# Spatial statistical analysis at the National Cancer Institute

Linda Williams Pickle, Ph.D.

ESRI Health GIS Conference

October 26, 2005

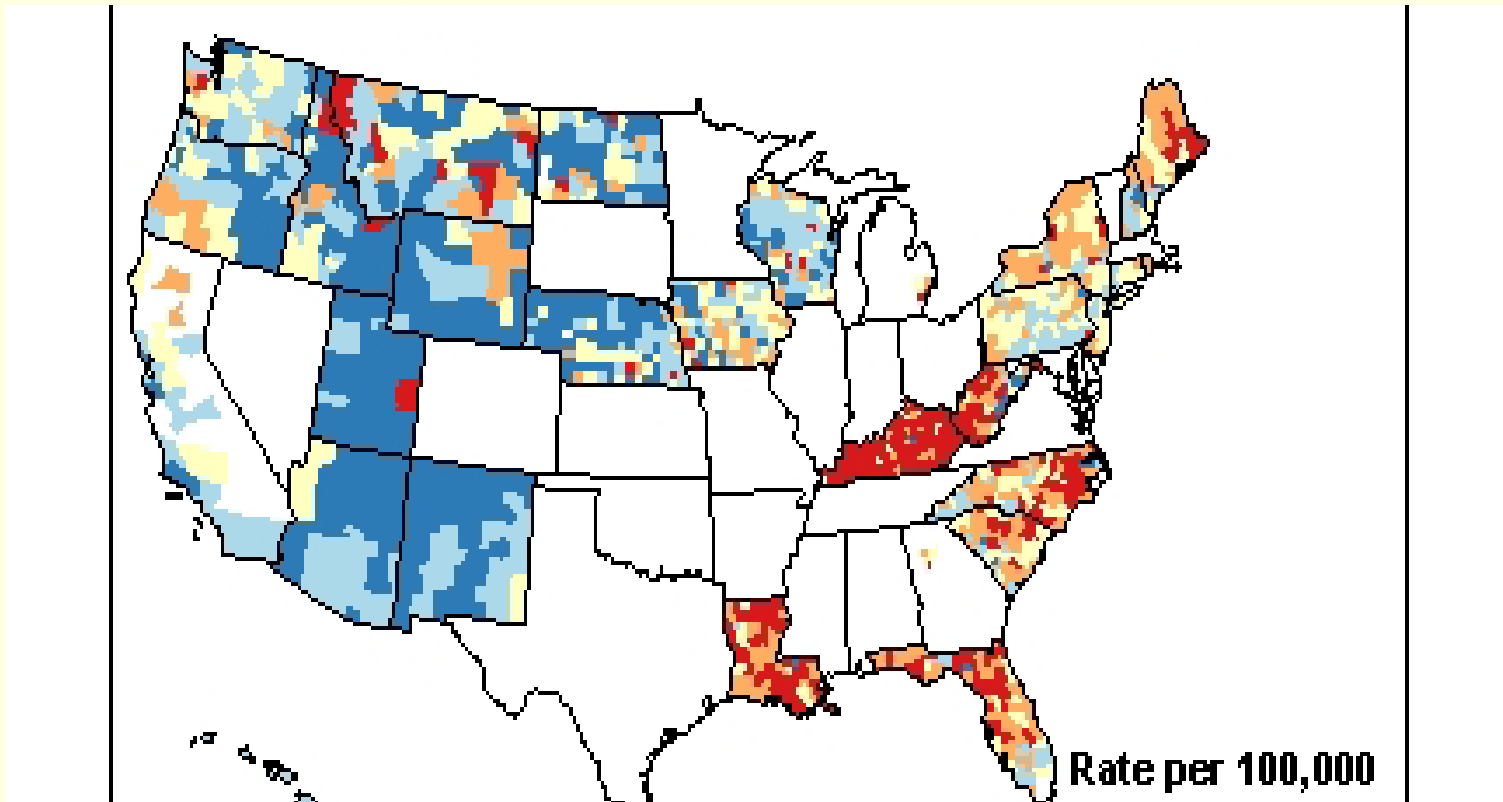# Steps in any statistical analysis project

- Data exploration
  - Quality control – errors in data?
  - Distribution of variables
  - Associations among predictor variables? between outcome & predictors?
  - Establish hypotheses to be tested, goals
- Statistical modeling (or hypothesis testing)
  - Verify model assumptions
  - Apply model; estimate parameters
  - Assess fit of the model, modify until satisfactory
- Communicate results to client, public, etc.

- Will illustrate application to spatial data using cancer rate examples

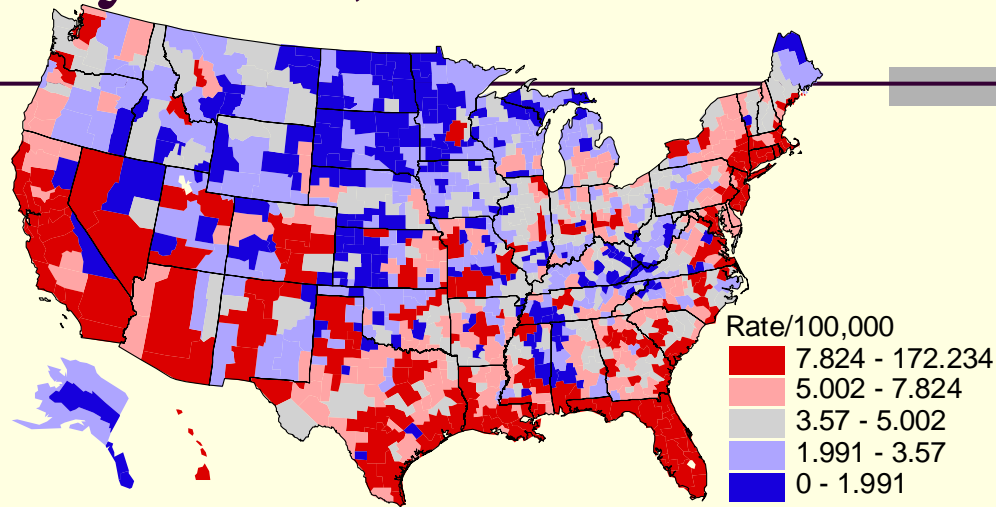# Spatial data exploration

- Visualization of patterns in original spatial data
  - QC: any anomalies? Holes in map, extreme outliers?
  - Evaluate general pattern
    - Smoothing
    - Clustering
    - Outliers
    - Measures of spatial correlation
- Hypothesis generation
  - What risk factor maps look similar to cancer rate map? (How to measure similarity of patterns on multiple maps?)
- Select potential predictors for model (covariates)
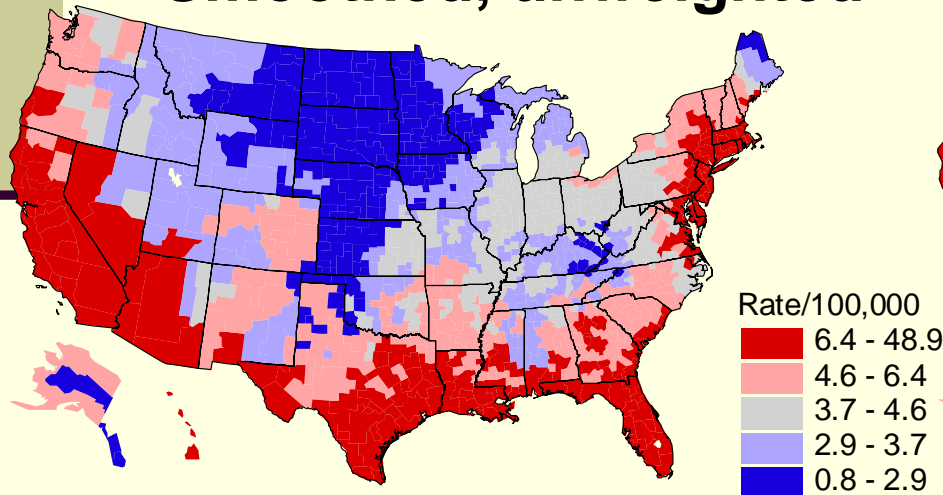
# Quality control: Missing some CA data



Rate per 100,000
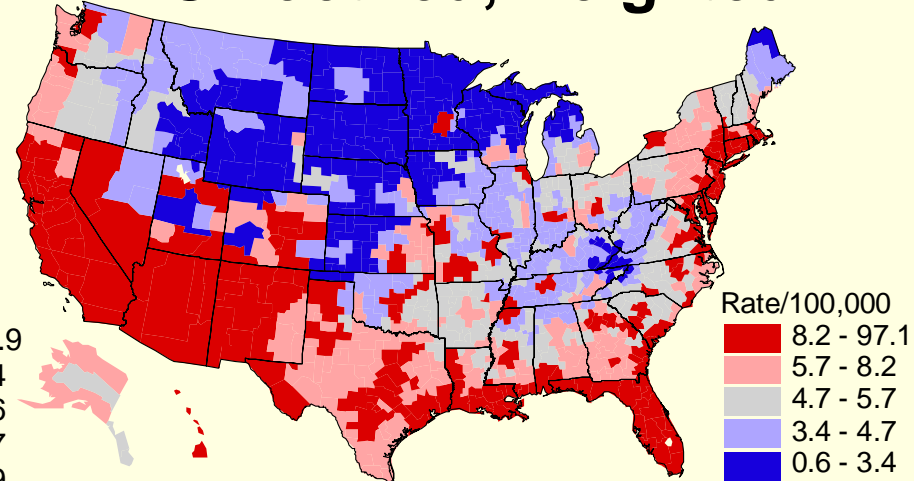
# Smoothing:
# HIV mortality rates, 1988-92

**Original data:**



Rate/100,000
- 7.824 - 172.234
- 5.002 - 7.824
- 3.57 - 5.002
- 1.991 - 3.57
- 0 - 1.991

**Smoothed, unweighted**



Rate/100,000
- 6.4 - 48.9
- 4.6 - 6.4
- 3.7 - 4.6
- 2.9 - 3.7
- 0.8 - 2.9

**Smoothed, weighted**



Rate/100,000
- 8.2 - 97.1
- 5.7 - 8.2
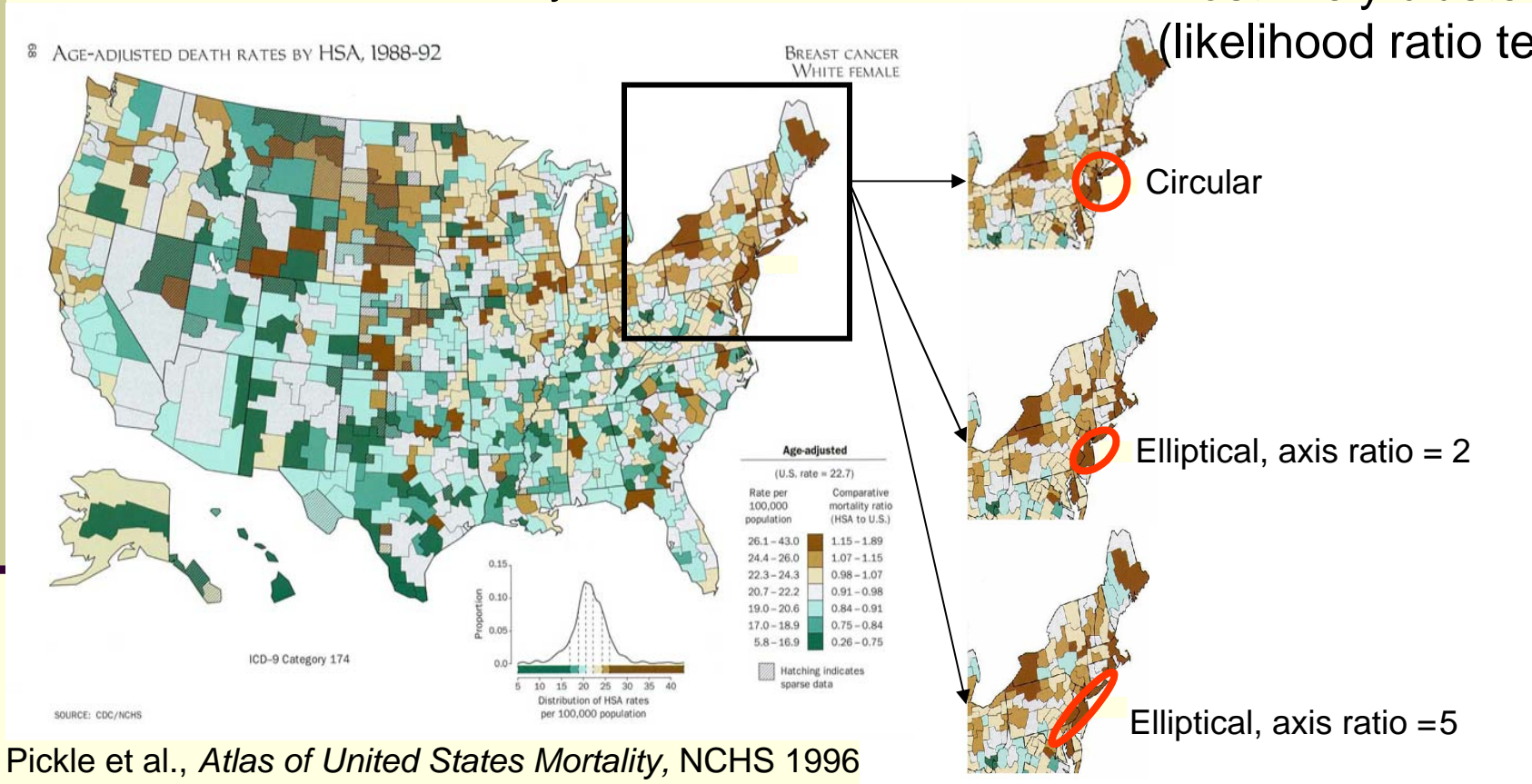- 4.7 - 5.7
- 3.4 - 4.7
- 0.6 - 3.4

**Source: Pickle et al., *Atlas of United States Mortality*, NCHS, 1996.**

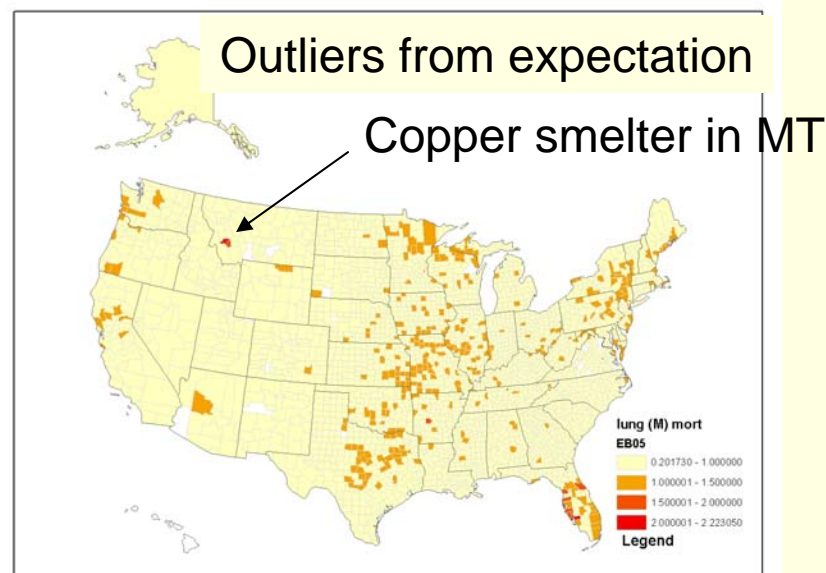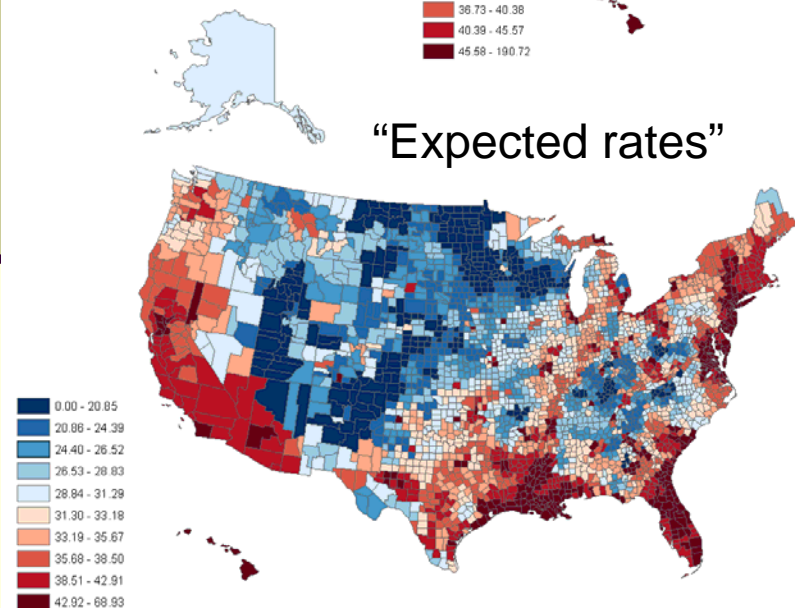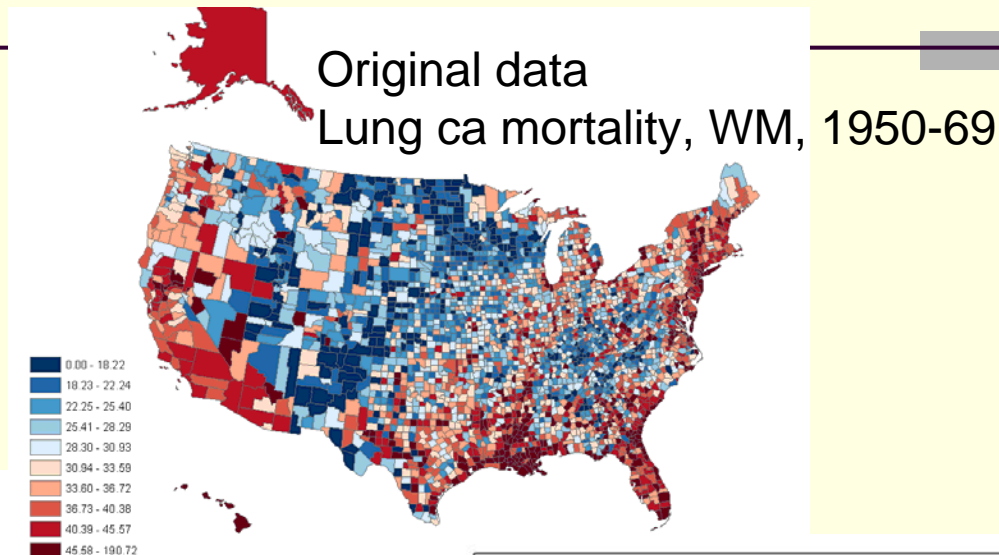# Cluster identification using SaTScan: Breast cancer mortality rates

**Breast cancer mortality rates**

**Most likely cluster (likelihood ratio test)**



Circular

Elliptical, axis ratio = 2

Elliptical, axis ratio = 5

Pickle et al., *Atlas of United States Mortality,* NCHS 1996
SaTScan by Martin Kulldorff, available at *www.satscan.org*

# Outlier identification:
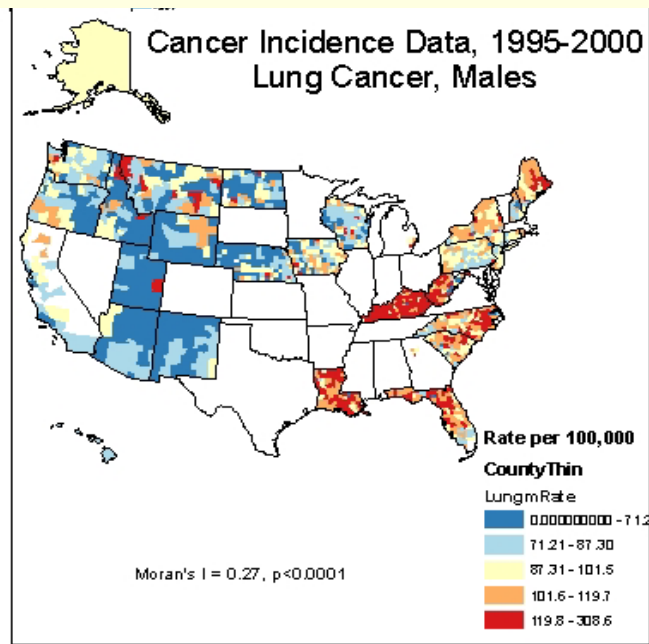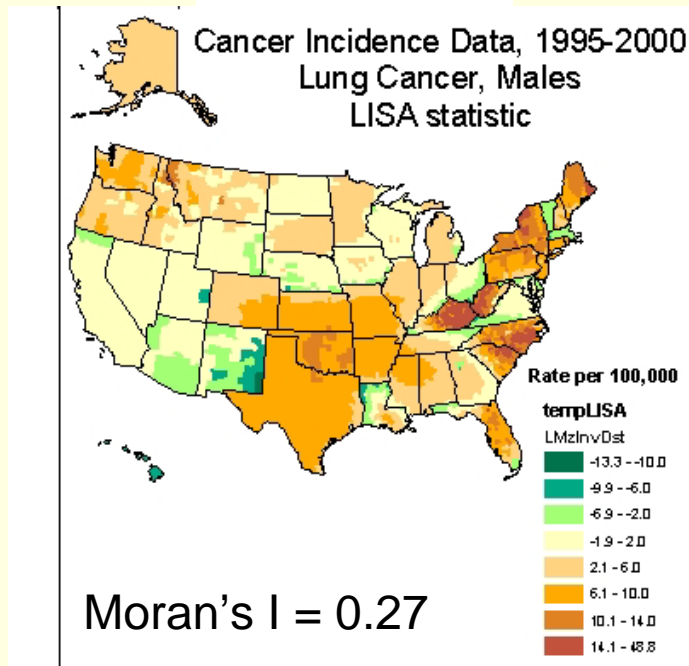# Values statistically different from expectation



Original data
Lung ca mortality, WM, 1950-69

"Expected rates"

Outliers from expectation

Copper smelter in MT

Method: *(DuMouchel & Pregibon, Proc KDD, 2001; Lincoln Technologies, Inc)*

# Cluster detection in ArcMap


Cancer Incidence Data, 1995-2000
Lung Cancer, Males
Rate per 100,000
CountyThin
LungmRate

| | |
|---|---|
| | 0.00000000 - 71.20 |
| | 71.21 - 87.30 |
| | 87.31 - 101.5 |
| | 101.6 - 119.7 |
| | 119.8 - 308.6 |

Moran's I = 0.27, p<0.0001

LISA statistic

Getis-Ord Gi* statistic


Cancer Incidence Data, 1995-2000
Lung Cancer, Males
LISA statistic
Rate per 100,000
tempLISA
LMzInvDst

| | |
|---|---|
| | -13.3 - -10.0 |
| | -9.9 - -6.0 |
| | -6.9 - -2.0 |
| | -1.9 - 2.0 |
| | 2.1 - 6.0 |
| | 6.1 - 10.0 |
| | 10.1 - 14.0 |
| | 14.1 - 48.8 |

Moran's I = 0.27

Color Breaks For Z:
min
-10
-6
-2
+2 (NS)
+6
+10
+14
max


Cancer Incidence Data, 1
Lung Cancer, Males
Getis-Ord Gi*
Rate per 100,000
tempGetisOrdG2
GiInvDst 2

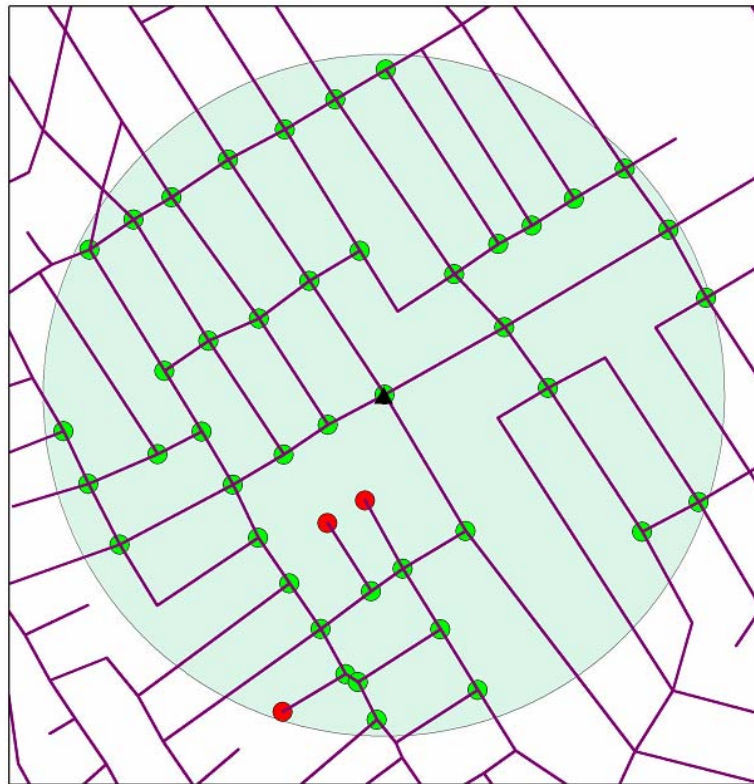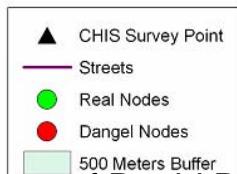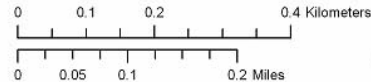| | |
|---|---|
| | -0.762228 - 0.000000 |
| | 0.000001 - 2.000000 |
| | 2.000001 - 5.214720 |

Moran's I = 0.27, p<0.0001

# Defining potential risk factors using a GIS:
# High and Low Connectivity Buffers in Los Angeles

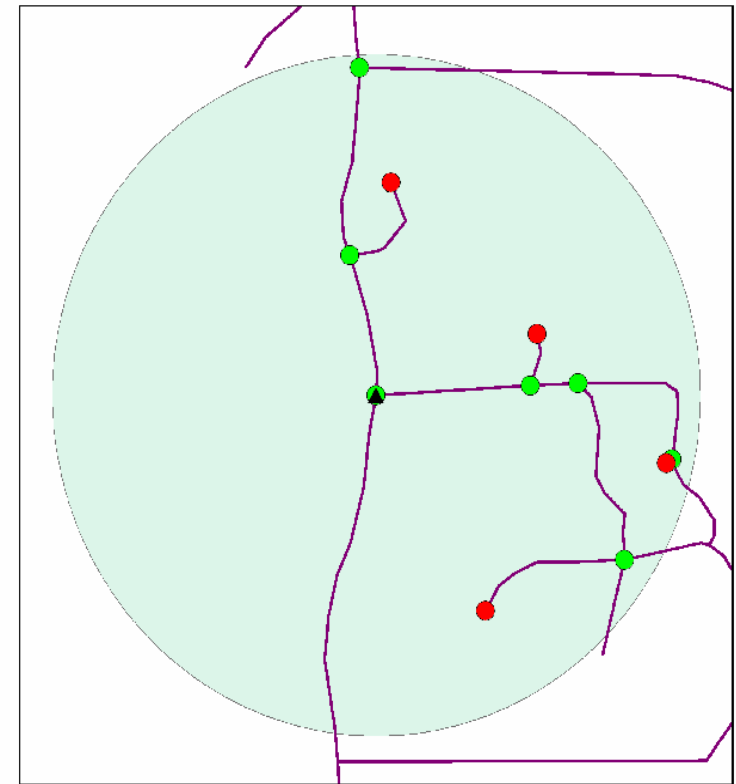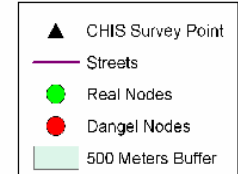(From The California Health Interview Survey, 2001)



Projection: NAD 1983 StatePlane California V FIPS 0405
Data Source: US Census Bureau
CHIS Survey

| Link Node Ratio | 2.02 |
| Intersection Density in Km | 54.73 |
| Connected Node Ratio | 0.93 |
| Street Network Density | 14.53 |
| Gamma Index | 0.46 |
| Alpha Index | 0.18 |
| Block Density | 37.98 |
| Median Block Length | 0.11 |
| Average Block Length | 0.16 |
| Population Density in Km | 3752.35 |
| Employment Density in Km | 647.18 |

▲ CHIS Survey Point
— Streets
● Real Nodes
● Dangel Nodes
▢ 500 Meters Buffer

Projection: NAD 1953 StatePlane California V FPS 0405
Data Source: US Census Bureau
CHIS Survey

| Link Node Ratio | 1.45 |
| Intersection Density in Km | 8.91 |
| Connected Node Ratio | 0.64 |
| Street Network Density | 3.53 |
| Gamma Index | 0.37 |
| Alpha Index | 0.00 |
| Block Density | 2.78 |
| Median Block Length | 0.20 |
| Average Block Length | 0.26 |
| Population Density in Km | 15.63 |
| Employment Density in Km | 1.14 |

▲ CHIS Survey Point
— Streets
● Real Nodes
● Dangel Nodes
▢ 500 Meters Buffer

# Choosing covariates for the analysis: Comparison of patterns of outcome and potential risk factors



Lung cancer mortality in Maryland White males, 1998-2002

**Age-Adjusted Annual Death Rate (Deaths per 100,000)**
Quantile Interval

- 96.4 to 120.5
- 86.7 to 96.3
- 81.2 to 86.6
- 76.7 to 81.1
- 67.6 to 76.6
- 41.3 to 67.5

**United States Rate (95% C.I.)** 75.2 (75.0 - 75.5)

**Maryland Rate (95% C.I.)** 76.1 (74.1 - 78.0)

**Healthy People 2010 Goal 03-02** 44.9

Somerset County 1998 - 2002 Age-adjusted death rate = 120.5 (90.0 - 160.

URL: statecancerprofiles.cancer.gov

% men ever smoked cigarettes

Data source: BRFSS, CDC

# Exploratory Spatio-Temporal Analysis Tool (ESTAT)
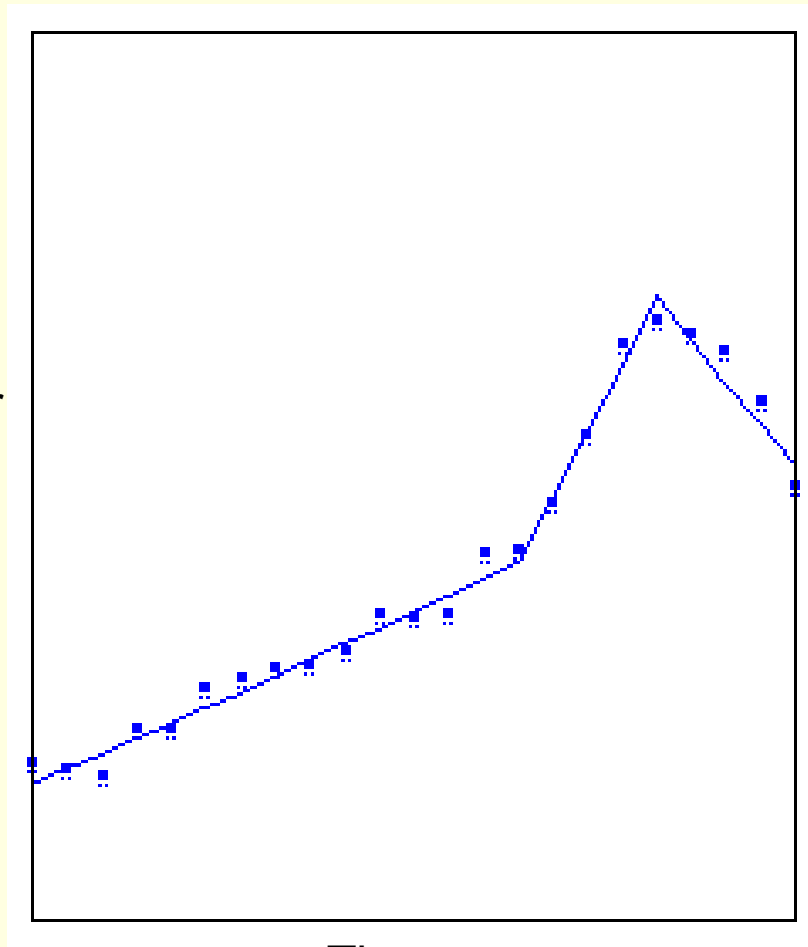


Map

Rate
Time
Series
Plot

Scatter
plot

Covariate
PCP plot

Developed by Alan MacEachren & GeoVista staff, Penn State University

# Cancer surveillance:
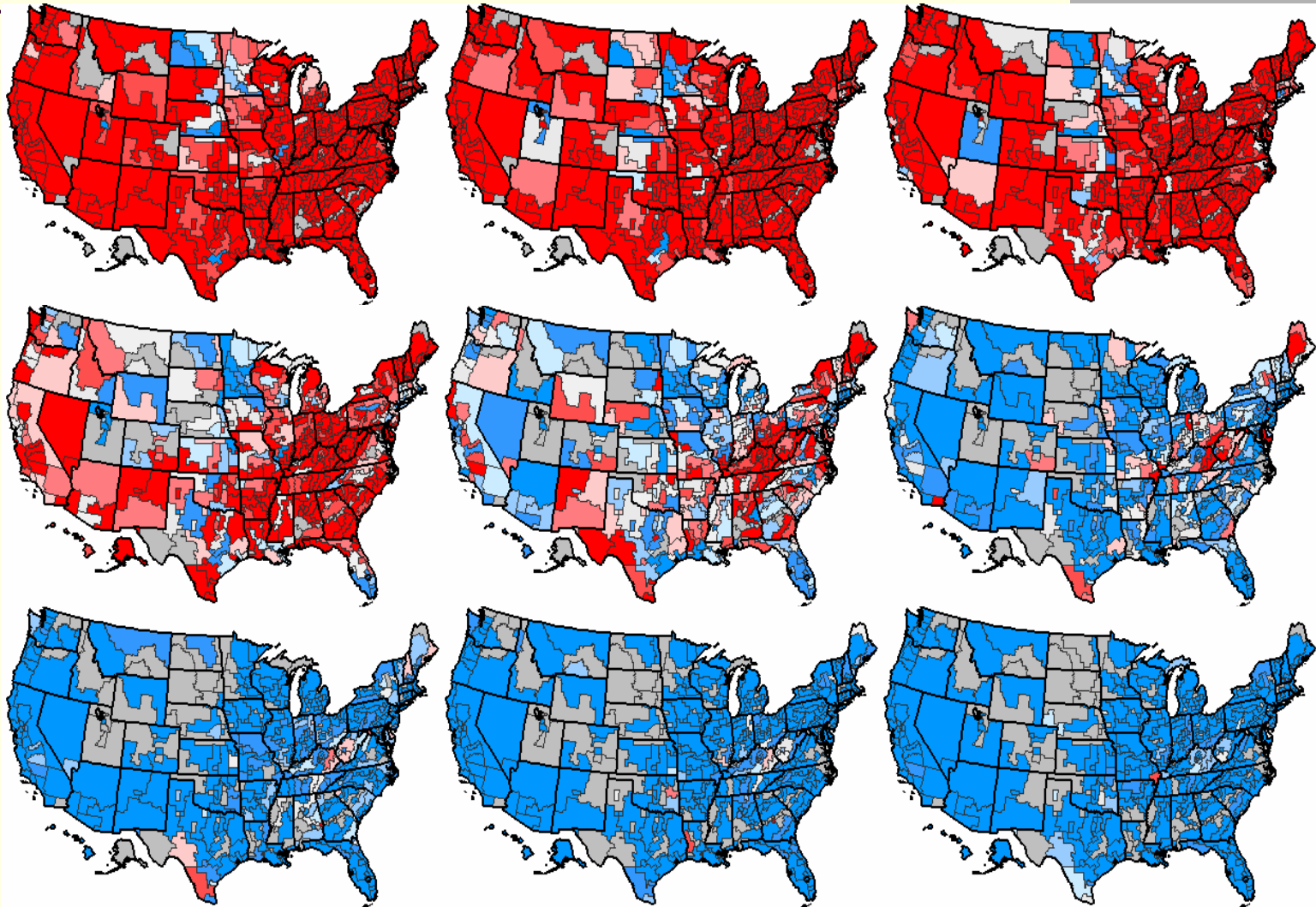# Has the cancer trend changed? If so, where?

Cancer rate

Joinpoint software determines when time trend changed significantly

Time

# Spatio-temporal patterns
## Cervical cancer mortality 1950-94 by 5 years

# Goals of statistical modeling

- Inference
  - Explain the patterns, calculate risks
  - Emphasis is on interpretation of model results
- Prediction
  - Either smooth observed data or fill in gaps
  - Emphasis is on fit of the model
  - Further complication: project ahead in time
- Approaches to modeling
  - Include as many covariates & interactions as possible to explain spatial patterns & correlation
  - Use a simpler model and let the spatial correlation "soak up" some of the spatial variation
- Spatial correlation must be included if necessary, otherwise variance estimates will be wrong
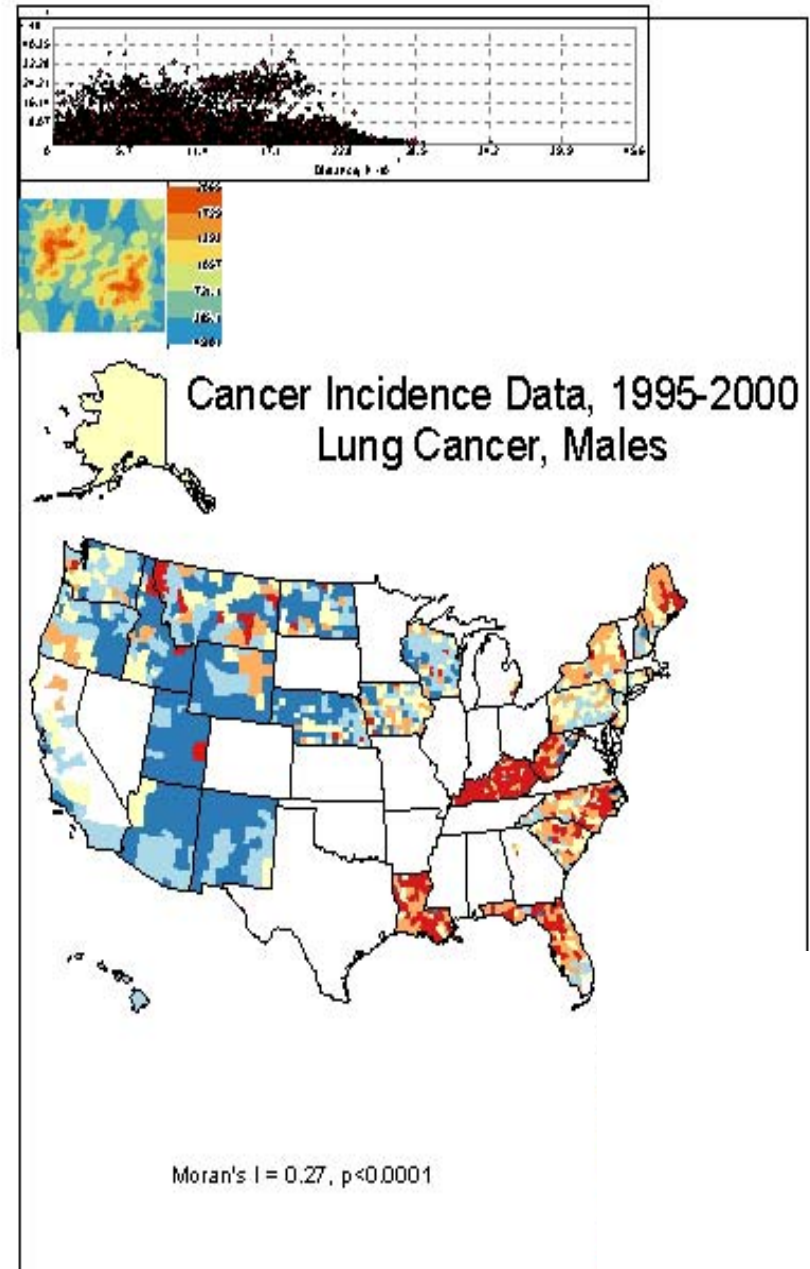
# The process of statistical modeling

- Verify model assumptions (statistical & spatial)
- Apply model to data, estimate parameters
- Assess fit of the model
    - Chi-square goodness-of-fit, deviance statistics, etc
    - Scatter plots of observed vs. predicted values, leverage (observations that most influence results),…
    - Map the results
        - Does predicted value map look like observed value map?
        - Do residuals appear to be spatially random?
        - Are residuals still spatially correlated? (variogram)
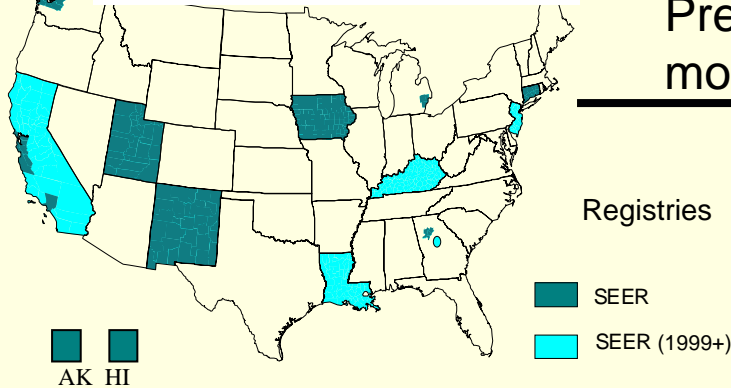- If necessary (almost always!), modify model & rerun

SAS, S+, Etc.

# Use of Geostatistical Analyst to check spatial assumptions

- Stationarity
- Isotropy
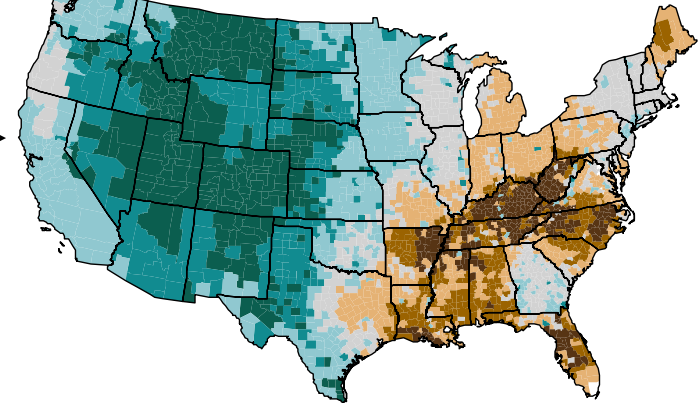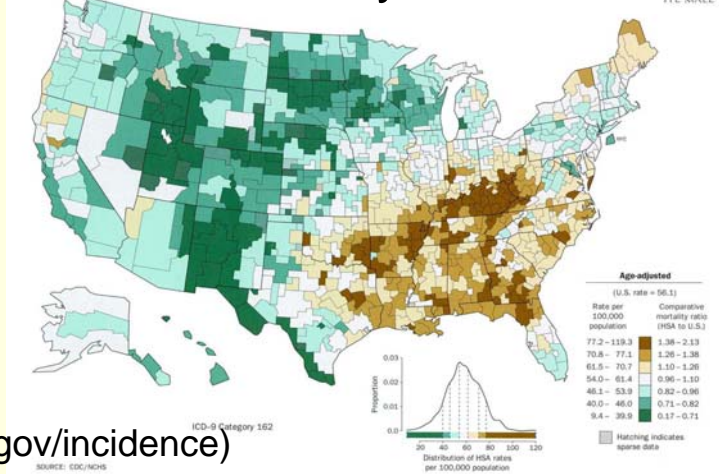- Functional form of spatial correlation (variogram models)



Cancer Incidence Data, 1995-2000
Lung Cancer, Males

Moran's I = 0.27, p<0.0001

# Map comparison as a goodness-of-fit tool



NCI cancer registries
Input data to model

Prediction model

Incidence, 1999

Registries

SEER

SEER (1999+)

AK  HI

Mortality, 1988-92

Top: Pickle et al., NCI monograph 2003 (URL: srab.cancer.gov/incidence)
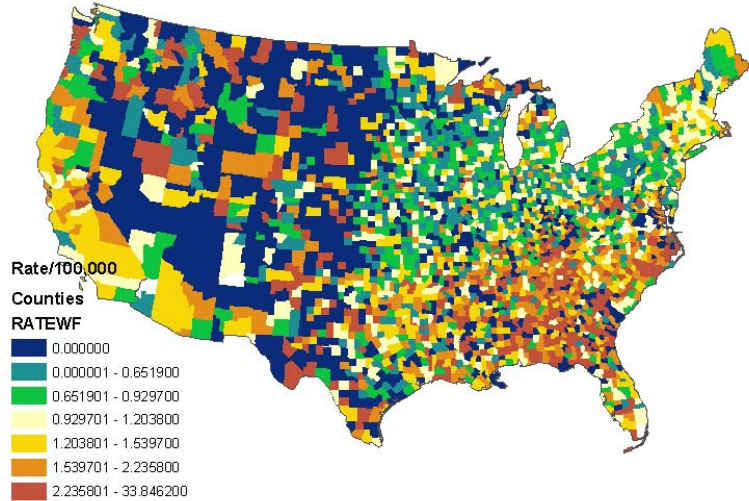Bottom: Pickle et al., *Atlas of United States Mortality,* NCHS 1996
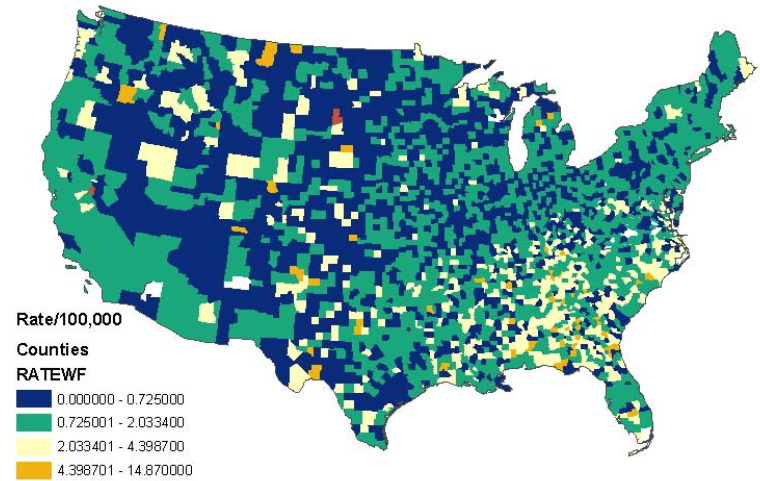
# Communicating results of spatial analysis

- Important for gov't agencies to disseminate information
  - Back to original data collectors (cancer registries, states)
  - To researchers in the subject area
  - To the public
- Information needs to be accurate and clear to diverse audience, dissemination tools need to be user friendly
- Accuracy on a map
  - Tension between precision (narrow rate categories) & readability (not too many categories)
  - Cartographic choices can impact visual impressions
- Uncertainty of statistic must be communicated
- Often need multiple maps for multiple purposes

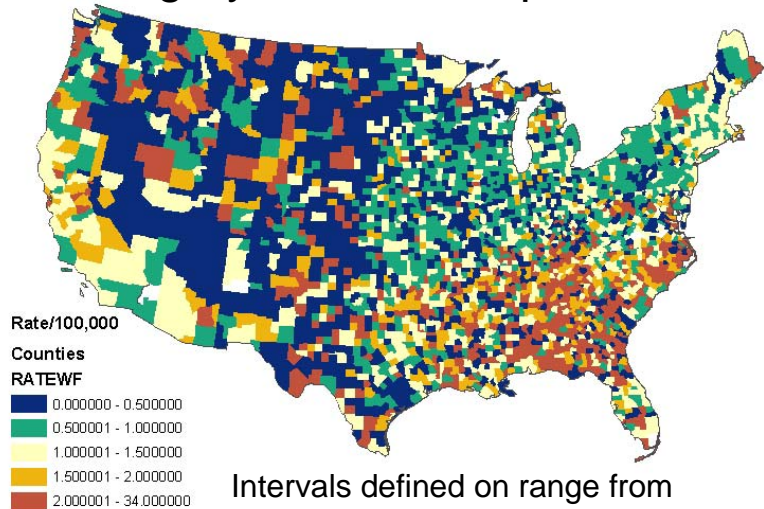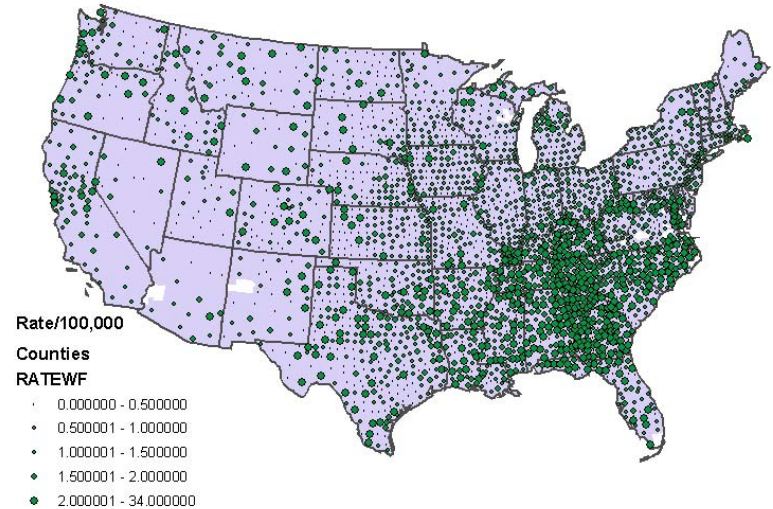# Oral cancer mortality, white females, 1950-69

## 7 category quantile

Rate/100,000
Counties
RATEWF
- 0.000000
- 0.000001 - 0.651900
- 0.651901 - 0.929700
- 0.929701 - 1.203800
- 1.203801 - 1.539700
- 1.539701 - 2.235800
- 2.235801 - 33.846200

## Jenks (4 categories)

Rate/100,000
Counties
RATEWF
- 0.000000 - 0.725000
- 0.725001 - 2.033400
- 2.033401 - 4.398700
- 4.398701 - 14.870000

## 5 category truncated equal interval

Rate/100,000
Counties
RATEWF
- 0.000000 - 0.500000
- 0.500001 - 1.000000
- 1.000001 - 1.500000
- 1.500001 - 2.000000
- 2.000001 - 34.000000

Intervals defined on range from
10%tile to 90%tile of distribution

## Categorical symbols

Rate/100,000
Counties
RATEWF
- 0.000000 - 0.500000
- 0.500001 - 1.000000
- 1.000001 - 1.500000
- 1.500001 - 2.000000
- 2.000001 - 34.000000

# Methods of communicating uncertainty by separation of value and variance information



All Cancer Sites
Year 2002
Latest Annual Death Rate
All Races
Both Sexes, All Ages

URL: statecancerprofiles.cancer.gov

Confidence intervals on Graphic linked to map

Separate maps
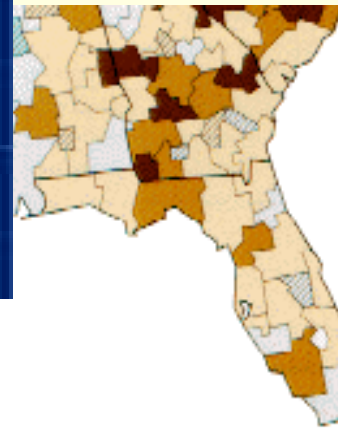
Cause 16   Age-adjusted Death Rates, 1988–1992

Cause 16
Reliability of
Death Rates

Source: MacEachren et al.,
Environment & Planning A,
1998.

# Methods of communicating uncertainty by superimposing uncertainty layer on map



Meteorologist's "cone of uncertainty"



Hatching indicates
Sparse data
(unreliable rates)

Source: Pickle et al., *Atlas of United States Mortality,* NCHS 1996

# Small multiple maps useful for comparisons: Relative mortality ratios (area:US), 1988-92 Selected causes of death, white males



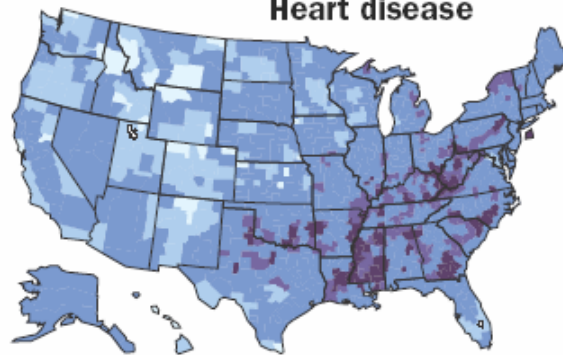Comparative mortality ratio (HSA to U.S.)
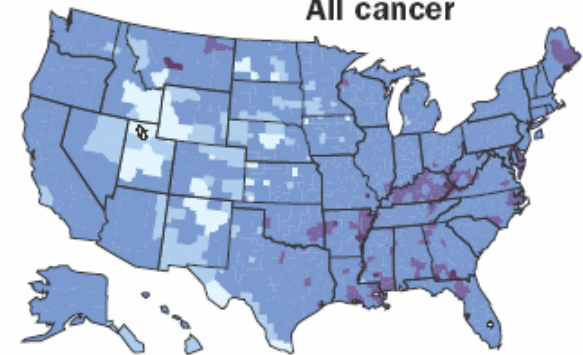
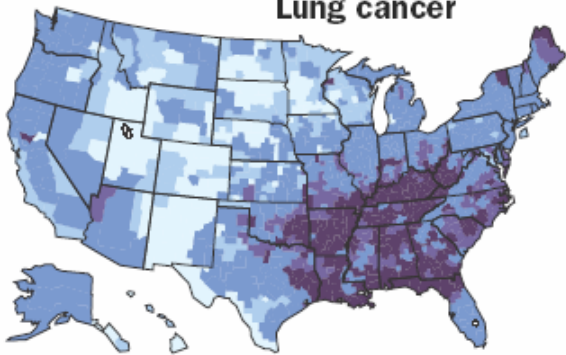- \> 1.25
- 1.16 – 1.25
- 0.85 – 1.15
- 0.75 – 0.84
- \< 0.75

Heart disease

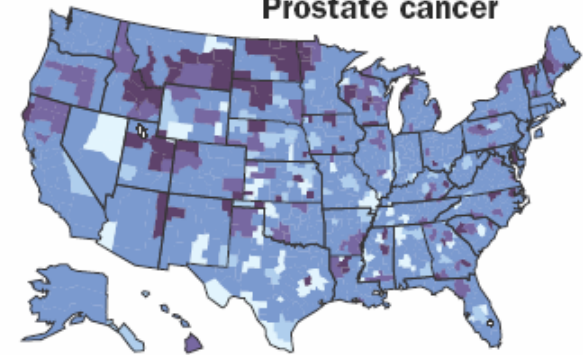All cancer
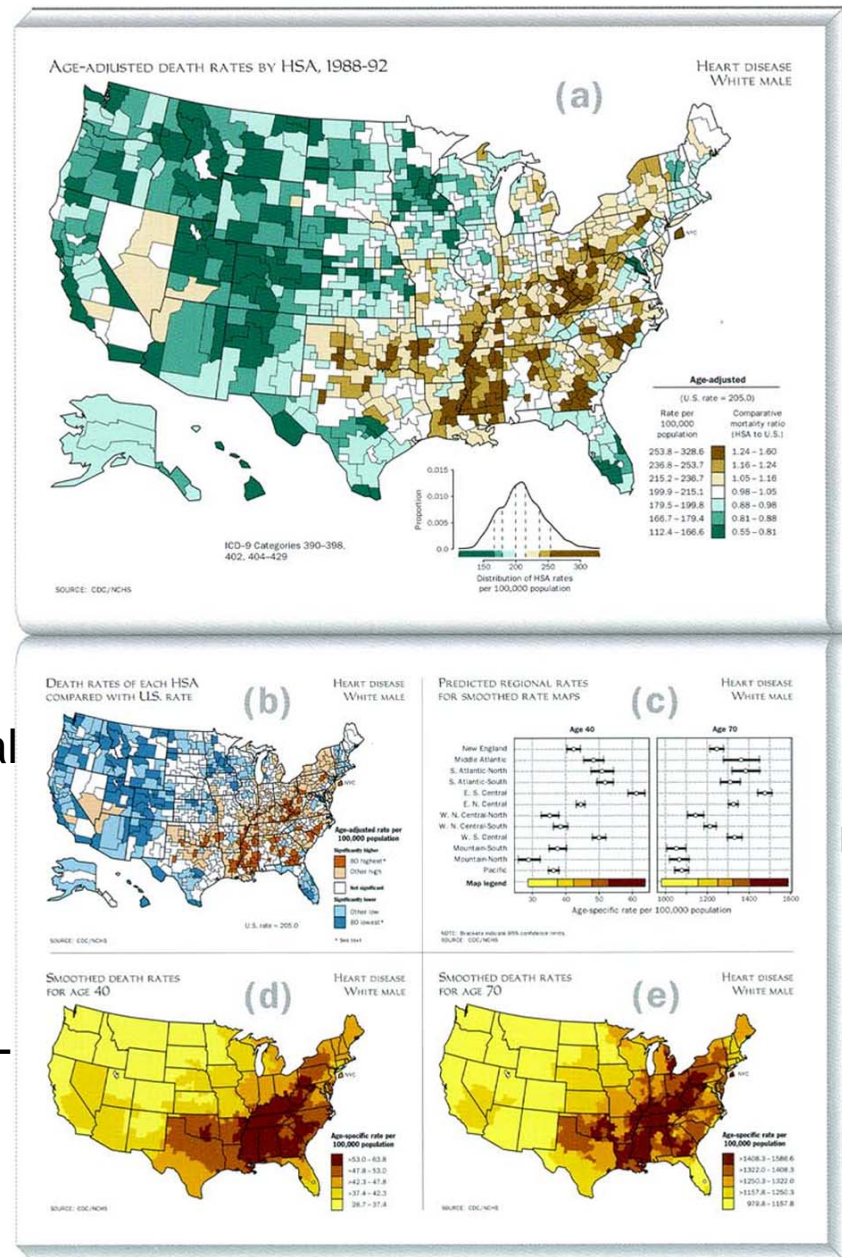
Lung cancer

Colorectal cancer

Prostate cancer

Pickle et al., *Atlas of United States Mortality,* NCHS 1996

# Multiple maps for multiple purposes



FIGURE 1. GRAPHICAL COMPONENTS OF THE TWO-PAGE ATLAS LAYOUT

Observed rates & reliability

Results of statistical significance test

Regional rates

Modeled age-specific rates

**Source: Pickle et al., *Atlas of United States Mortality*, NCHS, 1996.**

# Conclusions:
# Methods important for spatial statistical analysis

- Flexible smoothing methods to explore patterns
- Inclusion of weights for smoothing, spatial statistics in order to account for population heterogeneity
- Output must include a measure of reliability (variance)
- Methods to examine spatial pattern characteristics: spatial correlation, stationarity, isotropy
- Methods to identify significant patterns, e.g., clusters & outliers – eye can be fooled
- Link to commonly-used statistical software packages
- Cartographic choices that present least biased view of the data (colors, categorization, unreliability, symbology)