

Geospatial analysis

Jason G. Fleischer, Ph.D.

Asst. Teaching Professor

Department of Cognitive Science, UC San Diego

jfleischer@ucsd.edu

@jasongfleischer



<https://jgfleischer.com>

How a coastline 100 million years ago influences modern election results in Alabama

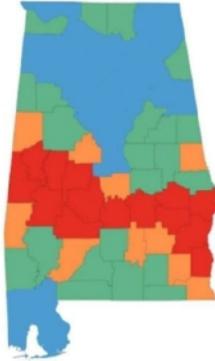
Cretaceous Sediments



Fertile Blackland Prairie Soil



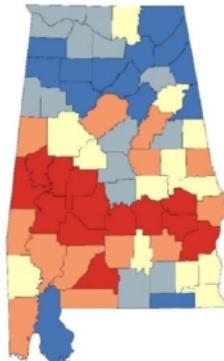
Average Farm Size, 1997



Slave Population, 1860



Black population, 2010



Election Results, 2020



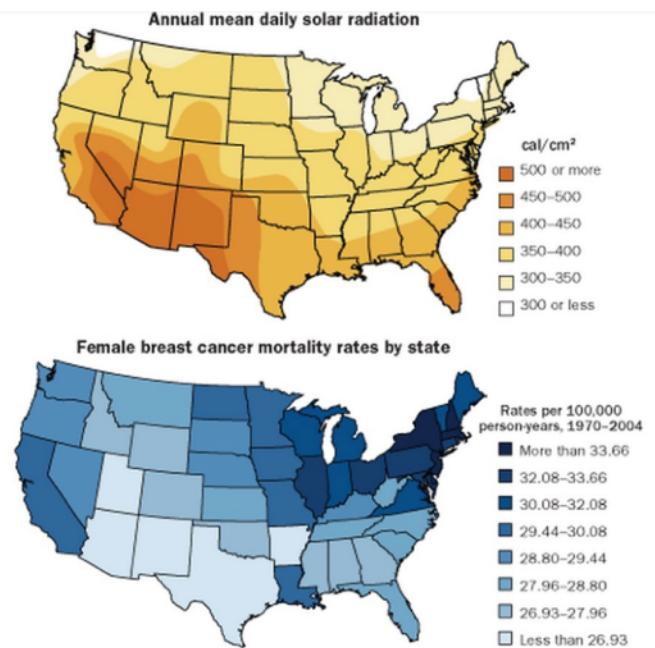
-Starkey Comics

Why Geospatial Analysis?

'Everything is related to everything else, but near things are more related than distant things.' -Tobler 1979

“...the purpose of geographic inquiry is to examine relationships between geographic features collectively and to use the relationships to describe the real-world phenomena that map features represent”
-Clarke 2001

Clearly visualizes
important differences in
disease distribution



ON THE MAP Scientists who study vitamin D can't help but notice that a host of diseases seem to vary with latitude. Type 1 diabetes, multiple sclerosis and even some cancers appear to be more common in areas that get less sun -- meaning less opportunity for the body to produce vitamin D. The maps above illustrate the apparent link between solar radiation and breast cancer mortality rates.

SOURCE, FROM TOP: D. M. HARRIS AND V.L.W. GO // J. OF NUTRITION 2004; NATIONAL CANCER INSTITUTE

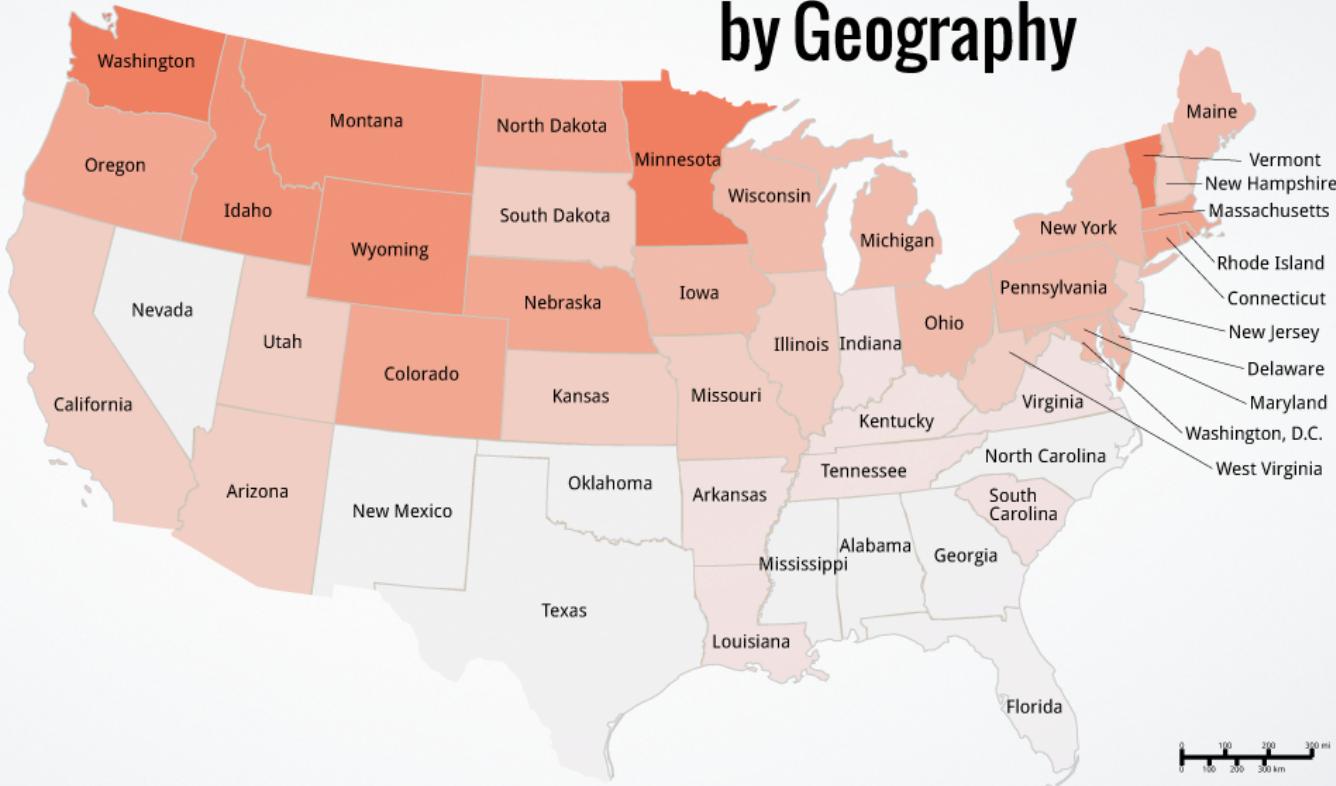
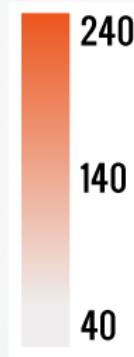
adapted from Brad Voytek

Visualizing Geospatial Data

Multiple Sclerosis by Geography

CASE-CONTROL RATIO OF MS

A higher ratio indicates greater prevalence

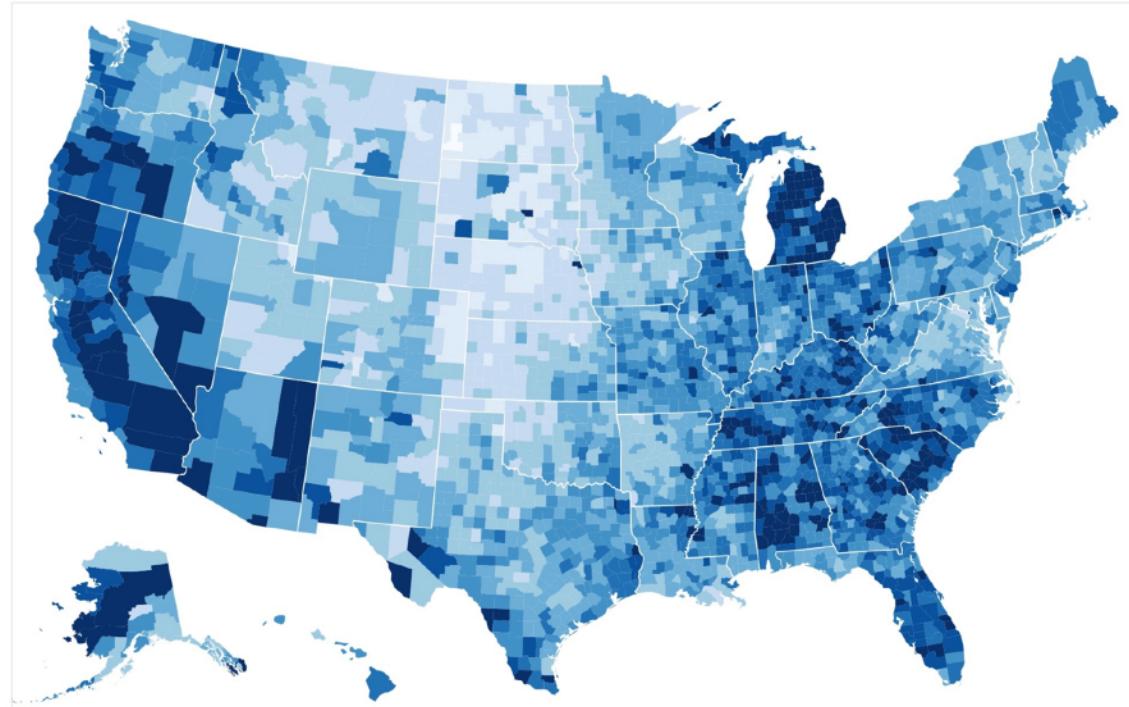


LEARN MORE AT
WWW.INVW.ORG/MS

PRODUCED BY: JASON ALCORN/INVESTIGATEWEST SOURCE: BERETICH AND BERETICH (2009)

adapted from Brad Voytek

Unemployment
rate by county
(August 2016)

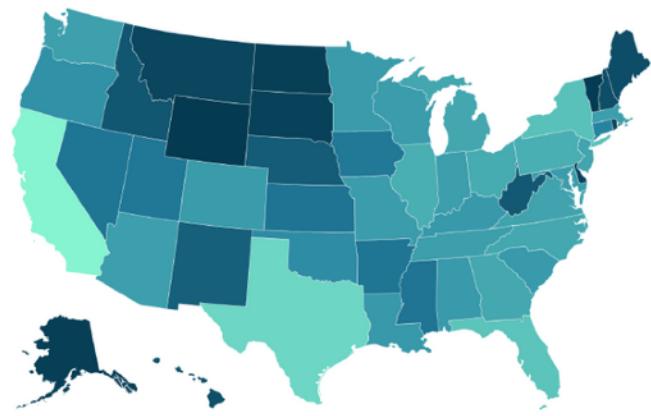


This choropleth encodes unemployment rates from 2008 with a [quantize scale](#) ranging from 0 to 15%. A [threshold scale](#) is a useful alternative for coloring arbitrary ranges.

[Open in a new window.](#)

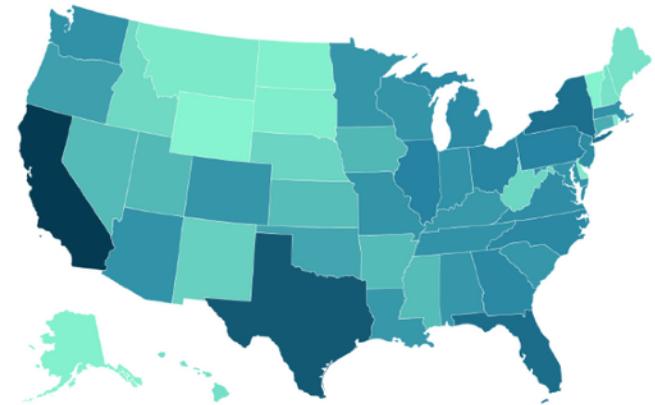
Choropleth maps are useful for visualizing *clear regional patterns* in the data

NOT IDEAL



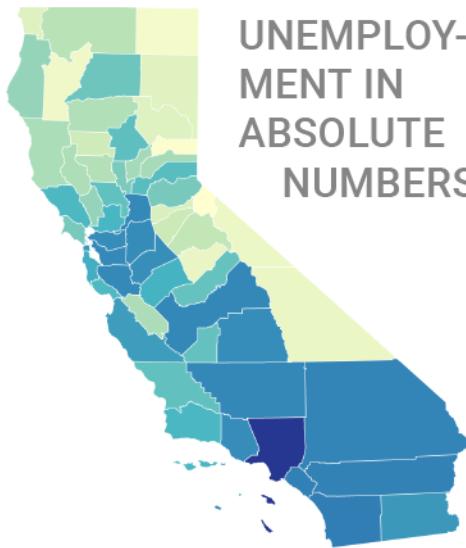
Use light colors for low values. Dark colors for high values.

BETTER



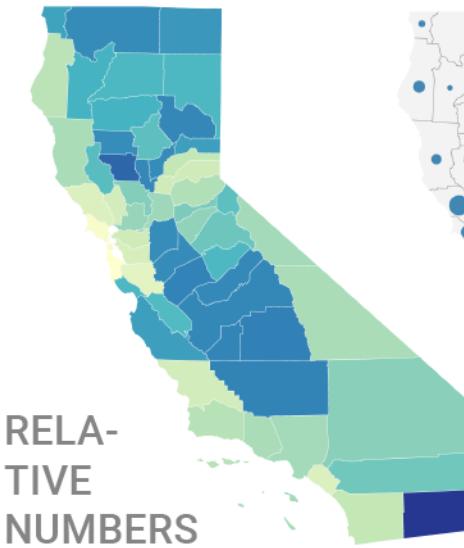
Choropleth should display relative differences, *not* absolute numbers

NOT IDEAL

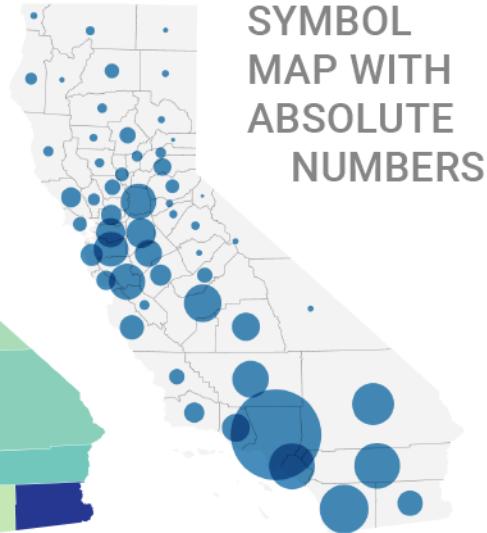


UNEMPLOY-
MENT IN
ABSOLUTE
NUMBERS

BETTER



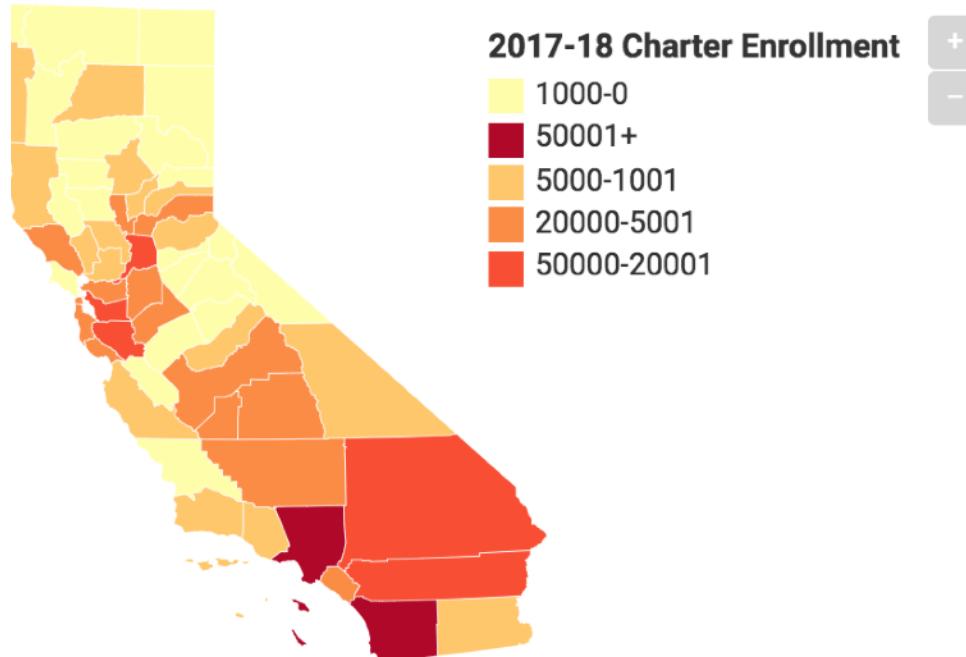
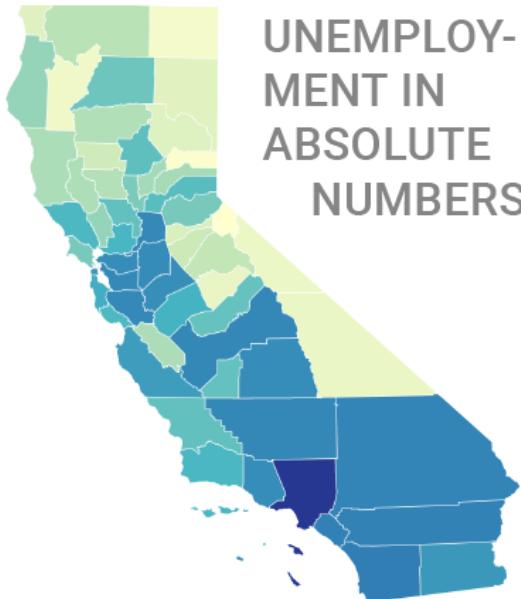
RELAT-
IVE
NUMBERS



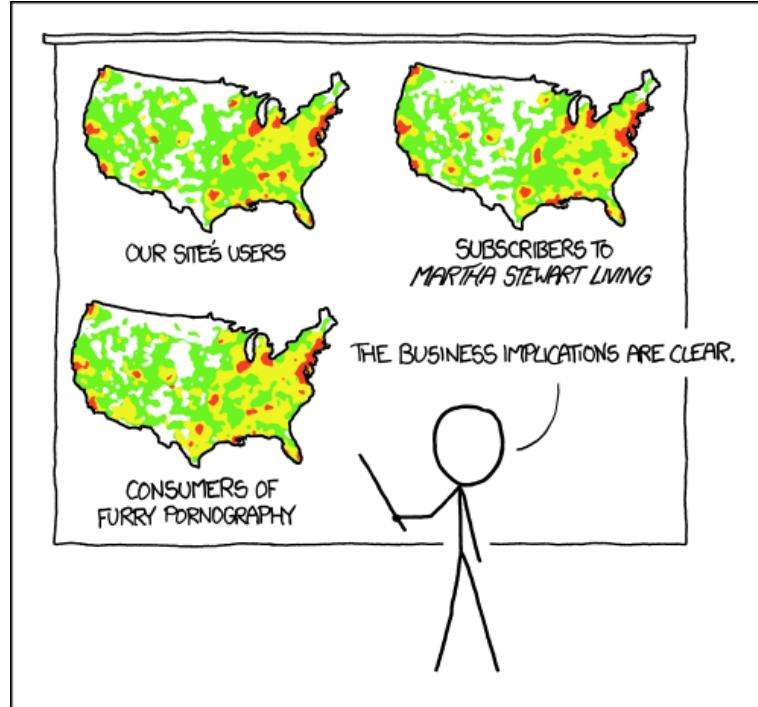
SYMBOL
MAP WITH
ABSOLUTE
NUMBERS

Map: Where Are Students Attending Charter Schools?

The majority of California's charter school student population is concentrated in Los Angeles, San Diego and Bay Area counties. Hover through the counties on each map for more information on their



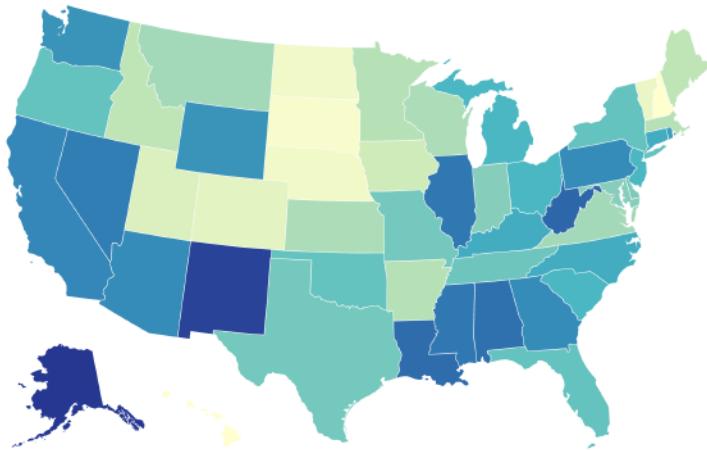
partment of Education • [Get the data](#) • Created with Datawrapper



PET PEEVE #208:
GEOGRAPHIC PROFILE MAPS WHICH ARE
BASICALLY JUST POPULATION MAPS

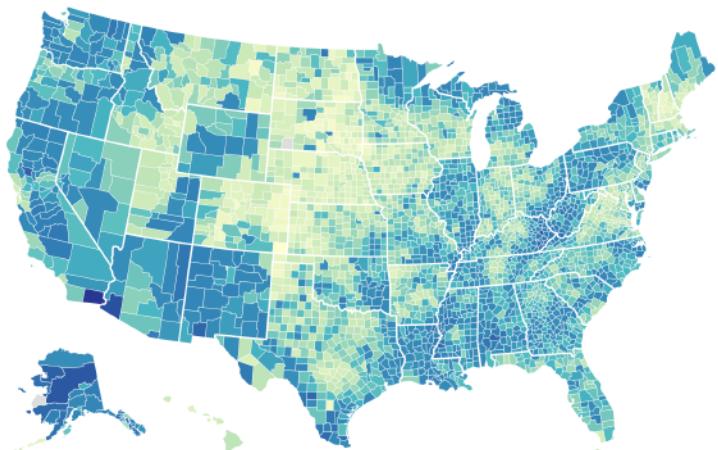
Choropleth maps can be misleading

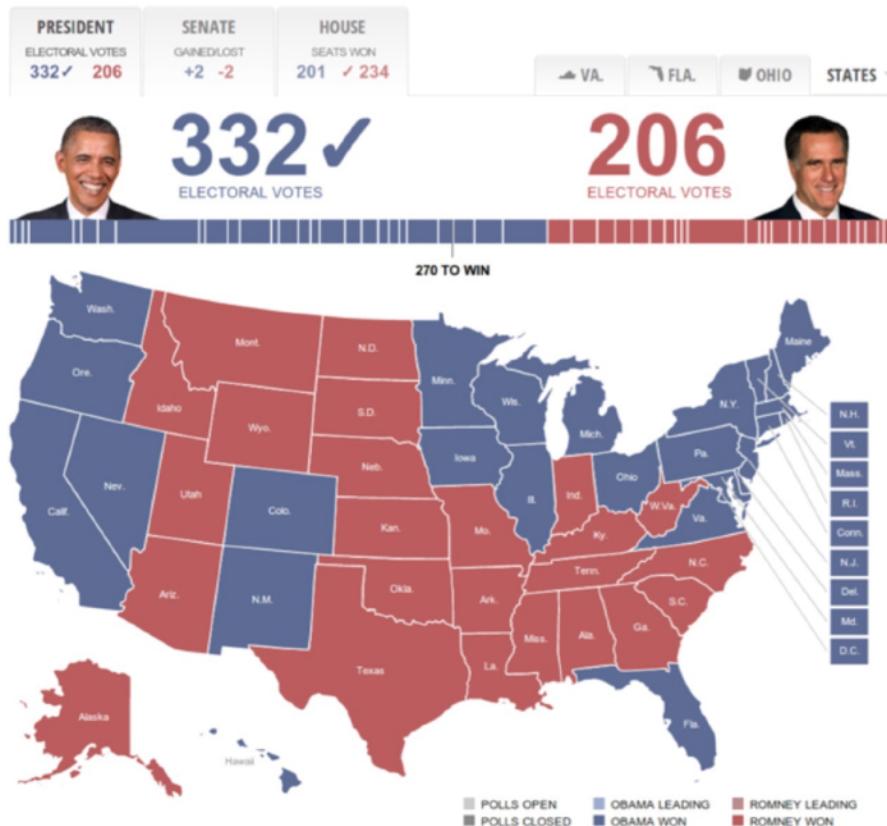
NOT IDEAL

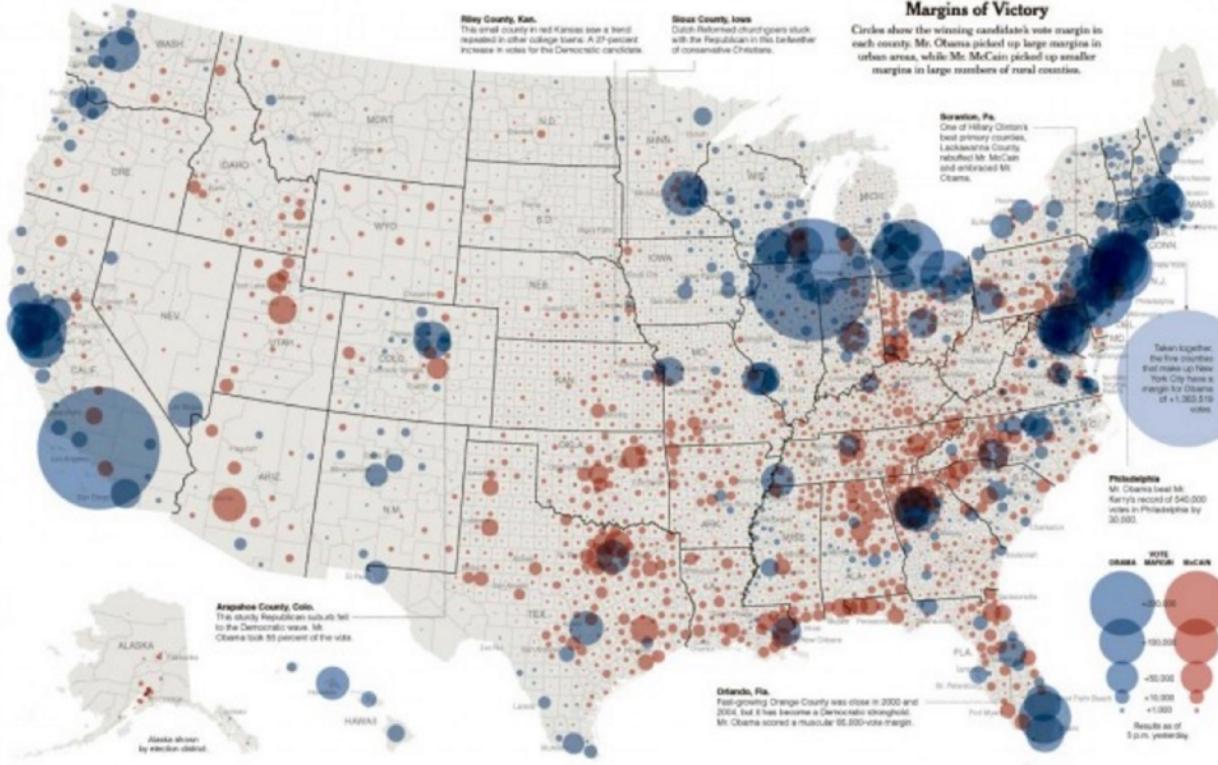


Consider using the smallest unit possible
(but there are exceptions!)

BETTER







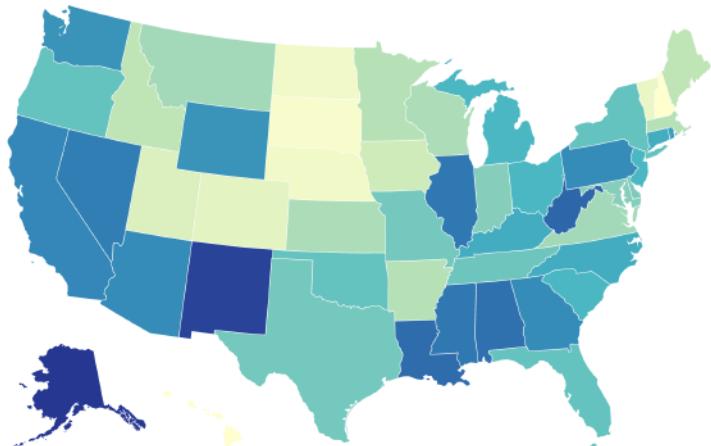
This **bubble graph** more accurately tells the full story, since the size of the bubbles is reflective of the population

...but same data can be displayed more effectively and informatively.

Visualization Choices

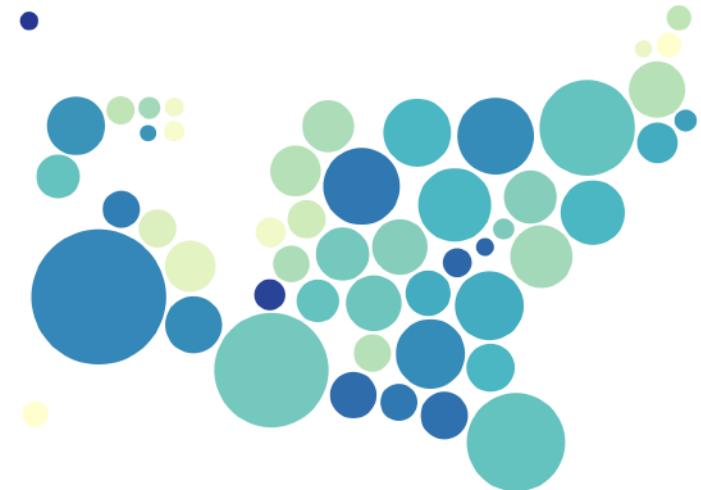
Cartograms should be considered when displaying how many people were affected

NOT IDEAL



Choropleths answer “How much area was affected?”

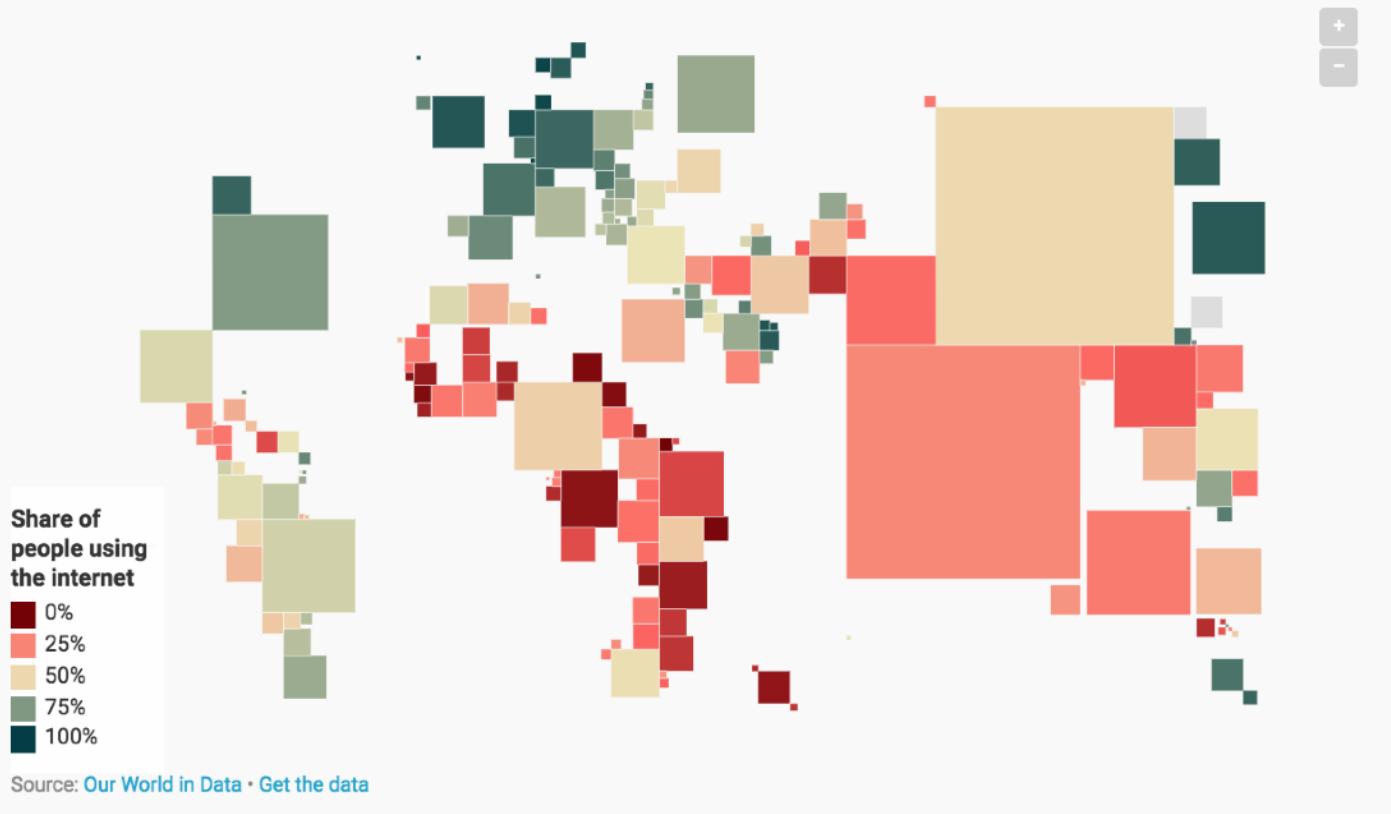
BETTER



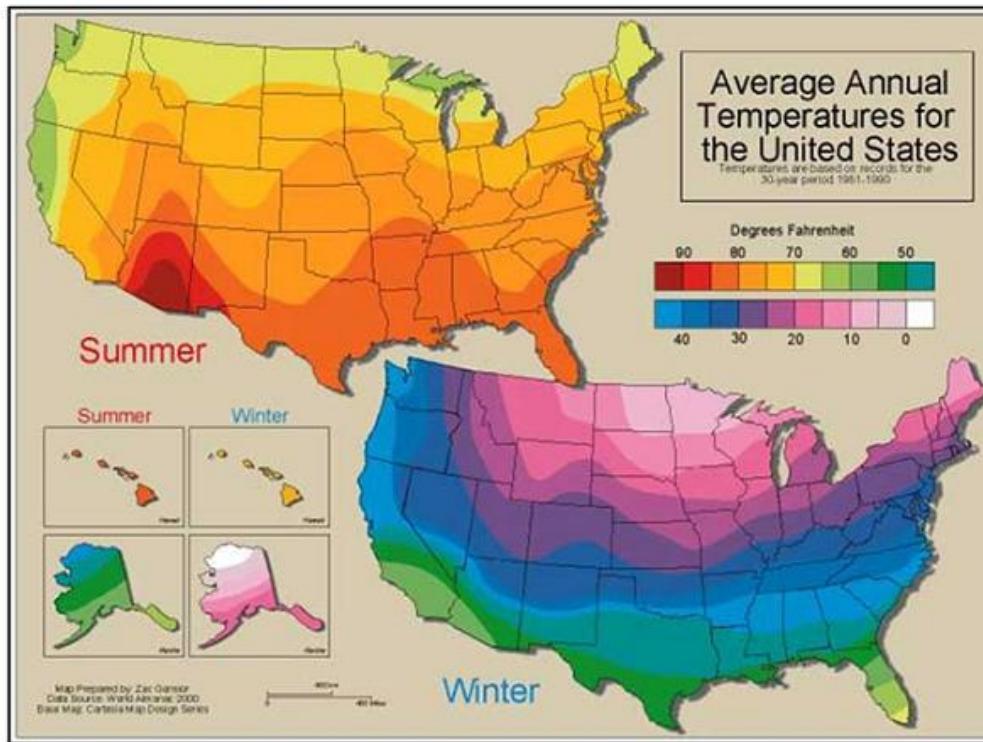
Cartograms answer “How many people were affected?”

Share of individuals using the internet, 2015

Share of individuals using the internet, measured as the percentage of the population. Internet users are individuals who have used the Internet (from any location) in the last 3 months. The Internet can be used via a computer, mobile phone, personal digital assistant, games machine, digital TV etc.



Isarithmic maps demonstrate smooth, continuous phenomena
(temperature, elevation, rainfall, etc.)

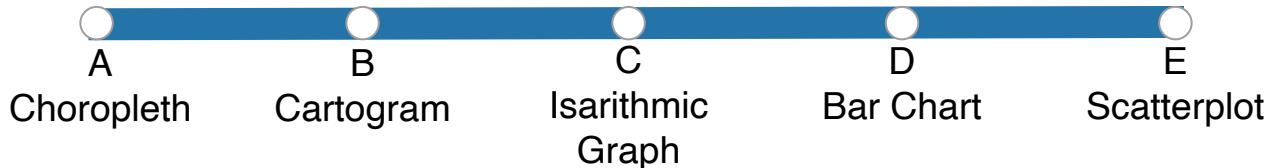




Visualizing Geospatial Data

You want to visualize how many people have been affected by COVID19 worldwide.

Best approach to visualize these data?



Spatial Statistics : The Why

Spatial Statistics

The statistical techniques we've discussed so far don't work well when considering spatial distributions...

Spatial Statistics

The statistical techniques we've discussed so far don't work well when considering spatial distributions...

...which means we have a chance to take a look at data and the relationship between the data in new and interesting ways (distance, adjacency, interaction, and neighbor)

Spatial data violate conventional statistics:

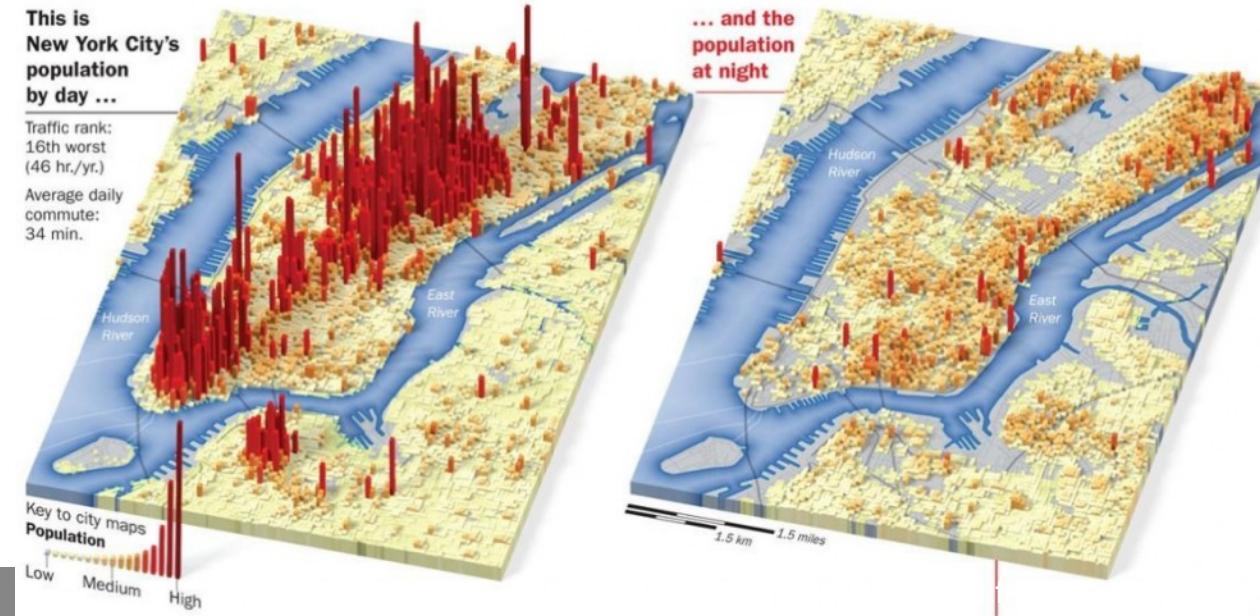
Violations of conventional statistics:

- Spatial autocorrelation
- Modifiable areal unit problem (MAUP)
- Edge effects (Boundary problem)
- Ecology fallacy
- Nonuniformity of space

Spatial Autocorrelation

Data from locations near one another in space are more likely to be similar than data from locations remote from one another:

- Housing market
- Elevation change
- Temperature



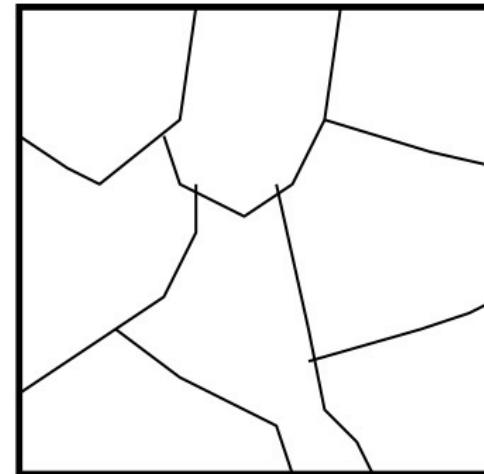
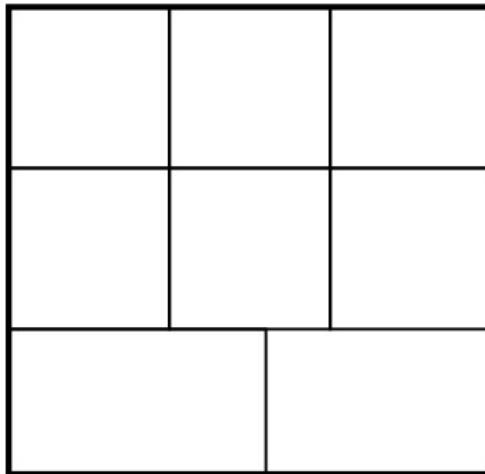
Modifiable Areal Unit Problem (MAUP)

The aggregation units used are arbitrary with respect to the phenomena under investigation, yet the aggregation units used will affect statistics determined on the basis of data reported in this way.

If the spatial units in a particular study were specified differently, we might observe very different patterns and relationships.

Modifiable Areal Unit Problem (MAUP)

modifiable area: Units are arbitrary defined and different organization of the units may create different analytical results.

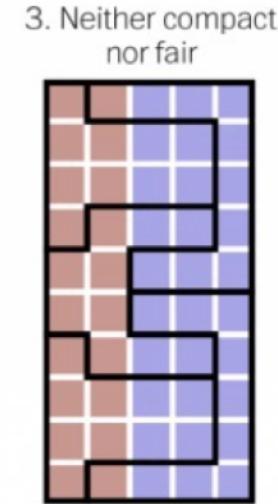
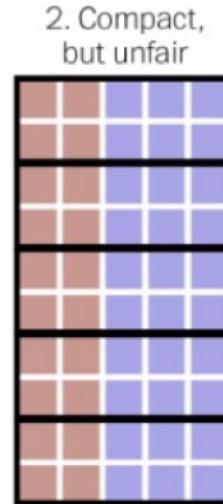
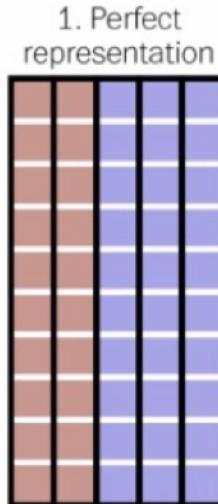
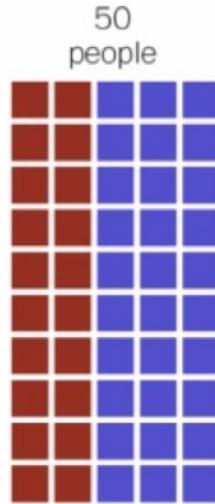


adapted from Brad Voytek

For example...gerrymandering

Gerrymandering, explained

Three different ways to divide 50 people into five districts



BLUE WINS

BLUE WINS

RED WINS

For example...gerrymandering

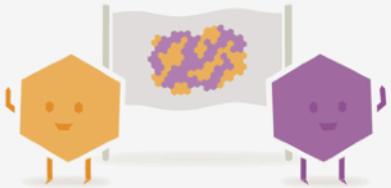
North Carolina

DISTRICTS REDRAWN TO OPTIMIZE COMPACTNESS



SOURCE: U.S. Census Bureau (top), Brian Olson (bottom)
GRAPHIC: The Washington Post. Published June 3, 2014

adapted from Brad Voytek



Welcome to Hexapolis



Every 10 years, Hexapolis redraws its congressional district lines — just like the United States does. But Hexapolis is a simpler place.



Lawmakers in either the **Purple Party** or **Yellow Party** control redistricting. To increase their advantage in upcoming elections, they have been known to gerrymander egregiously — even if it means leaving some voters disenfranchised.



Hexapolis has nine districts. Even though a majority of voters favor the Purple Party, that does not mean that the Yellow Party can't shift the state's partisan tilt.

<https://www.nytimes.com/interactive/2022/01/27/us/politics/congressional-gerrymandering-redistricting-game-2022.html>

Modifiable Areal Unit Problem (MAUP)

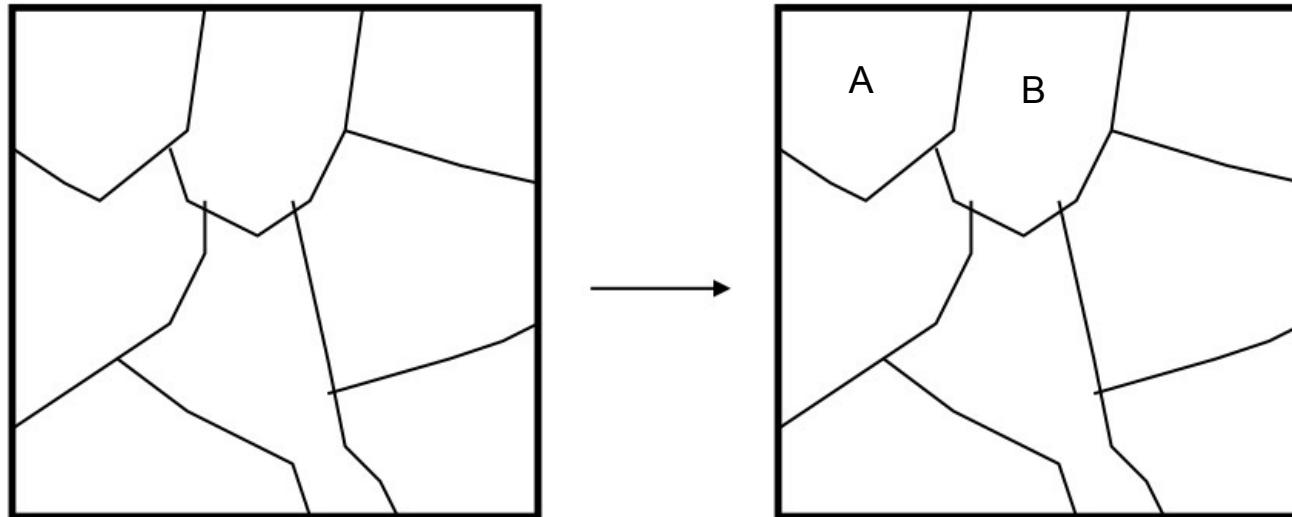
Potential problems in almost every field that utilizes spatial data.

In the 2016 U.S. presidential election, Hillary Clinton, with more of the population vote than Donald Trump, but failed to become president. (also true in Gore/Bush 2000)

A different aggregation of U.S. counties into states could have produced a different outcome.

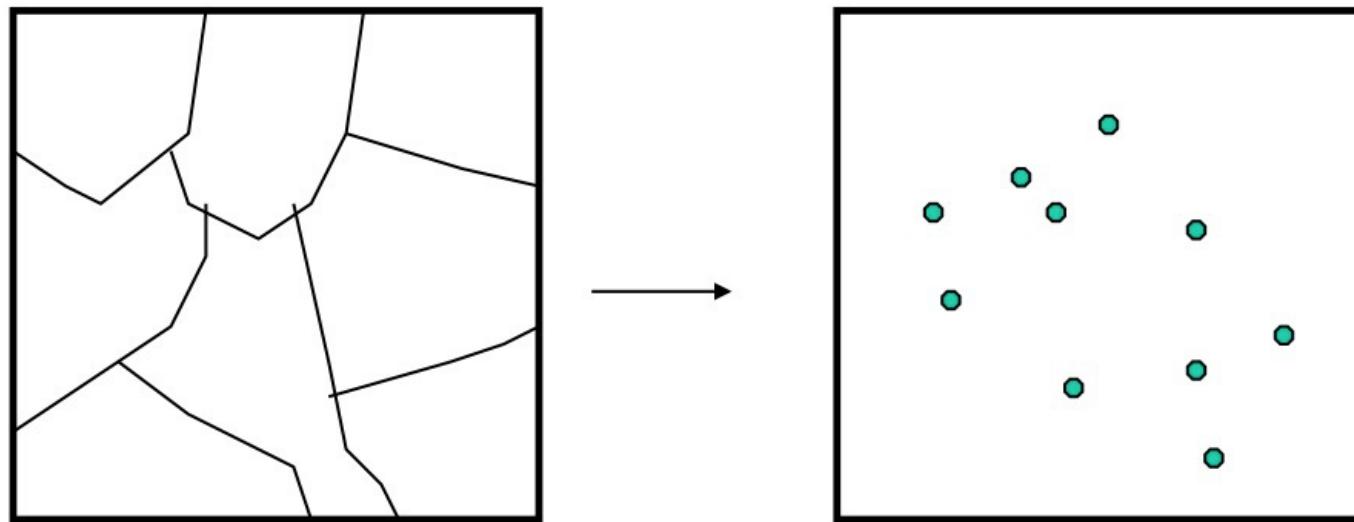
Edge Effects (The Boundary Problem)

Analyzing A vs B ignores
similarities between the two
based on their shared
boundary



Ecological Fallacy

The Ecological Fallacy is a situation that can occur when a researcher or analyst makes an inference about an individual based on aggregate data for a group.



adapted from Brad Voytek

Ecological Fallacy

Example: we might observe a *strong relationship between income and crime at the county level*, with lower-income areas being associated with higher crime rate.

Conclusions we might draw:

- Lower-income persons are more likely to commit crime
- Lower-income areas are associated with higher crime rates
- Lower-income counties tend to experience higher crime rates

The only valid conclusion!

Ecological Fallacy

Issues:

Inferences drawn about associations between the characteristics of an aggregate population and the characteristics of sub-units within the population are wrong. That is: *results from aggregated data (e.g. counties) cannot be applied to individual people*

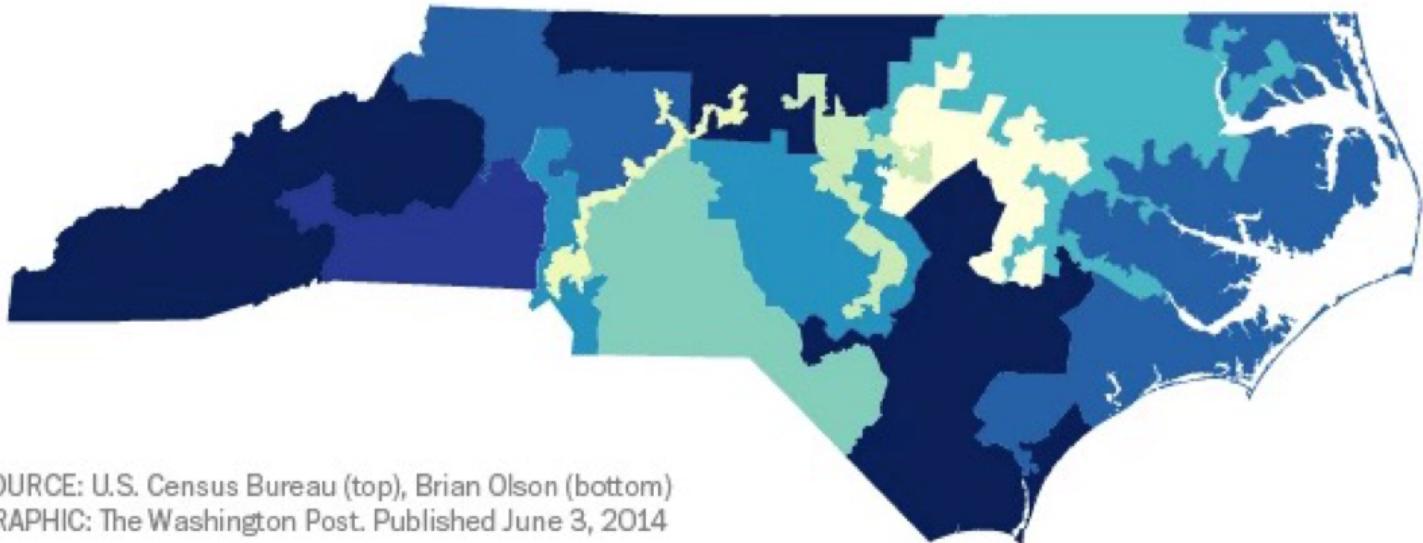
What should we do?

Be aware of the process of aggregating or disaggregating data may conceal the variations that are not visible at the larger aggregate level

For example...gerrymandering

North Carolina

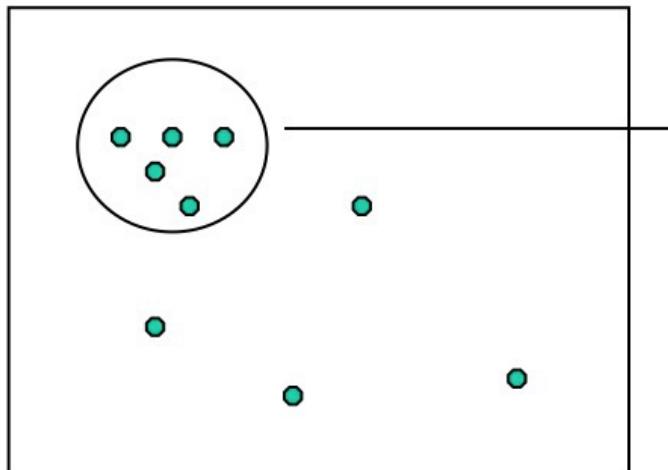
CURRENT CONGRESSIONAL DISTRICTS



SOURCE: U.S. Census Bureau (top), Brian Olson (bottom)
GRAPHIC: The Washington Post. Published June 3, 2014

adapted from Brad Voytek

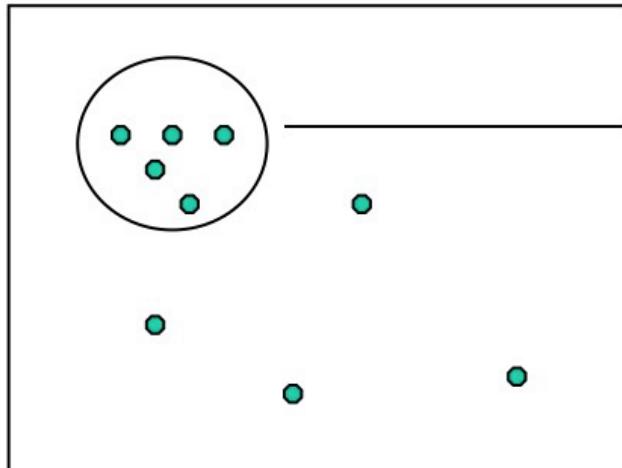
Nonuniformity



Area with high crime rates?

Crime locations

Nonuniformity



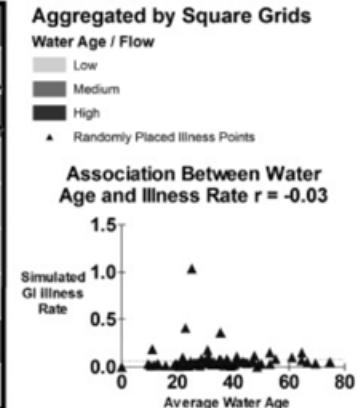
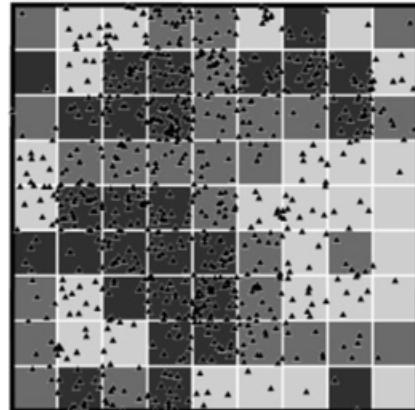
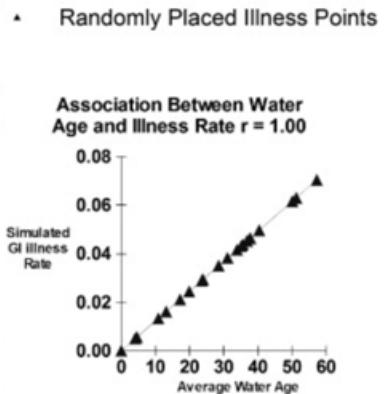
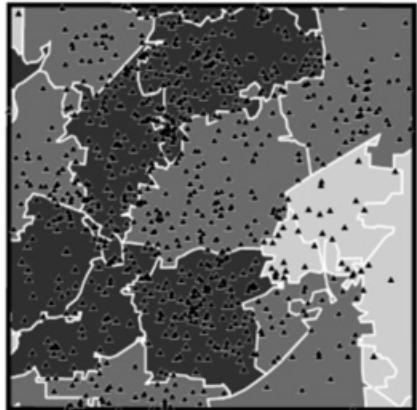
Area with high crime rates?

Crime locations

Conclusion: Bank robberies are clustered
....but only because banks are clustered!



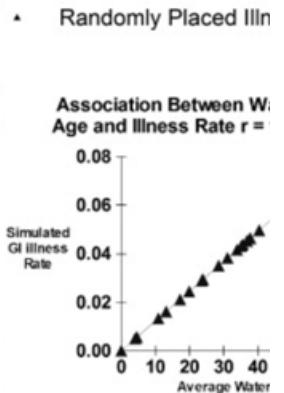
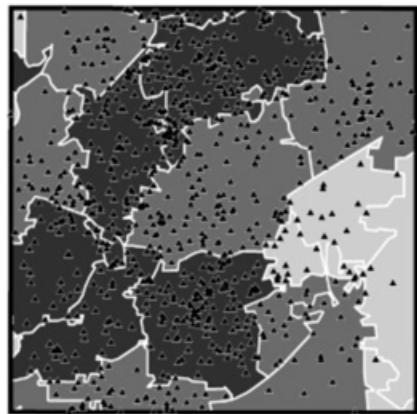
Spatial Statistics



What explains what's going on here?

- A Spatial Autocorrelation
- B MAUP
- C Edge Effects
- D Ecological Fallacy
- E Nonuniformity

Spatial Statistics



What explains variation?



A
Spatial
Autocorrelation



B
MAUP

