# School of Geography and Earth Sciences
## McMaster University

Advanced Topics in Spatial Statistics

# Area Data
# V & VI

# This session

- ○ Modeling area data
  - • Non-spatial regression models
- ○ Assumptions
  - • Non-constant variance
    - ○ Data transformations
  - • Error autocorrelation
    - ○ Moran's I
  - • Normality

# Example

○ <u>Voter turnout</u>
- What explains voter turnout in elections?

○ What do we know about the data?
- Variables follow a spatial pattern (Spatial autocorrelation)
- Variables are probably correlated

# Example

- <u>Voter turnout:</u> Variables
  - VOTE
  - Voting age population (POP)
  - Population with grade 12 or higher education (EDU)
  - Number of owner occupied houses (HOUSE)
  - Aggregate income (INCOME)
- Unit of analysis: county

# Example

○ Do we have prior expectations about the direction of the relationships?
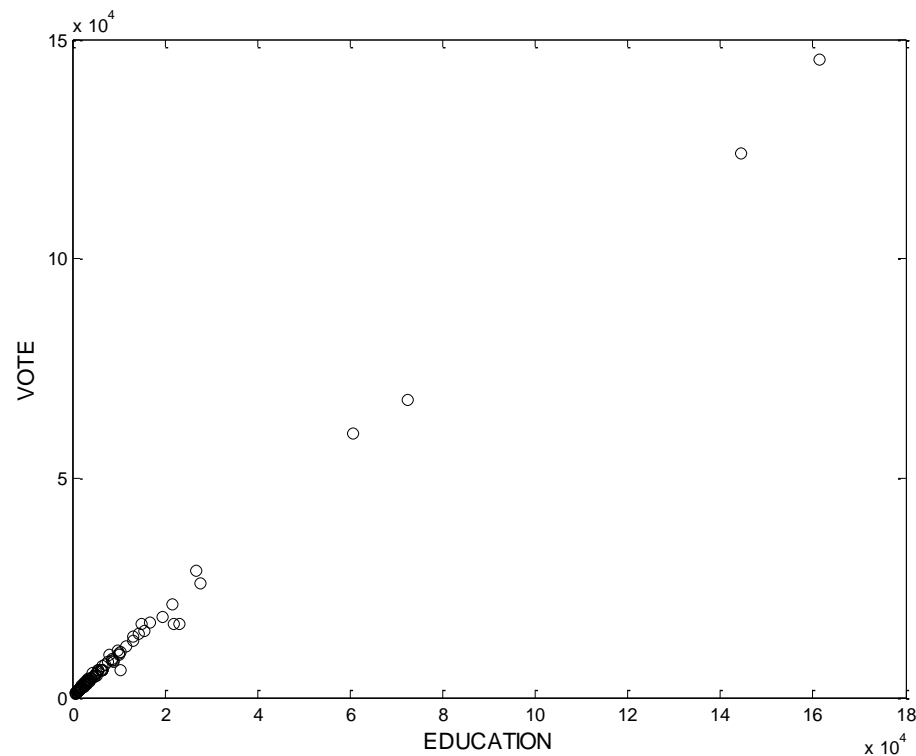
# Example

- Bivariate analysis
  - Scatterplots
  - Correlation coefficients

# Example

○ Scatterplot (EDU vs. VOTE)

# Example

○ Bivariate analysis

● Correlation coefficients

|  | VOTE | POP | EDU | HOUSE | INCOME |
|---|---|---|---|---|---|
| VOTE | 1.00 |  |  |  |  |
| POP | 0.993 | 1.00 |  |  |  |
| EDU | 0.998 | 0.990 | 1.00 |  |  |
| HOUSE | 0.997 | 0.994 | 0.994 | 1.00 |  |
| INCOME | 0.994 | 0.983 | 0.998 | 0.989 | 1.00 |

# Example

- Modeling area data
  - Non-spatial regression
    - Variables can be spatial
    - Variables may be autocorrelated
    - Assumes that there is no spatial autocorrelation
  - Multiple regression

$$Y = \mathbf{X}b + e$$

# Modeling Area Data

- Try different models
- Model evaluation guidelines

# Modeling Area Data

○ Model 1

|          | Model 1 (VOTE) |         |
| -------- | :------------: | :-----: |
| Variable | PARAMETER | t-value |
| CONST | 695.21 | 5.23 |
| EDU | 0.89 | 163.61 |

R^2=                          0.996
R^2(adj)=                     0.977
SIGMA^2=            1560013.858
SIGMA^2 (ML)=      1530299.308
SIGMA      =             1249.005
n=                              105
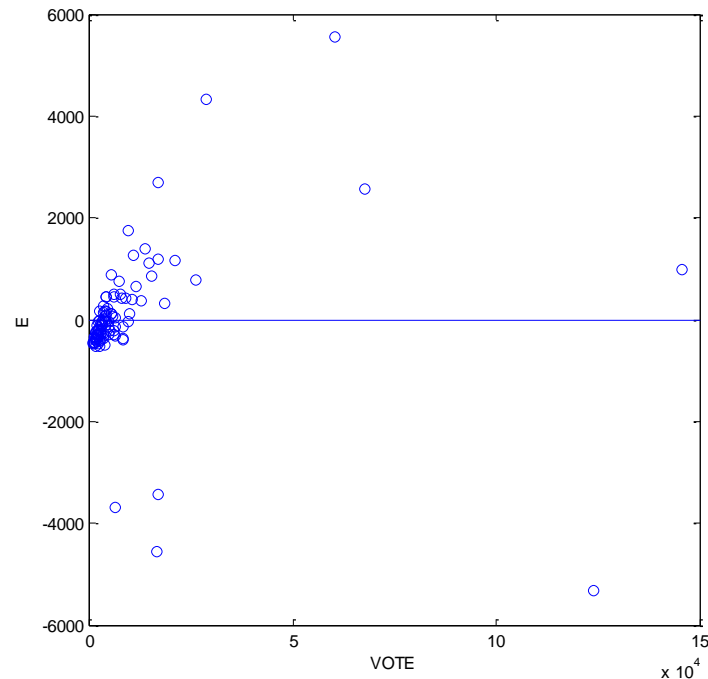
<< Normalized Moran's I >>
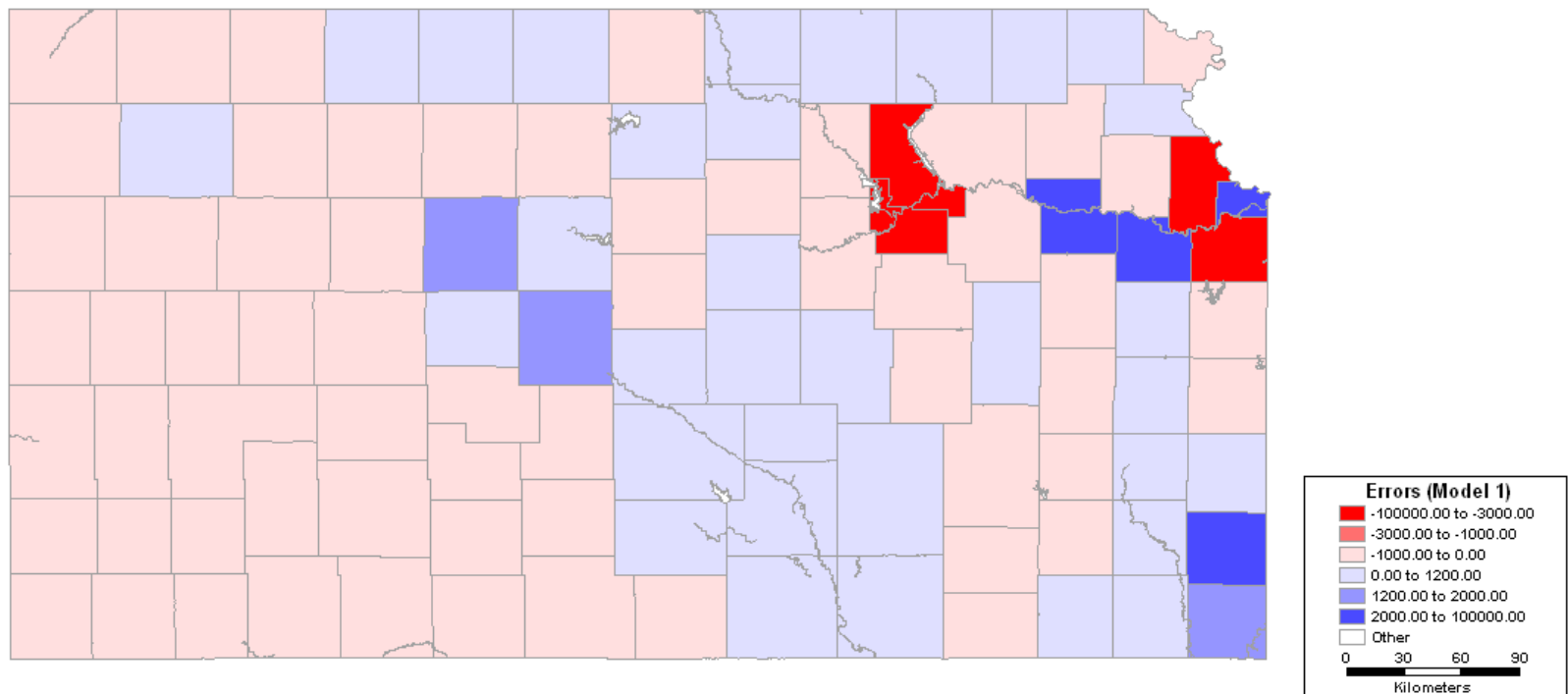Z ( I )   =                     -2.821

# Analysis of Residuals

○ Is variance constant?

- Scatterplot of errors vs. variables

# Analysis of Residuals

○ Are errors independent?

Produced by Academic TransCAD



Errors (Model 1)

| | |
|---|---|
| ▮ (red) | -100000.00 to -3000.00 |
| ▮ (light red) | -3000.00 to -1000.00 |
| ▮ (pink) | -1000.00 to 0.00 |
| ▮ (light blue) | 0.00 to 1200.00 |
| ▮ (blue) | 1200.00 to 2000.00 |
| ▮ (dark blue) | 2000.00 to 100000.00 |
| ▮ (white) | Other |

0   30   60   90
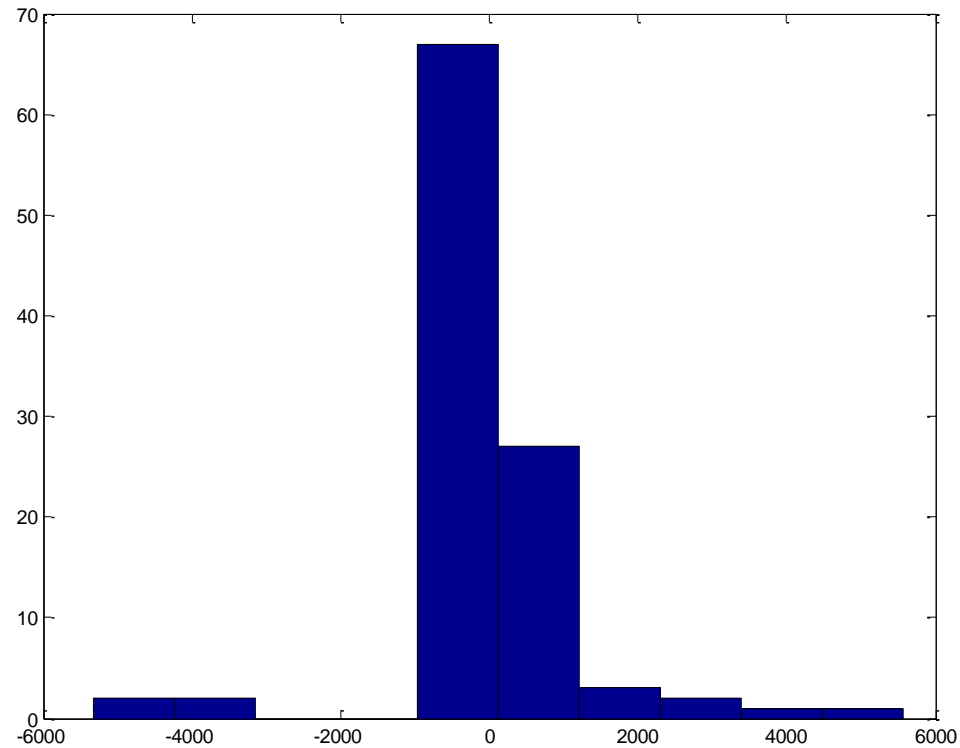Kilometers

# Analysis of Residuals

○ Are errors independent?
- Moran's I
- Expected (mean) value of $e$ is zero

# Analysis of Residuals

○ Are errors normally distributed?

- Histogram
- Probability plot
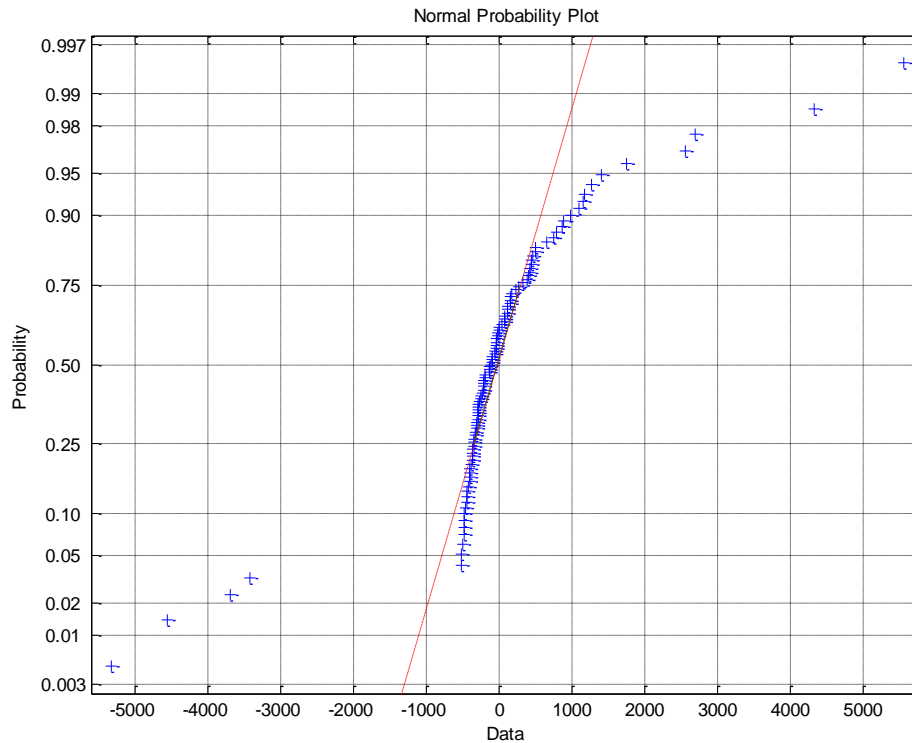- Other tests (Jarque-Bera, Kolmogorov-Smirnov,Lilliefors)

# Analysis of Residuals

○ Histogram

# Analysis of Residuals

○ Probability plot

# Modeling Area Data

- ○ What is the conclusion regarding this model?

# Modeling Area Data

○ Data transformations
- May help to reduce size effects
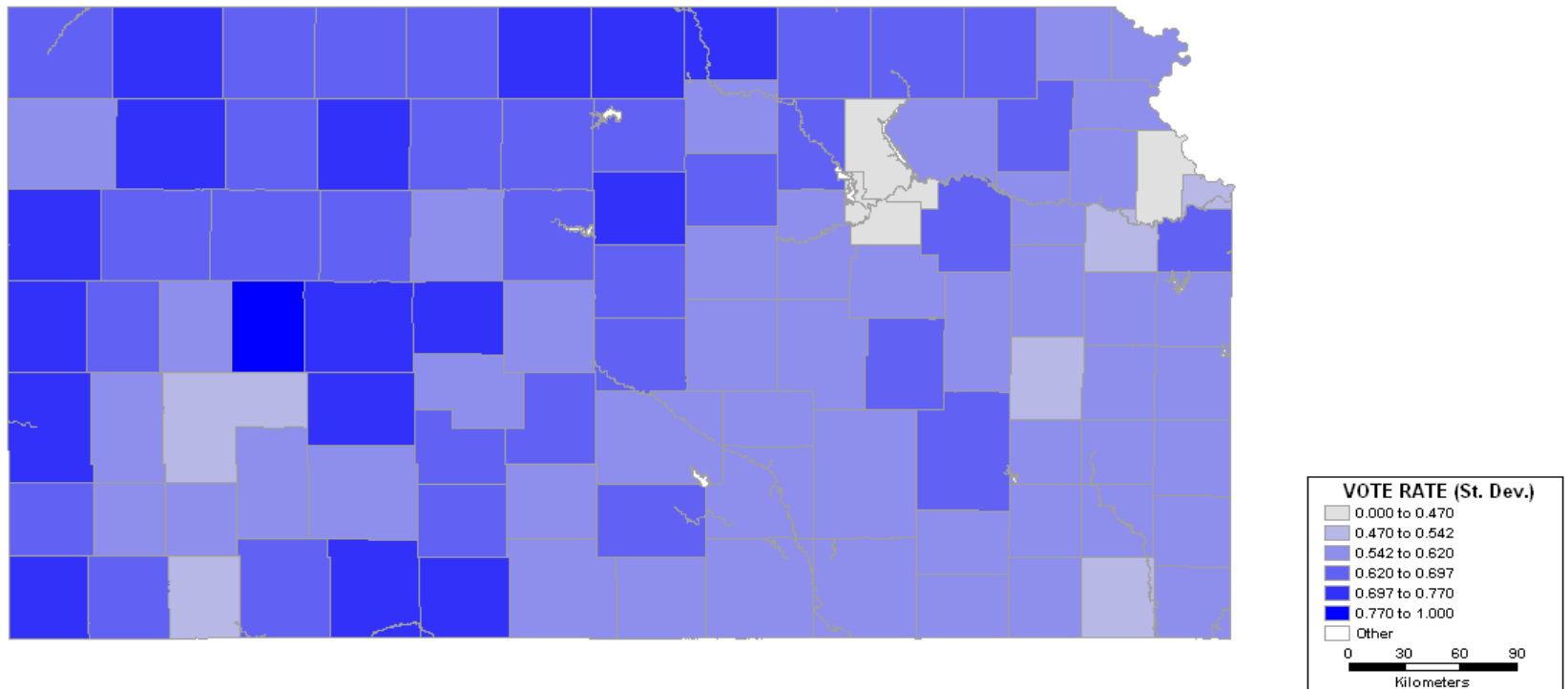- May increase the interpretability of the models

# Modeling Area Data

○ Data transformations

- Instead of modeling totals, model the proportions
- Proportional voter turnout: $VOTE/POP$

# Exploring the Proportions
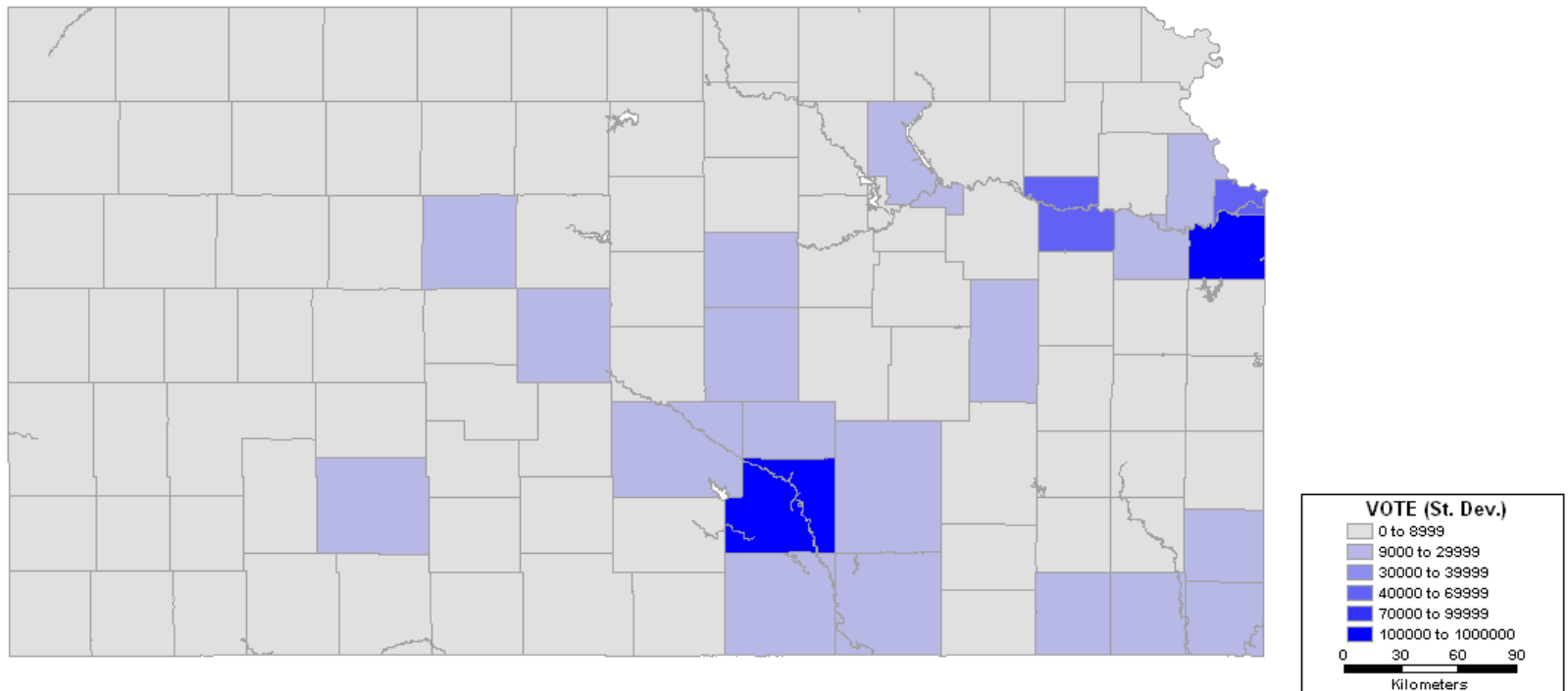
○ Choropleth map: PrVotes (Std. Dev.)



Produced by Academic TransCAD

**VOTE RATE (St. Dev.)**

| | |
|---|---|
| | 0.000 to 0.470 |
| | 0.470 to 0.542 |
| | 0.542 to 0.620 |
| | 0.620 to 0.697 |
| | 0.697 to 0.770 |
| | 0.770 to 1.000 |
| | Other |

0    30    60    90
Kilometers

# Exploring the Proportions

○ Choropleth map: Votes (Std. Dev.)
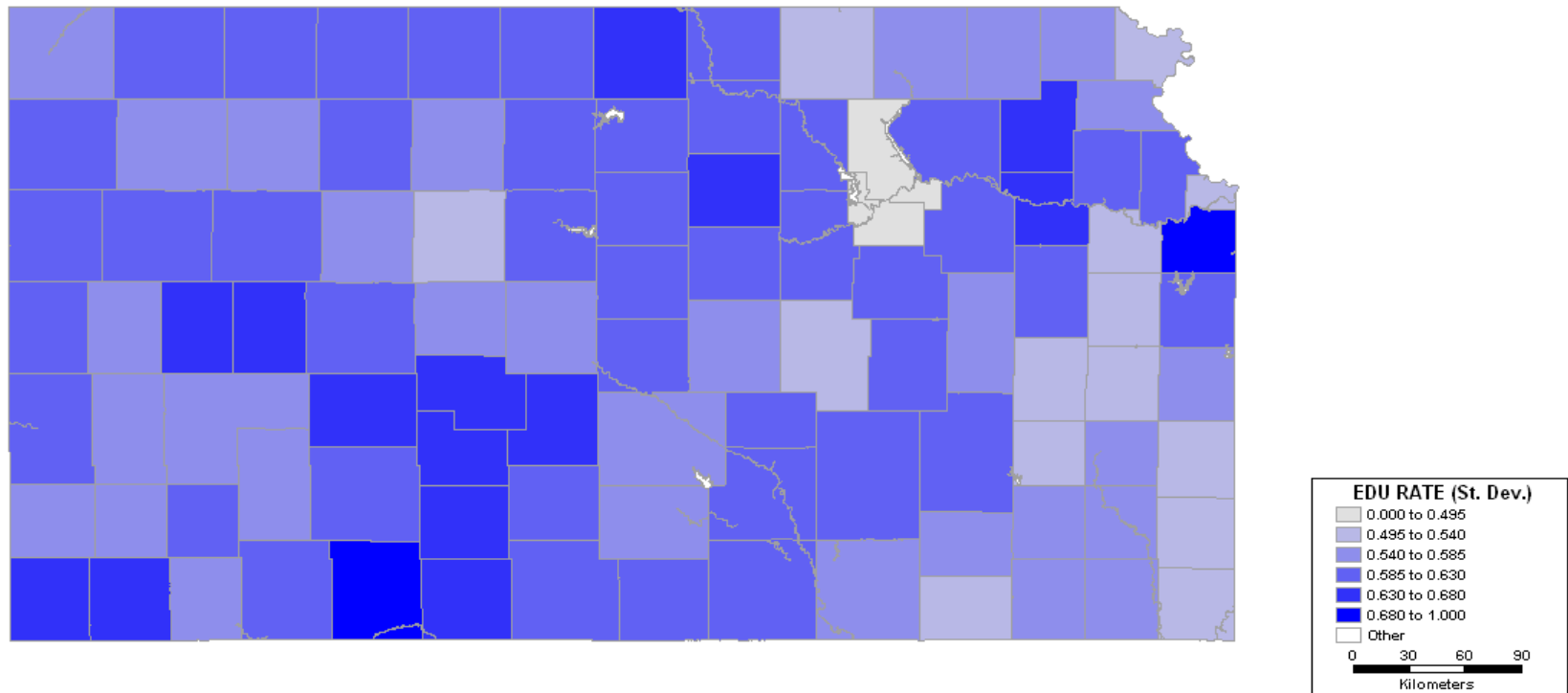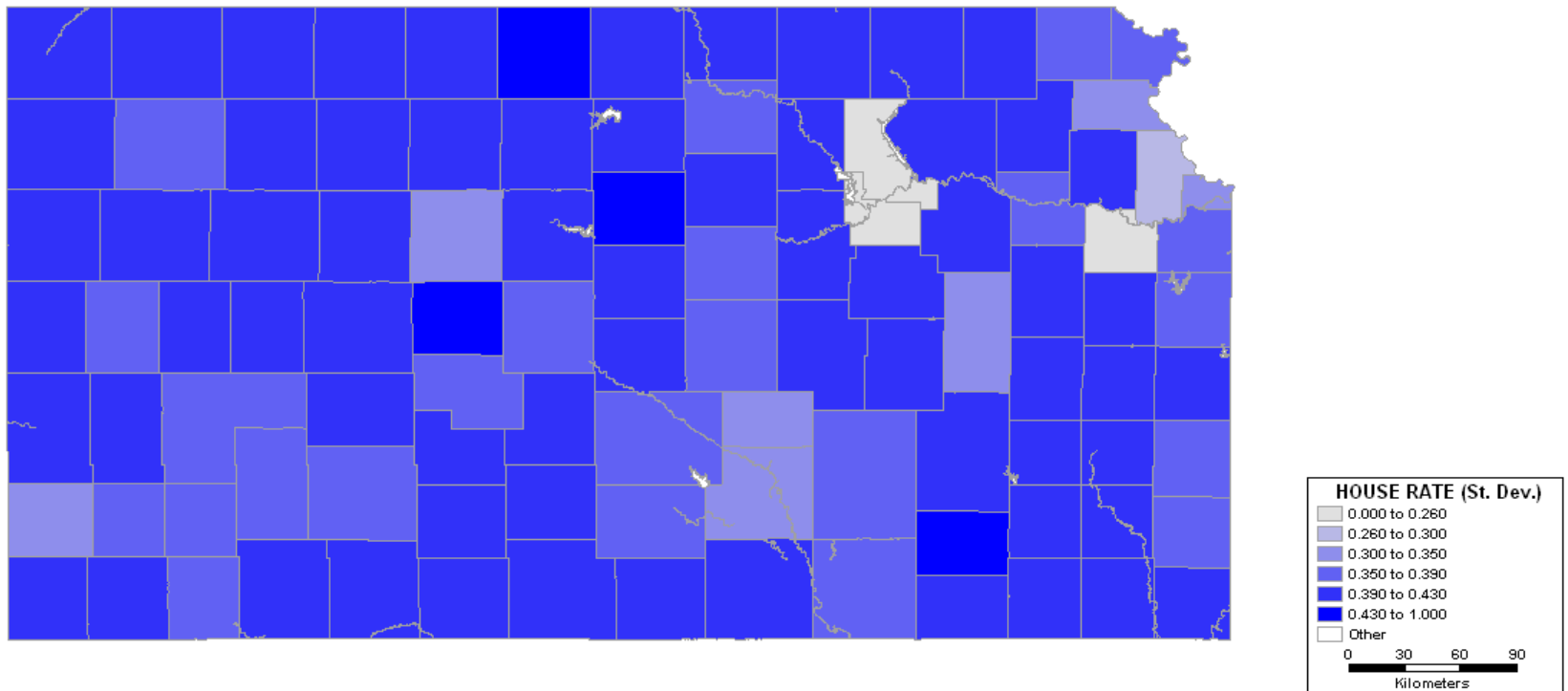
# Exploring the Proportions

○ PrEDU (Std. Dev.)



Produced by Academic TransCAD

# Exploring the Proportions

○ PrHOUSE (Std. Dev.)



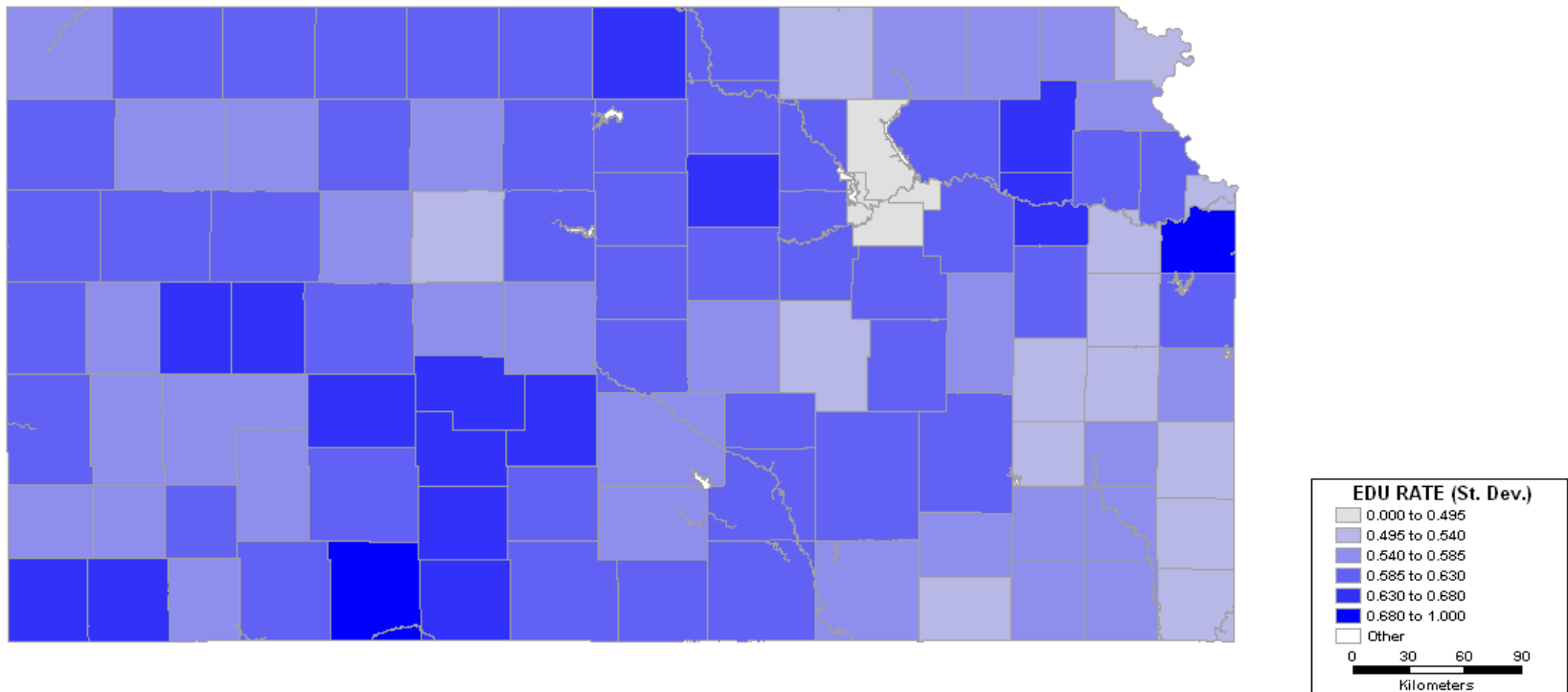Produced by Academic TransCAD

**HOUSE RATE (St. Dev.)**
- 0.000 to 0.260
- 0.260 to 0.300
- 0.300 to 0.350
- 0.350 to 0.390
- 0.390 to 0.430
- 0.430 to 1.000
- Other

0    30    60    90
Kilometers

# Exploring the Proportions

○ PrINCOME (Std. Dev.)



Produced by Academic TransCAD

EDU RATE (St. Dev.)
- 0.000 to 0.495
- 0.495 to 0.540
- 0.540 to 0.585
- 0.585 to 0.630
- 0.630 to 0.680
- 0.680 to 1.000
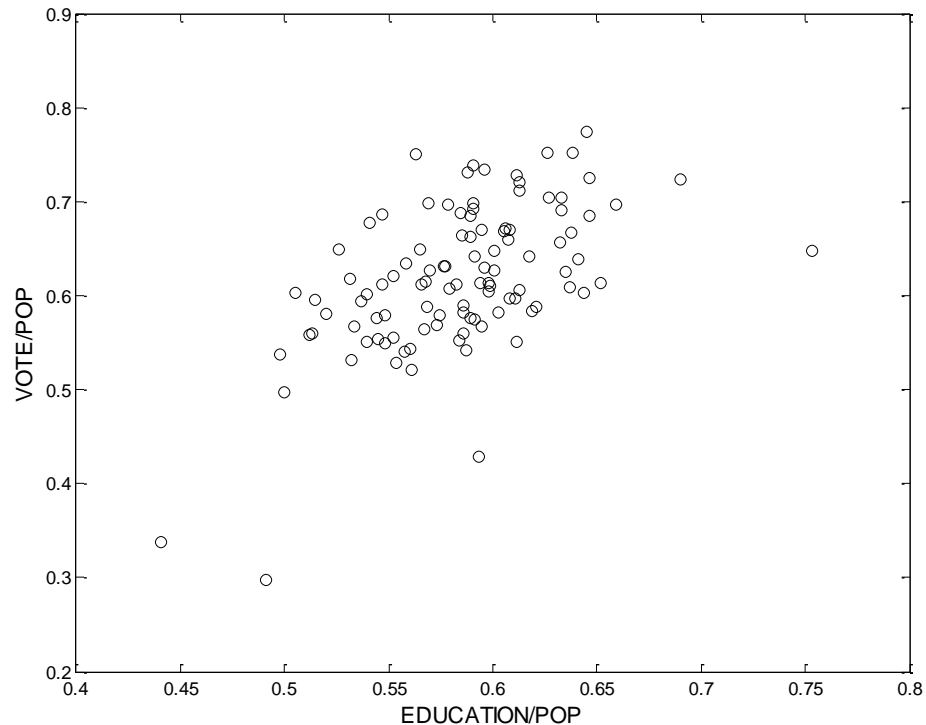- Other

0   30   60   90
Kilometers

# Example

- Bivariate analysis
  - Scatterplots
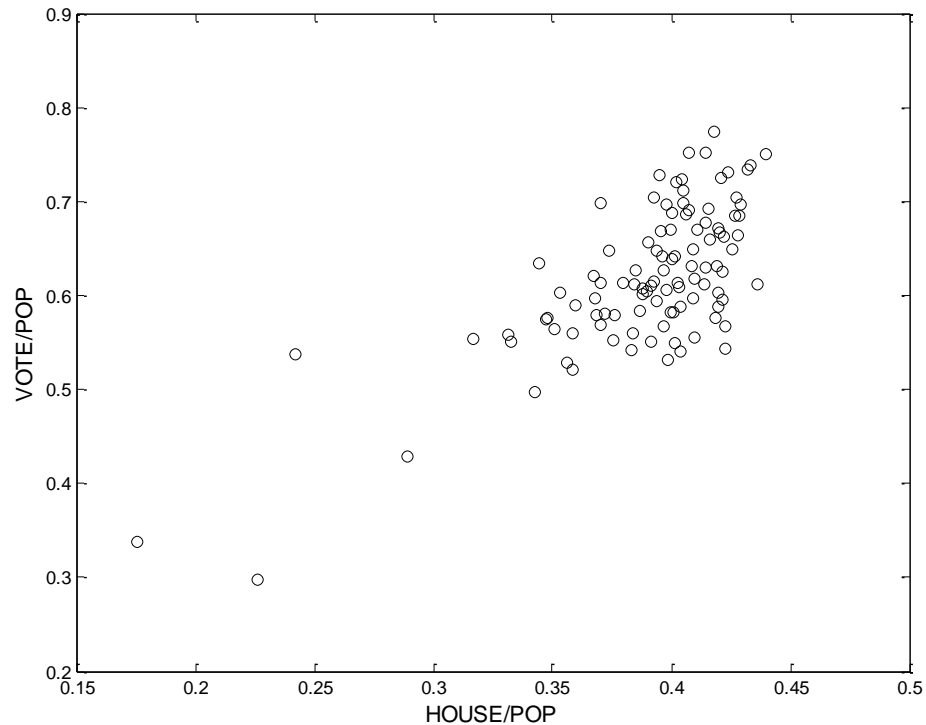  - Correlation coefficients

# Example

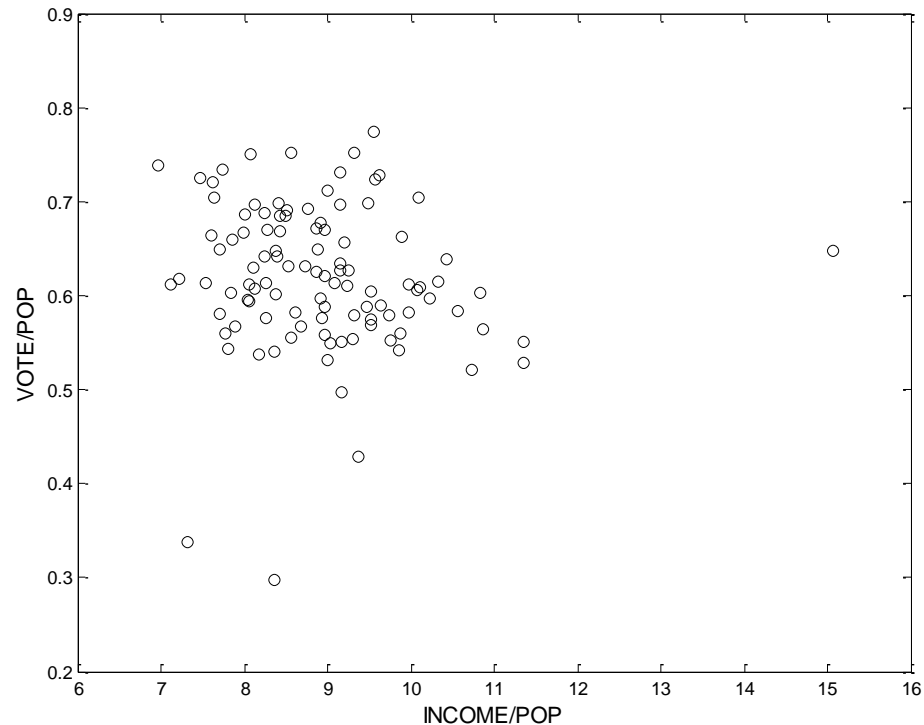○ Scatterplot (EDU/POP vs. VOTE/POP)

# Example

- Scatterplot (HOUSE/POP vs. VOTE/POP)

# Example

○ Scatterplot (INCOME/POP vs. VOTE/POP)

# Example

○ Bivariate analysis
  ● Correlation coefficients

|  | VOTE/POP | EDU/POP | HOUSE/POP | INCOME/POP |
|---|---|---|---|---|
| VOTE/POP | 1.00 |  |  |  |
| EDU/ POP | 0.563 | 1.00 |  |  |
| HOUSE/POP | 0.718 | 0.392 | 1.00 |  |
| INCOME/POP | -0.111 | 0.426 | -0.191 | 1.00 |

# Example

○ Model 2

<div align="center">Model 2 (VOTE/POP)</div>

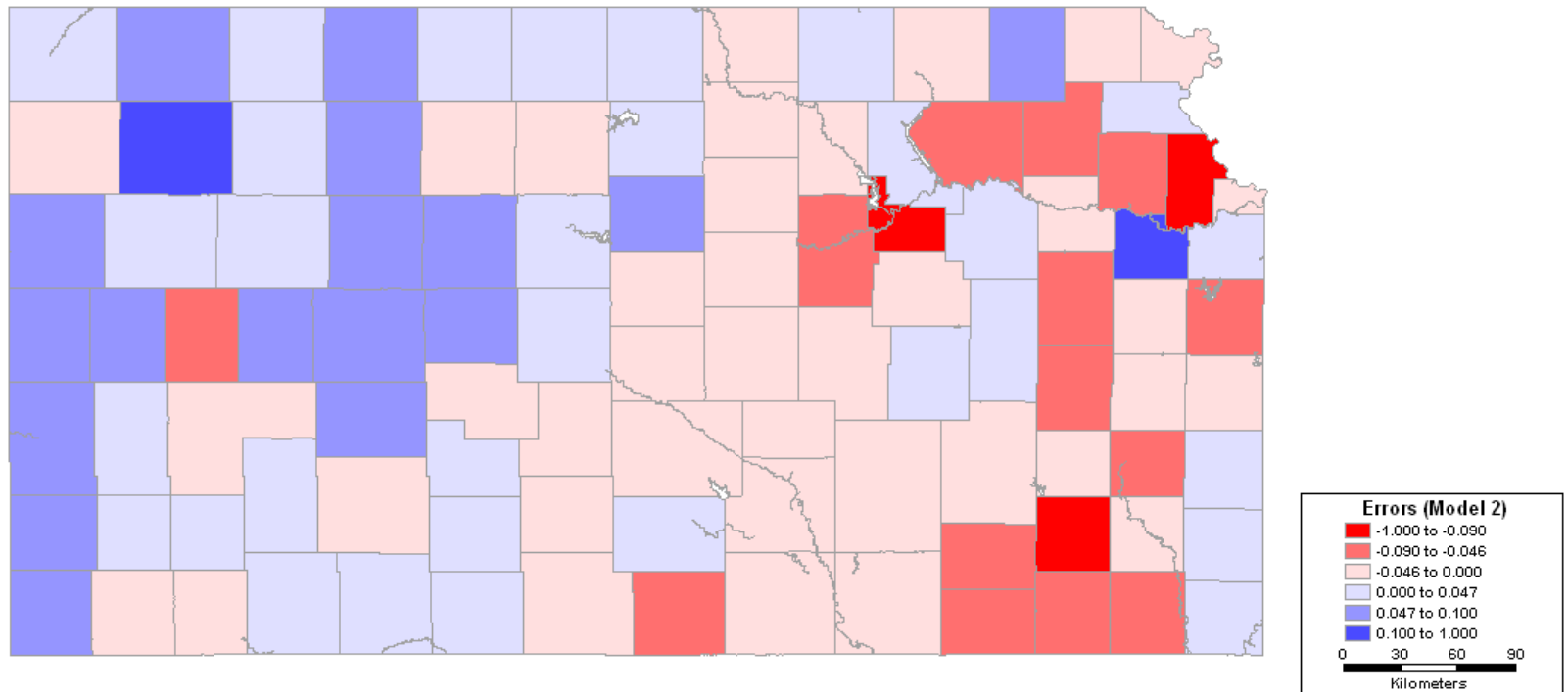| Variable | PARAMETER | t-value |
|---|---|---|
| CONST | -0.071 | -1.04 |
| EDU/POP | 0.782 | 5.76 |
| HOUSE/POP | 0.934 | 6.83 |
| INCOME/POP | -0.015 | -2.85 |
| R^2= | 0.636 | |
| R^2(adj)= | 0.611 | |
| SIGMA    = | 0.047 | |
| n= | 105 | |
| << Normalized Moran's I >> | | |
| Z ( I )  = | 4.023 | |

# Analysis of Residuals

○ Is variance constant?

- Scatterplot of errors vs. variables

# Analysis of Residuals

○ Are errors independent?

Produced by Academic TransCAD



**Errors (Model 2)**

| Color | Range |
|---|---|
| (red) | -1.000 to -0.090 |
| (light red) | -0.090 to -0.046 |
| (pink) | -0.046 to 0.000 |
| (light blue) | 0.000 to 0.047 |
| (blue) | 0.047 to 0.100 |
| (dark blue) | 0.100 to 1.000 |

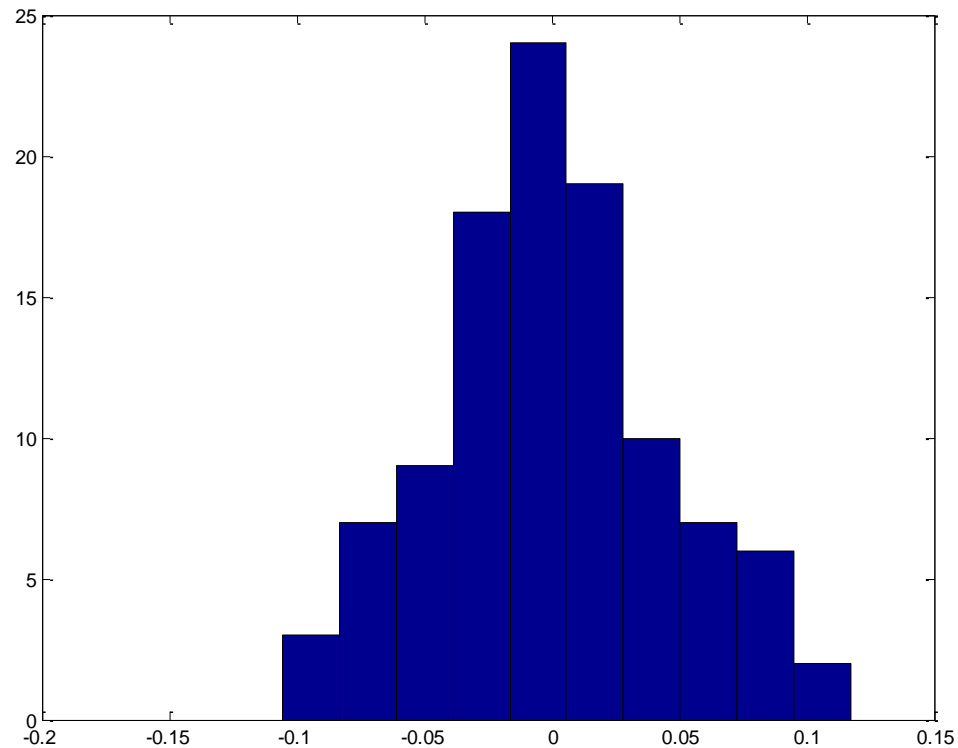0  30  60  90
Kilometers

# Analysis of Residuals

○ Are errors independent?

- Moran's I
- Expected (mean) value of $e$ is zero

# Analysis of Residuals

○ Are errors normally distributed?

- Histogram
- Probability plot
- Other tests (Jarque-Bera, Kolmogorov-Smirnov,Lilliefors)

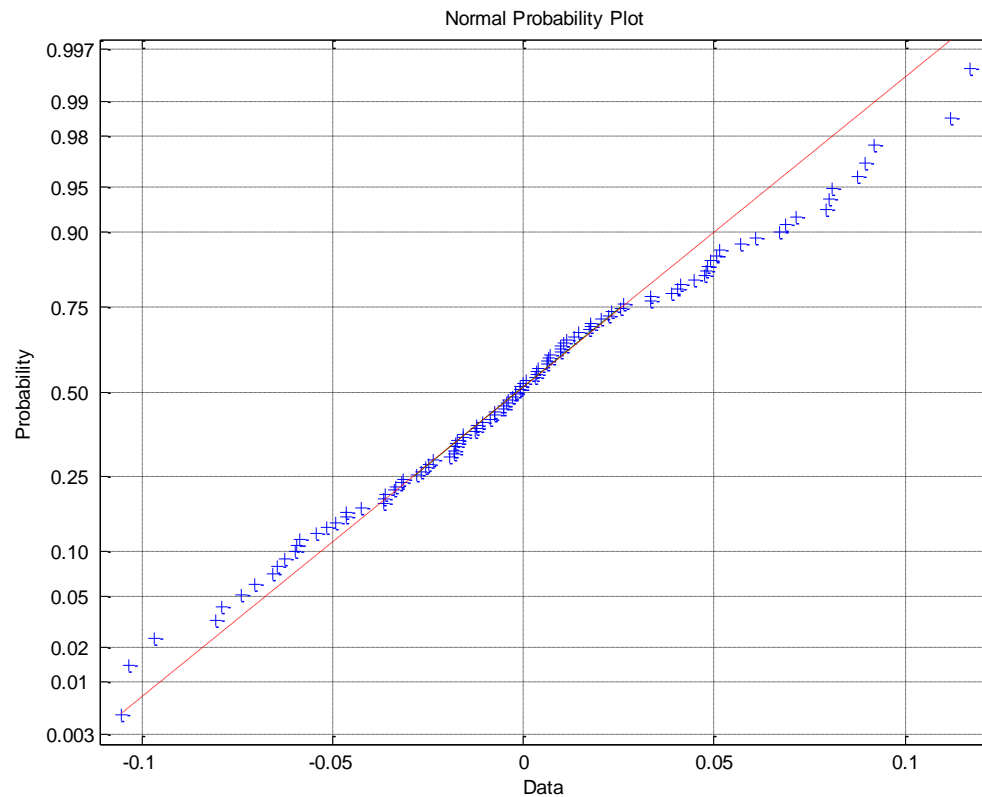# Analysis of Residuals

○ Histogram

# Analysis of Residuals

○ Probability plot

# Modeling Area Data

○ Data transformations
  ● Logarithm of proportional voter turnout:

$$\ln(PrVOTE)$$

  ● Correlation coefficients
  ● Model
  ● Analysis of residuals

# Variable Transformations

○ Taking logarithms:    $\ln(PrVOTE)$

Model 3 ln(VOTE/POP)

| Variable | PARAMETER | t-value |
|---|---|---|
| CONST | 0.926 | 4.70 |
| ln(EDU/POP) | 0.710 | 5.28 |
| ln(HOUSE/POP) | 0.614 | 8.21 |
| ln(INCOME/POP) | -0.205 | -2.57 |
| R^2= | 0.674 | |
| R^2(adj)= | 0.649 | |
| SIGMA    = | 0.081 | |
| n= | 105 | |
| << Normalized Moran's I >>' | | |
| Z ( I )   = | 3.138 | |

Normality? No
Constant variance? Yes

# Conclusion

- Model selection
  - Totals
  - Proportions
  - Logarithm of Proportions

# Conclusion

○ All three models still show residual error pattern (error autocorrelation)

# Non-spatial Regression

- Data transformations may help to reduce problems with the residuals + increase interpretability

- In the example: All three models show residual error pattern (error autocorrelation)

- Spatial regression models

# Spatial Regression Modeling

- Objective:
  - Testing for and estimating regression models that incorporate spatial effects
- Recall GLS

$$Y = \mathbf{X}\boldsymbol{\beta} + U$$

$$E[U] = \boldsymbol{0}$$

$$E[UU'] = \mathbf{C}$$

# Spatial Regression Modeling

○ Generalized Least Squares

$$Y = \mathbf{X}\boldsymbol{\beta} + U$$

- 2nd order effects?
- Questionable assumption of stationarity of 2nd order
- Distance measures?

# Spatial Regression Modeling

○ Autocorrelated errors model

$$Y = \mathbf{X}\boldsymbol{\beta} + U$$

Non-spatial part

$$U = \rho\mathbf{W}U + \varepsilon$$

Spatial part (residual pattern)

# Autocorrelated Errors Modeling

$$Y = \mathbf{X}\boldsymbol{\beta} + U$$

$$U = \rho\mathbf{W}U + \boldsymbol{\varepsilon}$$

$$E[\boldsymbol{\varepsilon}] = 0$$

$$E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'] = \sigma^2\mathbf{I}$$

- Rewriting:

$$Y = \mathbf{X}\boldsymbol{\beta} + \rho\mathbf{W}Y - \rho\mathbf{W}\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\mathbf{C} = \sigma^2\left(\left(\mathbf{I} - \rho\mathbf{W}\right)^T\left(\mathbf{I} - \rho\mathbf{W}\right)\right)^{-1}$$

# Estimation of the Spatial Model

○ OLS estimation not useful

○ Maximum likelihood estimation of $\beta$ and $\rho$

○ Assumptions:

$y_i$ are observations on $n$ random variables $Y$, jointly normally distributed with Mean = $\mathbf{X}\boldsymbol{\beta}$ and Covariance matrix $\mathbf{C}$, where

$$\mathbf{C} = \sigma^2 \left( \left( \mathbf{I} - \rho \mathbf{W} \right)^T \left( \mathbf{I} - \rho \mathbf{W} \right) \right)^{-1}$$

# Estimation of a Spatial Model

○ Maximum likelihood estimation of $\beta$ and $\rho$

○ What is the maximum likelihood function?

# Alternative Spatial Models

○ Lag models

$$Y = \mathbf{X}\boldsymbol{\beta} + U$$

$$U = \rho_1 \mathbf{W}^{(1)} U + \rho_1 \mathbf{W}^{(2)} U + \ldots + \varepsilon$$

# Alternative Spatial Models

○ Mixed regressive autoregressive

$$Y = \rho \mathbf{W} Y + \mathbf{X} \boldsymbol{\beta} + \varepsilon$$

Spatial part    Non-spatial part

# Alternative Spatial Models

○ Pure autoregressive

$$Y = \rho \mathbf{W} Y + \varepsilon$$

Spatial part

# Alternative Spatial Models

○ Moving average model (map pattern)

$$Y = (\mathbf{I} - \rho\mathbf{W})\varepsilon$$

Spatial part

# Spatial Regression

○ Fit autocorrelated errors model

### Model 4 VOTE/POP

| Variable | PARAMETER | t-value |
|---|---|---|
| CONST | 0.032 | 0.52 |
| EDU/POP | 0.645 | 4.80 |
| HOUSE/POP | 0.898 | 6.92 |
| INCOME/POP | -0.015 | -3.09 |
| RHO | 0.468 | 4.14 |
| SIGMA    = | 0.002 | |
| n= | 105 | |

<< Lagrange Multiplier >>
Omitted Sp. Lag=          1.038 --> Chi2 (1 DF)
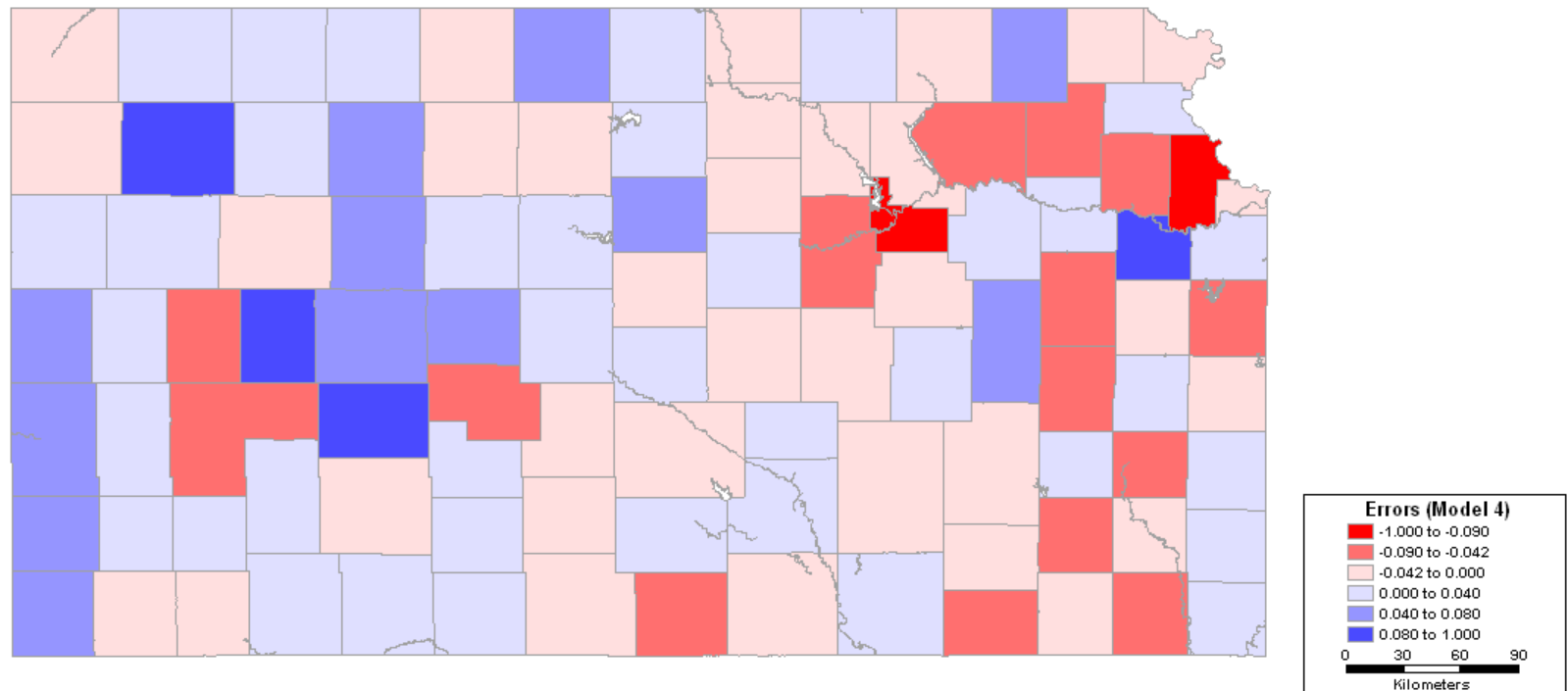
# Example

- Model 2

<div align="center">

Model 2 (VOTE/POP)

| Variable | PARAMETER | t-value |
|---|---|---|
| CONST | -0.071 | -1.04 |
| EDU/POP | 0.782 | 5.76 |
| HOUSE/POP | 0.934 | 6.83 |
| INCOME/POP | -0.015 | -2.85 |
| $R^2$= | 0.636 | |
| $R^2$(adj)= | 0.611 | |
| SIGMA     = | 0.047 | |
| n= | 105 | |
| << Normalized Moran's I >> | | |
| Z ( I )   = | 4.023 | |

</div>

# Analysis of Residuals

○ Are errors independent? Model 4

Produced by Academic TransCAD



**Errors (Model 4)**
- -1.000 to -0.090
- -0.090 to -0.042
- -0.042 to 0.000
- 0.000 to 0.040
- 0.040 to 0.080
- 0.080 to 1.000

0　30　60　90
Kilometers

# Analysis of Residuals

- Are errors independent? Model 2



Produced by Academic TransCAD

Errors (Model 2)
- -1.000 to -0.090
- -0.090 to -0.046
- -0.046 to 0.000
- 0.000 to 0.047
- 0.047 to 0.100
- 0.100 to 1.000

0  30  60  90
Kilometers

# Analysis of Residuals

○ Is variance constant?
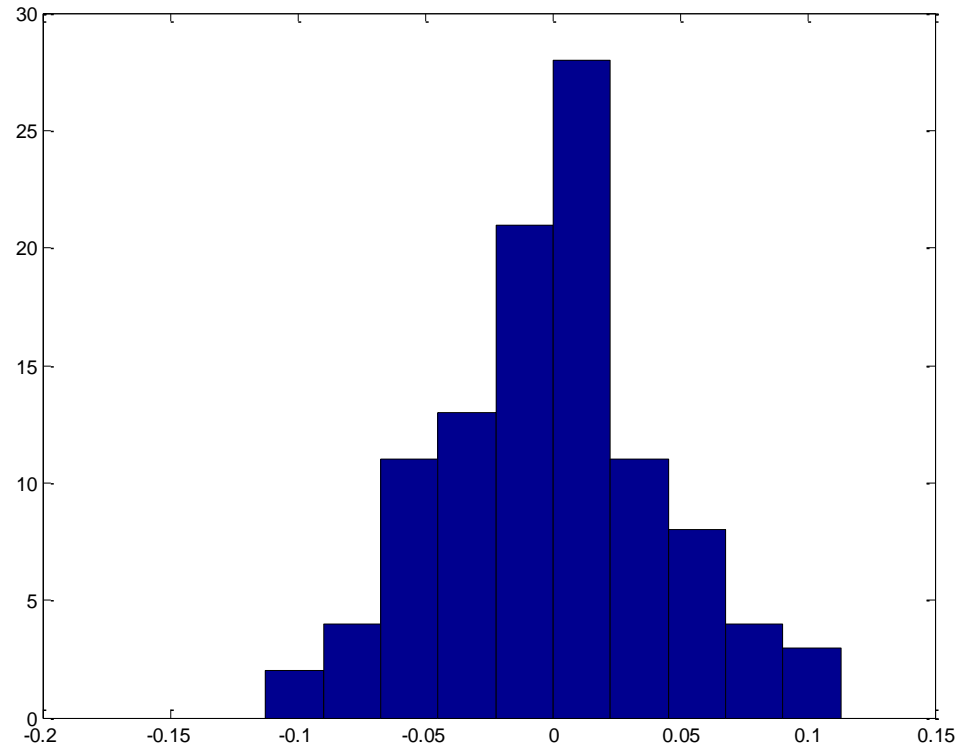- Scatterplot of errors vs. variables

# Analysis of Residuals

○ Are errors normally distributed?

- Histogram
- Probability plot
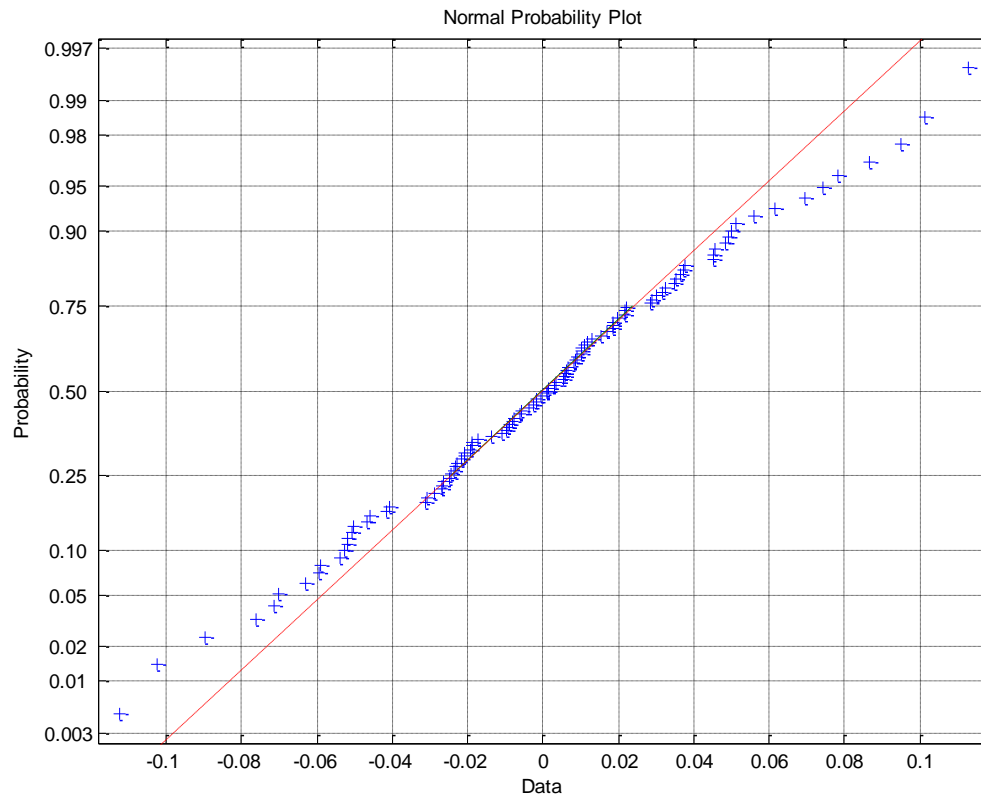- Other tests (Jaques-Bera, Kolmogorov-Smirnov,Lilliefors)

# Analysis of Residuals

○ Histogram

# Analysis of Residuals

○ Probability plot

# Analysis of Residuals

○ Other tests

| Test | Normal? |
|------|---------|
| Jarque-Bera | YES |
| Kolmogorov-Smirnov | YES |
| Lilliefors | YES |

# Conclusion

- The autocorrelated errors model is the best alternative
  - Satisfies all major assumptions
  - Accounts for all systematic spatial variation

# Summary of modeling approach

- Exploration
- Models – non-spatial models, check assumptions
- Models – spatial models, check assumptions and statistical fit
- Think Theory!