

---

# Geographical Analysis

## Introduction to Geographic Analysis and Visualization

# Objectives

---

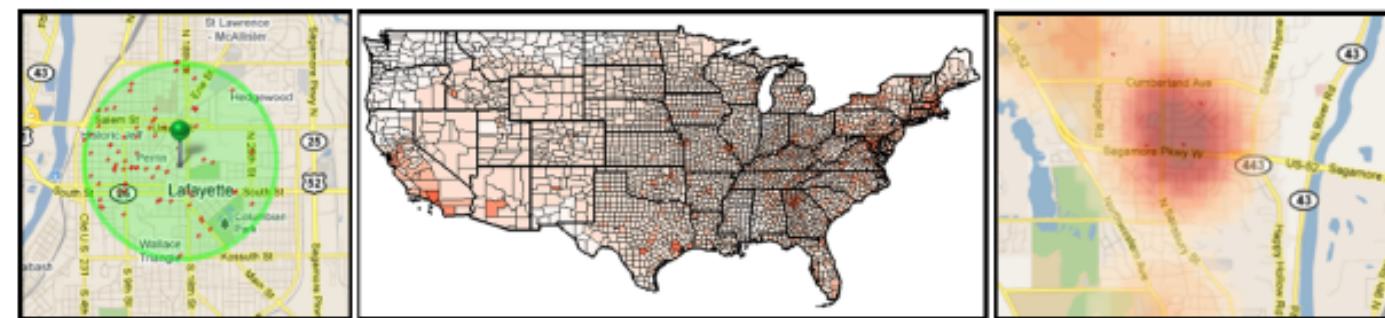


Objective

Describe tools and techniques that are designed to support analyses that focus on datasets with a geographic component

# Geovisualization

- | Primarily denotes tools and techniques that are designed to support analyses that focus on datasets with a geographic component
- | Visual representations are designed and built utilizing cartographic principles
- | Look for trends over geographic regions



# Geographic Visualization

---

Utilizes sophisticated,  
interactive maps to explore  
information

Recently, focuses on  
incorporating temporal  
components into analysis



Moving towards a  
combination of interactive  
maps and statistical analysis  
methods

# Tobler's First Law of Geography



“Everything is related to everything else, but near things are more related than distant things”

---

# Geographical Analysis

## Thematic Maps

# Objectives

---



Objective

Describe tools and techniques that are designed to support analyses that focus on datasets with a geographic component

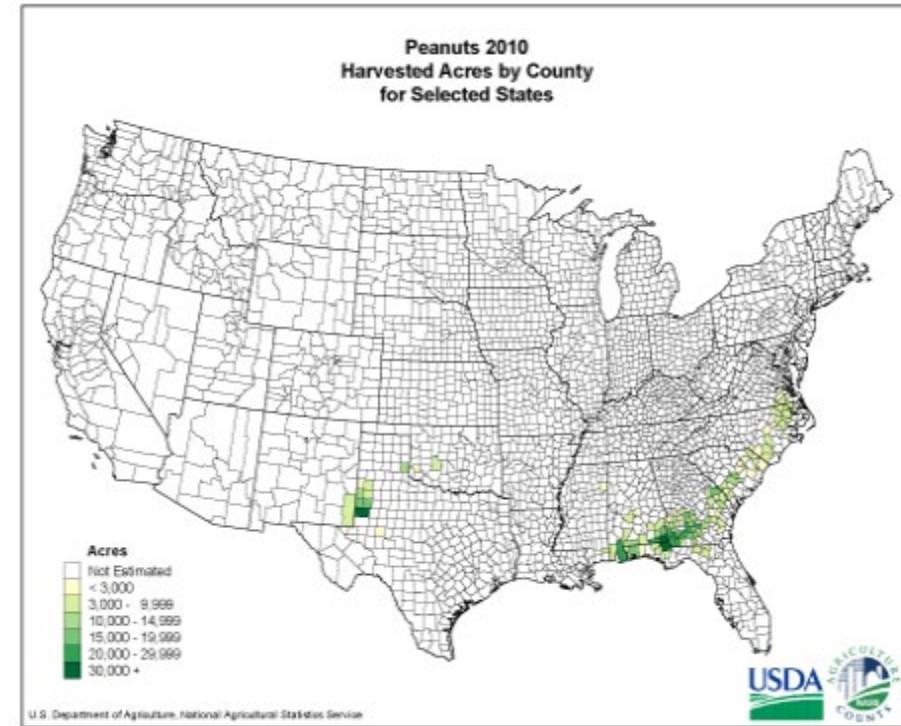
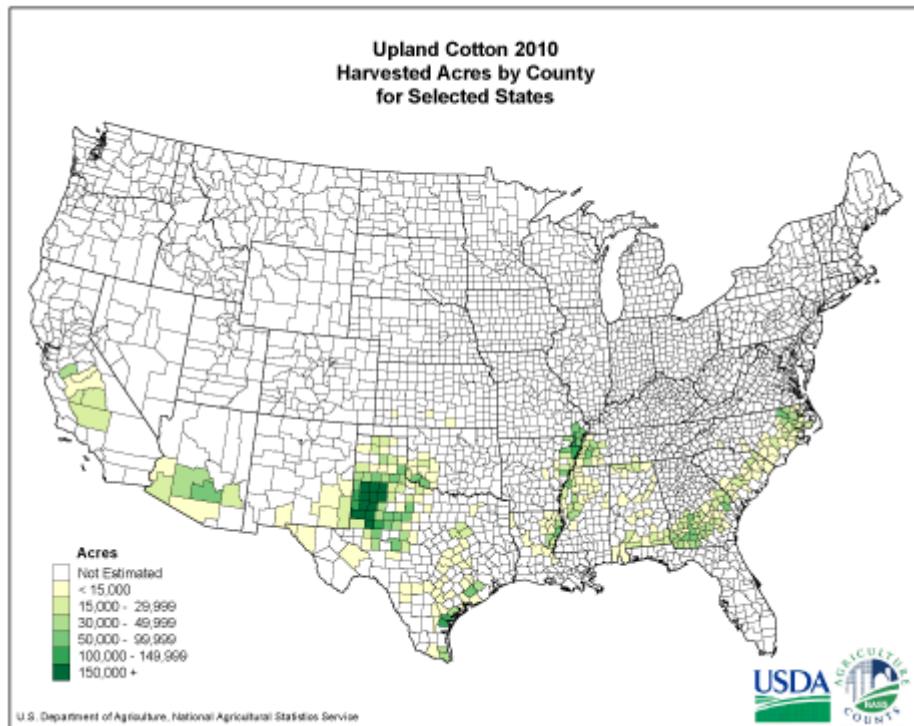
# Thematic Maps



- | A thematic map (or statistical map) is used to display the spatial pattern of a theme or attribute
- | Focus is on spatial pattern as opposed to a general-reference map which focuses on location

- | Goal of a thematic map is to emphasize spatial patterns of geographic attributes (e.g., population density)

# Examples



# How are Thematic Maps Used?



Locations

- To provide specific information about particular locations

Patterns

- To provide general information about spatial patterns

Comparisons

- To compare patterns on multiple maps



# Geographical Analysis

## Map Design: Coordinate System

# Objectives

---



Objective

Describe tools and techniques that are designed to support analyses that focus on datasets with a geographic component

# Basic Steps for Communicating Map Information



- 1. Consider what the **real-world distribution** of the phenomenon might look like**
- 2. Determine the purpose of the map and the intended audience**
- 3. Collect data appropriate for the map's purpose**
- 4. Design and construct the map**
- 5. Determine whether users find the map **useful/informative****

# 5 Types of Map Phenomena

Spatial Dimension

Point  
Phenomena

Linear  
Phenomena

Areal  
Phenomena

2.5D  
Phenomena

3D  
Phenomena

# Coordinate Systems

---

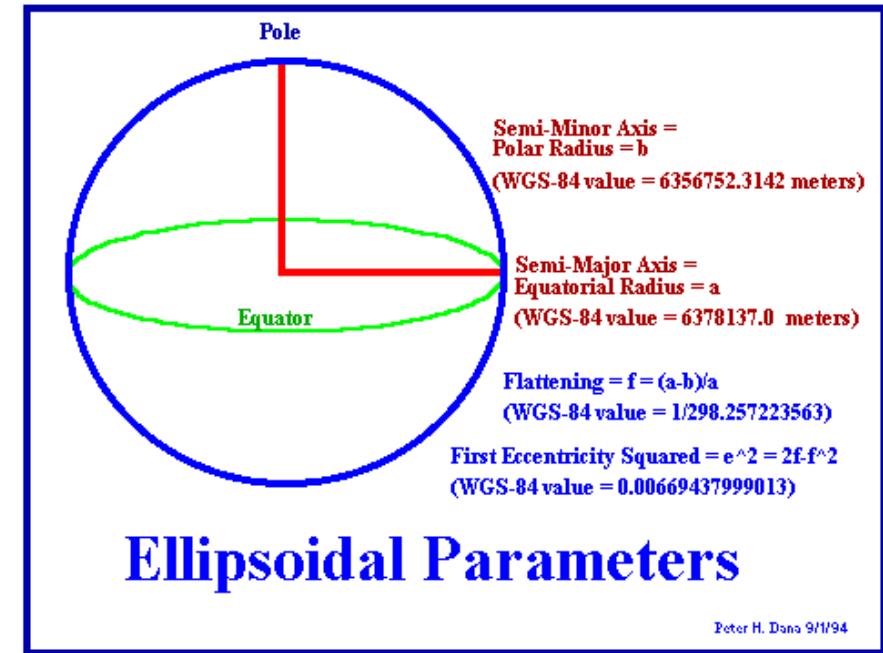


Types of Map Coordinate  
Systems:

- Latitude/Longitude
- Universal Transverse Mercator
- State Plane
- Metes and bounds

# Geodetic Datums

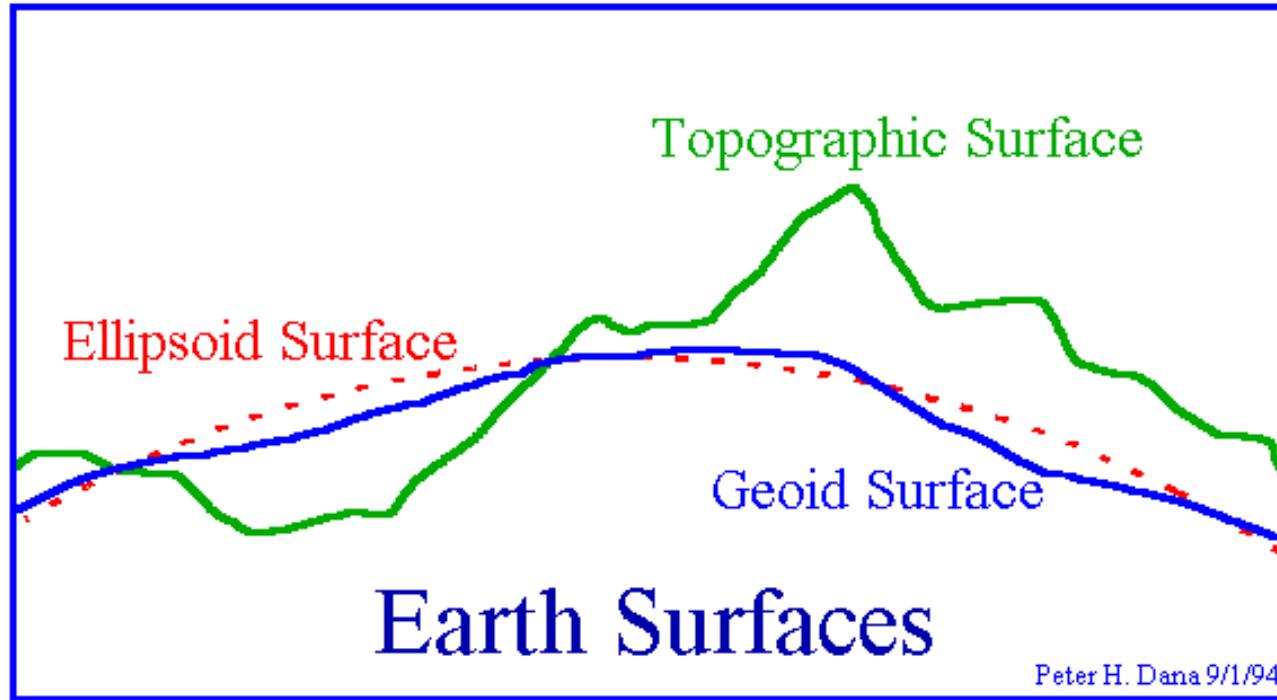
- | Datums define the size and shape of the earth
- | The initial point of the coordinate system is determined by the projection, ellipse model and datum



Common datums in the US:

- North American Datum 1927
- North American Datum 1983
- World Geodetic System

# Types of Geodetic Datums



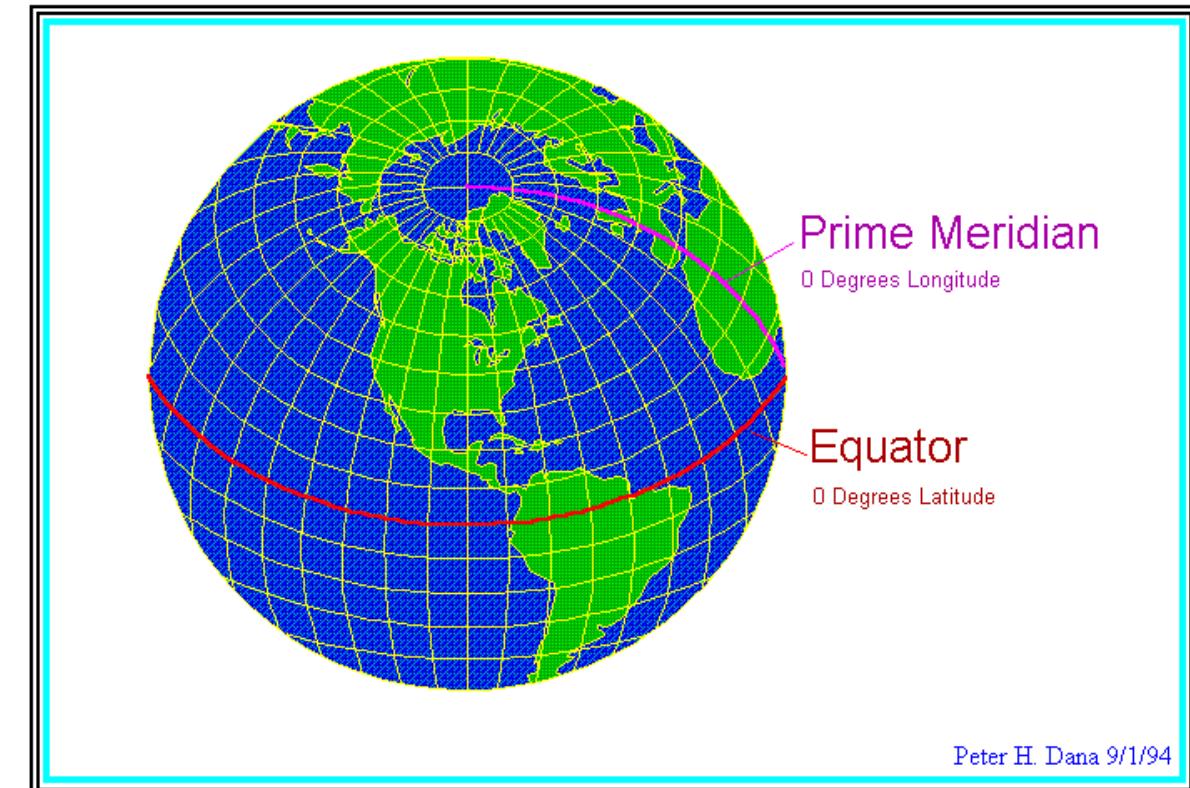
Topographical Surface

Sea Level

Gravity Models

# Coordinate Systems: Latitude and Longitude

- | Most commonly used coordinate system
- | No transformations necessary between areas
- | Scale, shape and direction distortion all increase with increasing area of interest



Peter H. Dana 9/1/94

---

# Geographical Analysis

## Map Projections

# Objectives

---



Objective

Describe tools and techniques that are designed to support analyses that focus on datasets with a geographic component

# Map Projections

## | Area/ Shape Distortion

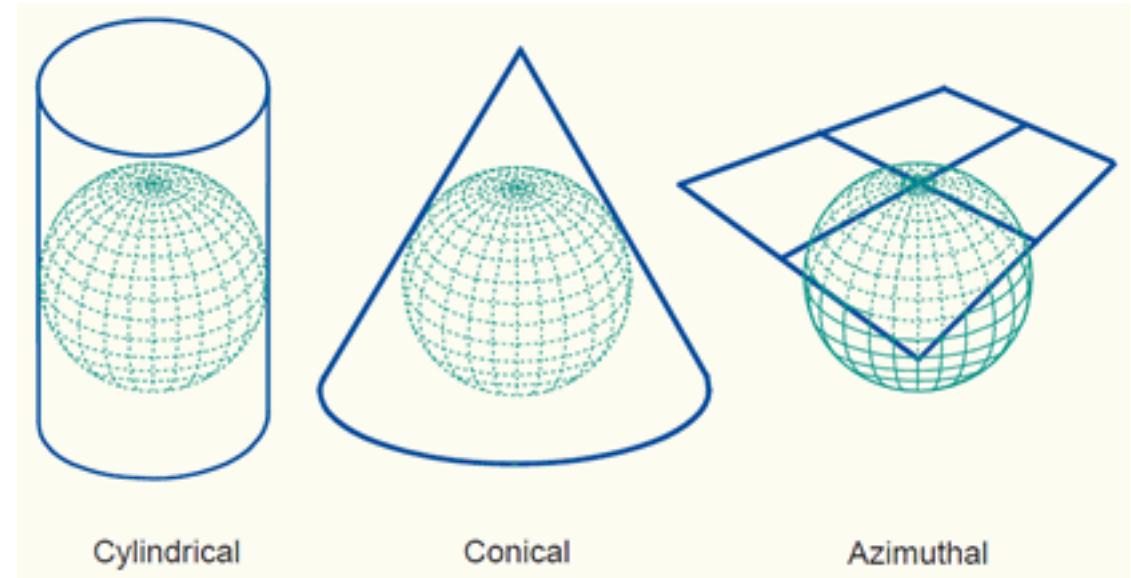
- **Equivalent** – Area is similar on globe and flat map, but shape is not
- **Conformal** – Shape is similar on globe and flat map, but the area is not

| <do we have a pic to illustrate this? >

# Map Projections

## | Shape of Projections

- **Cylindrical** – projection of sphere onto a cylinder
- **Conic** – projection of sphere onto a cone
- **Azimuthal** – projection of a sphere onto a plane
- **Miscellaneous**

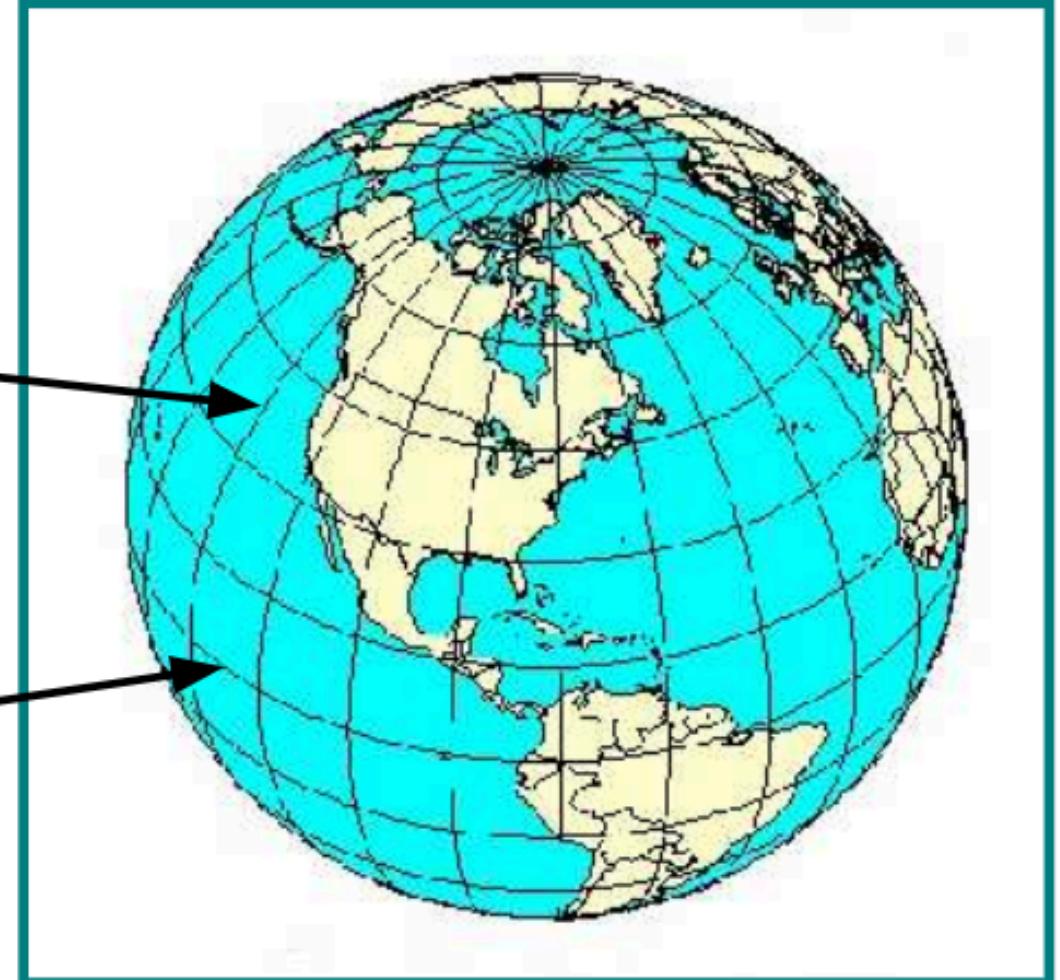


# Map Projections

---

**Meridians**  
North-South lines (such as Longitude)

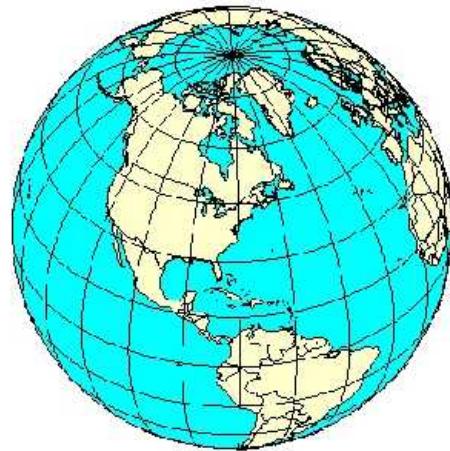
**Parallels**  
East-West lines (such as Latitude)



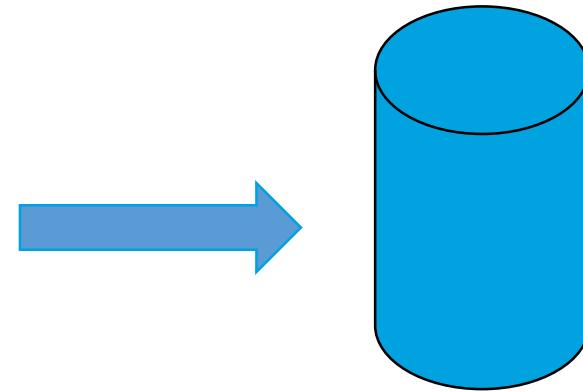
# Cylindrical Projections

## Characteristics:

- Straight meridians and parallels
- Meridians equally space
- Parallels unequally spaces

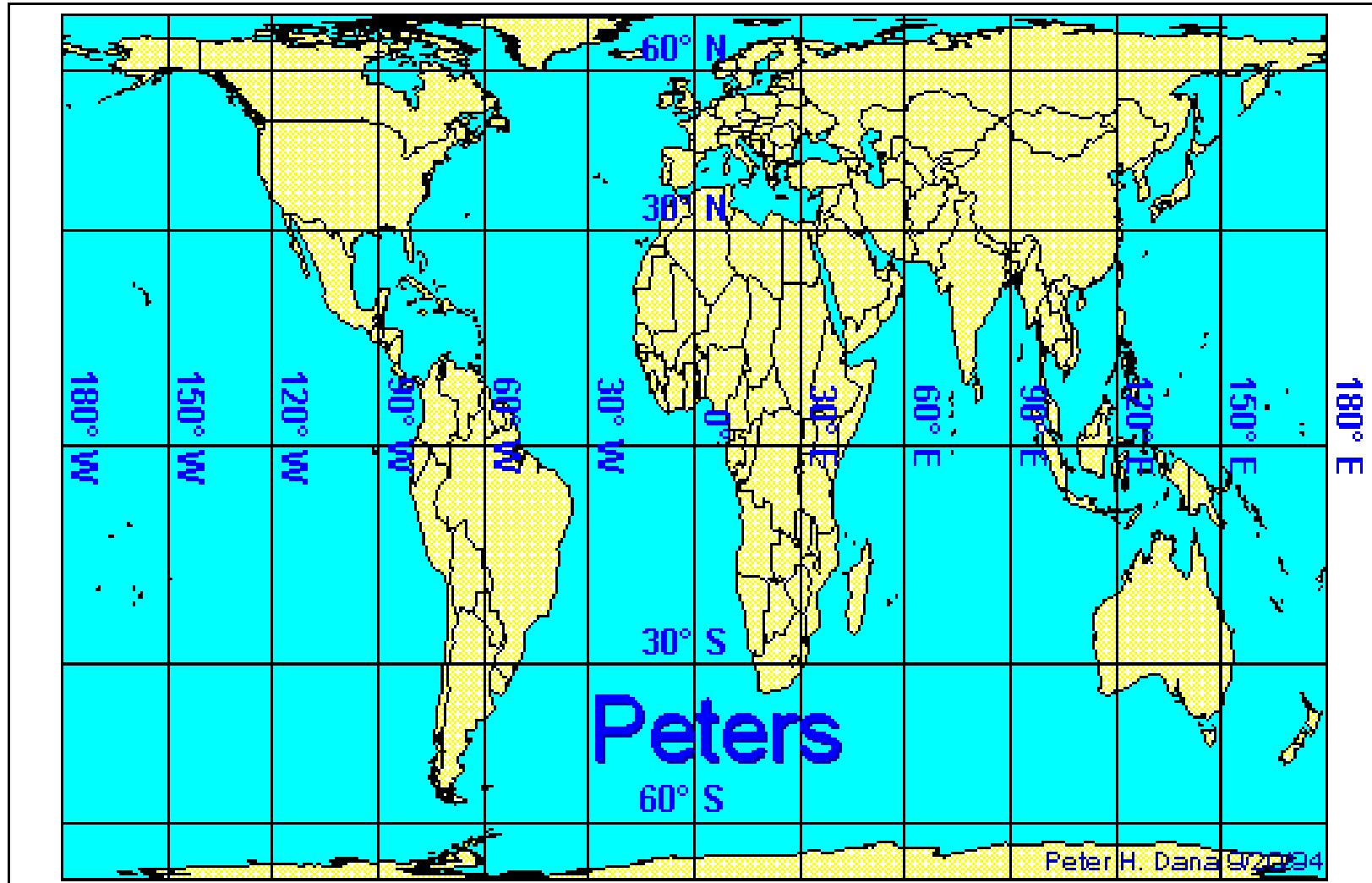


**Mercator**

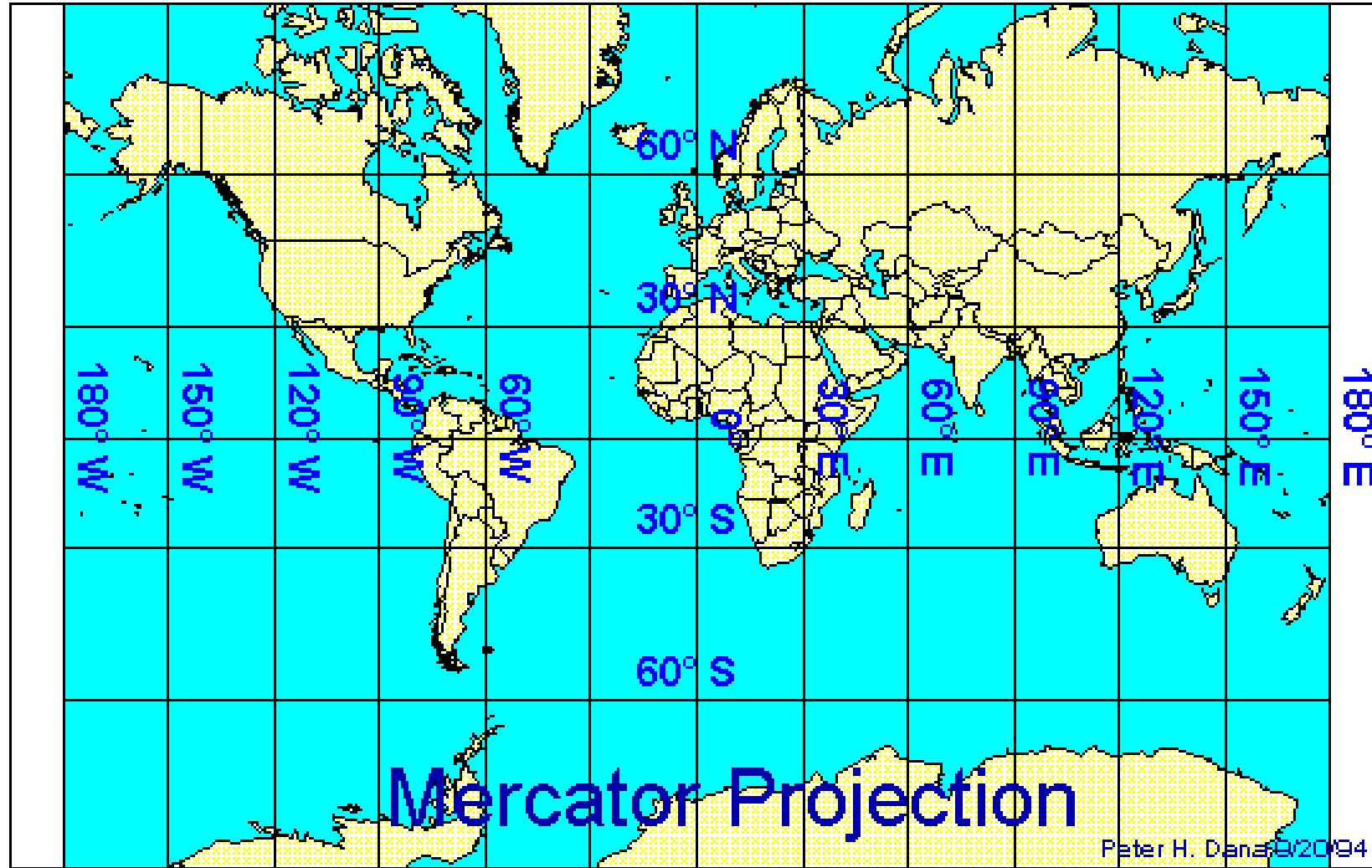


**Universal Mercator**

# Cylindrical Projections - Peters



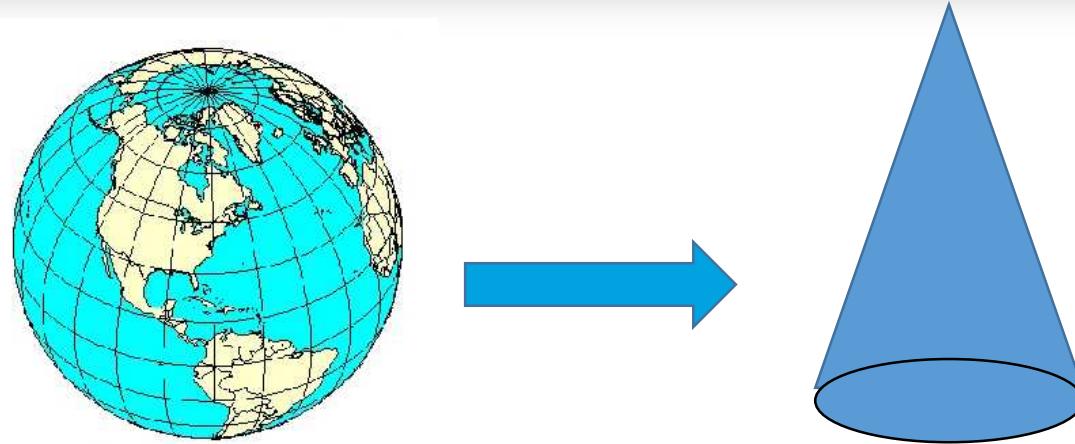
# Cylindrical Projections - Mercator



# Conic Projections

## Characteristics:

- Straight meridians, curved parallels
- Meridians radiate from poles
- Parallels may be equally spaces

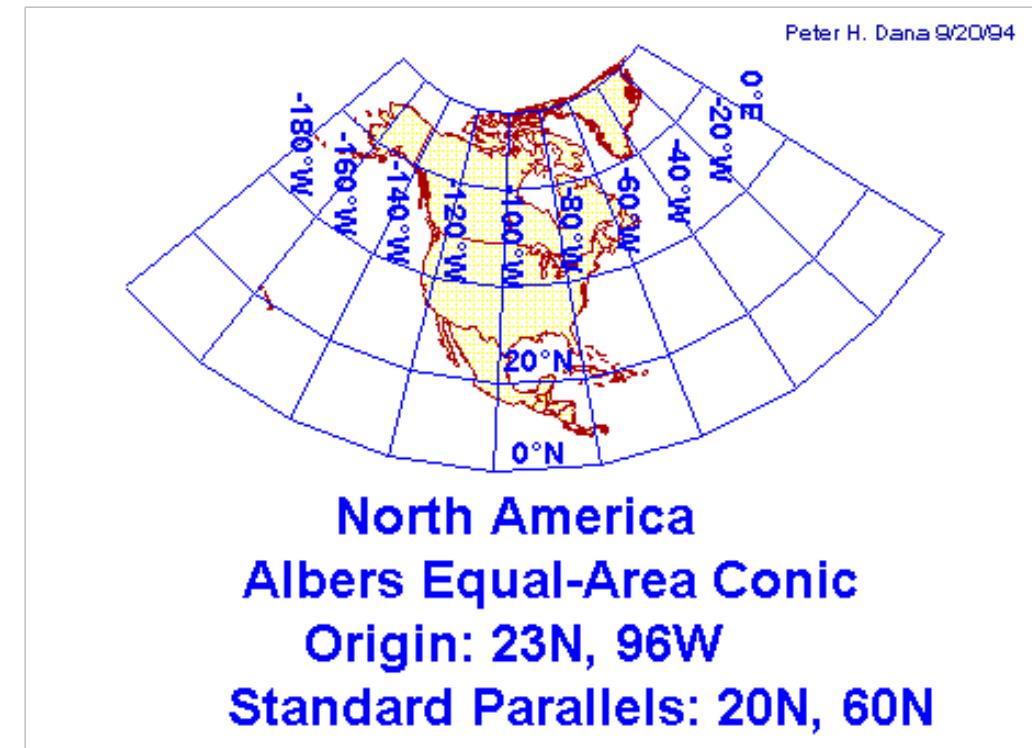


## Common conic projections

- Albers
- Lambert
- Polyconic

# Conic Projections – Albers Equal Area

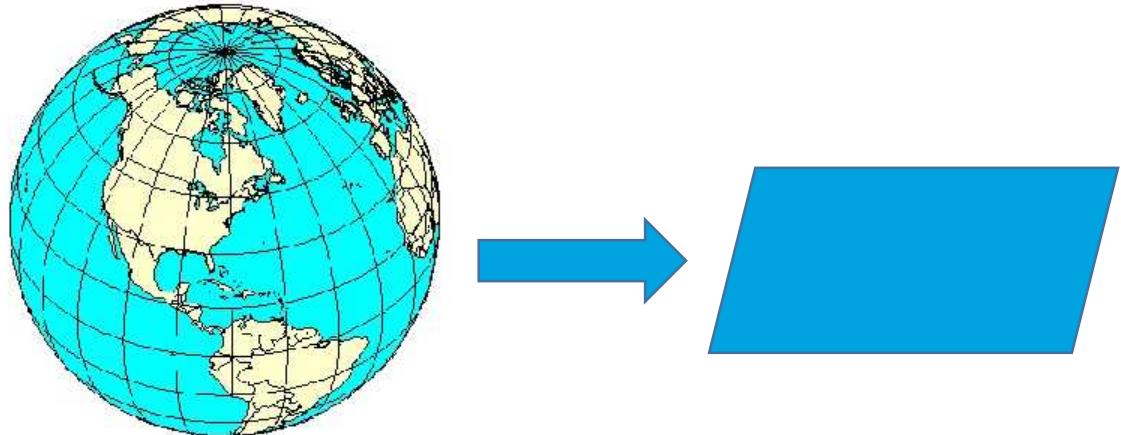
- | Direction, area and shape are distorted away from the standard parallels
- | Area and directions are true only in limited portions of a map



# Azimuthal Projections

## Characteristics:

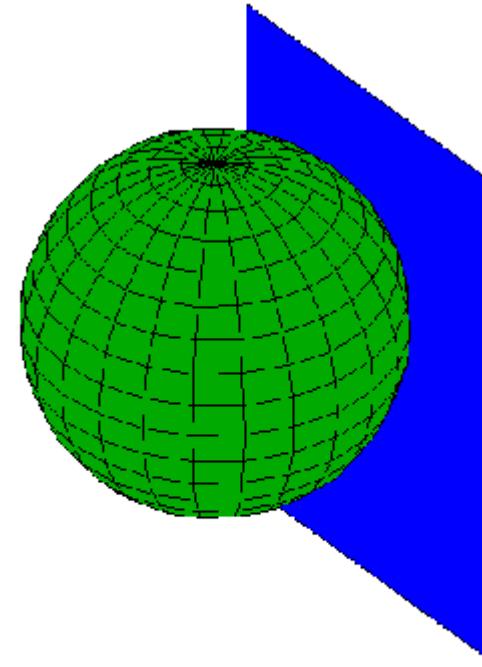
- Straight or curved meridians, curved parallels
- Meridians radiate from poles
- Parallels may be equally spaced



# Azimuthal Projections - Orthographic

- | Simplest form is orthogonal projection
- | Adequate only for very small areas
- | Scale and area distortion increases as distance from the tangent center increases

Peter H. Dana 9/20/94



**Planar Projection Surface**



---

# Geographical Analysis

## Map Design: Map Elements and Typography

# Objectives

---



Objective

Describe tools and techniques that are designed to support analyses that focus on datasets with a geographic component

# Map Elements and Typography



## Map Elements

Frame line and  
neat line

Legend

Mapped area

Data source

Inset

Scale

Title and  
subtitle

Orientation

Legend

Title



# Legend



| Map element that **defines all** of the thematic symbols of the map

## Symbols

- if self explanatory or not directly related to the map's theme can be omitted
- should be vertically centered with their definition

| Textual definitions should be horizontally centered

| Legend heading can be used to further explain the maps theme

| Scale is added to indicate the distances

# Typography

## | Use

- Use bold and italic type sparingly
- A realistic lower limit for all type size
- Type size should correspond with the size or importance of map features

## | Do not use

- decorative type families
- Script, cursive and ornate style
- Two type families on a given map

## | Limited Use

- Reserve italics for water features or to identify the data source

# Specific Guidelines for Map Elements



| Orient type horizontally

| Avoid overprinting

| Ensure type is placed so that it is clearly associated with the feature it is representing

---

# Geographical Analysis

Introduction to Choropleth Maps and Color Schemes

# Objectives

---



Objective

Differentiate types of  
geographic  
visualizations

# Introduction to Choropleth Maps

---

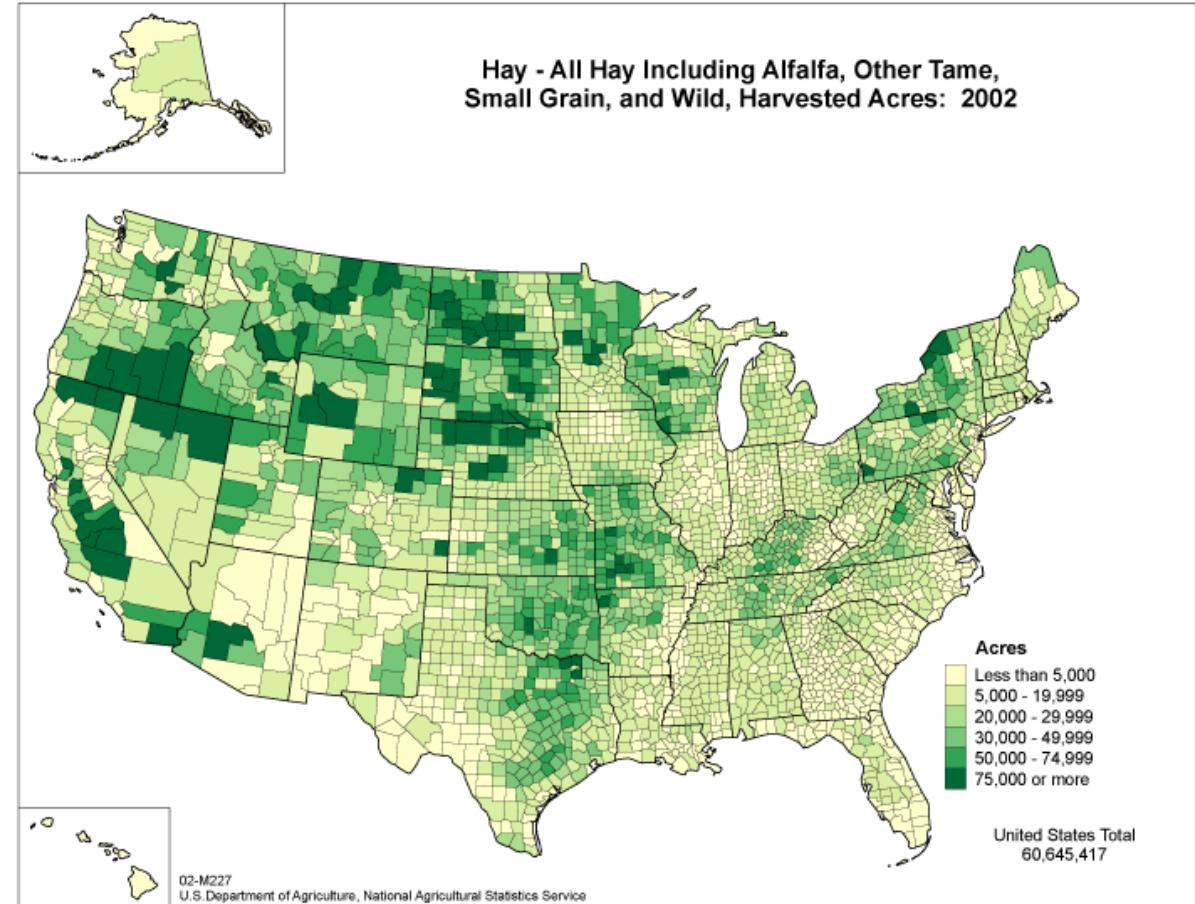
| Earliest known choropleth map was created in 1826 by Baron Pierre Charles Dupin

| Term choropleth map was coined in 1938 by John Kirtland Wright

| Choropleth maps are based on statistical data aggregated over defined geographical regions

# Choropleth Maps

- | Areas of the map are shaded in proportion to a measured variable
- | Coloring is based on a classification (histogram binning) of the distribution of the measured variable



# Coloring Choropleth Maps



Relates to number of classes

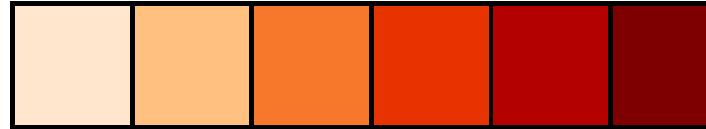
Cartographic rule =  
5-7 classes

Colors

Color Schemes:  
sequential,  
divergent, qualitative

Choose carefully to  
allow viewers to see  
trends

# Color Schemes: Sequential



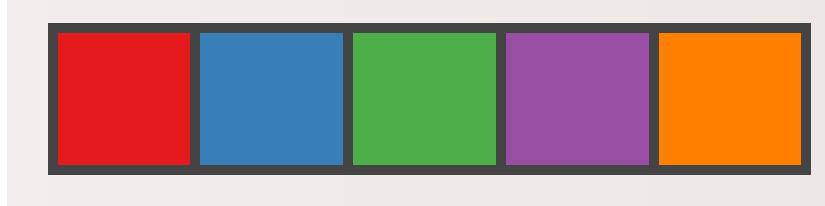
- | Suited for ordered data
- | Lightness steps dominate the look of the scheme
- | Light values are low data values, dark are high
- | Good for Ordinal, interval and ratio data types

# Color Schemes: Diverging



- | Puts an emphasis on critical midrange values
- | Color change represents deviation from a meaningful midrange critical value
- | Good for ratio data types where looking at data above and below a 'zero' point

# Color Schemes: Qualitative



- | Does not imply magnitude difference
- | Used to show differences between classes
- | Good for Nominal data types

---

# Geographical Analysis

## Data Classifications

# Objectives

---



Objective

Differentiate types of  
geographic  
visualizations

# Class Interval Selection



| Choices for optimizing the class interval selection are highly dependent on the underlying data distribution

| Similar to the concept of histogram binning

| Popular choices for class interval selection include

- Equal interval selection
- Jenks' Natural Breaks
- Minimum boundary error

# Equal Interval



| Classifies data such that each case occupies an equal interval along the number line

## Advantage

- Easy to compute

$$\frac{\text{range}}{\text{NumClasses}} = \frac{\text{High} - \text{Low}}{\text{NumClasses}}$$

## Disadvantage

- Fails to consider how data are distributed

# Quantiles



| The number of color bins (classes) will determine the number of quantiles

## Advantages

- Easy to compute
- Percentage of observations in each class will be the same
- Class assignment is based on rank order

$$\text{Number In Class} = \frac{\text{TotalSamples}}{\text{NumClasses}}$$

## Disadvantages

- Fails to consider data distribution
- Dissimilar data can be placed into same class

# Mean-Standard Deviation



| Classes are formed by adding or subtracting some number of standard deviations from the mean

## | Advantages:

- If data are normally (or near normally) distributed, the mean serves as a useful dividing point
- Legend will contain no gaps

## | Disadvantage:

- Only works well only with data that are normally distributed

# Maximum Breaks



| Goal is to consider individual data values and group those that are similar

- Order data from low to high and differences between adjacent values are computed
- The largest differences (“breaks”) are used for the class divisions

| Advantages:

- Easy to compute

| Disadvantage:

- May miss natural clusters

# Natural Breaks



| Data values are examined visually to determine logical breaks within the data

| Goal is to **minimize** differences between data values in the same class and **maximize** differences between different classes

## | Advantage:

- Tries to take into account natural underlying structure of the data

## | Disadvantage:

- Subjective, different mapmakers may choose other values

# Optimal Classification



| Same as natural breaks, but minimizes an objective function

| Goal is to place values into groups

| Measure an “error” – common is to use the sum of absolute deviations about class medians

| Calculate each class median and then add the resulting sums of absolute deviations

| Many computer based algorithms have been developed to do this:

- Jenks-Caspall
- Fisher-Jenks

# Optimal Classification



## Advantages

- Good empirical method for grouping data
- Can assist in determining appropriate number of classes

## Disadvantages

- Hard to explain to novice users
- May leave gaps in the map legend

# Utilizing Spatial Context



## Optimal Map (with spatial constraint)

Raw Data

1, 2, 3, 4, 5, 6, 9, 10, 11, 12, 13, 16, 17, 18, 19, 20

1	9	5	6
20	3	16	12
19	4	13	11
18	10	17	2

$$GADF = 0.70$$

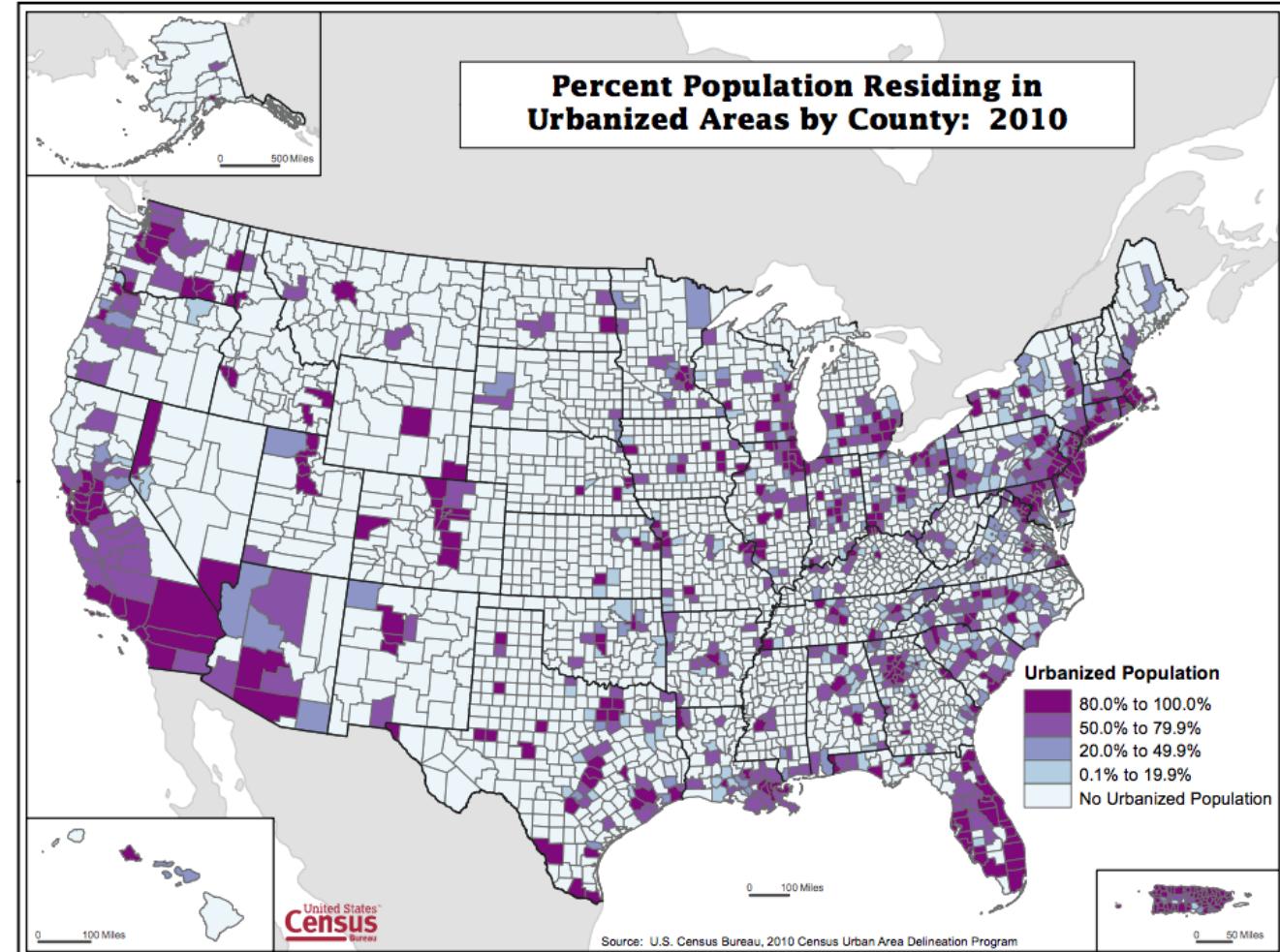
$$C_F = 6/16 = 0.38$$

# Choropleth Maps

| Need to standardize the statistics

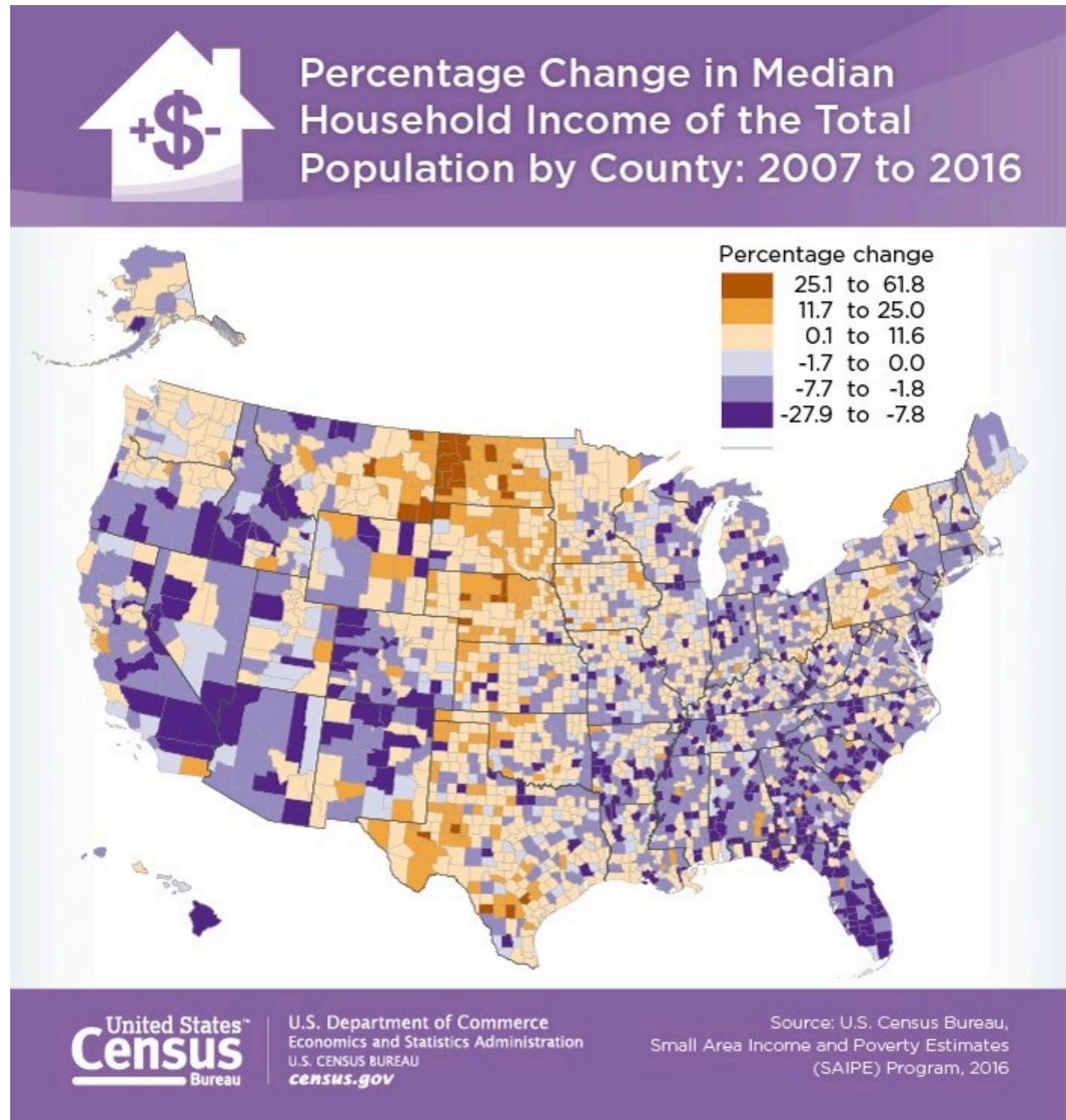
- Divide by population
- Divide by area

| Modifiable areal unit problem – source of statistical bias occurring when data is aggregated into districts



# Ecological Fallacy

- Inferences about individuals are based solely upon aggregate statistics collected for the group to which those individuals belong
  - Assumes that individual members of a group have the average characteristics of the group at large
  - Group characteristics do not necessarily apply to individuals within that group



---

# Geographical Analysis

## Spatial Statistics

# Objectives

---



Objective

Explain spatial statistics

# Tobler's First Law of Geography



“Everything is related to everything else, but near things are more related than distant things”

# Spatial Statistics



| Spatial dependency is the co-variation of properties within a geo-space

| If correlation exists (either positively or negatively), there are at least three possible explanations:

- There is a simple spatial correlation relationship
- Spatial causality
- Spatial interaction

# Covariance



$$Cov(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

**Covariance:** The pattern of common variation observed in collection of two (or more) datasets, or partitions of a single dataset

# Pearson's Correlation Coefficient

---

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

**Pearson's Correlation Coefficient:** A measure of similarity between two or more paired datasets

# Entropy



$$I = - \sum_{i=1}^k p_i \log(p_i)$$

**Entropy:** A measure of the amount of pattern, disorder, or information in a set  $\{x_i\}$  where  $p_i$  is the proportion of events or values occurring in the  $i^{\text{th}}$  class or range.

# Diversity

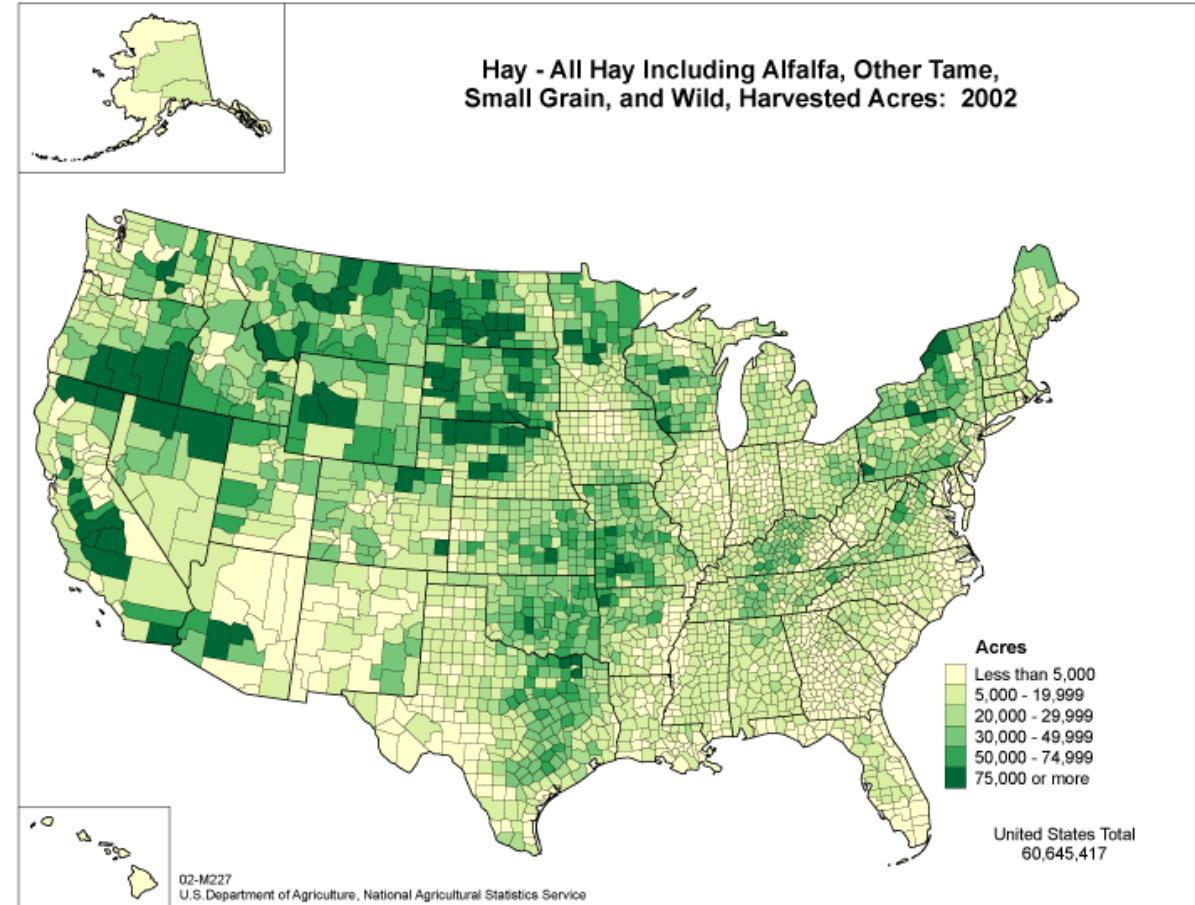


$$Div = \frac{-\sum_{i=1}^k p_i \log(p_i)}{\log(k)}$$

**Diversity:** entropy standardized by the number of classes, k

# Spatial Statistics

- | For spatial data, we want to find related regions
- | Can do analysis prior to the visualization to find areas that are statistically correlated and focus the visual representation on these areas
- | One method of doing this is using spatial autocorrelation



# Issues in Spatial Statistics



Scaling

Sampling

Logical Fallacy

Ecological Fallacy

# Distance and Direction

---

| Knowledge of location also allows the analyst to determine the distance and direction between objects



| Many types of spatial analysis require the calculation of a table expressing the relative proximity of pairs of places

# Elements of Matrix W

---

## | Elements of Matrix W:

- **1** if the places share a **common boundary**, else **0**
- The **length of any common boundary between places**, else **0**
- A decreasing function of the distance between the places, or between their representative points



# Geographical Analysis

## Spatial Autocorrelation

# Objectives

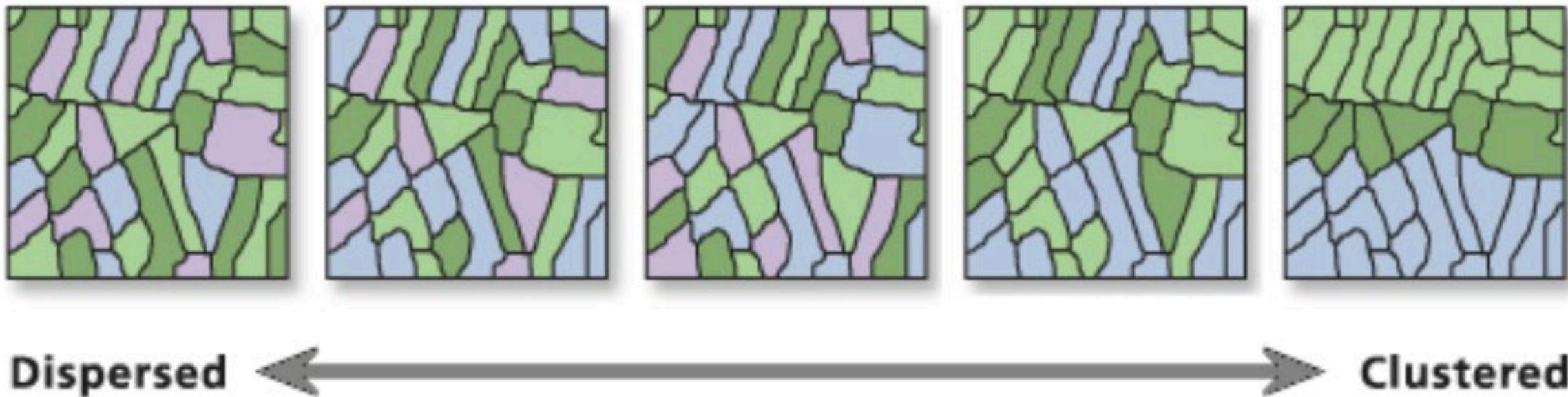
---



Objective

Explain spatial statistics

# Spatial Autocorrelation



# Moran's I



$$I = \frac{N}{\sum_i \sum_j w_{ij}} \frac{\sum_i \sum_j w_{ij}(X_i - \bar{X})(X_j - \bar{X})}{\sum_i (X_i - \bar{X})^2}$$

- | Given a set of features and an associated attribute, Global Moran's I evaluates whether the pattern expressed is clustered, dispersed or random
- | Values near +1.0 indicate clustering while values near -1.0 indicate dispersion
- | Global Moran's I function also calculates a Z score value that indicates whether we can reject the null hypothesis of “there is no spatial clustering”

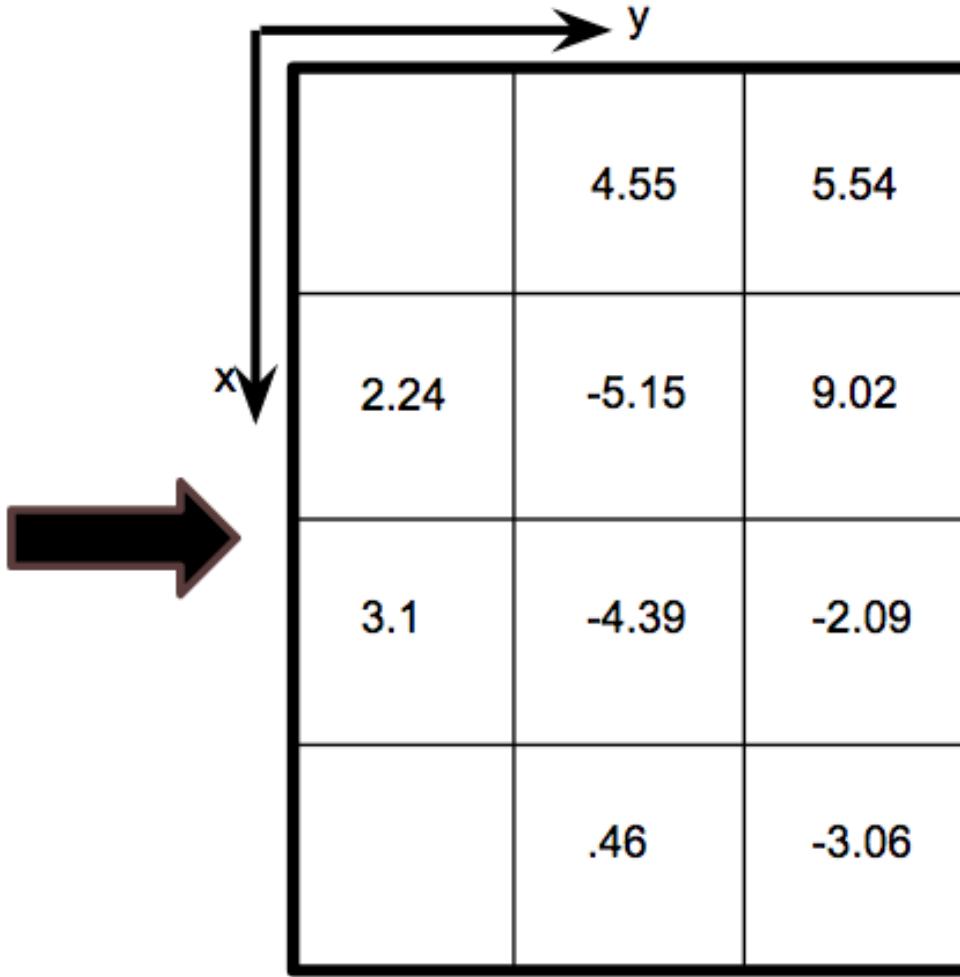
# What is a Z Score?



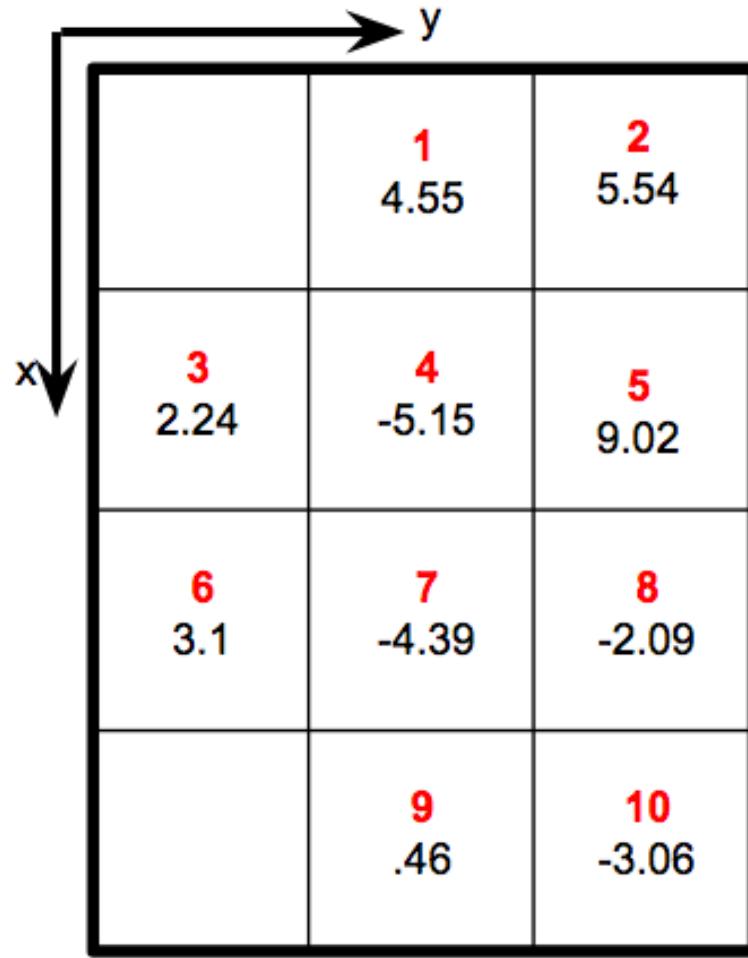
- | Most statistical tests begin by identifying a null hypothesis
- | The Z score is a measure of **standard deviations**
- | The Z score is associated with a **normal distribution**
- | Critical Z score values when using a 95% confidence level are **-1.96 and +1.96 standard deviations.**

# Calculating Moran's I

X	Y	Z
1	2	4.55
1	3	5.54
2	1	2.24
2	2	-5.15
2	3	9.02
3	1	3.1
3	2	-4.39
3	3	-2.09
4	2	.46
4	3	-3.06



# Adjacency Matrix

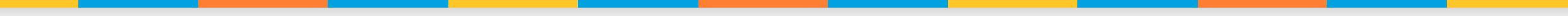


A 3x3 matrix with row and column indices labeled 1 through 10. The matrix has red numbers and black numbers. Red numbers are at positions (1,1), (1,2), (1,3), (2,1), (2,2), (2,3), (3,1), (3,2), and (3,3). Black numbers are at positions (1,4), (1,5), (1,6), (1,7), (1,8), (1,9), (1,10), (2,4), (2,5), (2,6), (2,7), (2,8), (2,9), (2,10), (3,4), (3,5), (3,6), (3,7), (3,8), (3,9), and (3,10).

			1 4.55	2 5.54					
	3 2.24	4 -5.15	5 9.02						
6 3.1	7 -4.39	8 -2.09							
		9 .46	10 -3.06						

	1	2	3	4	5	6	7	8	9	10
1	0	1	0	1	0	0	0	0	0	0
2	1	0	0	0	1	0	0	0	0	0
3	0	0	0	1	0	1	0	0	0	0
4	1	0	1	0	1	0	1	0	0	0
5	0	1	0	1	0	0	0	1	0	0
6	0	0	1	0	0	0	1	0	0	0
7	0	0	0	1	0	1	0	1	1	0
8	0	0	0	0	1	0	1	0	0	1
9	0	0	0	0	0	0	1	0	0	1
10	0	0	0	0	0	0	0	1	1	0

# Geary's C



$$C = \frac{(N - 1) \sum_i \sum_j w_{ij} (X_i - X_j)^2}{2W \sum_i (X_i - \bar{X})^2}$$

| Value of Geary's C lies between 0 and 2

- Values of 1 means no spatial autocorrelation
- Smaller than one means positive spatial autocorrelation
- Larger than one means negative spatial autocorrelation

| Geary's C is more sensitive to local spatial autocorrelation

# Local Indicators of Spatial Association



| Moran's I is a **global** measure of correlation

| The individual components can be mapped and tested for significance to provide an indication of clustering patterns within the study region

# Getis-Ord Statistic



$$G_i^* = \frac{\sum_{j=1}^N w_{ij}x_j - \bar{x}\sum_{j=1}^N w_{ij}}{S\sqrt{\frac{[N\sum_{j=1}^N w_{ij}^2 - (\sum_{j=1}^N w_{ij})^2]}{N-1}}}$$

Unlike Moran's I, the Getis and Ord statistic identifies the degree to which high or low values cluster together

# Significance Tests for Autocorrelation Indices

- Autocorrelation coefficients can be tested for statistical significance under two different model assumptions.
  - Classical statistical assumption of Normality
  - Assume that values are **independent** and **identically** distributed drawings from a Normal distribution
  - Observed pattern of a set of values is **assumed to be just one realization** from all possible random permutations
  - Utilizes Monte Carlo testing

---

# Geographical Analysis

## Spatial Scan Statistics

# Objectives

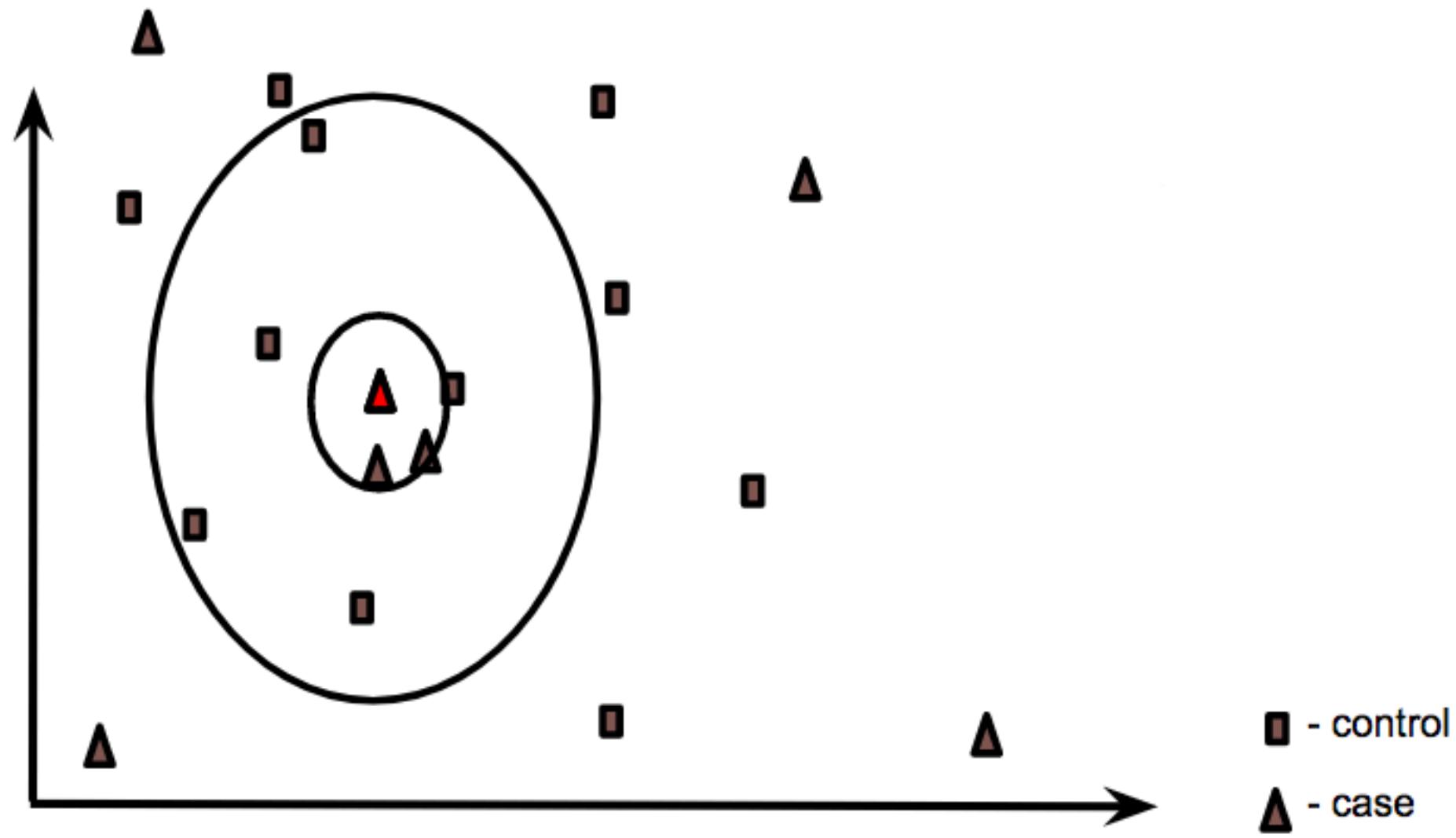
---



Objective

Explain spatial statistics

# Scan Statistics



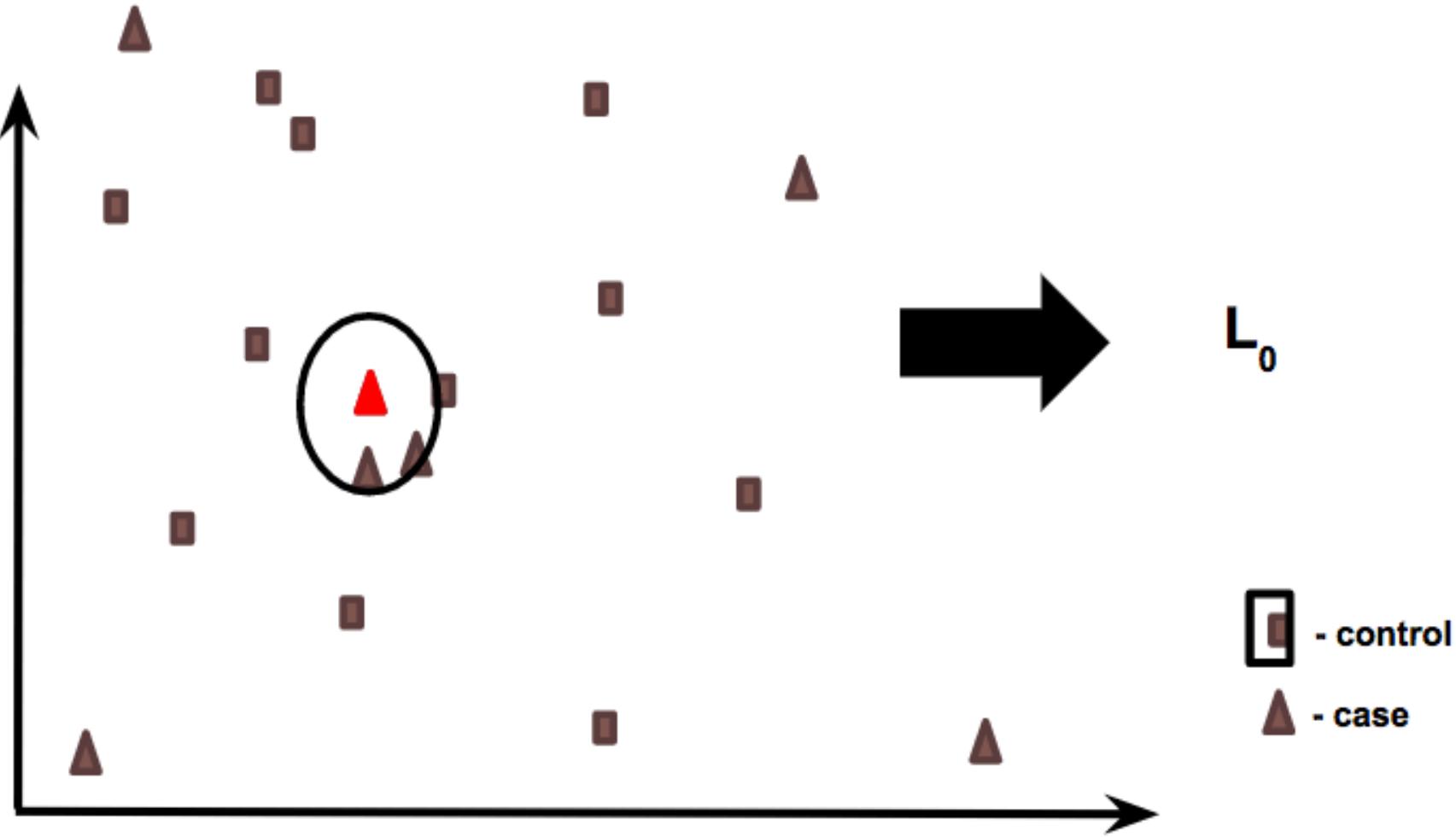
# Scan Statistics

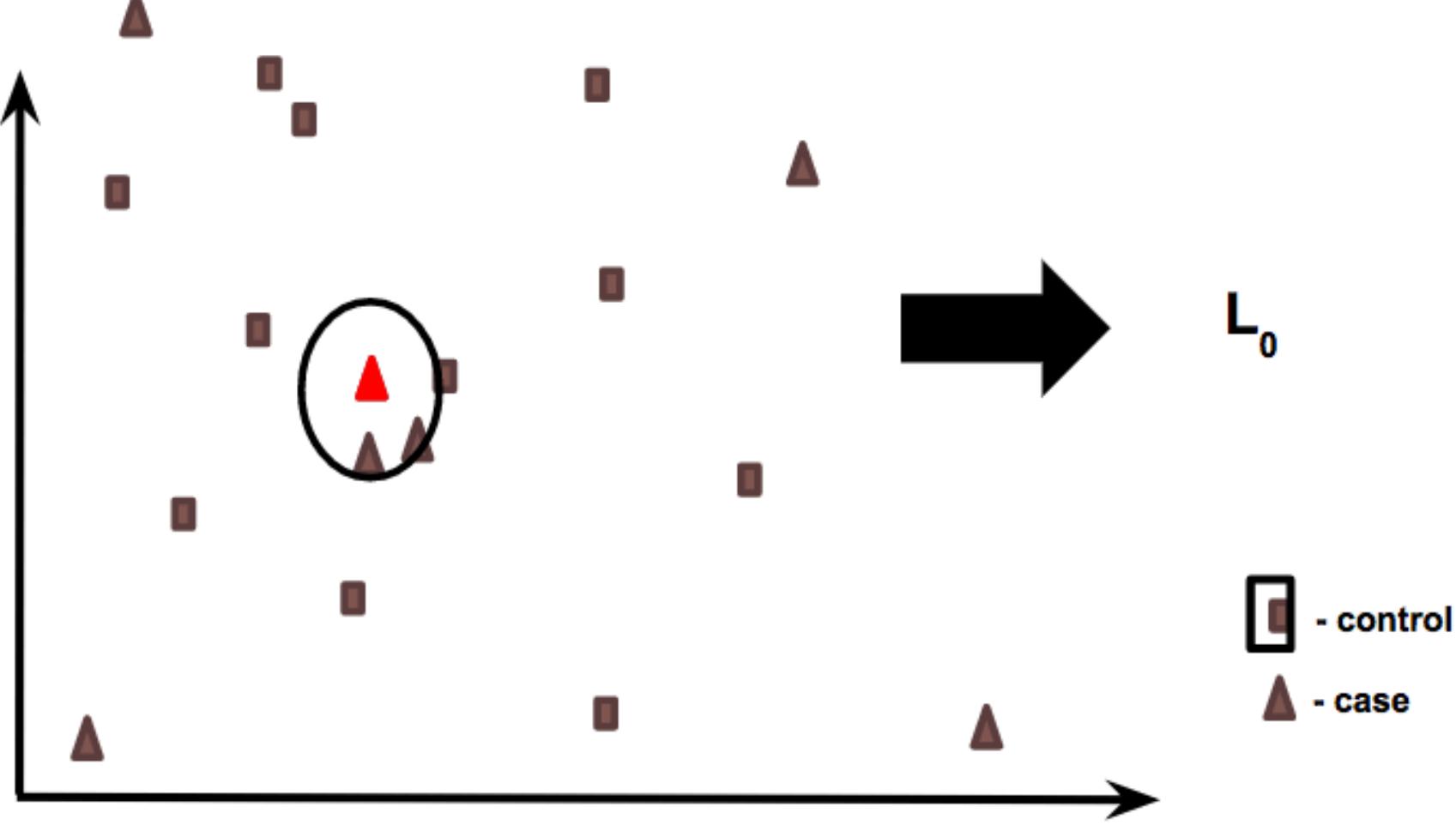


For each circle (window) compute the likelihood function

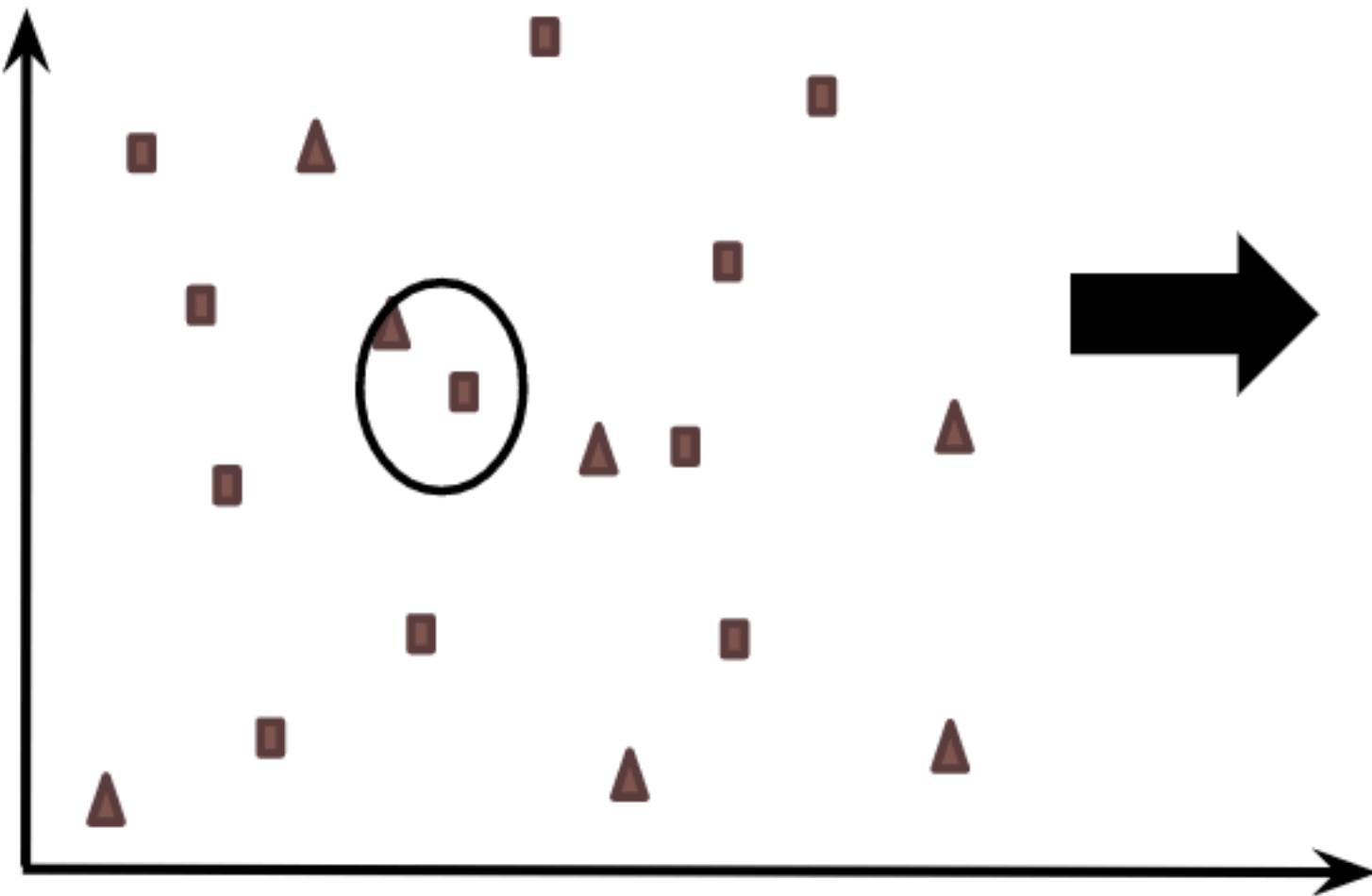
For Bernoulli distributions, the function is:

$$\left(\frac{c}{n}\right)^c \left(\frac{n-c}{n}\right)^{n-c} \left(\frac{C-c}{N-n}\right)^{C-c} \left(\frac{(N-n)-(C-c)}{N-n}\right)^{(N-n)-(C-c)} \text{IO}$$



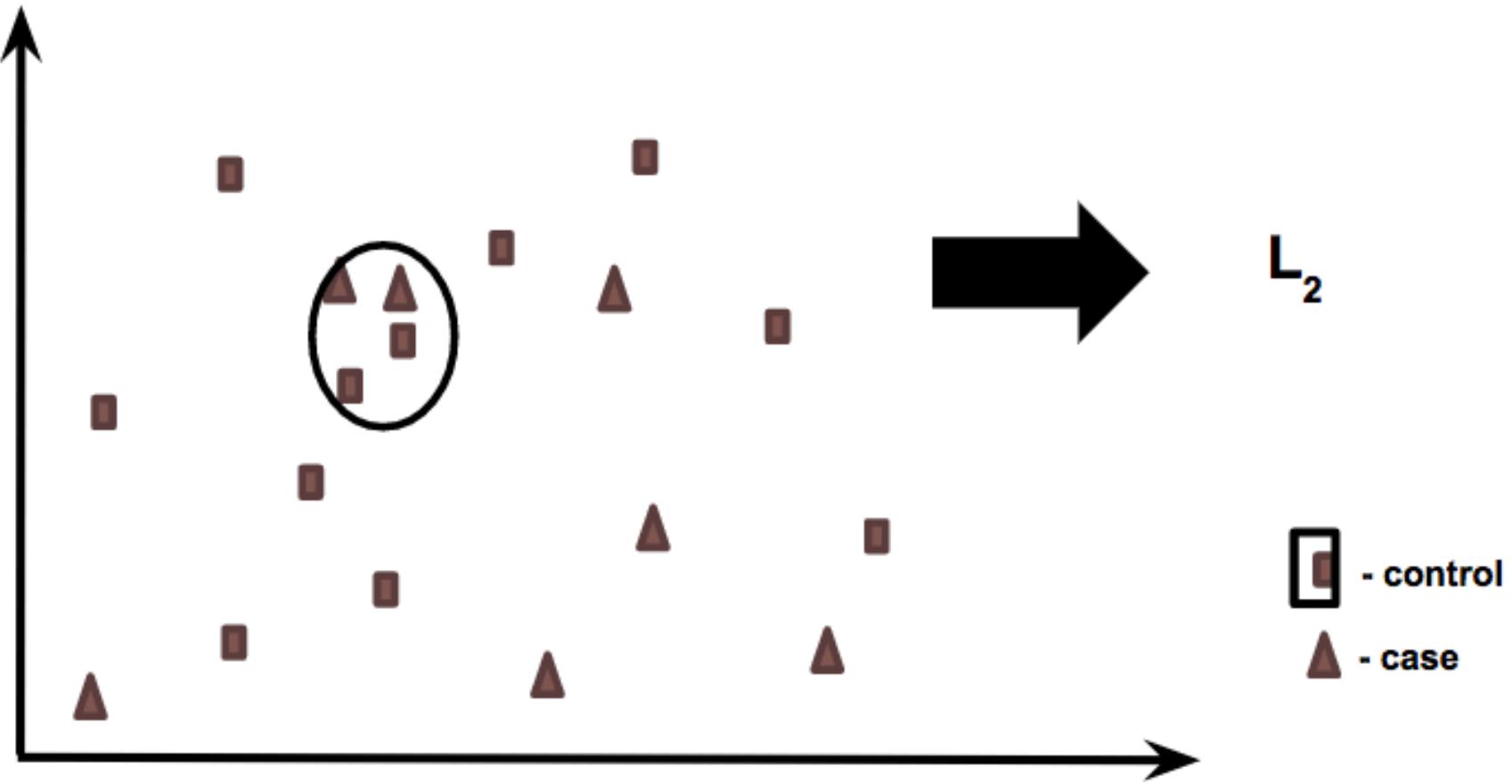


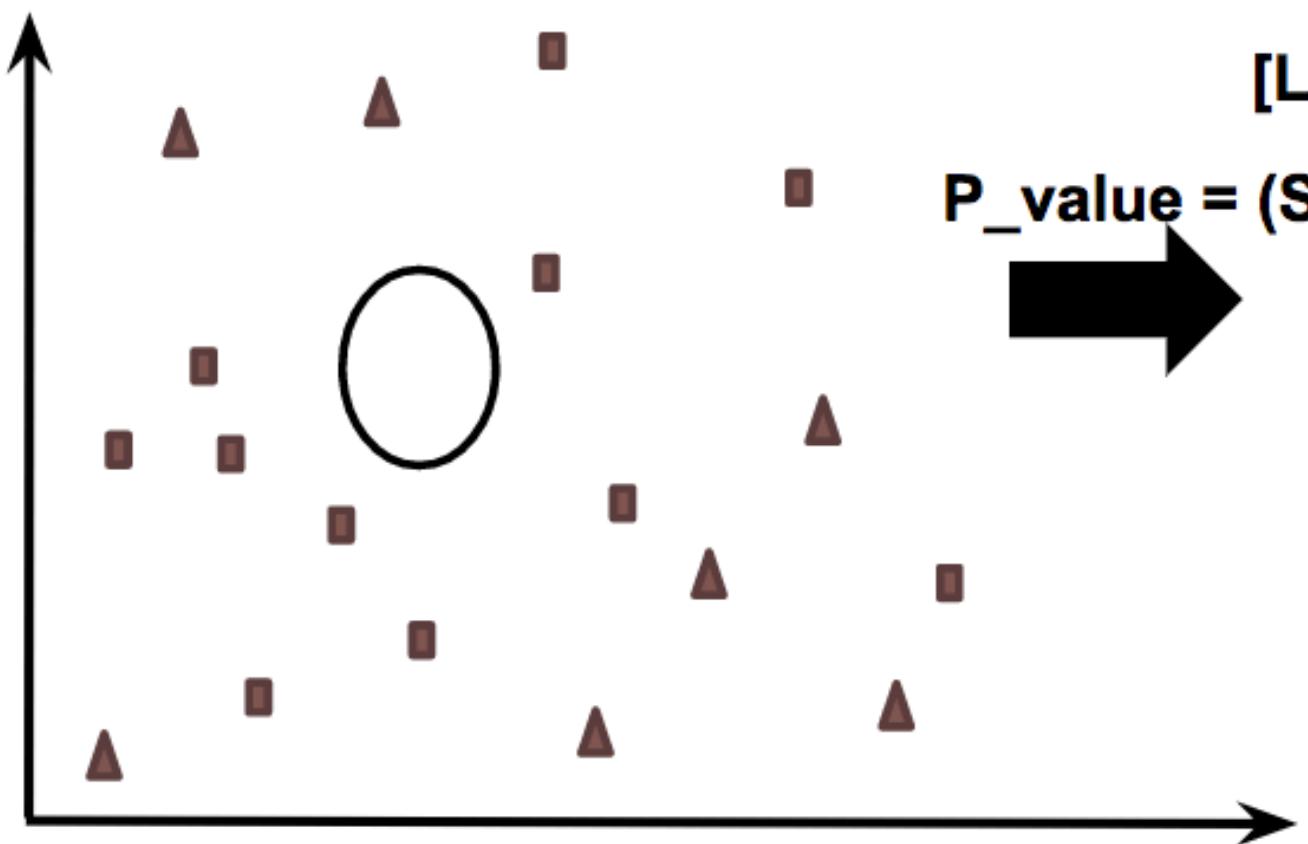
[ $L_0$ ]



- control  
 - case

$[L_0, L_1]$





$[L_0, L_1, L_2, \dots, L_{M-1}]$

$[L_5, L_2, L_{88}, L_0, \dots, L_{M-22}]$

$P\_value = (\text{Sorted position of } L_0 \text{ in list})/M$

$L_M$

 - control  
 - case

---

# Geographical Analysis

## Geovisual Analytics Systems

# Objectives

---

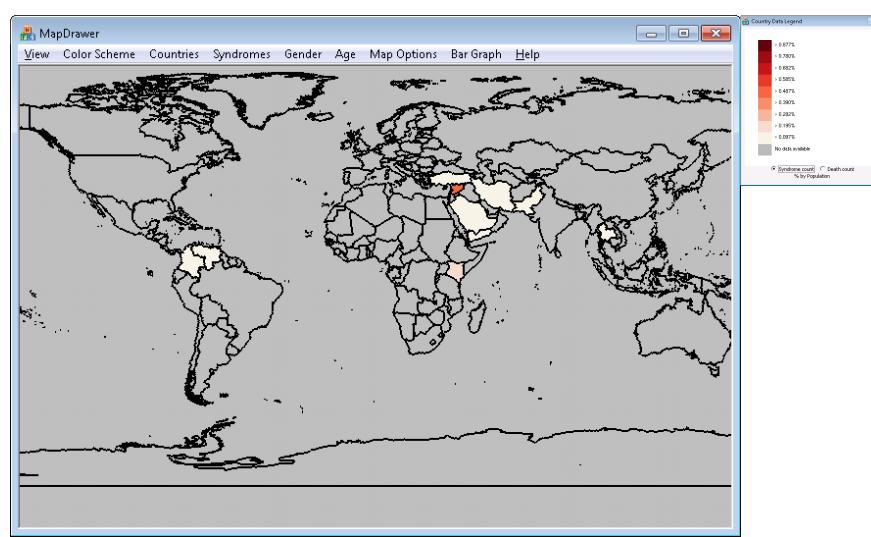


Objective

Apply methods of spatial  
analysis

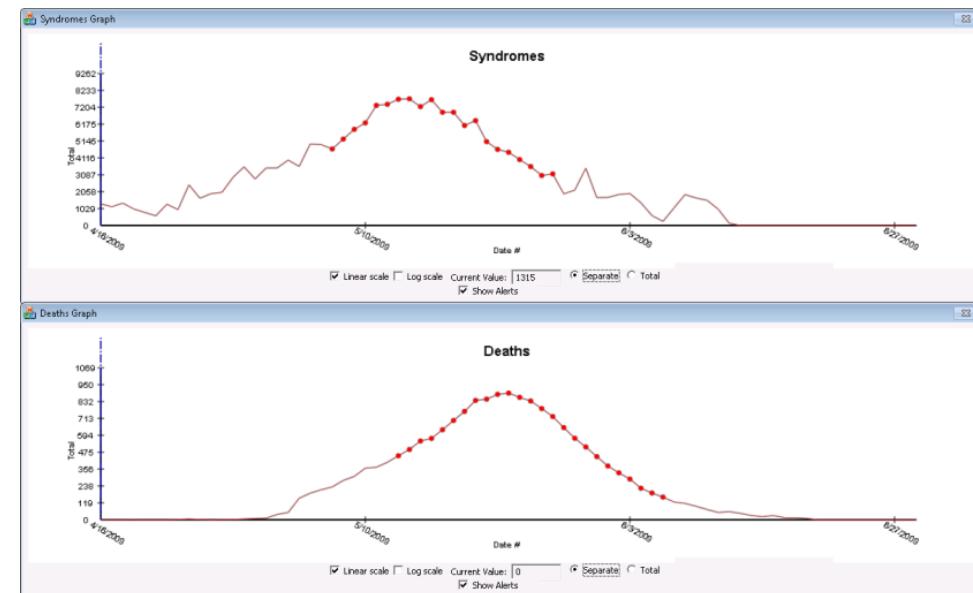
Chief Complaints	
Admittance records of patients:	
abd pain,vomiting, percent = 2.633%	
abdominal pain acute, percent = 2.558%	
abd pain, percent = 2.525%	
abdominal pain/abdominal pain, percent = 2.475%	
vomiting alone, percent = 2.465%	
nose bleeds, percent = 19.140%	
nose bleed, percent = 9.376%	
vomiting/blood/vomiting blood, percent = 9.235%	
nose/bleed/nose bleed, percent = 9.200%	
vomiting blood, percent = 9.094%	
vomiting/diarrhea, percent = 3.351%	
vomiting, diarrhea, percent = 2.703%	
abd pain,vomiting, percent = 2.541%	
abd pain, percent = 2.432%	
abd cramping, percent = 2.270%	
Admittance records of patients that died:	
abd pain,vomiting, percent = 3.821%	
abdominal pain, percent = 3.252%	
vomiting, abd pain, percent = 3.252%	
vomiting & diarrhea, percent = 3.171%	
abd pain fever, percent = 3.099%	
nose bleed, percent = 15.522%	
nose bleed, percent = 15.510%	
vomiting blood, percent = 15.510%	
vomiting/blood/vomiting blood, percent = 15.102%	
vomiting blood, percent = 13.469%	
abd pain fever, percent = 4.678%	
abd pain,vomiting, percent = 4.678%	
abdominal pain, percent = 4.678%	
abdominal pain, percent = 4.094%	
abd pain/fever, percent = 4.094%	
vomiting blood, percent = 27.586%	
vomiting/blood/vomiting blood, percent = 20.690%	

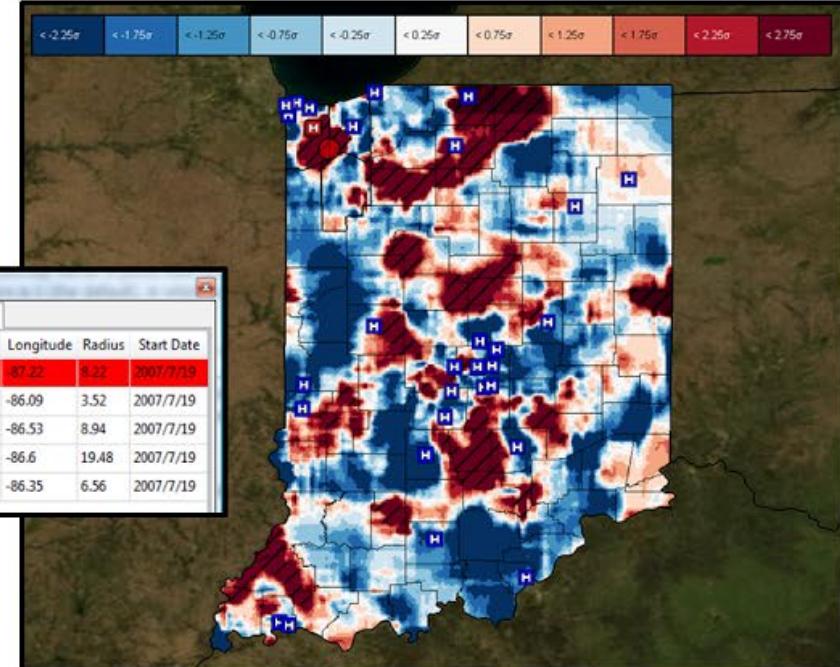
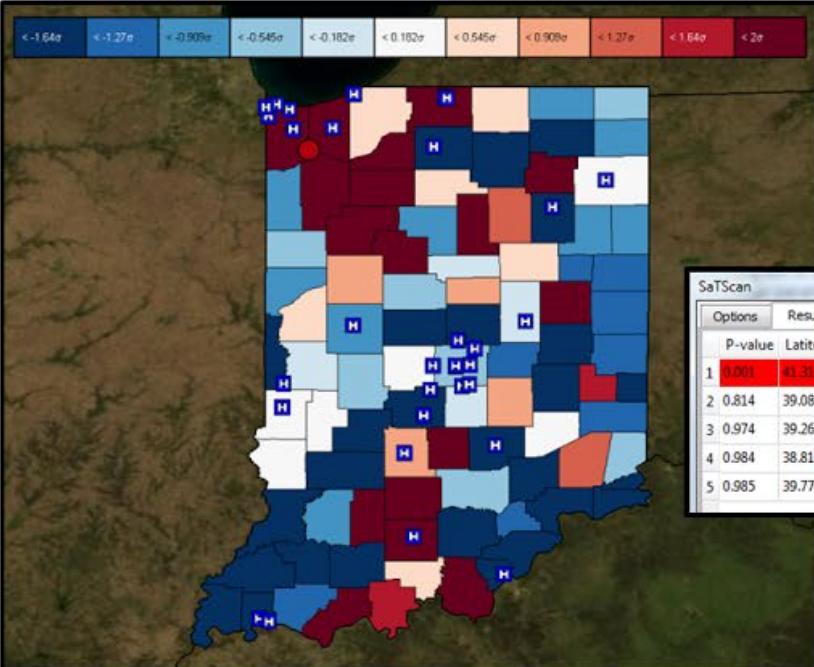
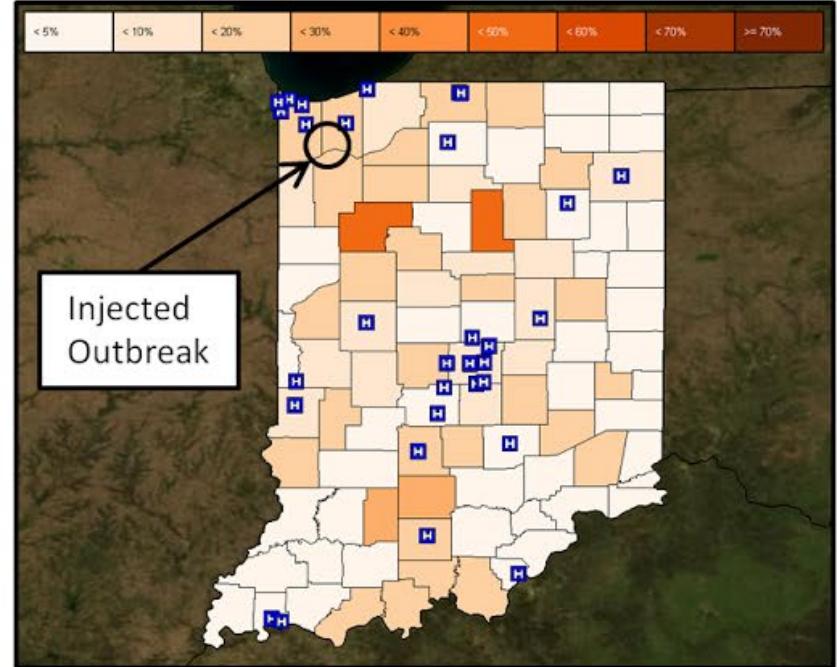
Summary statistics	
Summary for day 5/10/2009:	
Total Patients =	14839
Male percentage =	50.037%
Female percentage =	49.963%
Age range infected:	
< 5 =	0.013%
5 - 19 =	1.361%
20 - 39 =	37.193%
40 - 59 =	51.971%
> 60 =	9.462%
Summary for deaths for day 5/10/2009:	
Total deaths =	1098
Male percentage =	49.180%
Female percentage =	50.820%
Botulic percentage =	2.277%
Constitutional percentage =	8.470%
Gastrointestinal percentage =	39.435%
Hemorrhagic percentage =	7.286%
Neurological percentage =	2.914%
Rash percentage =	1.457%
Respiratory percentage =	1.393%
Other percentage =	36.521%



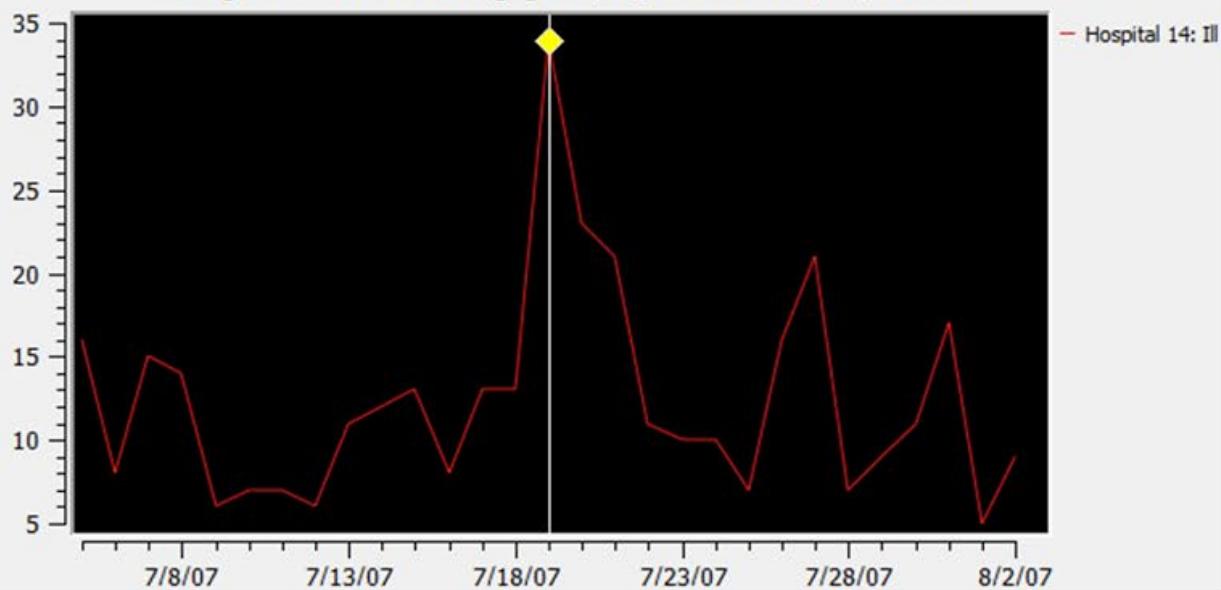
Patient ID	Date	Chief Complaints	Location
9398	4/16/09	Ear pain	Karachi
10816	4/16/09	Stuffy nose	Lebanon
1491	4/16/09	Fever	Allepo
16237	4/16/09	Head bleed	Yemen

CoCo  
Classifier

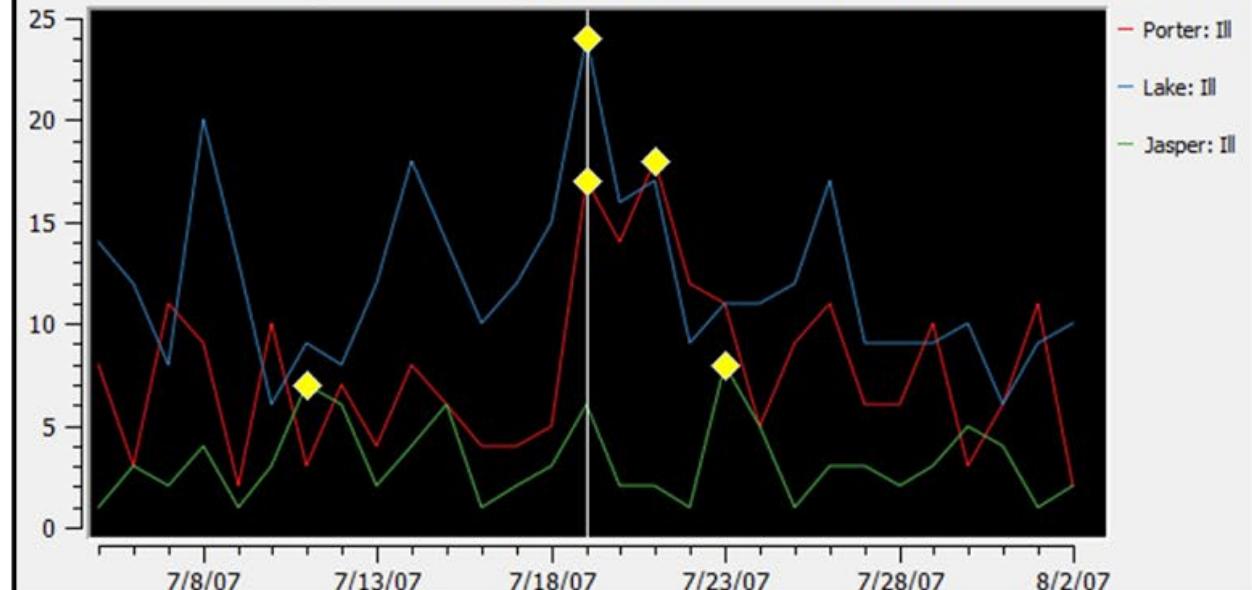




**Hospital Selection(s): 7/5/2007 - 8/2/2007**



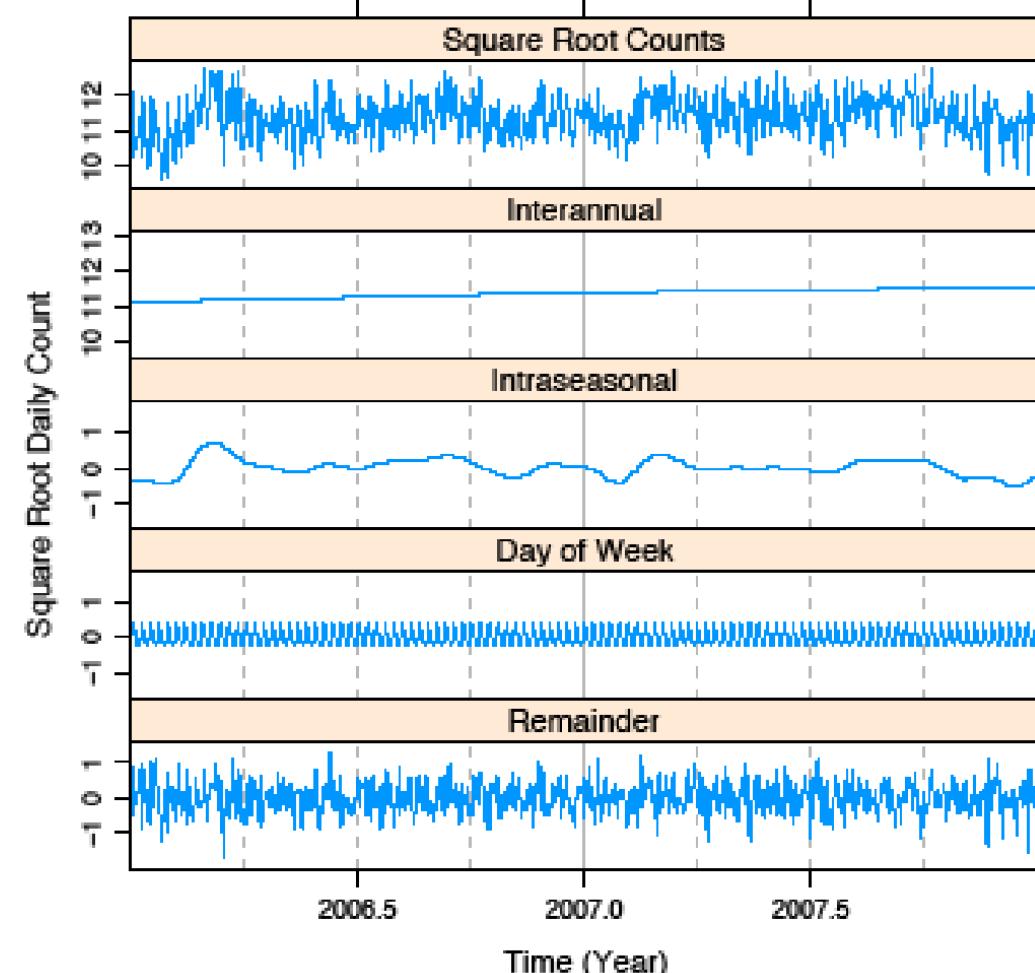
**County Selection(s): 7/5/2007 - 8/2/2007**



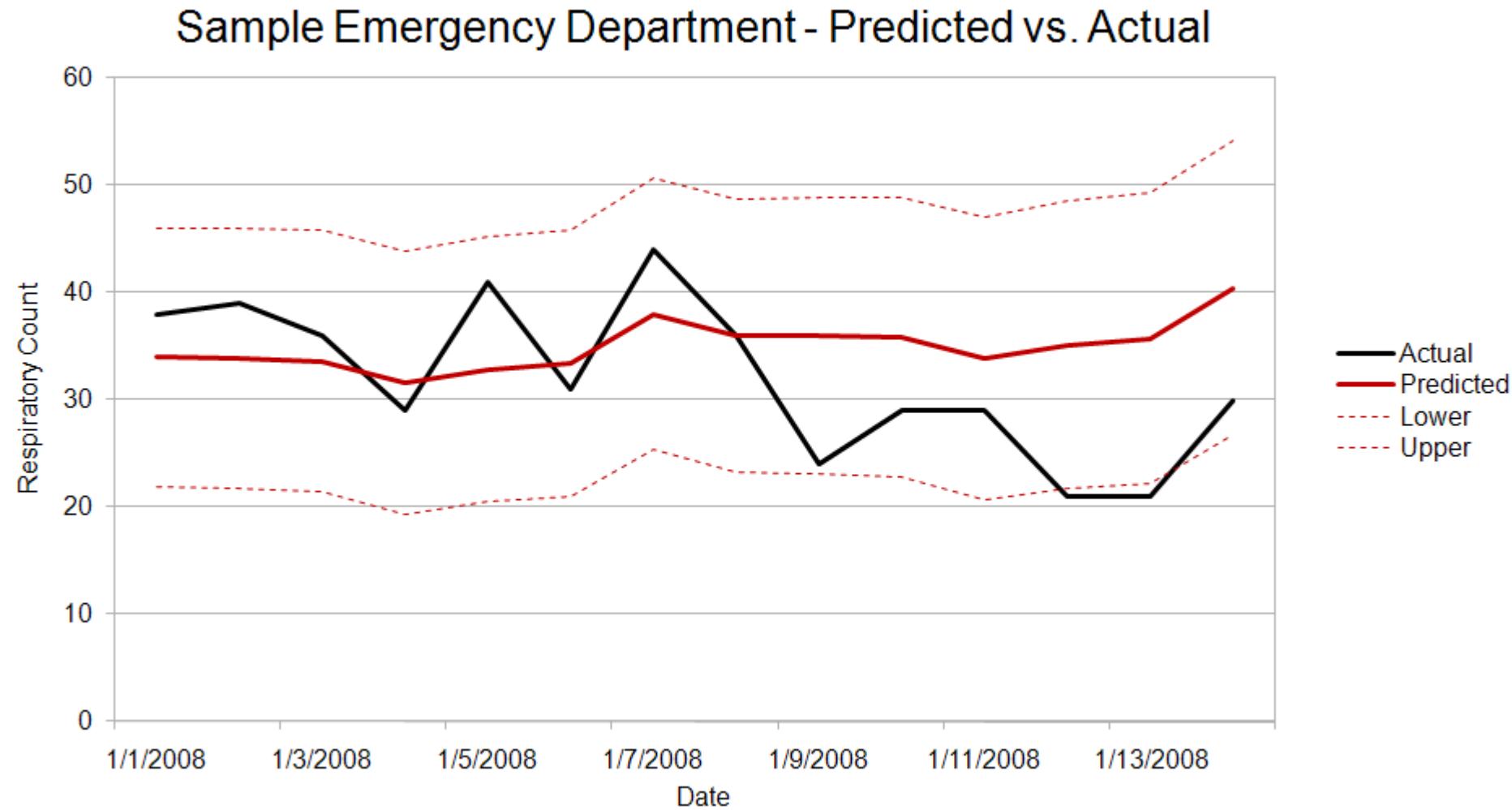
# Time Series Modeling

## Seasonal-Trend Decomposition Based on Loess

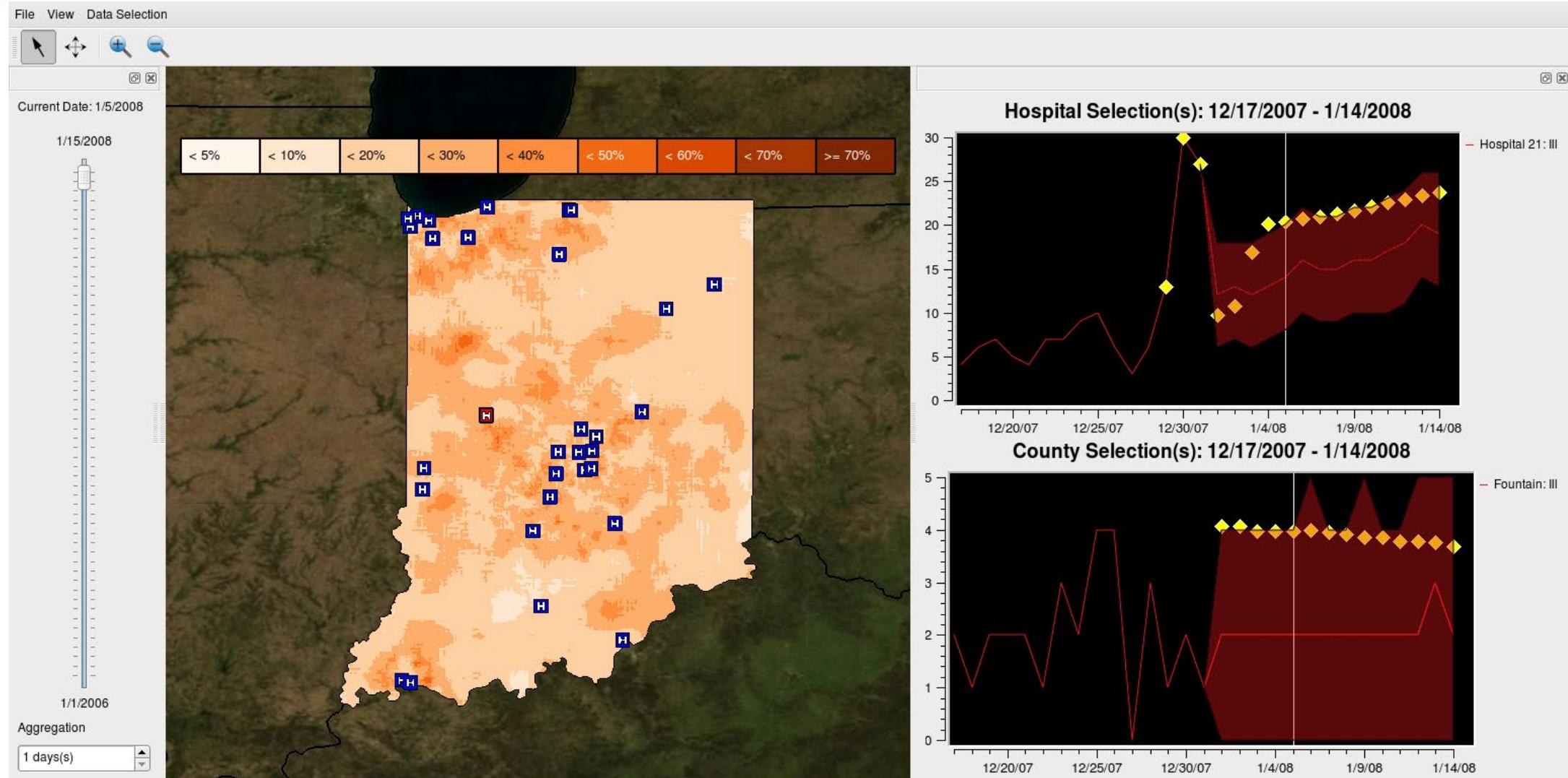
- Time series can be viewed as the sum of multiple trend components
- For each data signal, components are extracted
- Can then analyze correlation between components

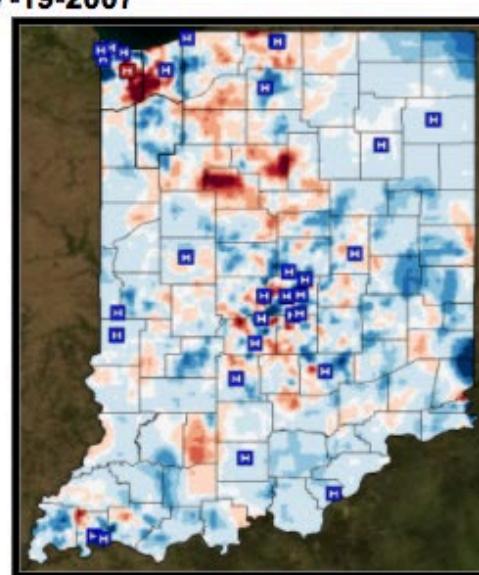
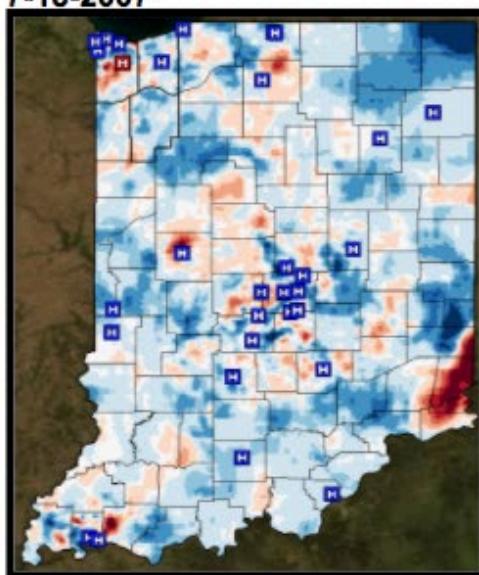
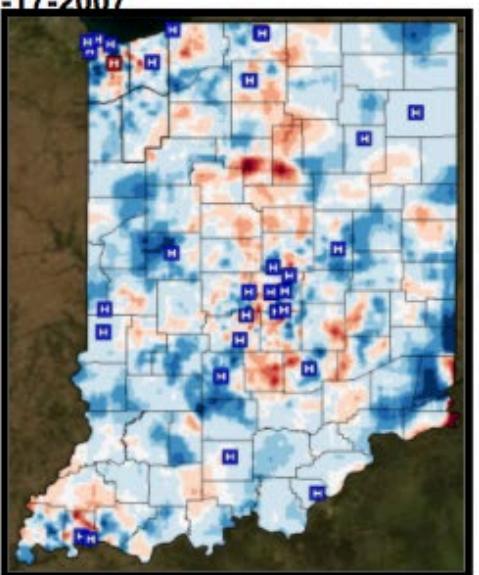
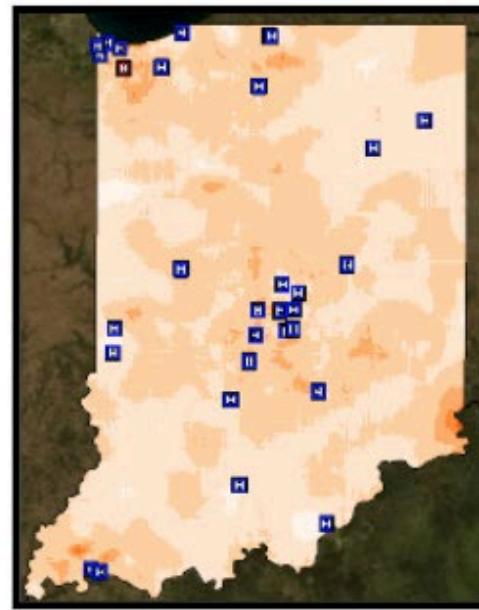
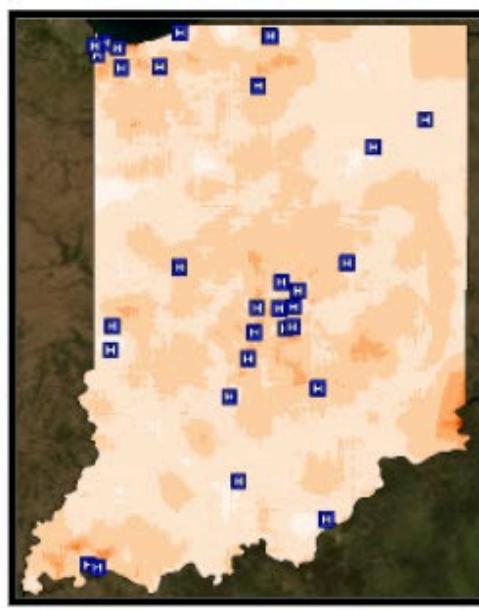
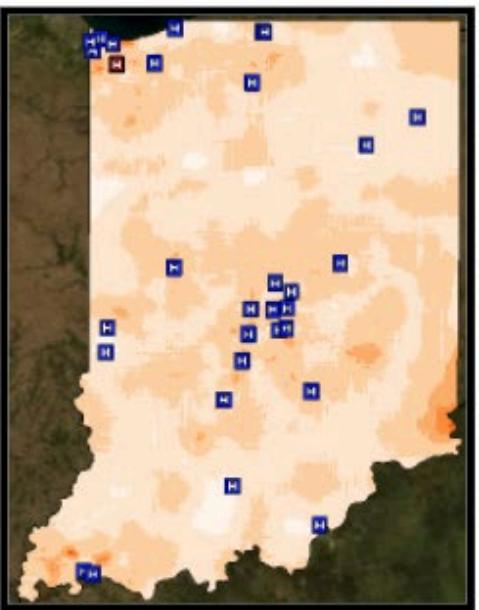


# Predictive Visual Analytics



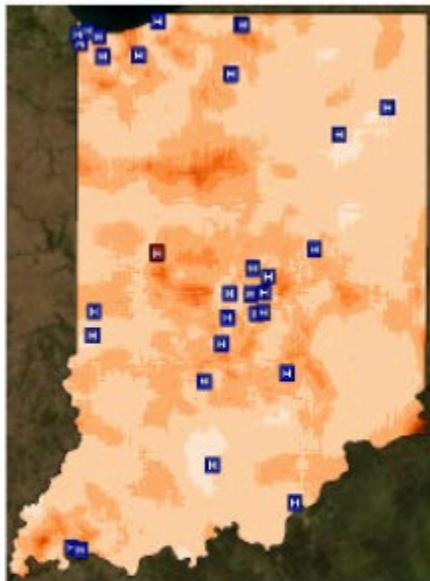
# Predictive Visual Analytics



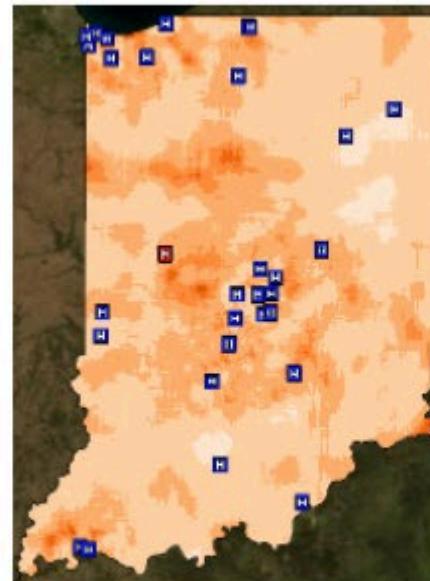


# Spatiotemporal

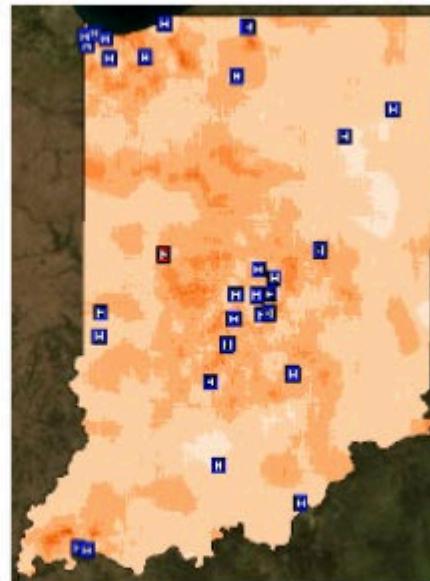
| The probability distribution of tomorrow can be based on past probability distributions



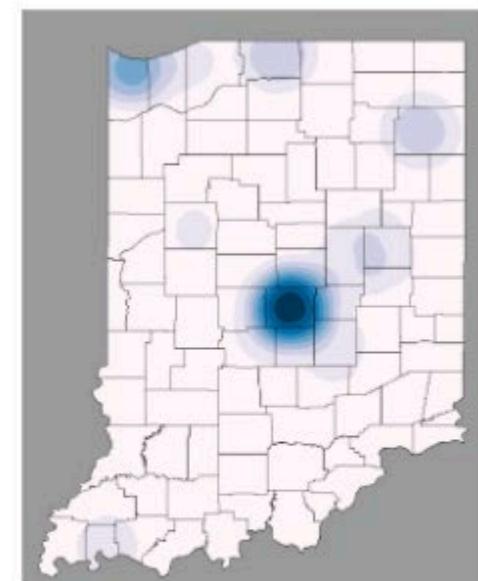
Day: t-3



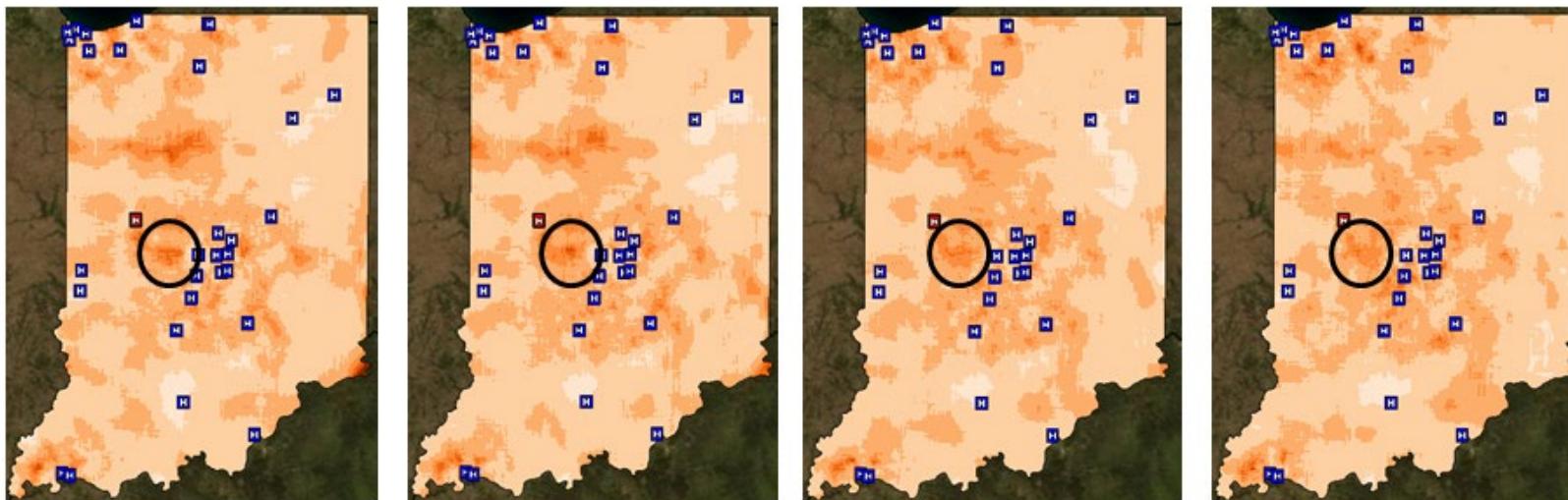
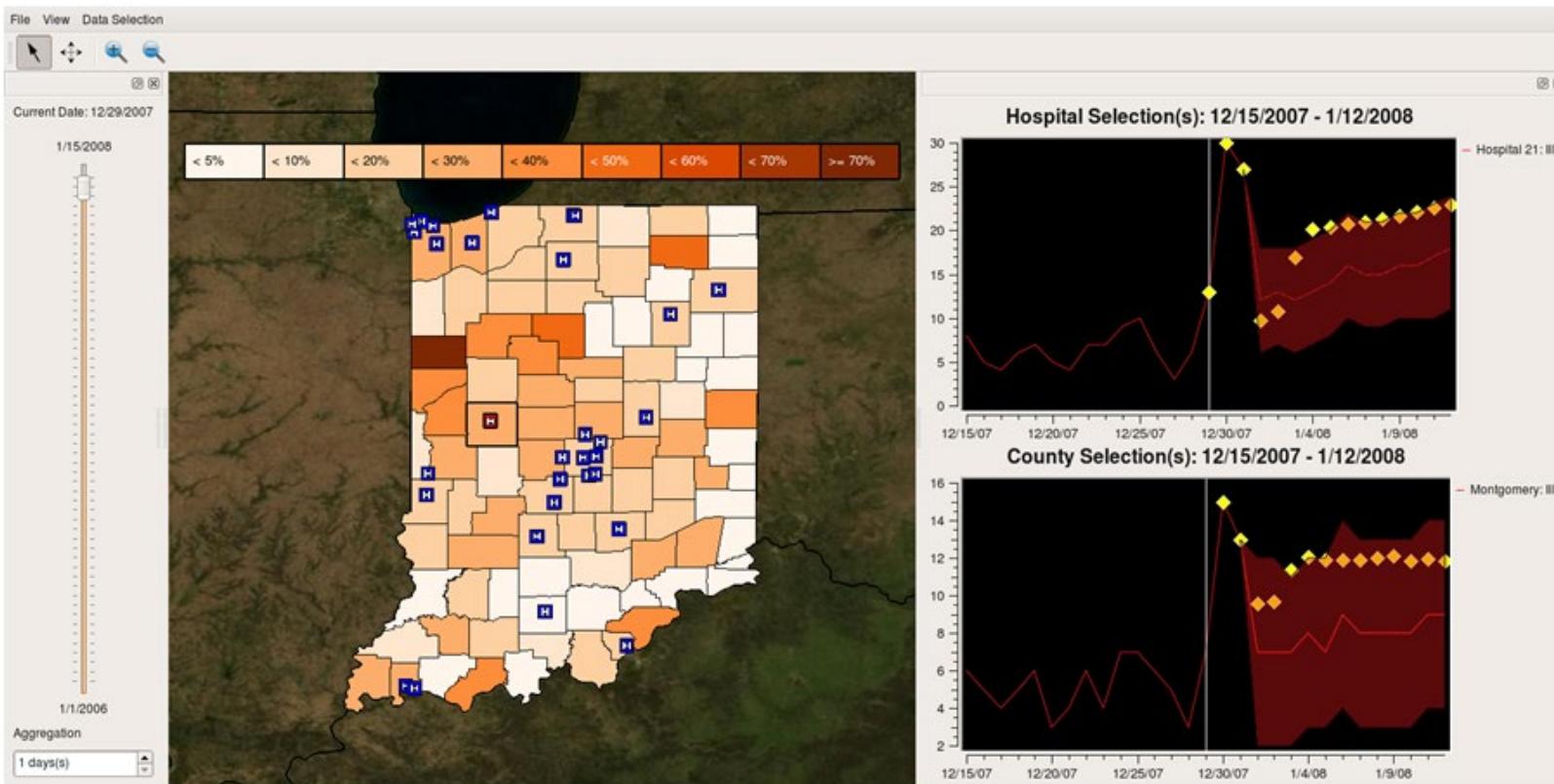
Day: t-2

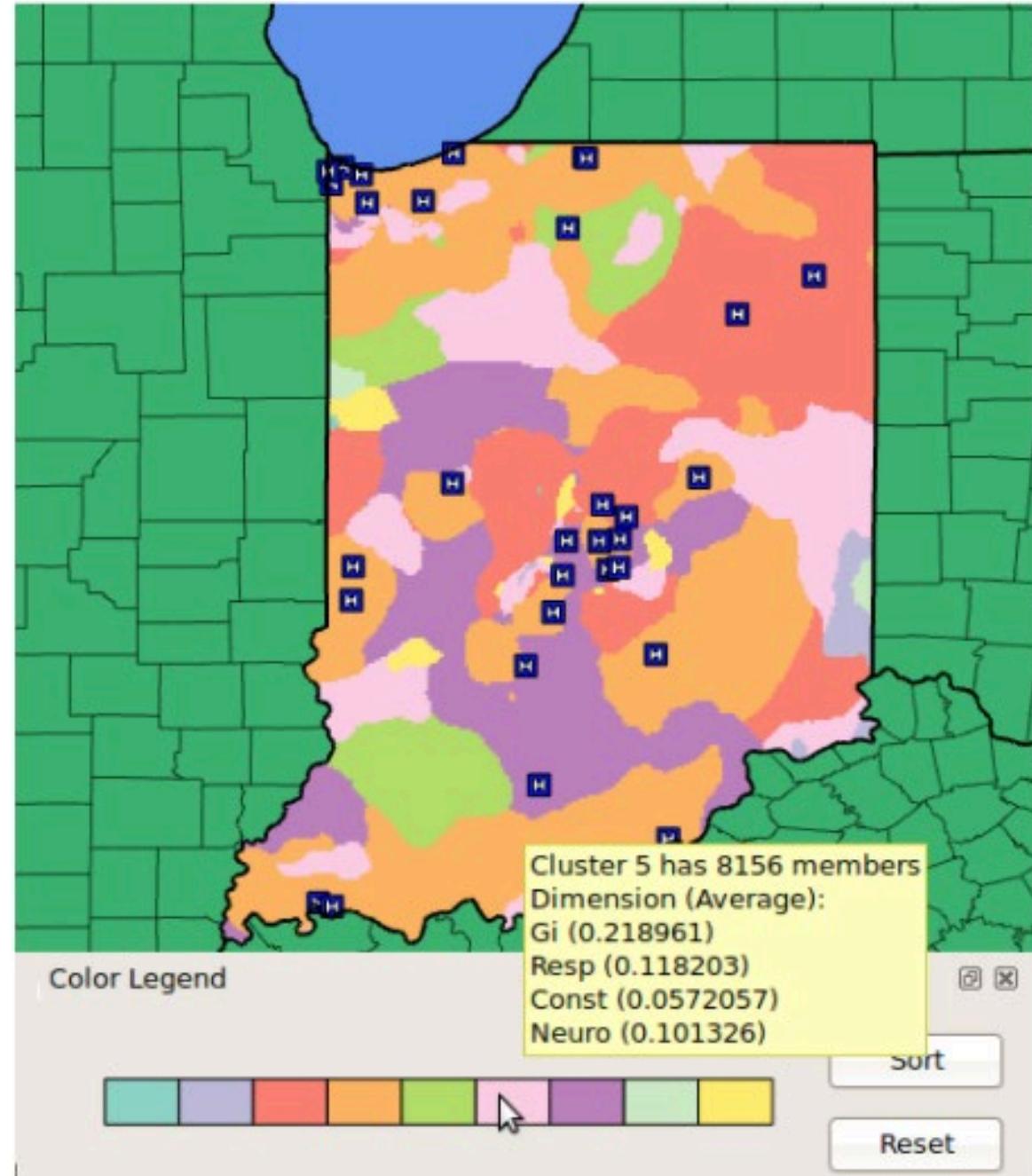
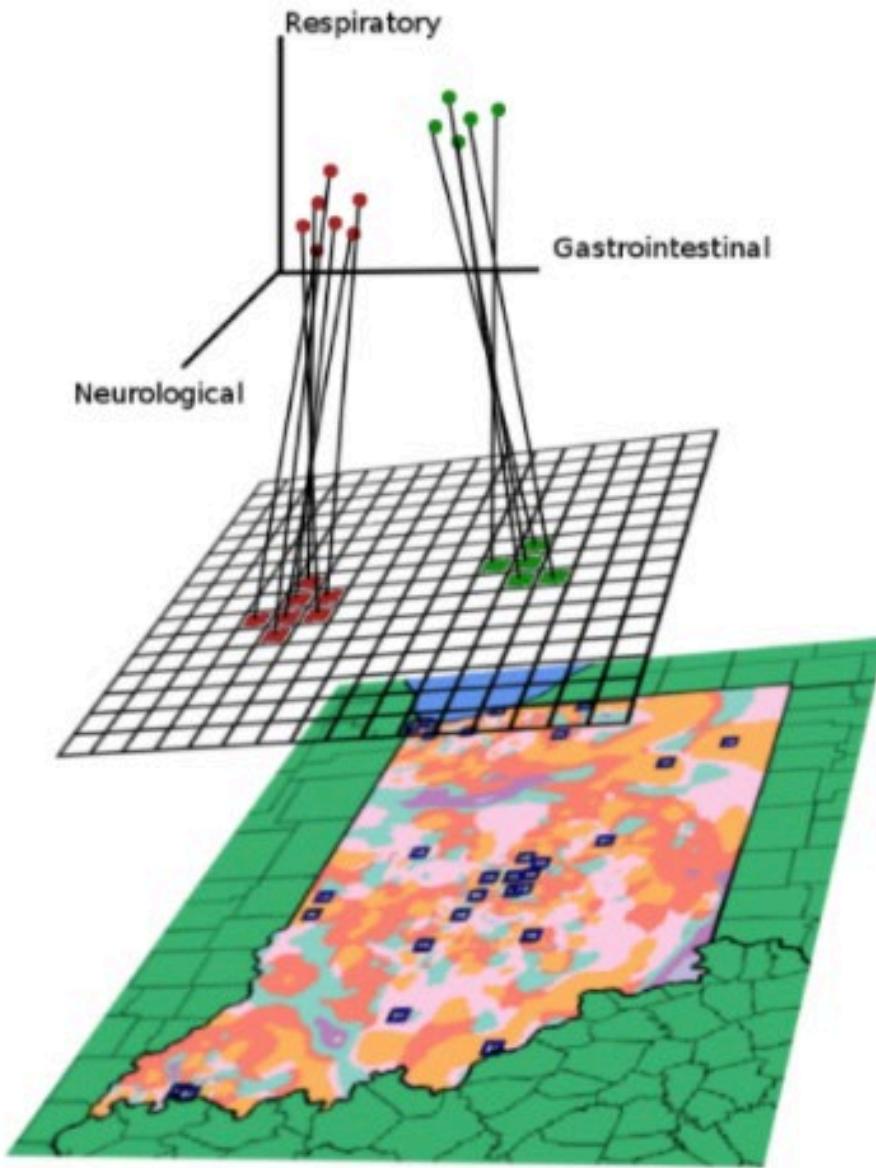


Day: t-1



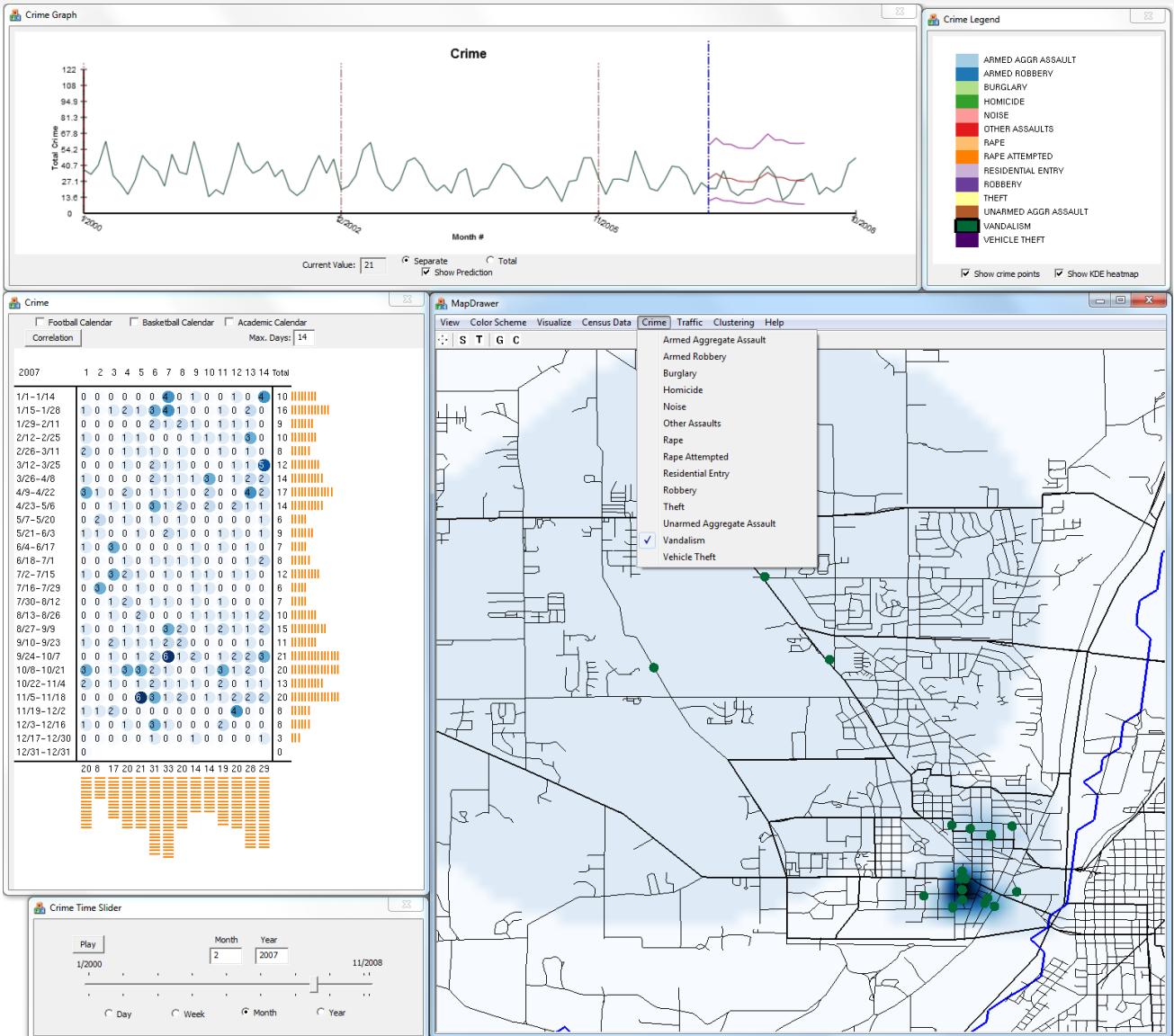
Expected Distribution  
using all past records

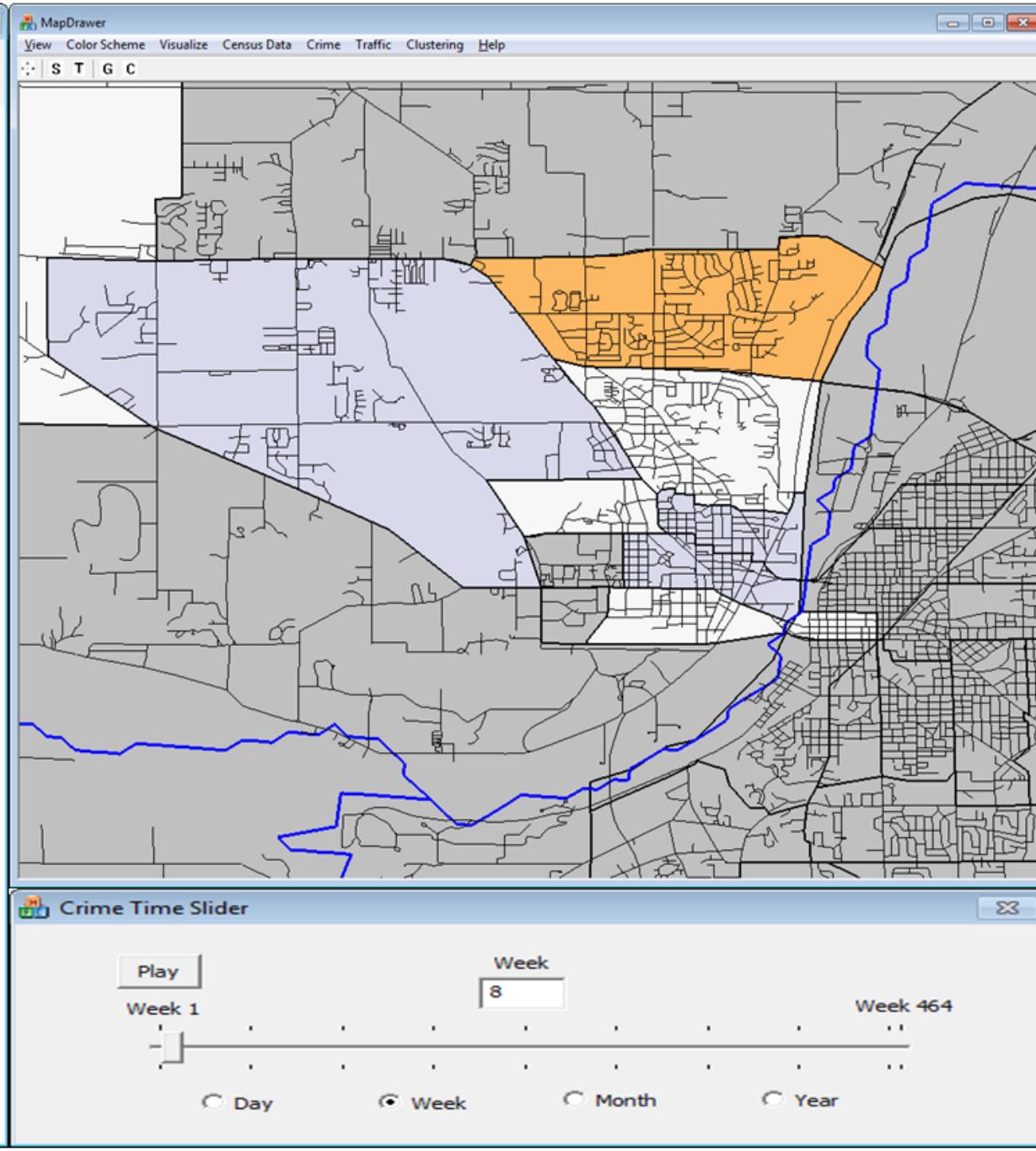
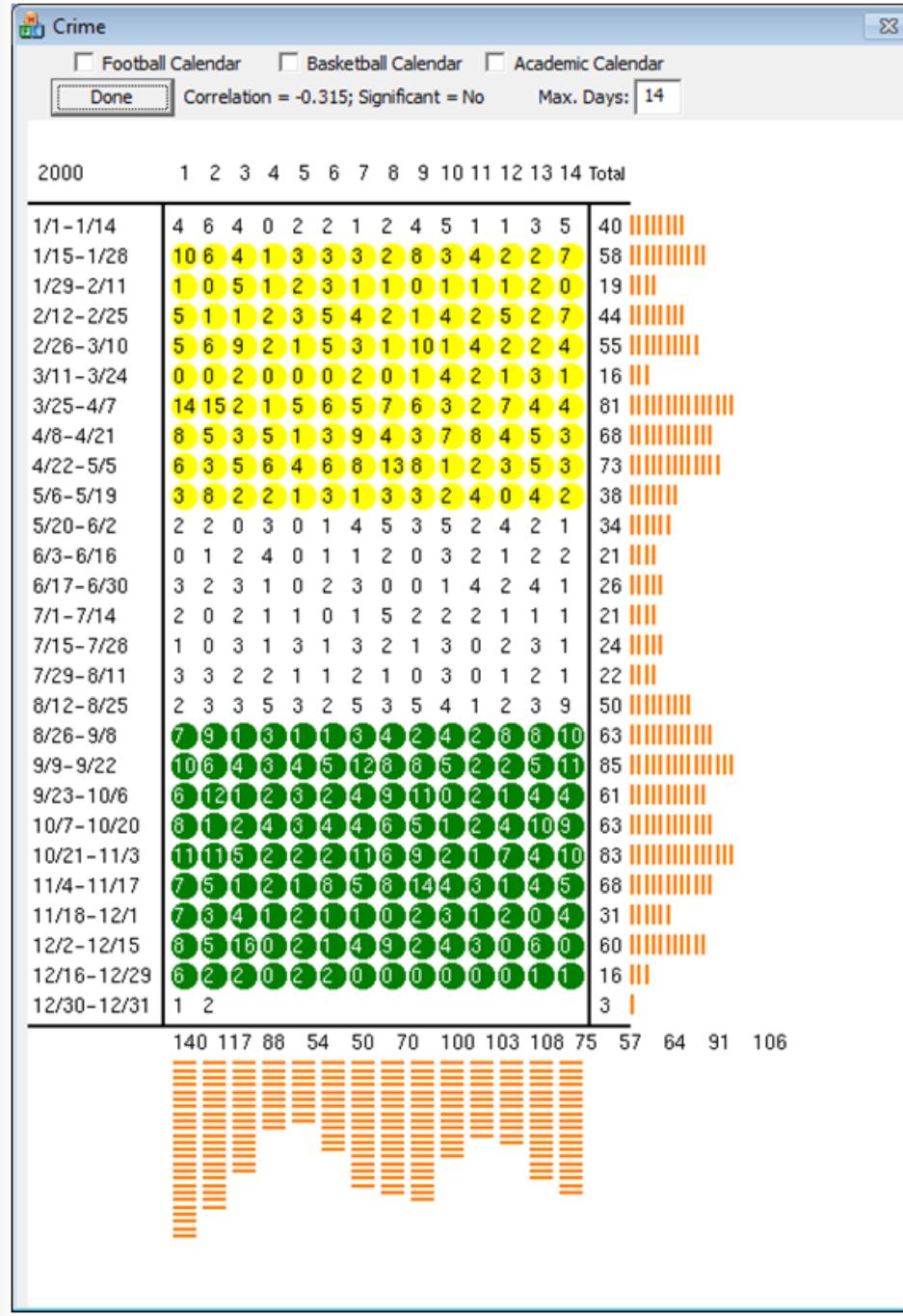


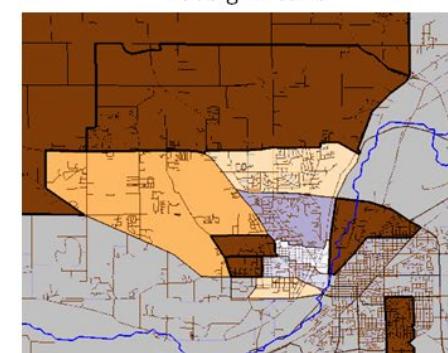
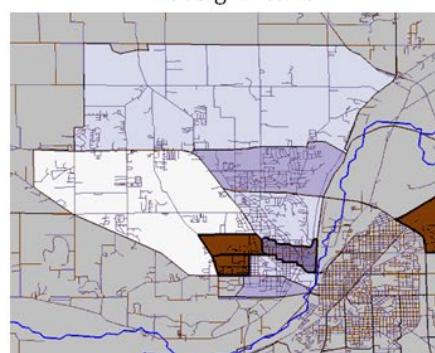
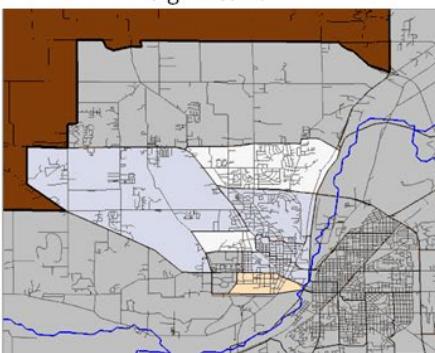
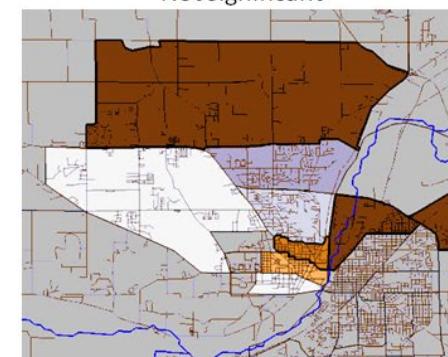
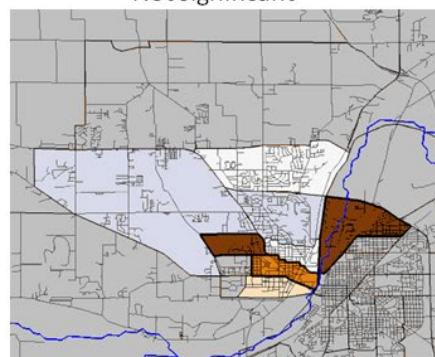
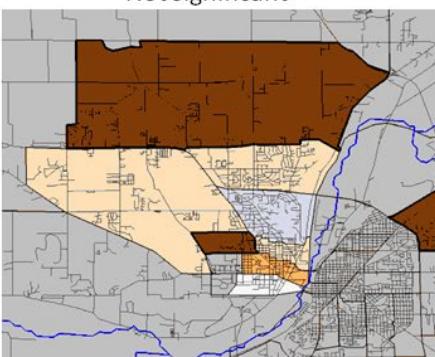
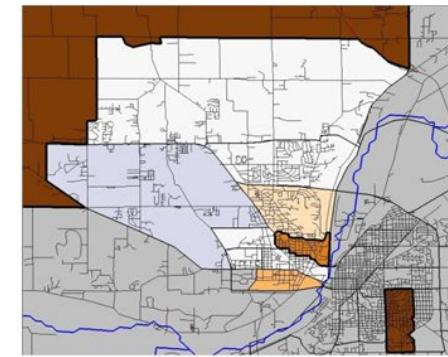
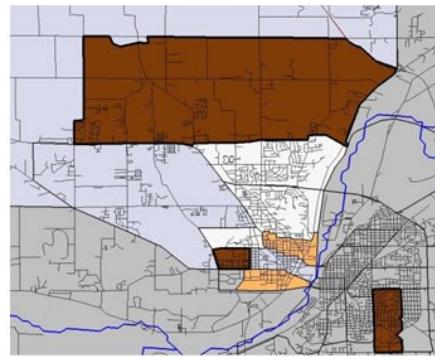
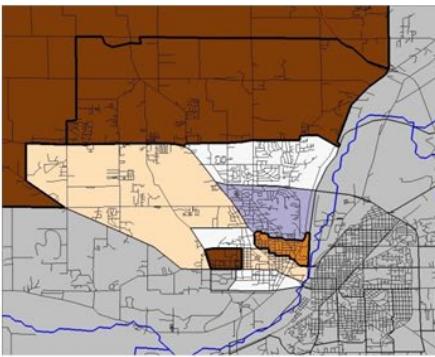


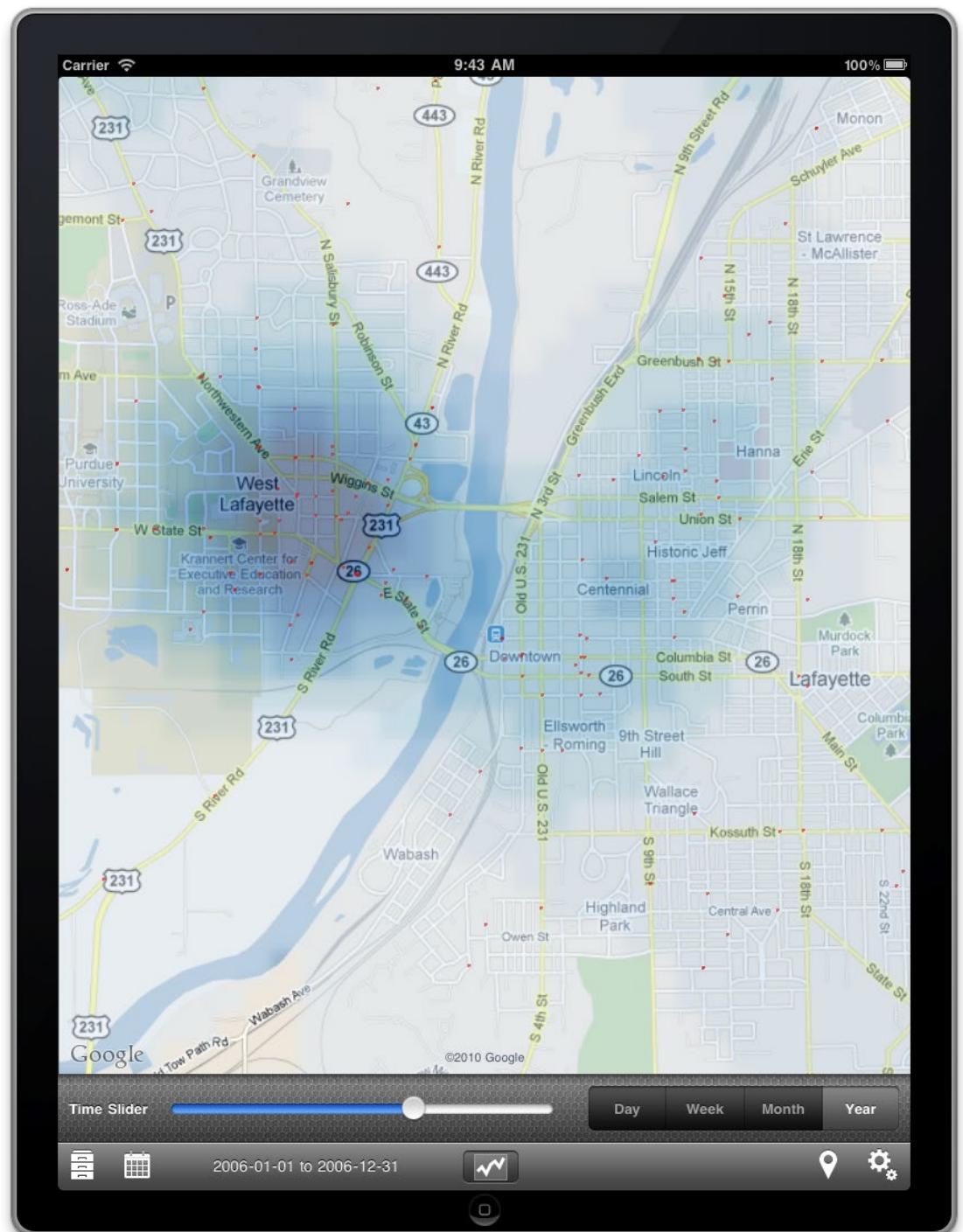
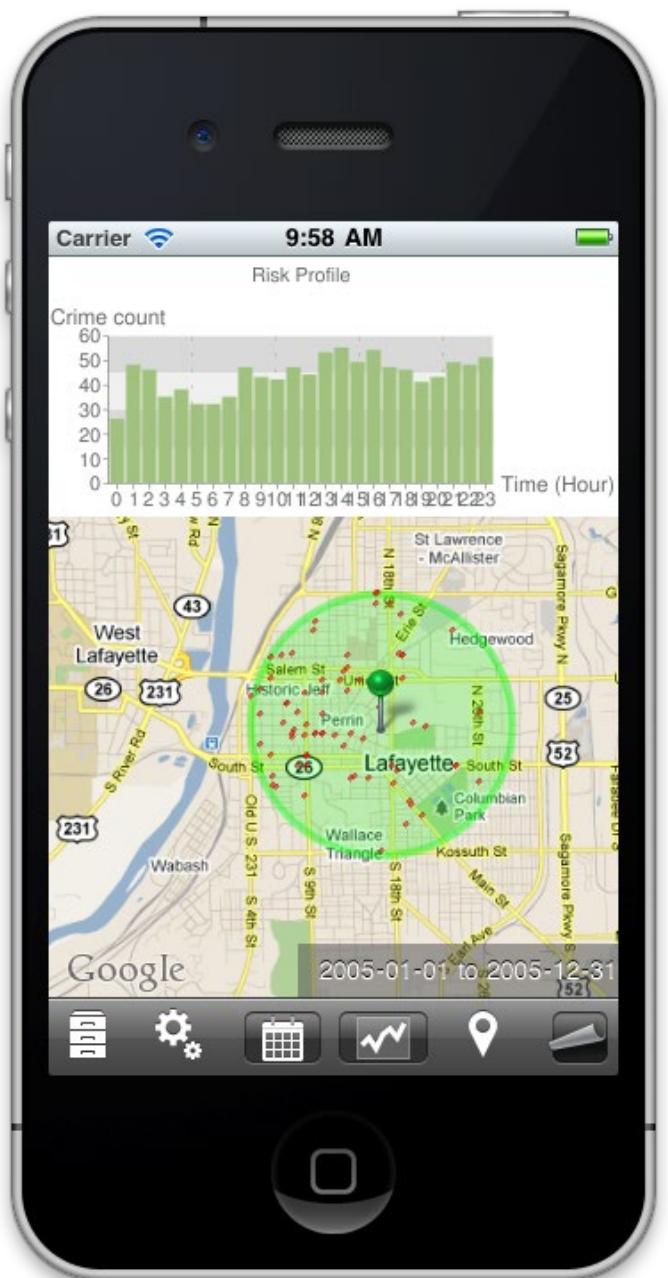
# Crime Mapping

- | The problem is analogous to syndromic surveillance
- | Instead of patient addresses, now there is criminal incident reports
- | More data (who, what, when, where, how)

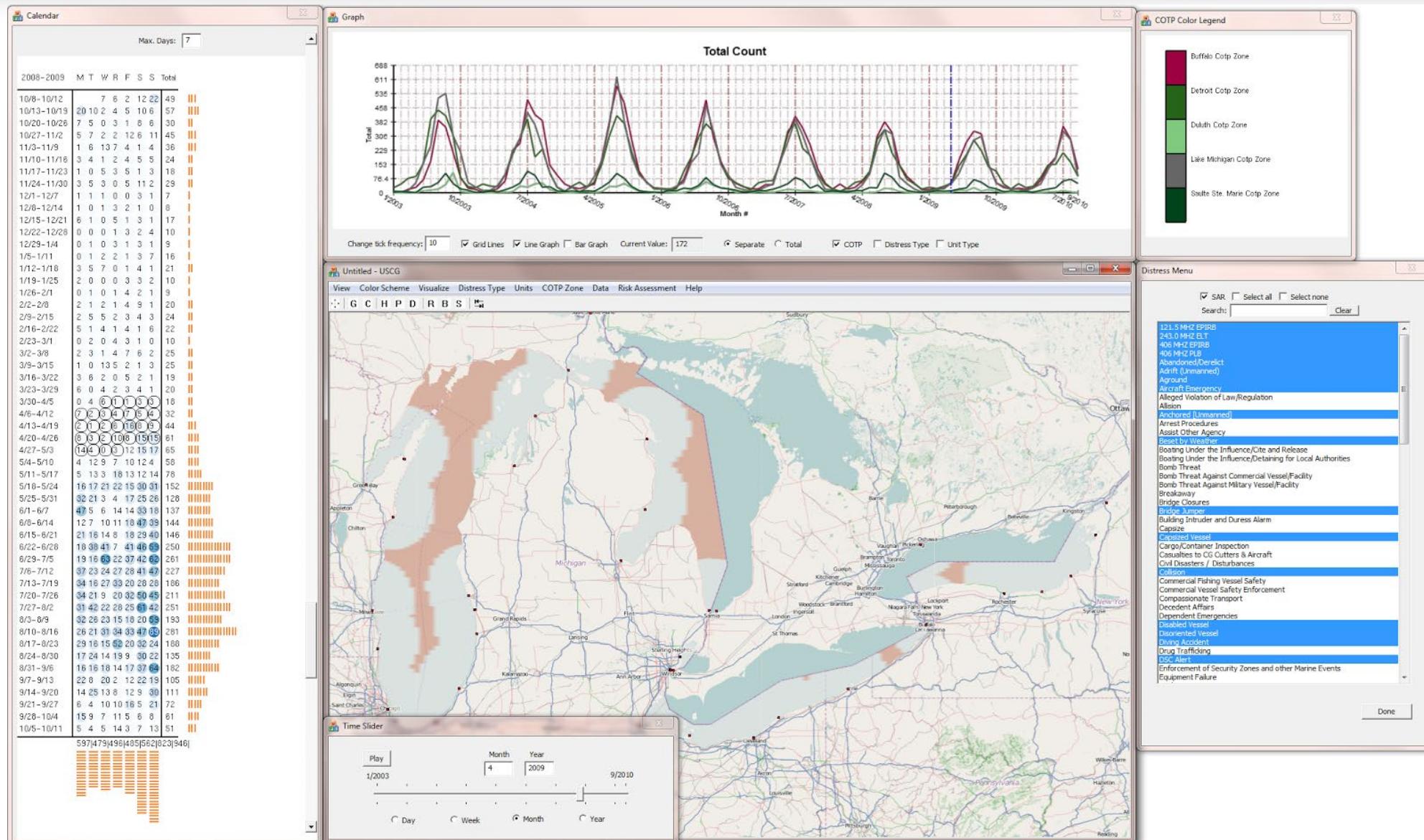








# Coast Guard Search and Rescue Visual Analytics (CGSARVA)



# Assessing Risk in the Great Lakes

