastern	12 (39.30976
0149 044 031616	149	44	310 N BRU	3/16/2016	FRANKLIN SQUARE	WESTERN	9 (39.29235
0244 074 031516	244	74	1708 LEMI	3/15/2016	UNION SQUARE	SOUTHERN	9 (39.28584
1114 000 031418	1114	0	1215 E FER	3/14/2018	Oliver	Eastern	12 (39.30845



# When you look at a spreadsheet...



1		
2		
3		
4		
5		

# You ask for more

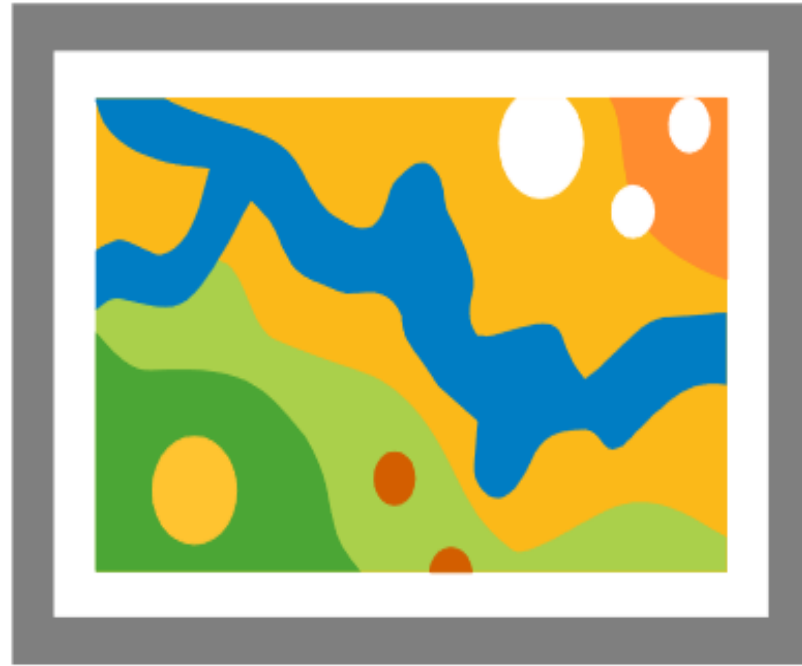
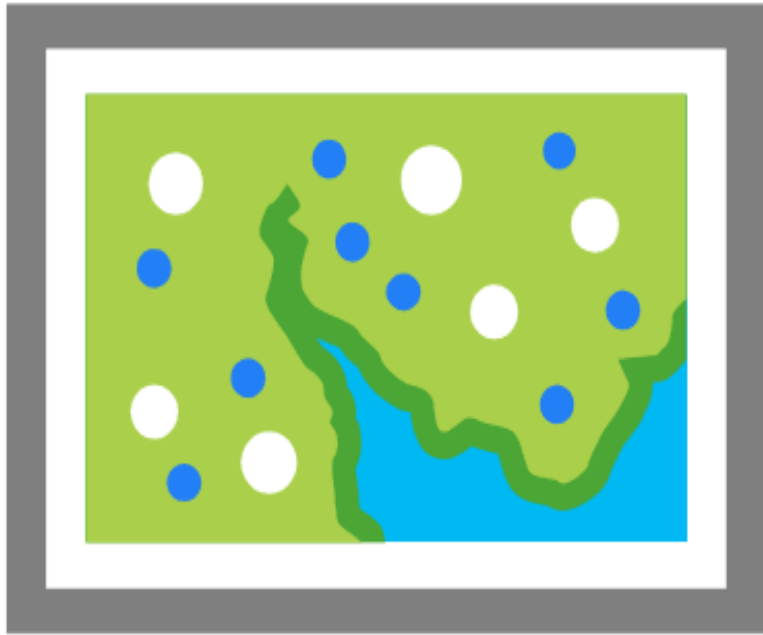
1			
2			
3			
4			
5			



- Mean
- Standard Deviations
- Min and Max
- ...

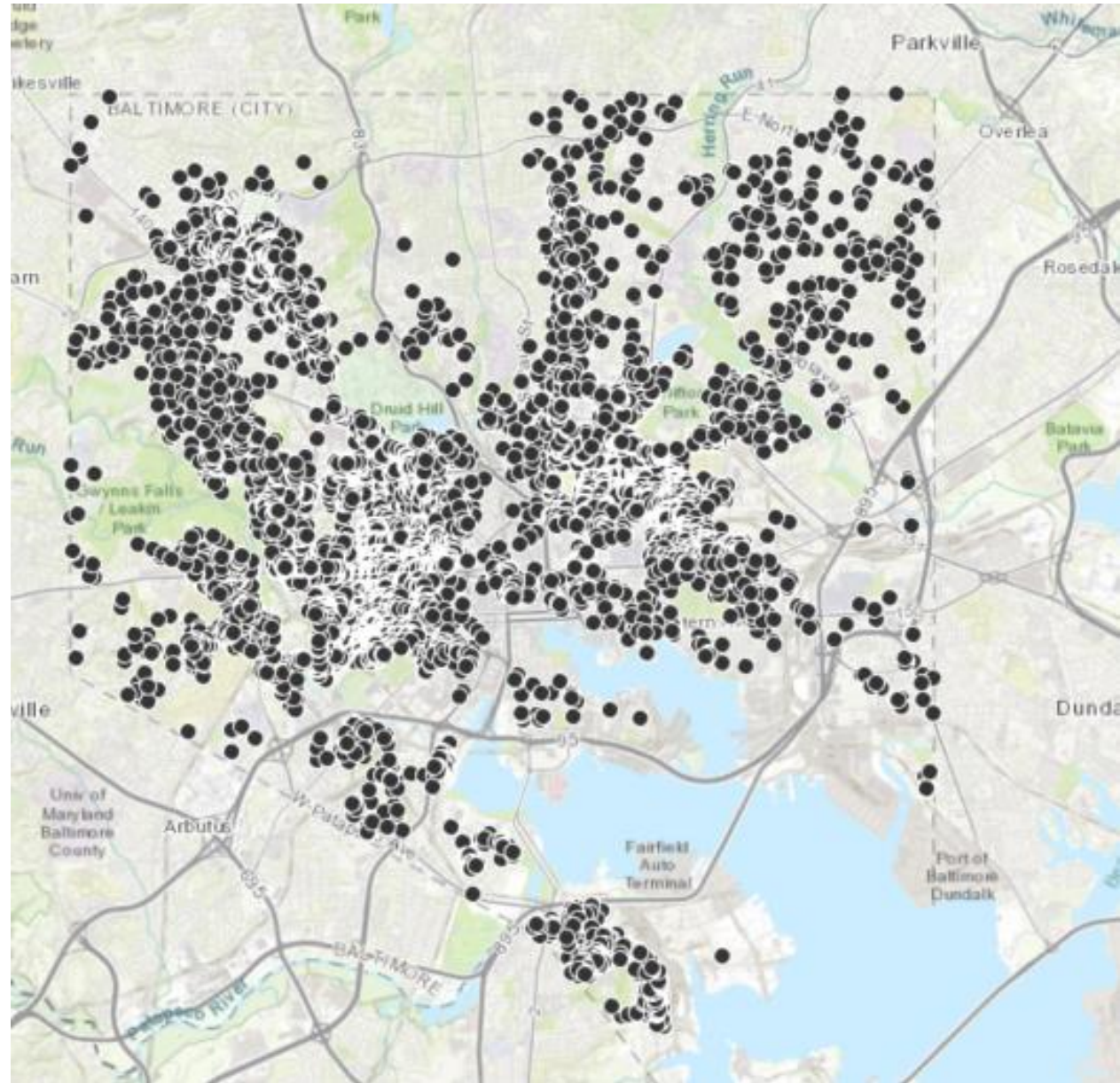
# Maps

## Data or Information?

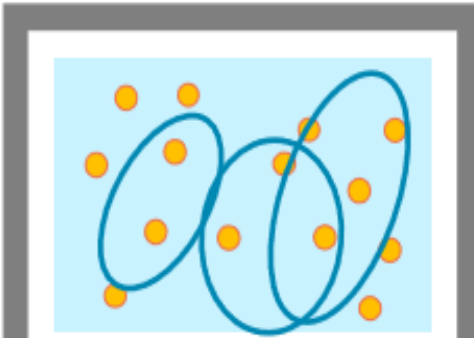
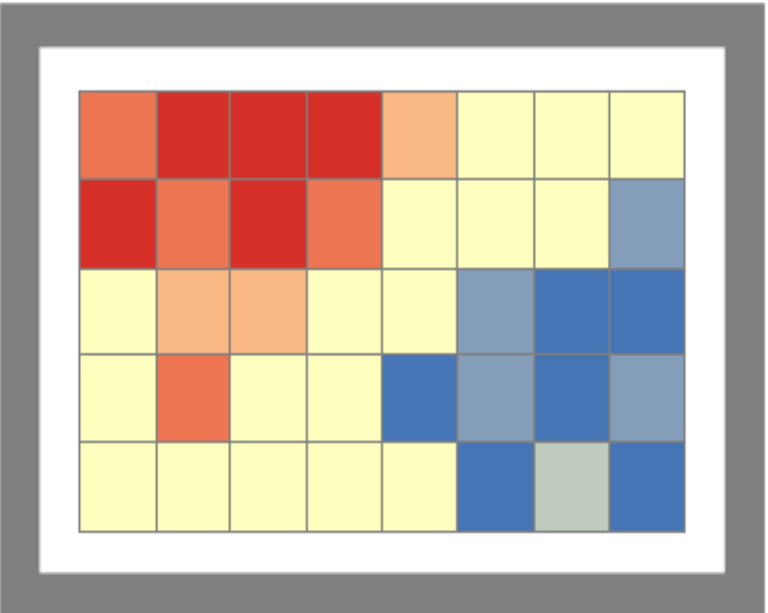


Same goes for maps!





We can do more





# Reconnecting Milwaukee

A BikeAble™ Study of Opportunity, Equity and Connectivity

# ArcGIS Connection with Open Source Libraries

The background features a dark blue gradient with abstract, layered geometric shapes in various shades of blue and red. These shapes, including lines and polygons, are arranged in a way that suggests depth and movement, particularly towards the right side of the frame. The overall aesthetic is modern and technical.



# Navigating the Python Ecosystem

Enterprise & Online



## Sample Use Case:

**Developing a Weather Repository**



- R: A widely used statistical programming language
- More than 10,000 Packages
- Expand Workflows



United States<sup>™</sup>  
**Census**  
Bureau



# Sample Use Cases:

## Predicting Seagrass Locations

- **Ingesting / Analyzing Multidimensional Scientific Data**
- **Modeling Vacancy Rates in DC -- from Baltimore Vacancy Data**

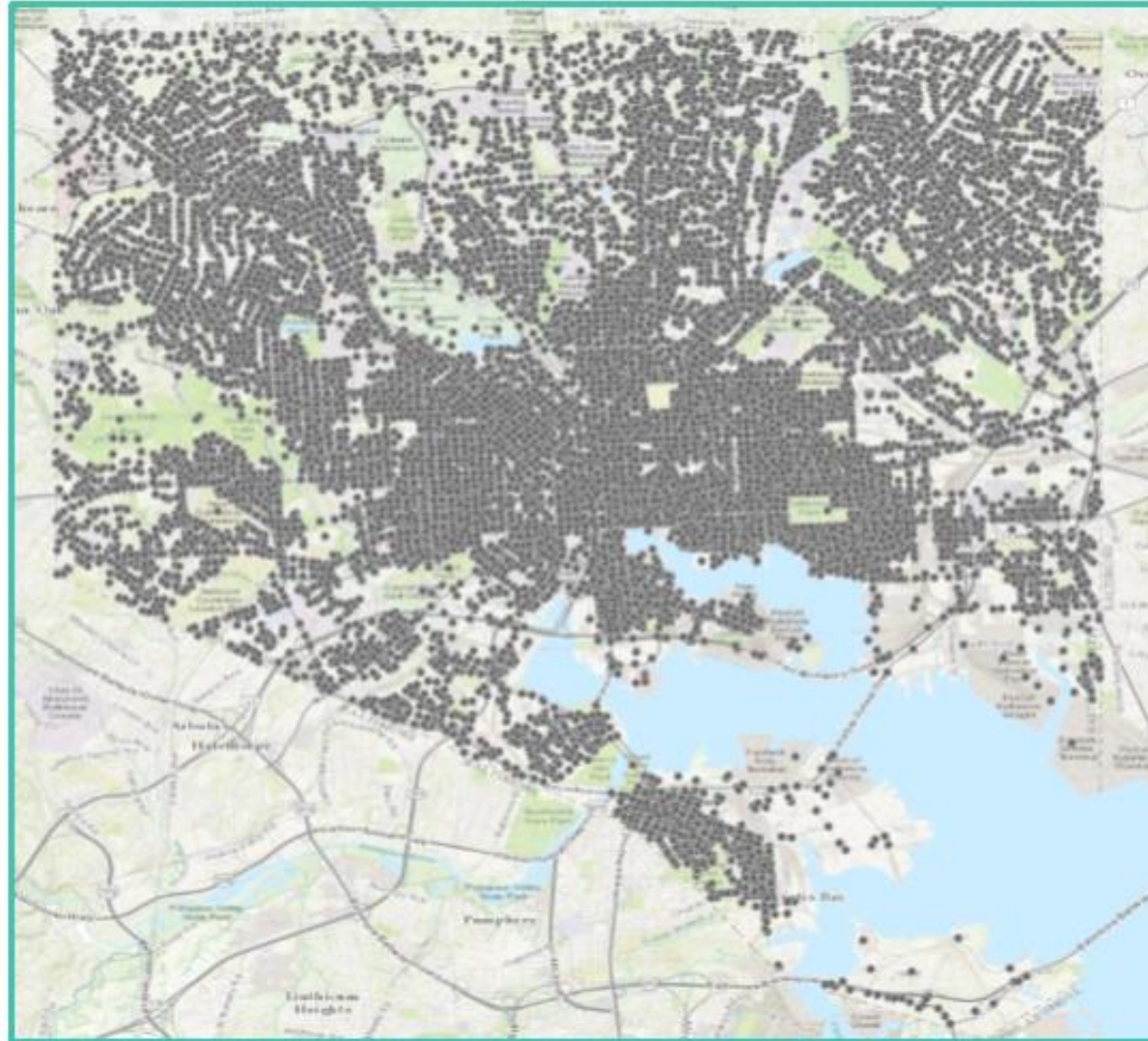
# Clustering, Prediction & Classification

Kristen Hocutt

# Subjectivity of Maps

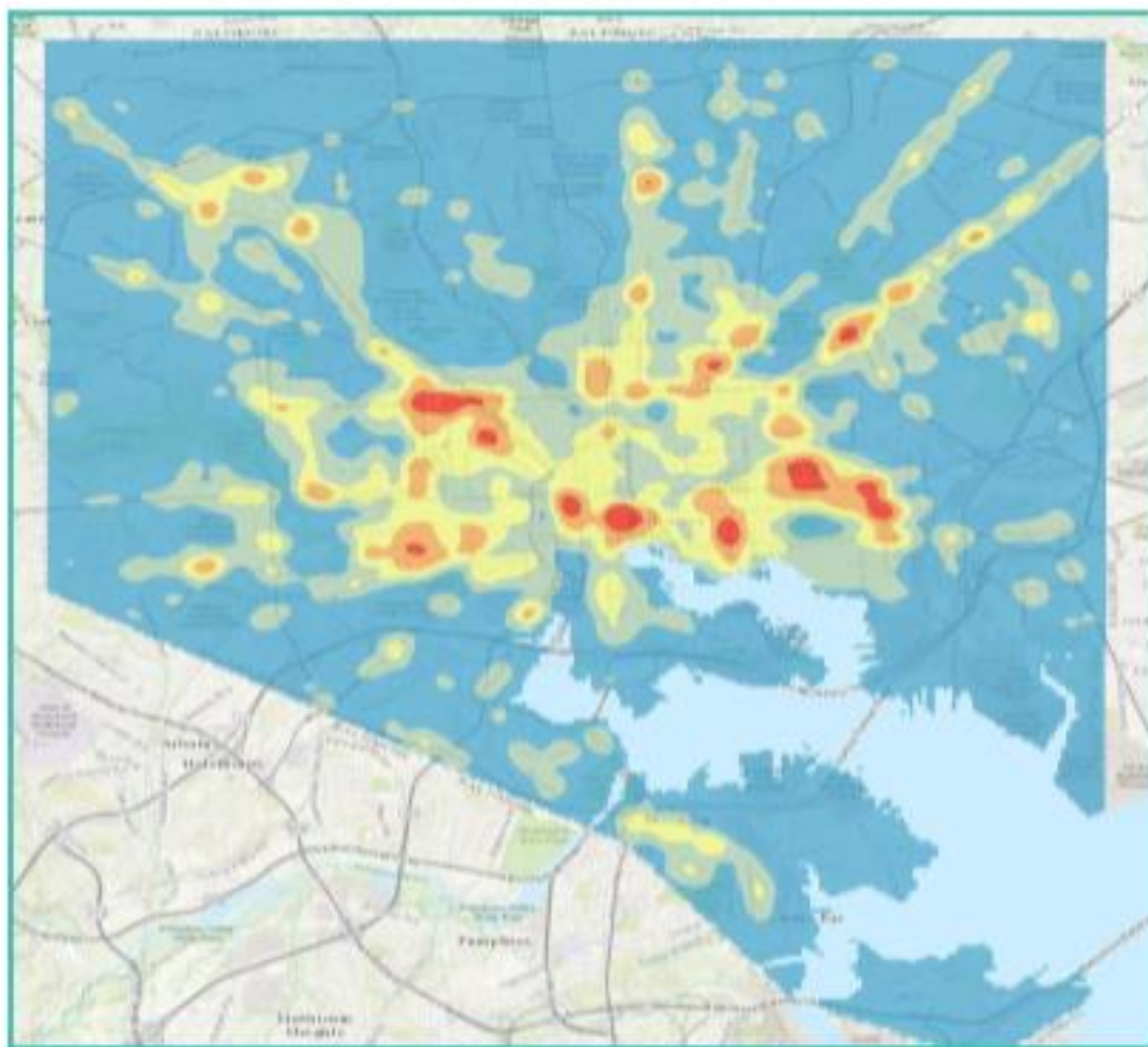
# The map as data

High Priority 911 Calls in Baltimore



# The map as data

High Priority 911 Calls in Baltimore

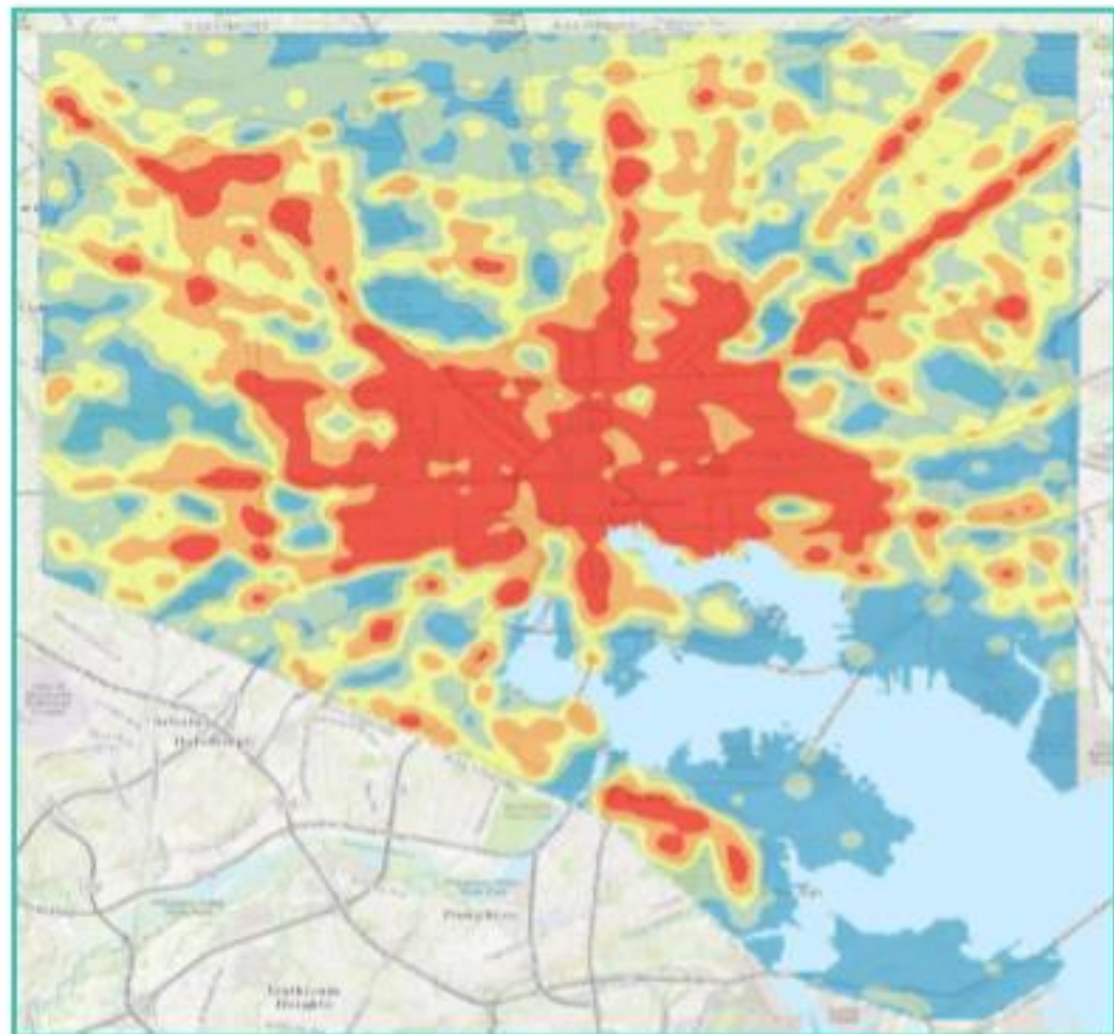


Where are the hot spots? Where is the variation greater?



# The map as data

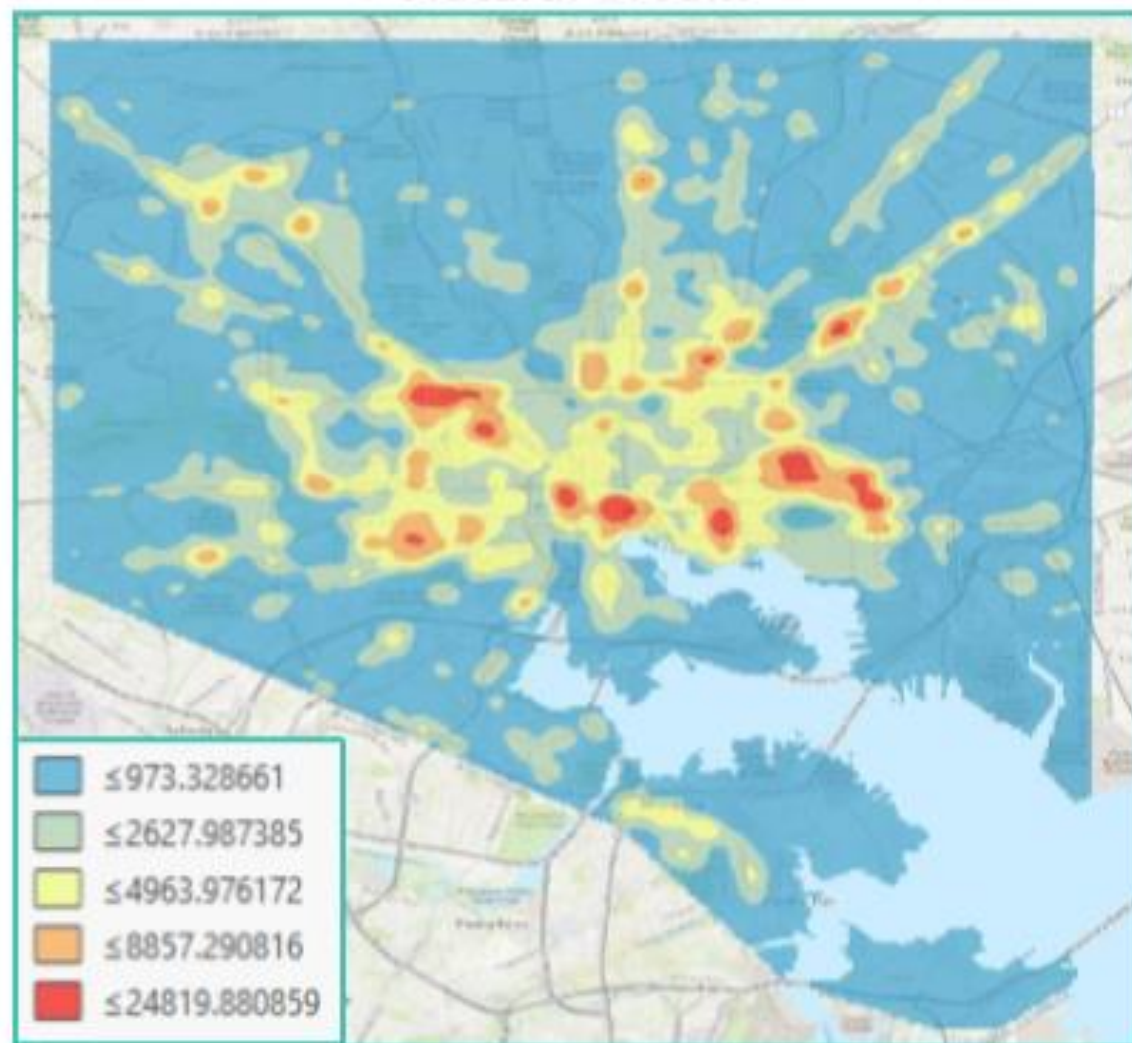
High Priority 911 Calls in Baltimore



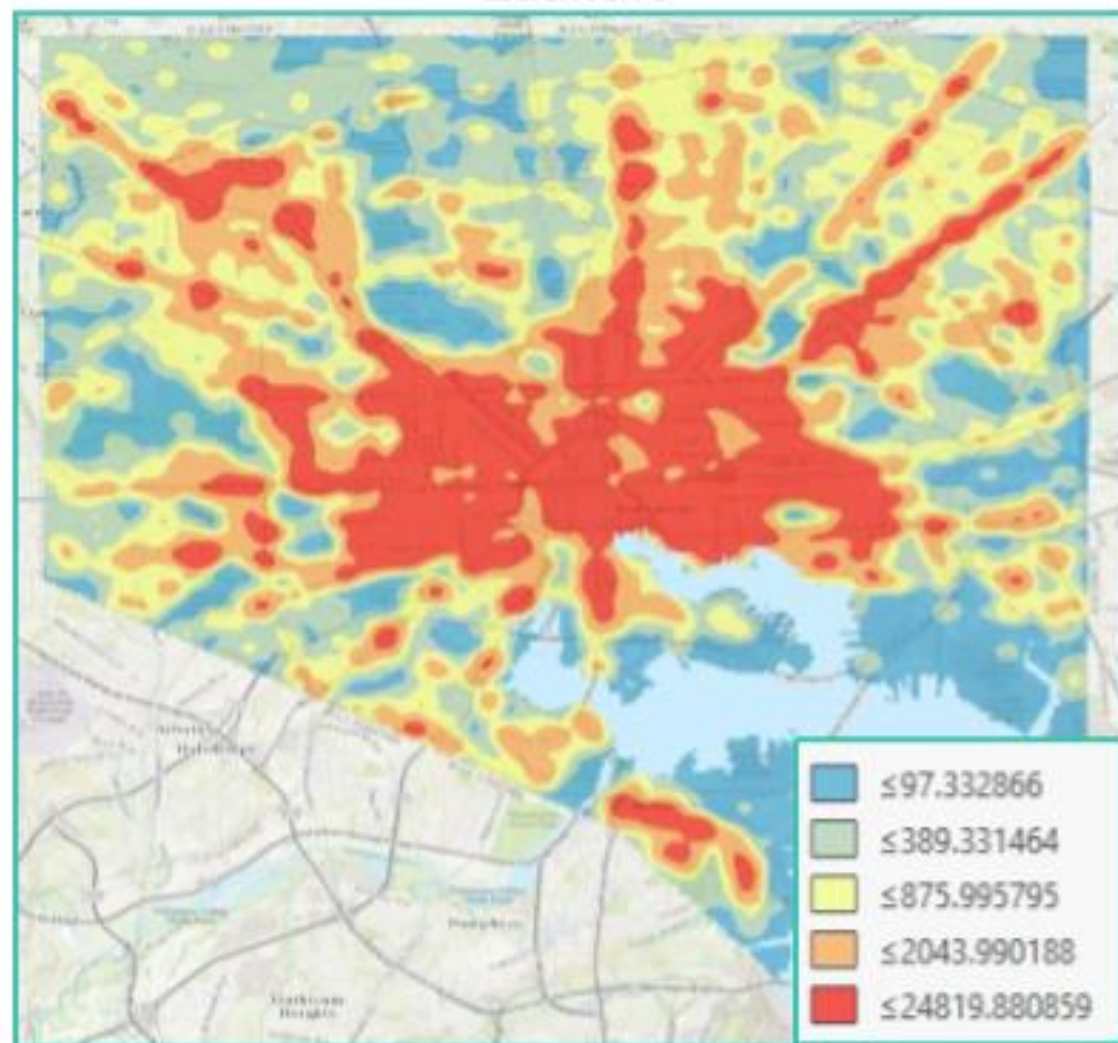
Where are the hot spots? Where is the variation greater?

# The **subjectivity** of visual pattern analysis

Natural Breaks



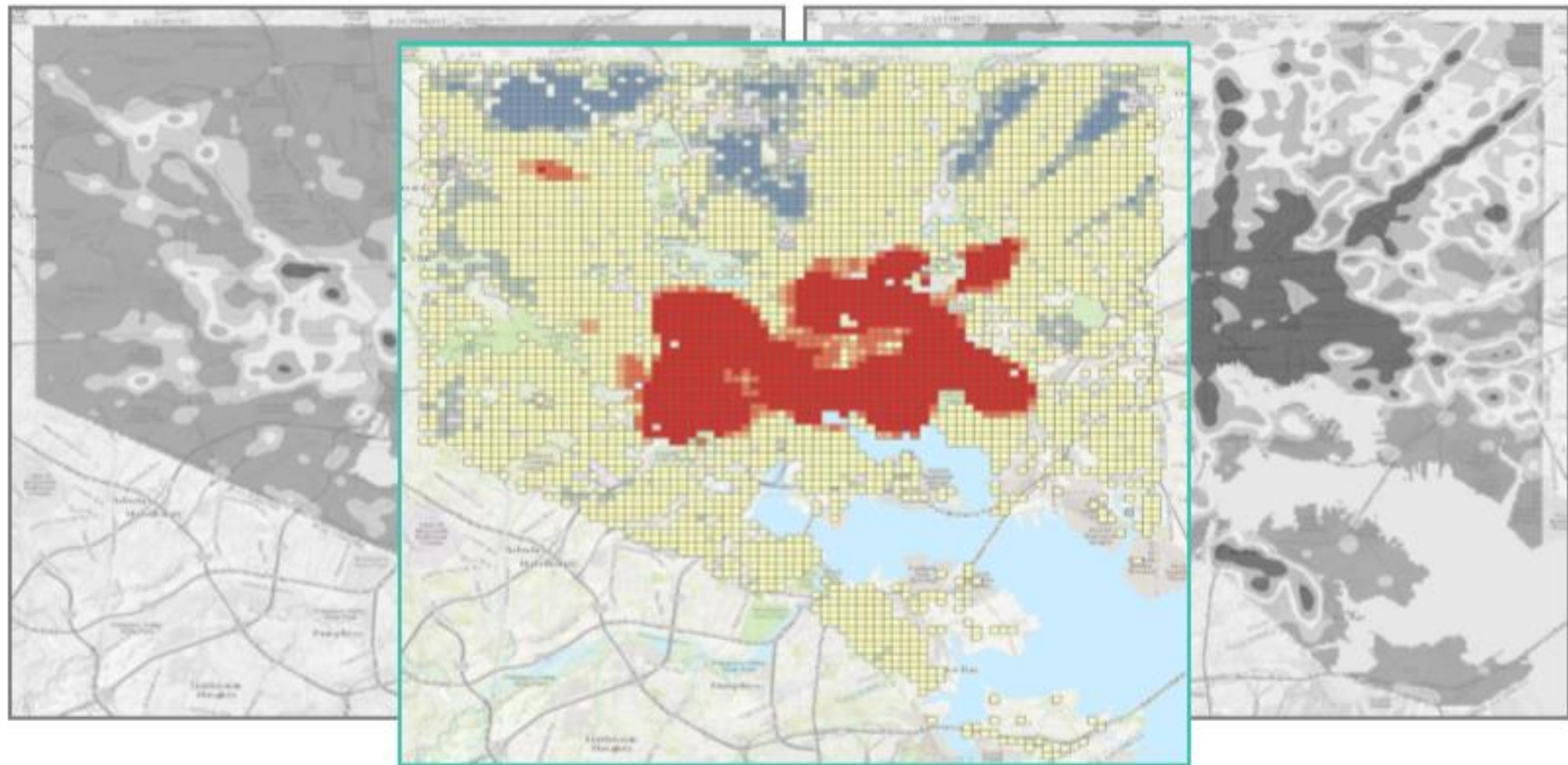
Quantile



Where are the hot spots? Where is the variation greater?

# Minimizing the subjectivity

Turning the map into **information**

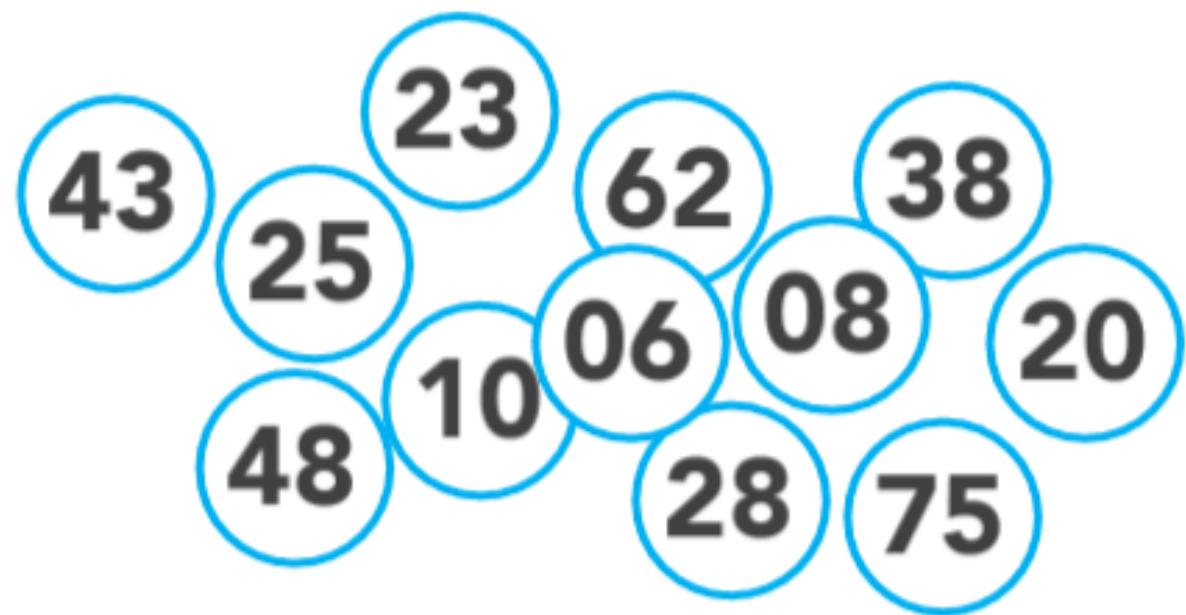
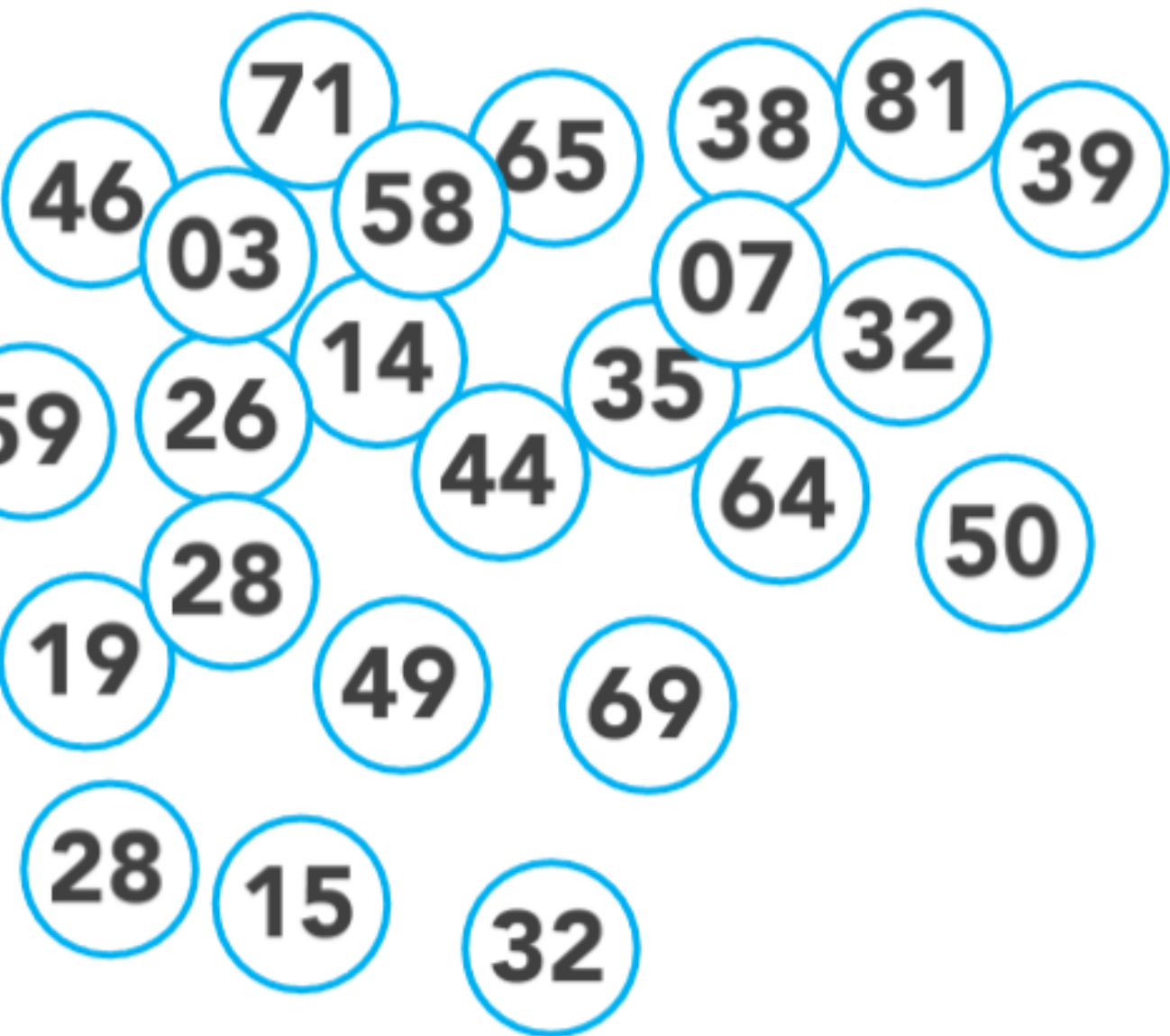


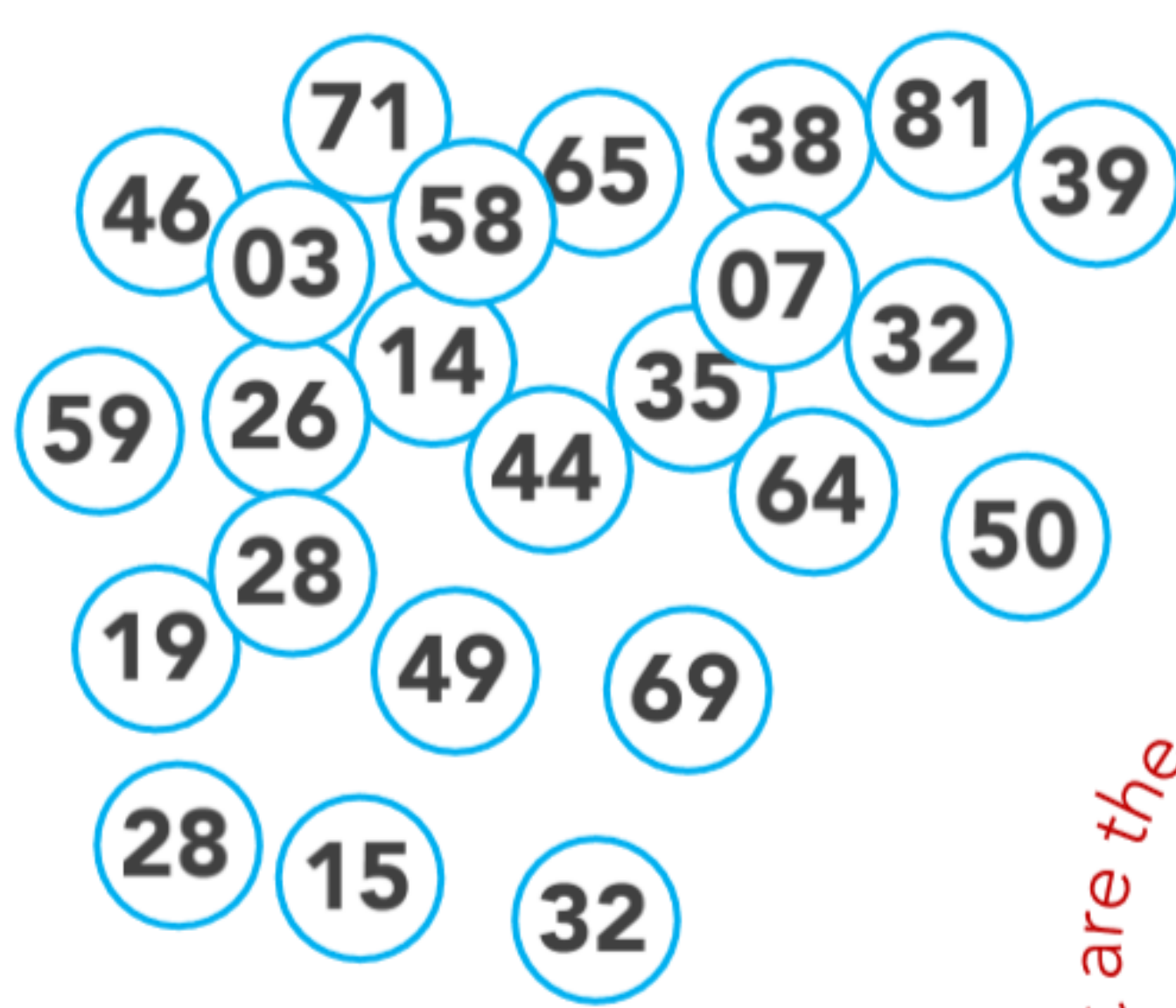


# Complete Spatial RANDOMNESS

Is there a **PATTERN**?

14	15	92	65	35	89	79	32	38
46	26	43	38	32	79	50	28	84
19	71	69	39	93	75	10	58	20
97	49	44	59	23	07	81	64	06
28	62	08	99	86	28	03	48	25





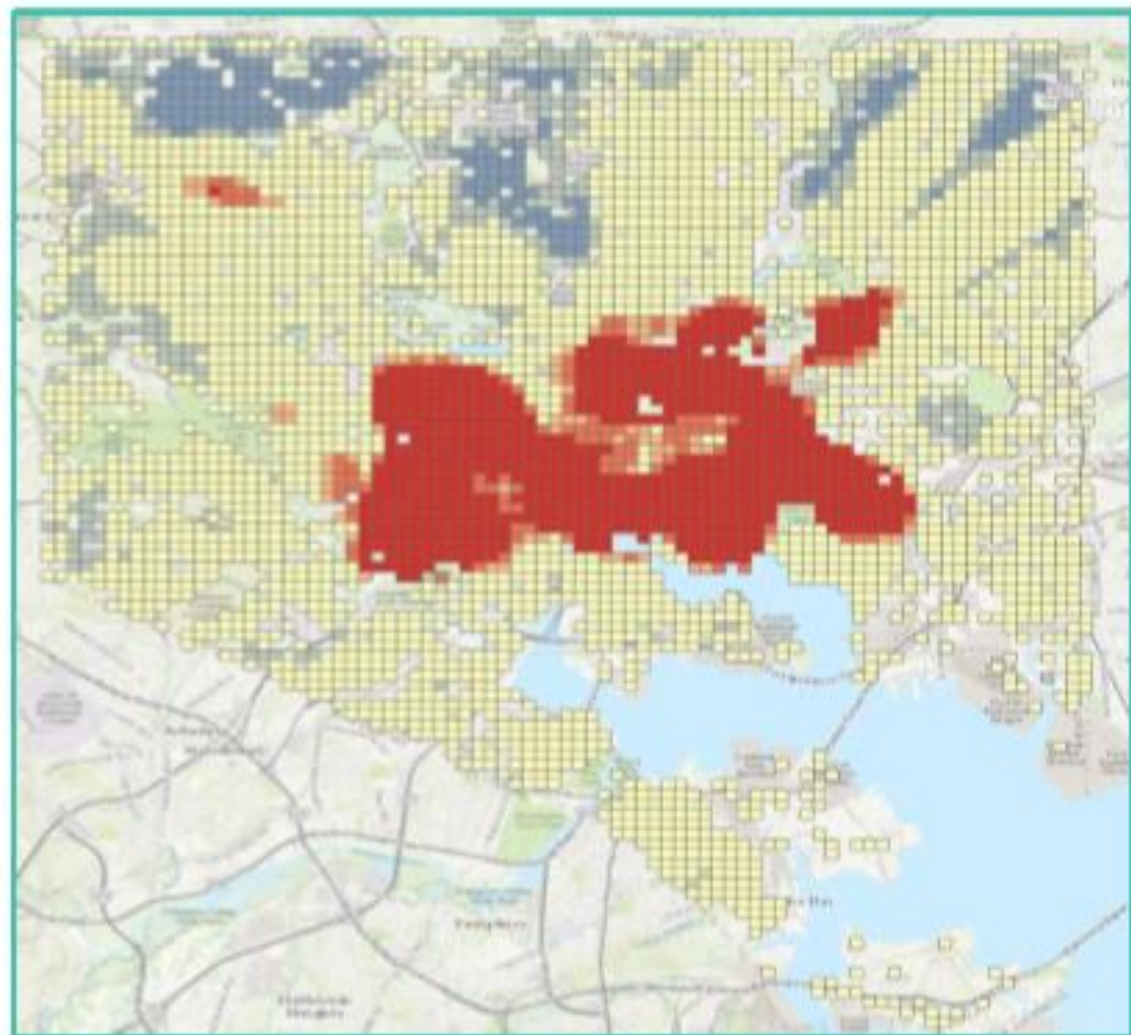
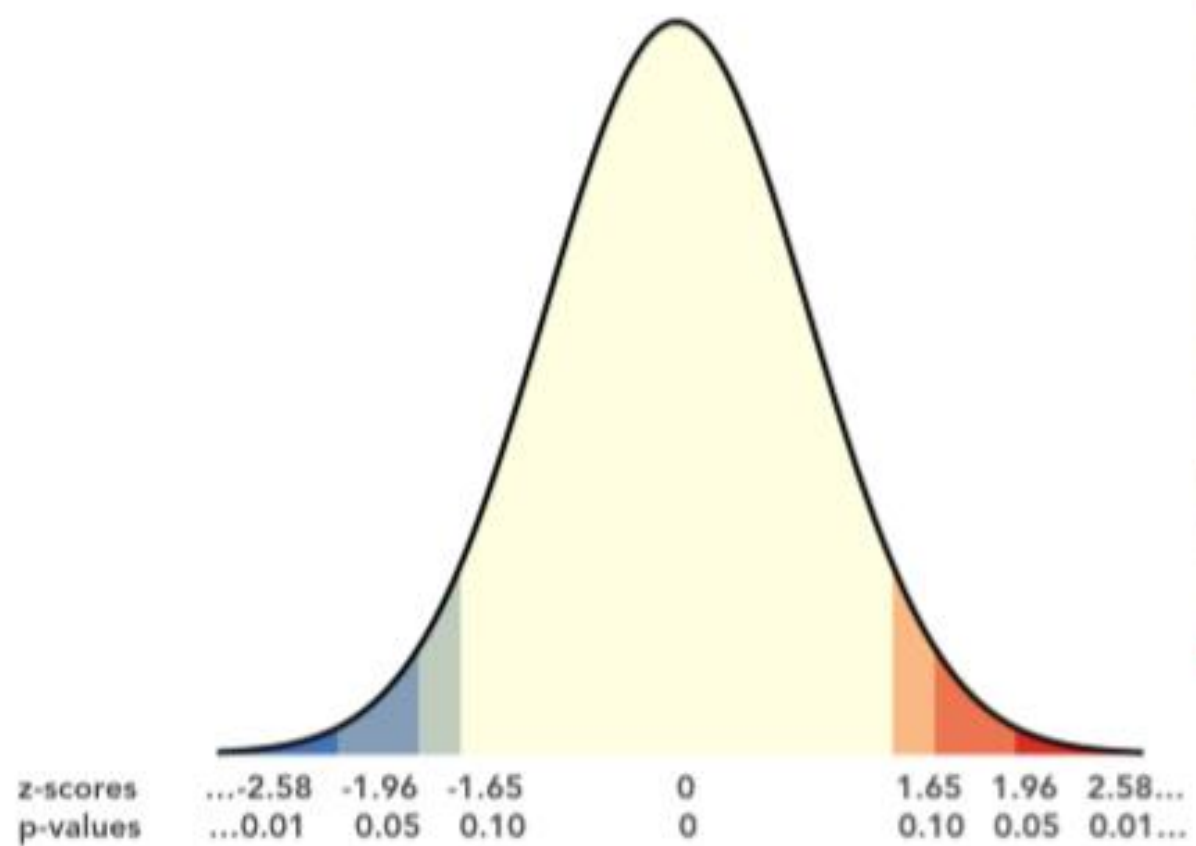
What are the chances this happened RANDOMLY???





**z-scores**

**p-values**



# Hot Spot Analysis

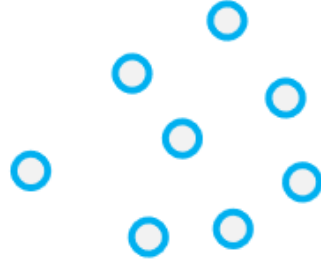
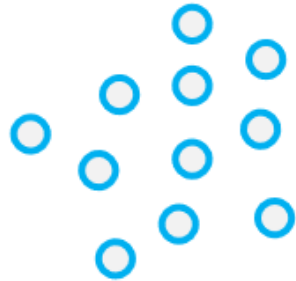
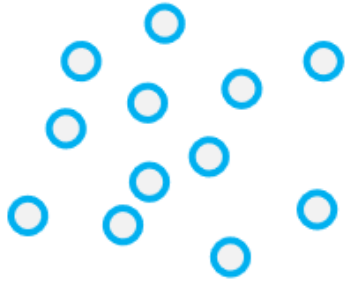
given a set of weighted features, identifies statistically significant hot spots and cold spots using the Getis-Ord  $G_i^*$  statistic

...how do we know if it's  
**SIGNIFICANTLY**  
different???

**МАТН!**

# Density-based Clustering

finds clusters based on feature locations







**DBSCAN** – defined distance

**HDBSCAN** – self adjusting

**OPTICS** – multi-scale

# Spatial Analysis of Traffic Accidents

A Density Based Clustering Approach



# Models

Representative  
generalizations used for  
**prediction**



# Why model

Use information we have  
to **predict** information we  
don't have

Which areas  
are most  
contaminated?

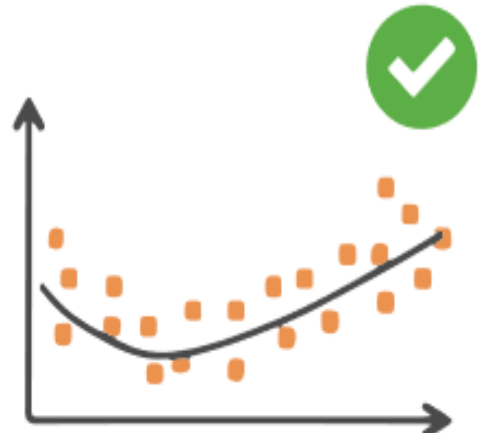
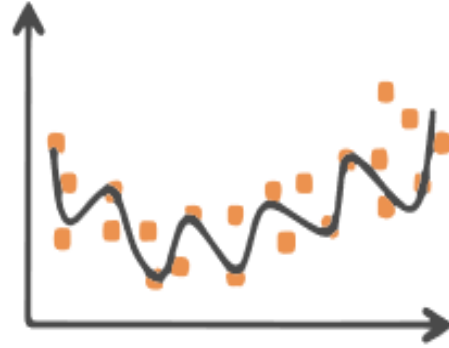
What drives  
sales?

Which  
buildings will  
fail inspection?

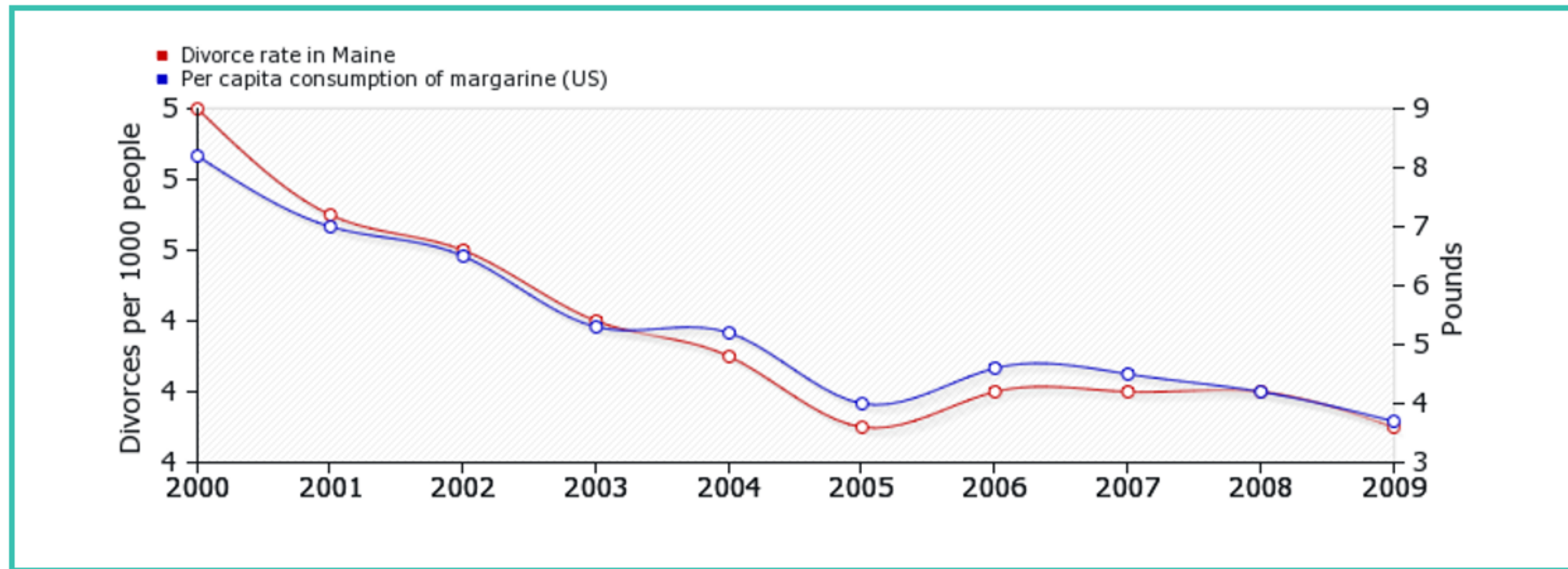
What will the  
weather be like  
tomorrow?

# When we can't trust a model

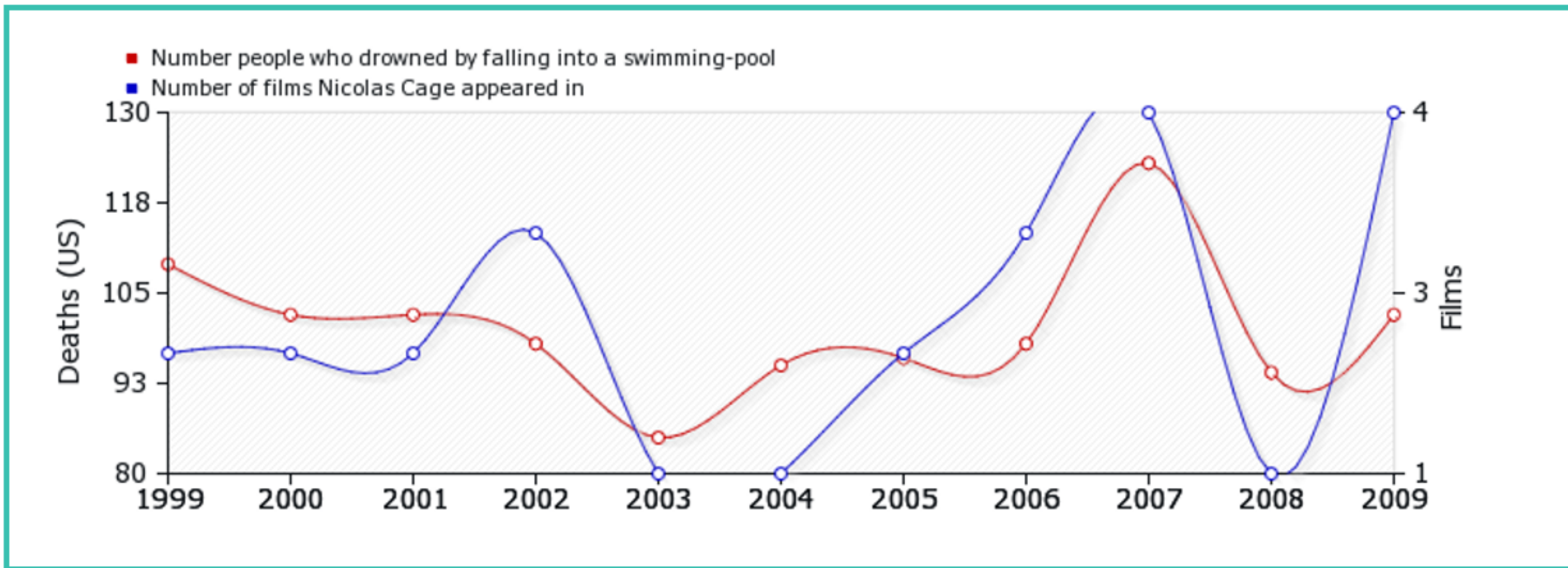
Mimics training dataset and models **noise** instead of generalizing a trend



# Divorce Rate in Maine vs Per Capita Consumption of Margarine



# Number People Who Drowned by Falling into a Swimming-Pool vs Number of Nicolas Cage Films



**Correlation: 0.666004**

# Many ways to model

Generalized Linear Regression

Geographically Weighted Regression



Forest-based Classification and Regression





Generalized

Linear

Regression

Modeling linear relationships

# Three model types

- **Gaussian – continuous**
- **Logistic – binary**
- **Poisson – count**

Gaussian

Ordinary Least Squares

# Dependent Variable

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \varepsilon$$

What are you trying to predict or understand?

# Exploratory Regression



Geographically

Weighted

Regression

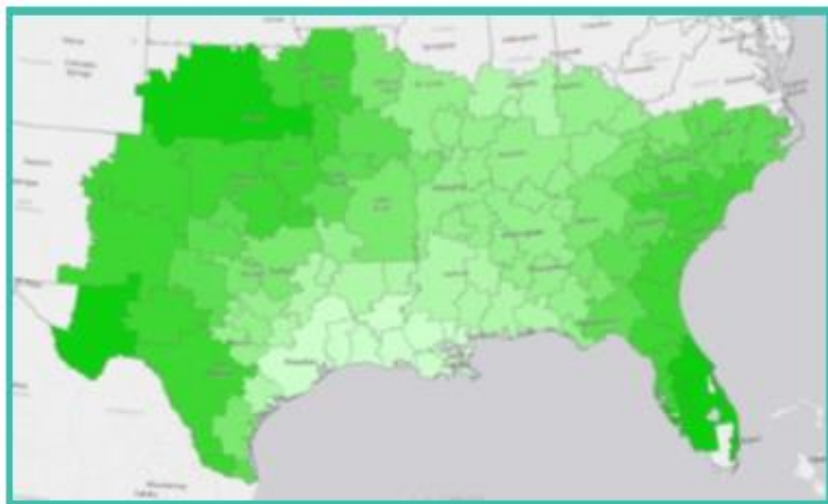
Exploring spatial variation



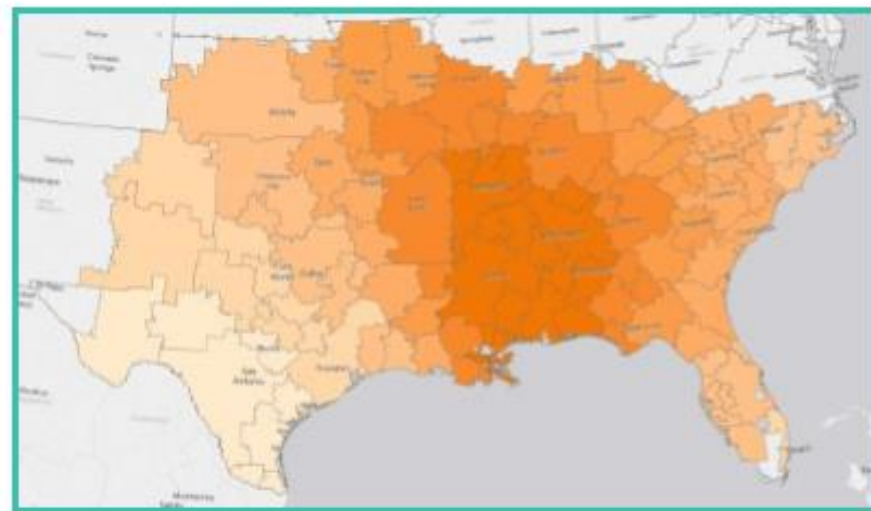
each  
feature  
gets a  
separate  
equation



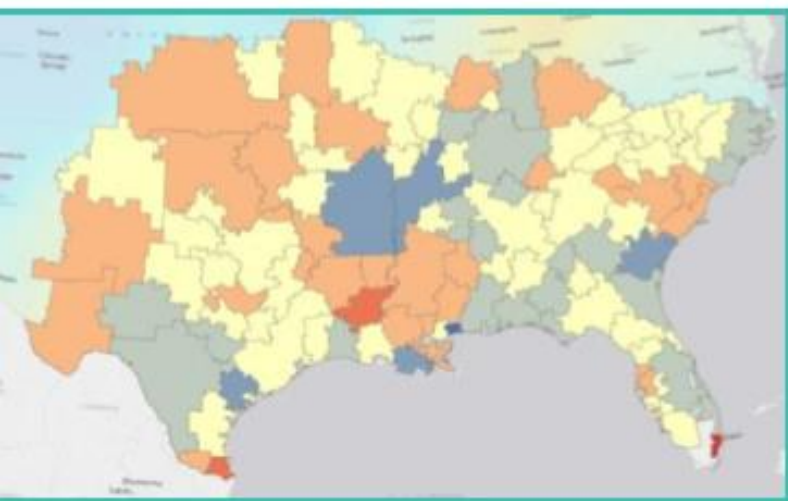




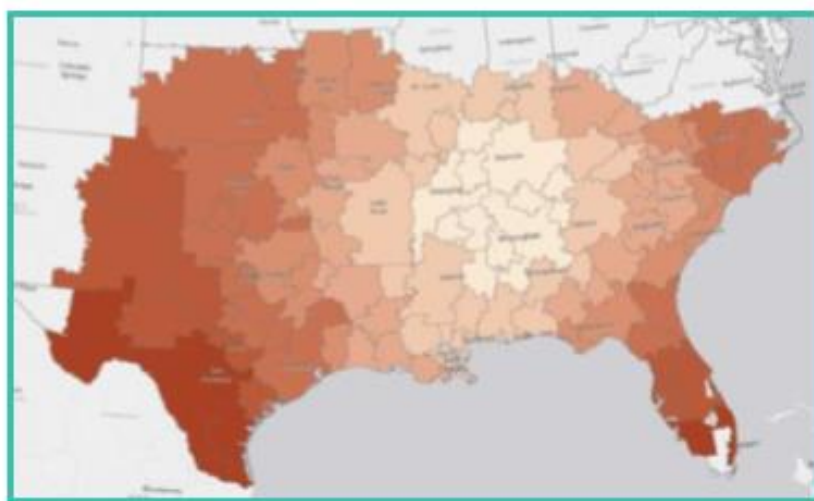
Local R-Squared



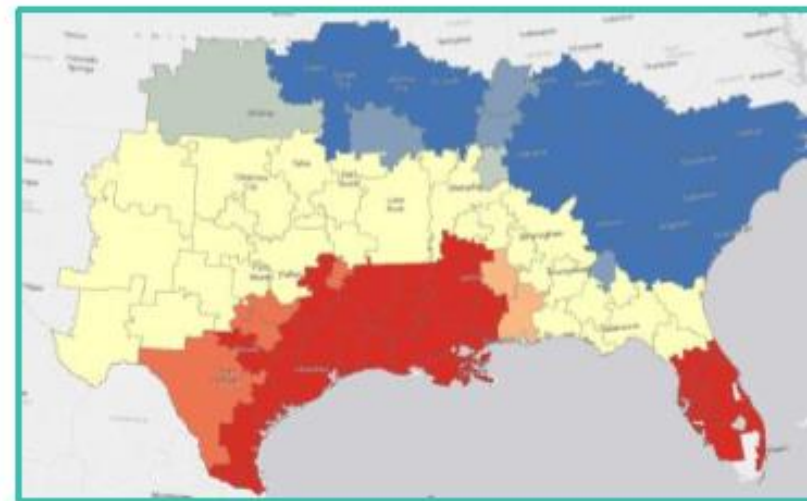
Coefficients



Residuals



Condition Number



Predictions

**Forest-based**

**Classification &**

**Regression**

**Predicting using machine learning**



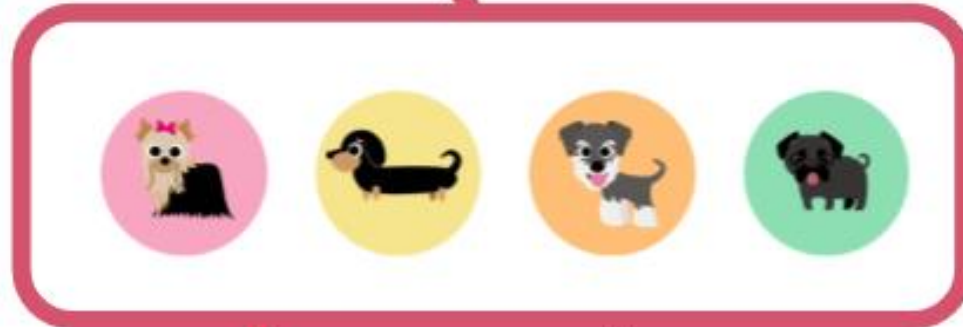
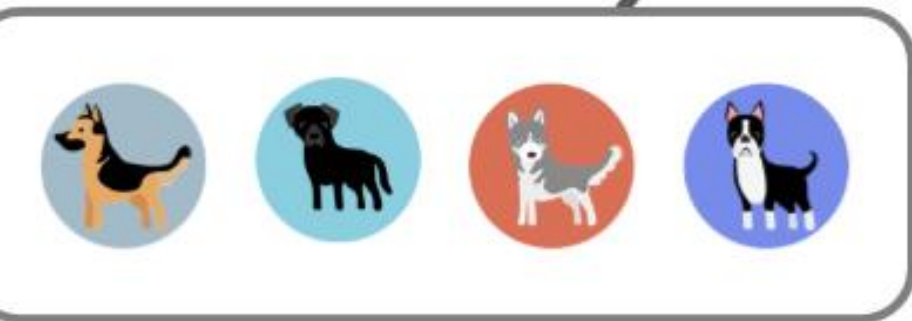
# Training

variable to predict

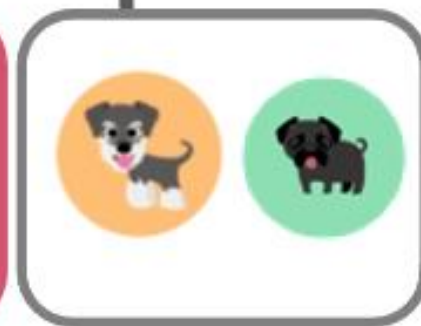
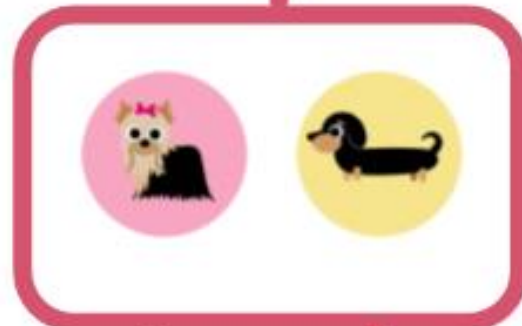
**Breed**

Size  
Color  
Fur  
Ears  
Tail  
Age  
Weight

# Decision Tree

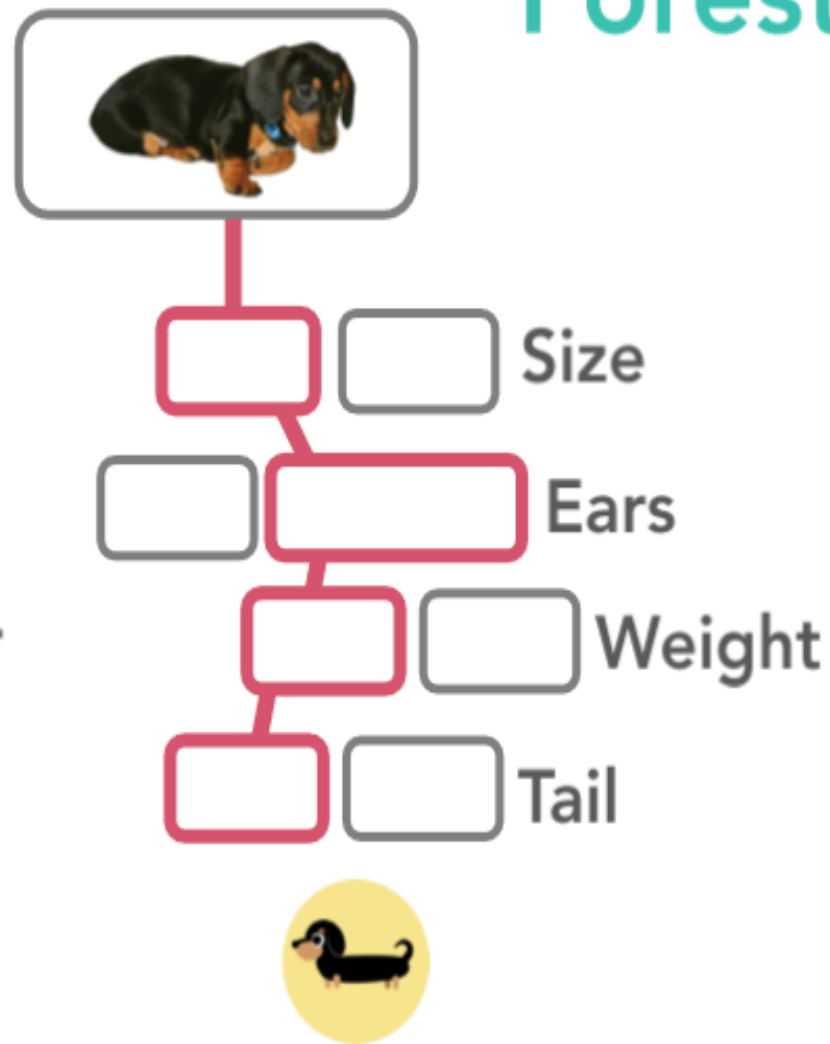
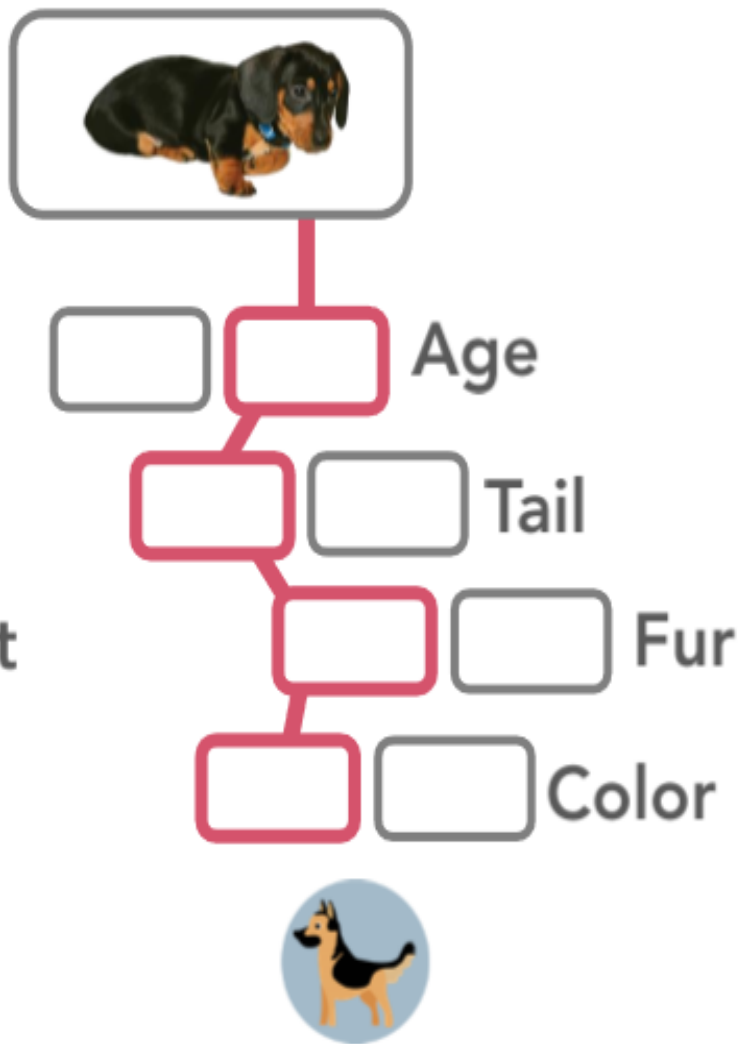
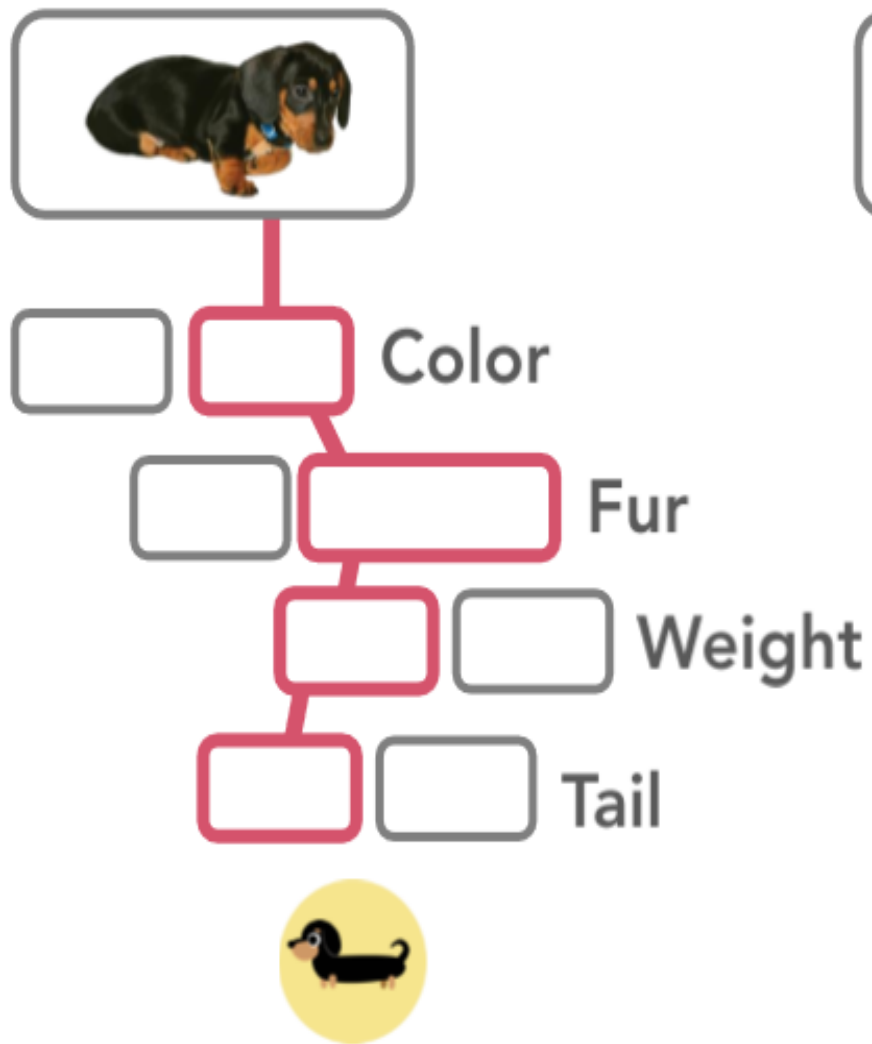


Color

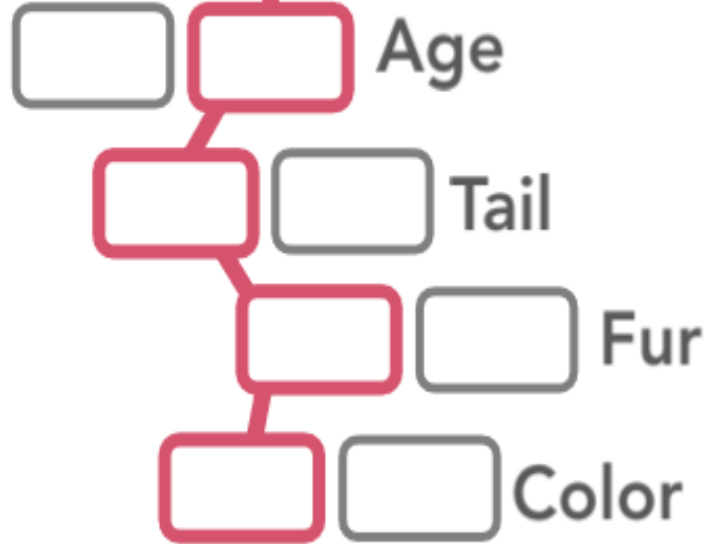
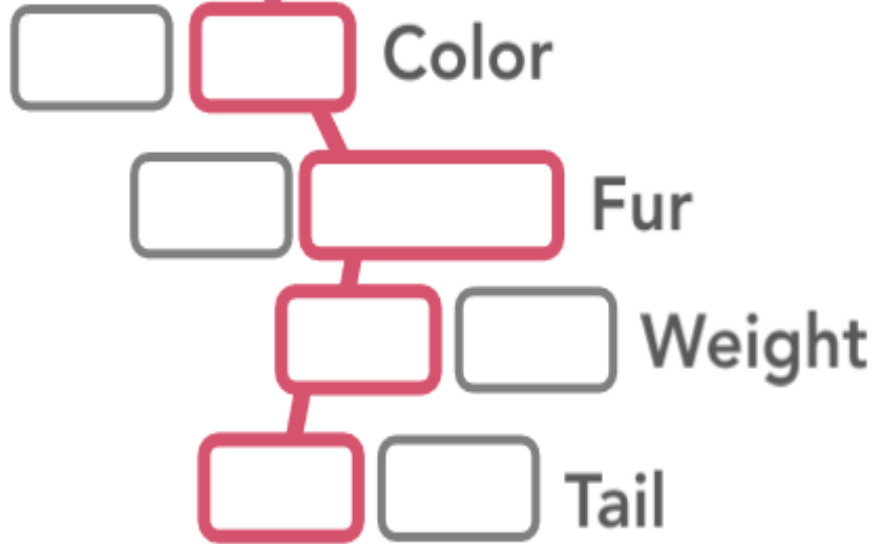


Ears





**Random** subset of data and variables used in each tree



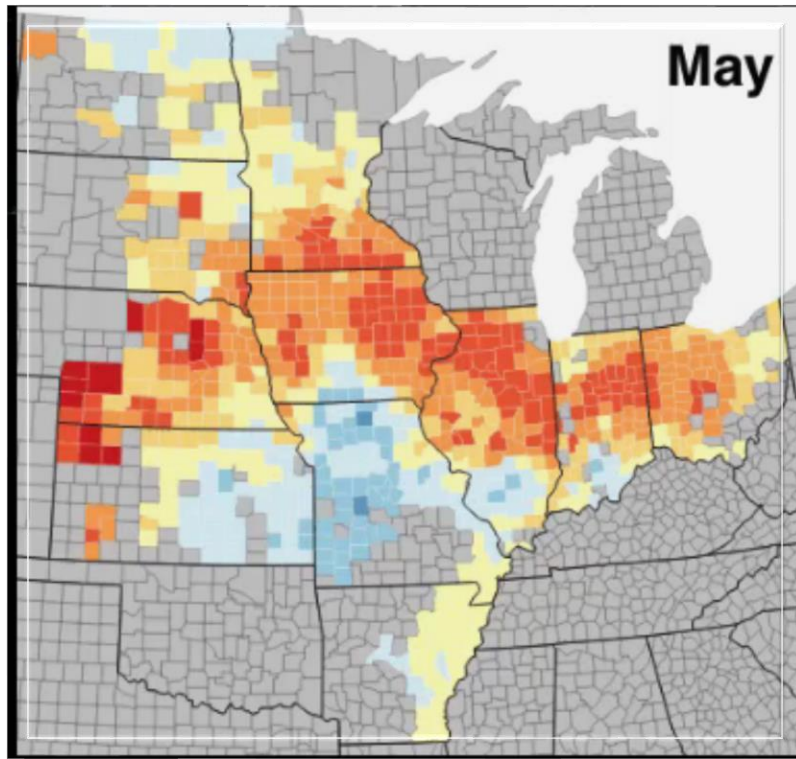
Majority vote wins



=

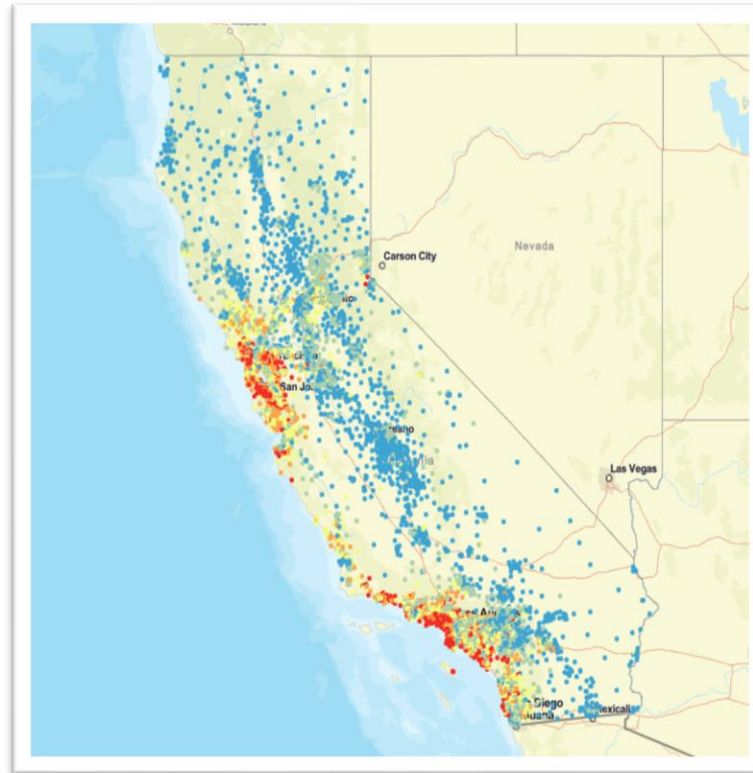


# Potential Applications/Where it's used Today

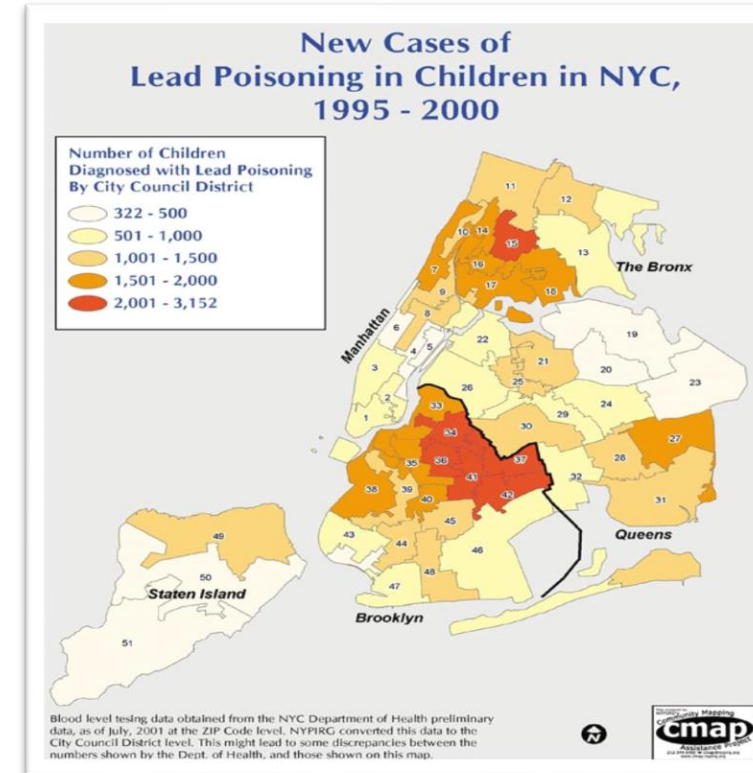


## Crop Yield

JIAXUAN YOU, XIAOCHENG LI, MELVIN  
LOW, DAVID B. LOBELL, STEFANO ERMON-  
STANFORD UNIVERSITY



## Housing Values



## Lead Poisoning

# Data Visualization

Kristen Hocutt

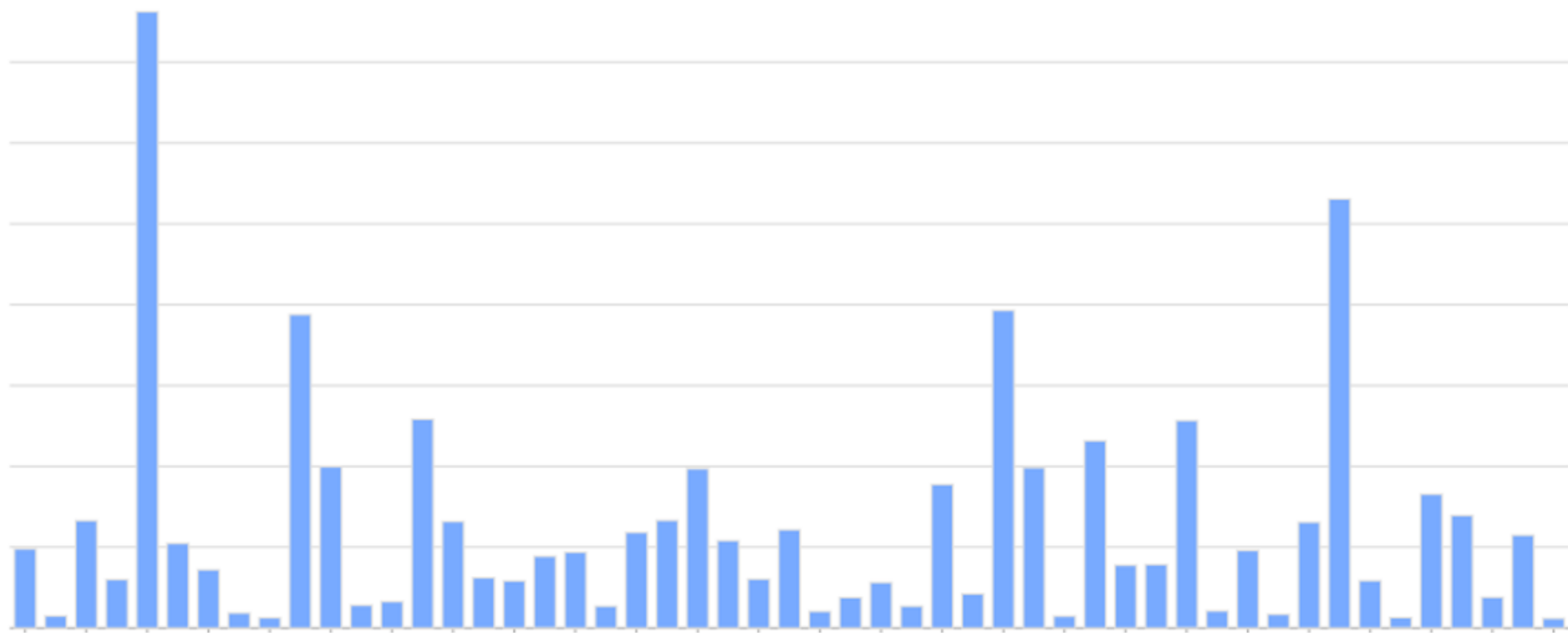


# Why visualize data?

Convert slow reasoning tasks  
into fast perception tasks



1399683	635200		2908933	2106392
6962578	5401748	8875318	8275961	4885854
	3922722	6581982	6074504	3018484
1019462	1336265	2815039		
1345609	3111844	2917750	6636256	9978939
	6656872	38120066	3880520	4796559
724027	1877879	11575704	9913774	2995330
844322	19631599	12914651	6539407	4681639
			26538203	19383475
587106	12817894	629120		9853722
5731663	3601157	927030	5903286	
		1885932	5235100	742404
1628760	1050228		4419036	

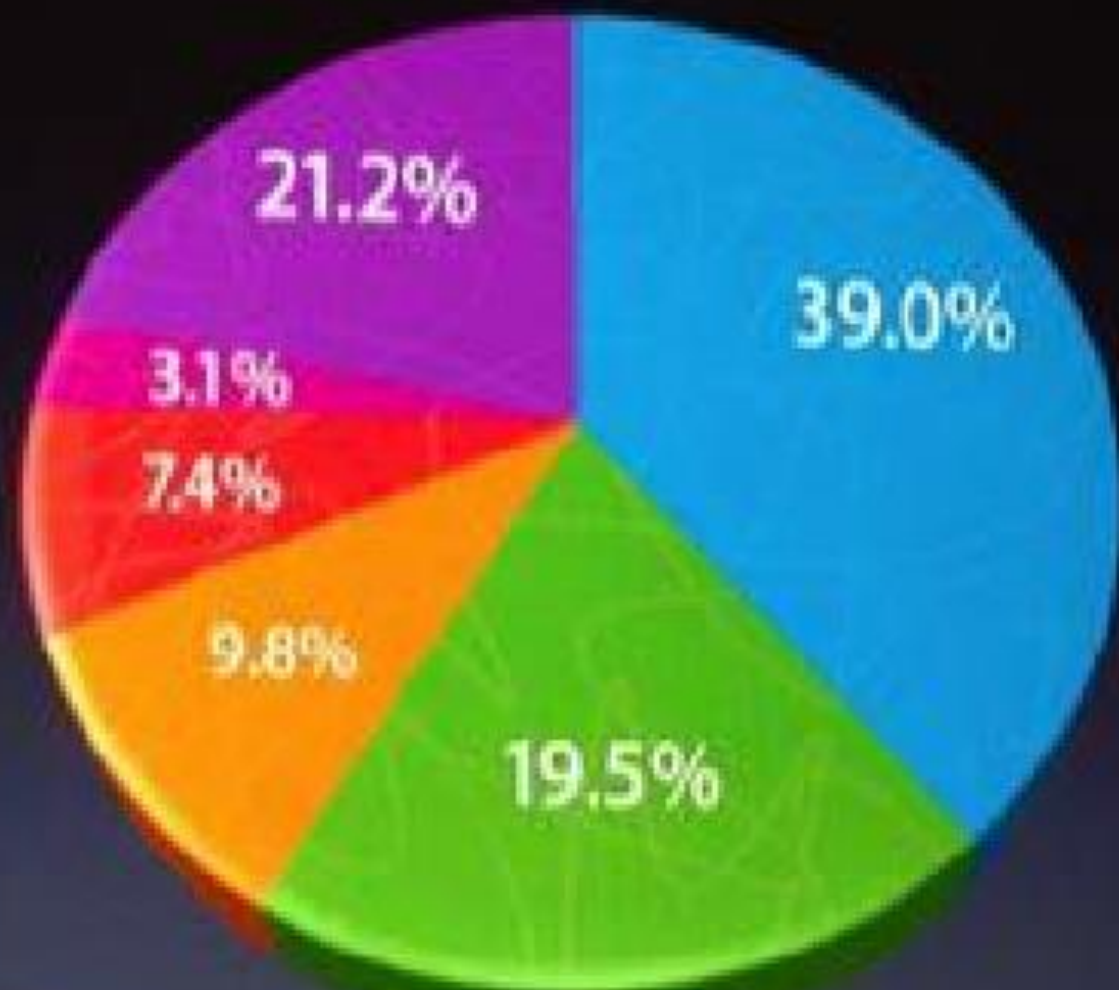


1399683	635200		2908933	2106392
6962578	5401748	8875318	8275961	4885854
1019462	3922722	6581982	6074504	3018484
1345609	1336265	2815039	6636256	9978939
	3111844	2917750		
724027	6656872	<b>38120066</b>	3880520	4796559
	1877879	11575704	9913774	2995330
844322			6539407	4681639
	19631599	12914651		
587106		629120	26538203	19383475
	12817894			9853722
5731663	3601157	927030	5903286	
1628760	1050228	1885932	5235100	742404
			4419036	

# Good Viz vs Bad Viz



# U.S. SmartPhone Marketshare



# Visualizations to support spatial analysis

Distributions and frequency

Category comparisons

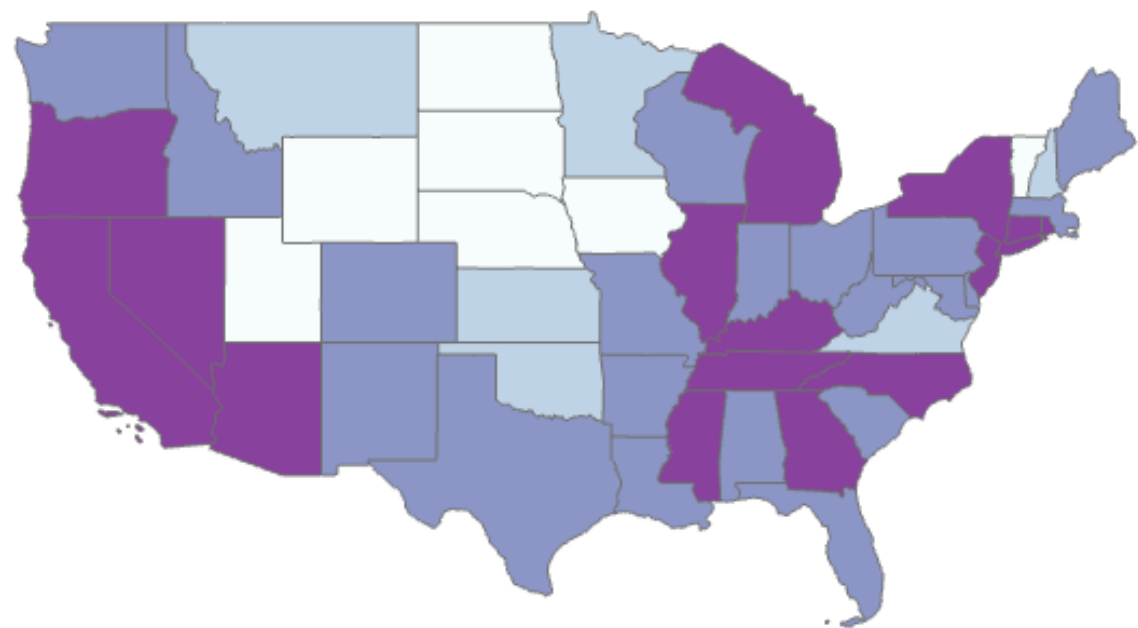
Relationships and correlations

Change over time or distance

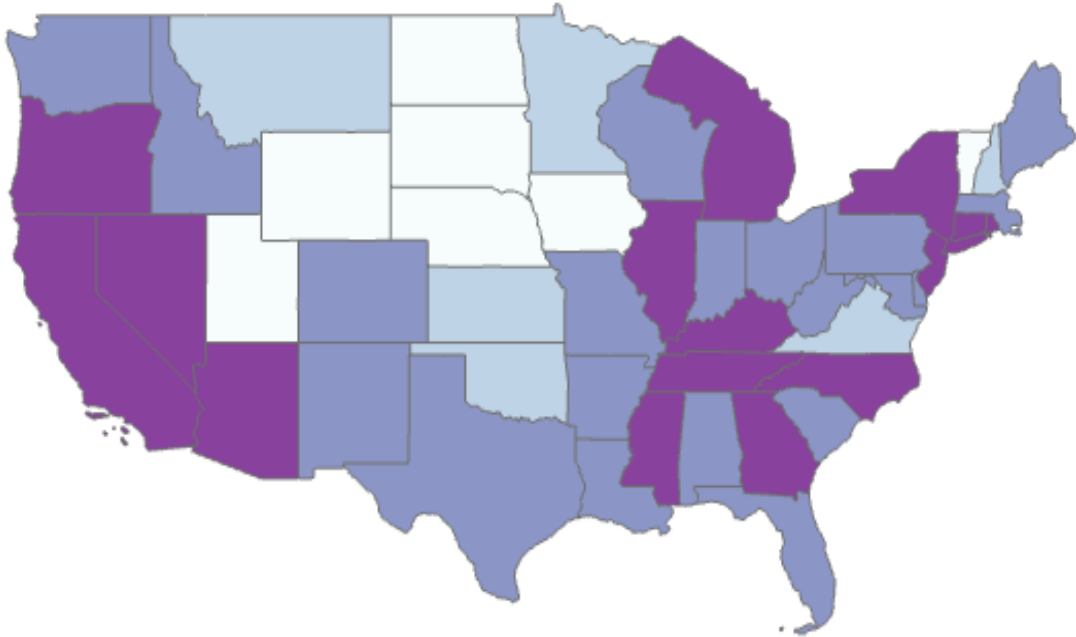
**When a map (alone)  
isn't the best option...**



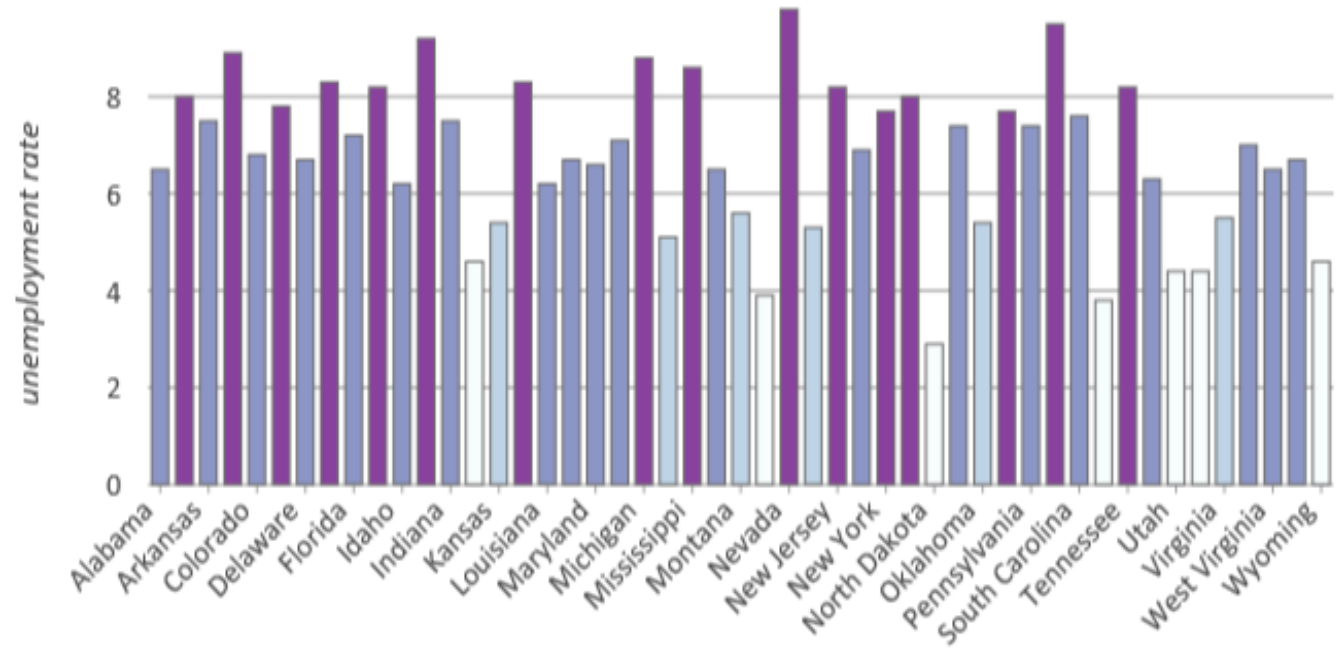
When a map (alone)  
isn't the best option...



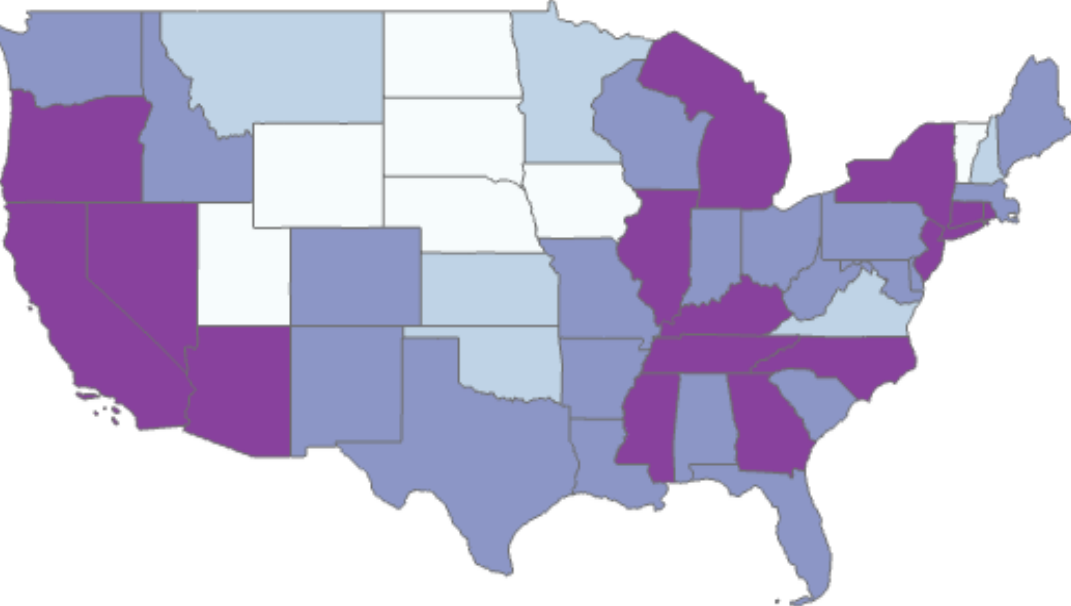
When a map (alone)  
isn't the best option...



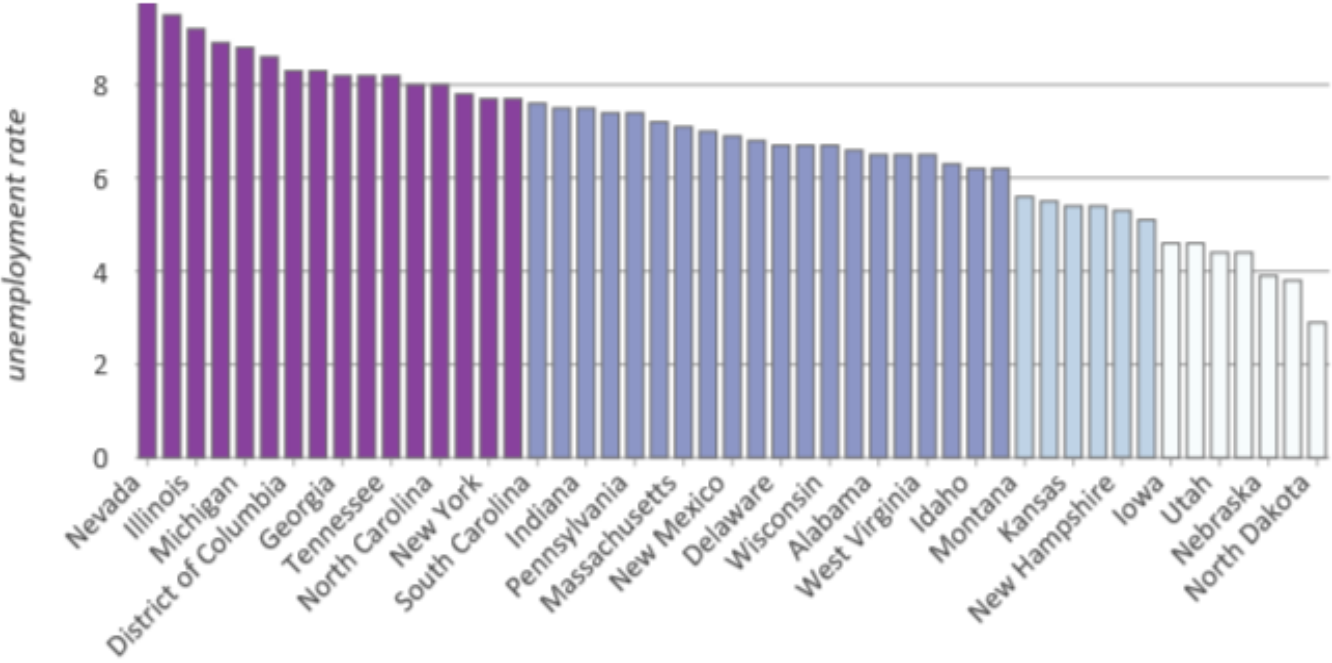
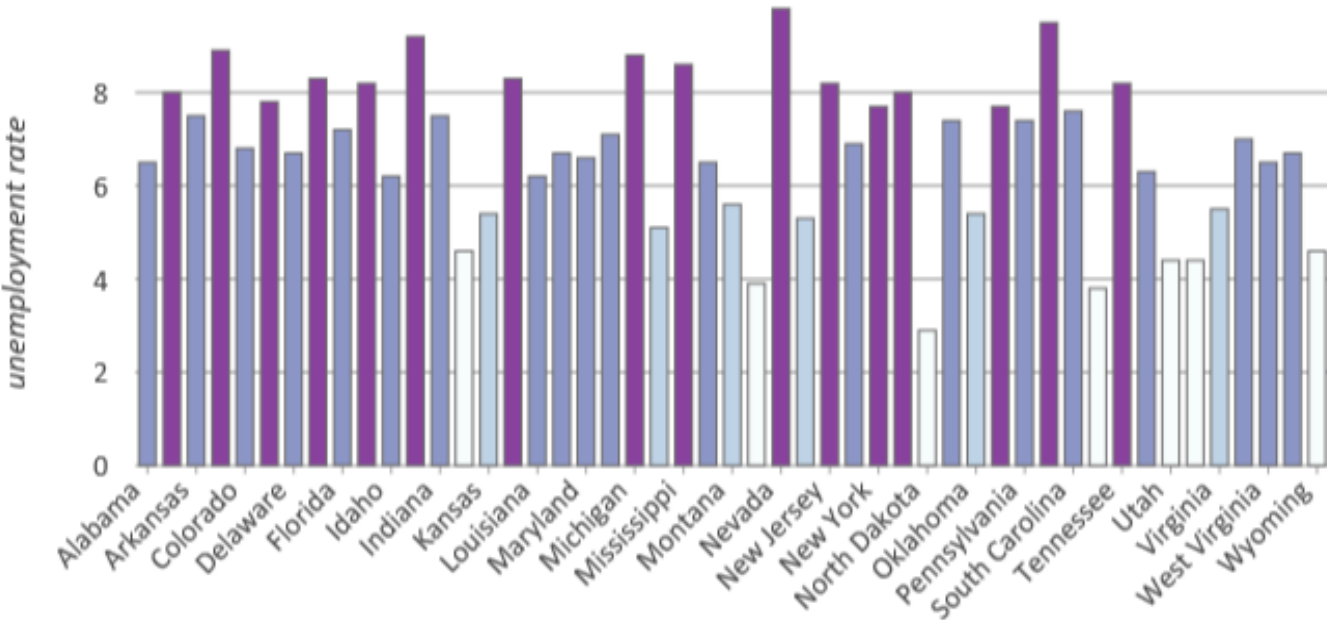
Comparison of unemployment rate by State



# When a map (alone) isn't the best option...



Comparison of unemployment rate by State



# Key Takeaways

Something happens  
Everywhere

Open Source paired  
with COTS can  
enhance workflows

Data Visualization &  
Messaging drive  
quick decisions

# Resources

- [Learn ArcGIS](#)
- [R Bridge](#)
- [ArcGIS Pro](#)
- [Spatial Statistics Page](#)

# Questions?

Kristen Hocutt- [khocutt@esri.com](mailto:khocutt@esri.com)

Lu Zhang- [l.zhang@esri.com](mailto:l.zhang@esri.com)