

Geographical Information Systems Institute
Center for Geographic Analysis, Harvard University

Spatial Regression with GeoDa

This lab includes discussion of two types of models of spatial dependence

- 1) Spatial lag
- 2) Spatial error

It then shows how to estimate simple spatial regression models in GeoDa

- 1) Classical OLS regression with diagnostics
- 2) Spatial lag model vs. Spatial error model

Data and Shapefiles:

We will use a US 2000 Census data for Manhattan in New York City. The zipped data is in the GeoDa folder you downloaded earlier this week.

For this lab we will use the shapefile **newyork.shp**, it includes the borough of Manhattan in New York City with Census 2000 data from summary file 3. These are socioeconomic attributes for 297 Census tracts. It includes the following variables:

Vraiable name	Label
POLYID	Polygon ID
STATE	State FIPS
COUNTY	County FIPS
TRACT	Census Tract ID
sctrct00	FIPSID
hvalue	Median housing value
t0_pop	Total population
pctnhw	Percent non-Hispanic white persons
pctnhb	Percent non-Hispanic black persons
pcthsp	Percent Hispanic persons
pctasn	Percent Asian persons
t0p_own	Percent homeowners
t0p_coll	Percent college educated
t0p_prf	Percent of people employed in professional/managerial occupations
t0p_uemp	Percent of people unemployed
t0p_for	Percent foreign born persons
t0p_rec	Percent recent immigrants
t0_minc	Median household income
t0p_poor	Percent total population below poverty

1. Standard linear regression - *Ordinary Least Squares (OLS)*

The general purpose of linear regression analysis is to find a (linear) relationship between a dependent variable and a set of explanatory variables:

$$y = X\beta + \varepsilon$$

The method of ordinary least squares (OLS) estimation is referred to as a Best Linear Unbiased Estimator (BLUE). The OLS estimates β by minimizing the sum of squared prediction errors, hence, least squares. In order to obtain the BLUE property and make statistical inferences about the population regression coefficients from the estimated b , certain assumptions about the random error of the regression equation need to be made. These include:

- a) the random errors have a mean of zero (there is no systematic misspecification or bias in the population regression equation);
- b) the random errors have a constant variance (homoskedasticity) and are uncorrelated;
- c) the random errors have a normal distribution.

2. Spatial dependence

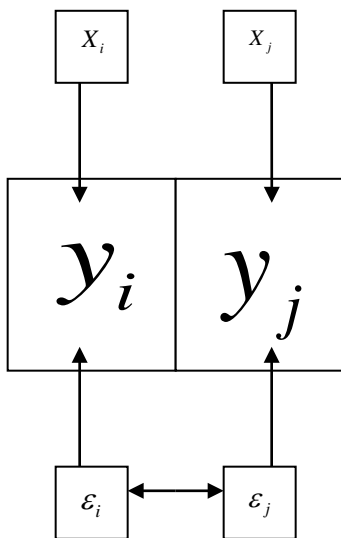
These assumptions may not be always satisfied in practice. When a value observed in one location depends on the values observed at neighboring locations, there is a spatial dependence. And spatial data may show spatial dependence in the variables and error terms.

Why should spatial dependence occur? There are two reasons commonly given. First, data collection of observations associated with spatial units may reflect measurement error. This happens when the boundaries for which information is collected do not accurately reflect the nature of the underlying process generating the sample data.

A second reason for spatial dependence is that the spatial dimension of a social or economic characteristic may be an important aspect of the phenomenon. For example, based on the premise that location and distance are important forces at work, regional science theory relies on notions of spatial interaction and diffusion effects, hierarchies of place and spatial spillovers.

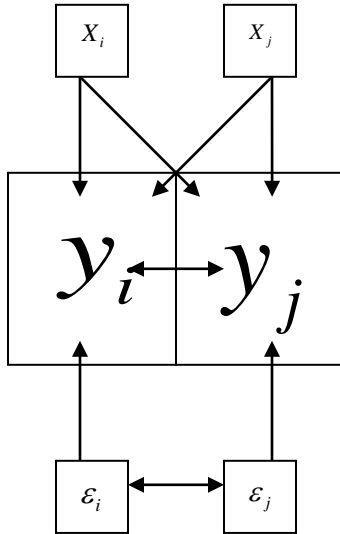
There are two primary types of spatial dependence:

- a. Spatial error -the error terms across different spatial units are correlated



With spatial error in OLS regression, the assumption of uncorrelated error terms is violated. As a result, the estimates are inefficient. Spatial error is indicative of omitted (spatially correlated) covariates that if left unattended would affect inference.

b. Spatial lag- the dependent variable y in place i is affected by the independent variables in both place i and j .



With spatial lag in OLS regression, the assumption of uncorrelated error terms is violated; in addition, the assumption of independent observations is also violated. As a result, the estimates are biased and inefficient. Spatial lag is suggestive of a possible diffusion process – events in one place predict an increased likelihood of similar events in neighboring places.

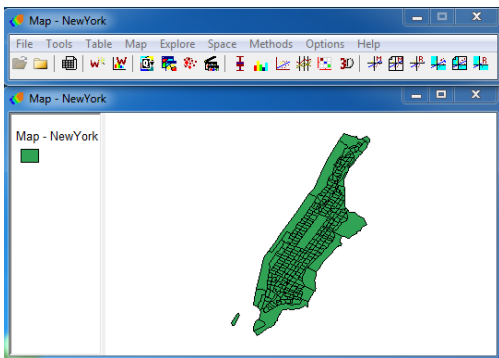
GeoDa provides a range of diagnostics to detect spatial dependence. It also provides unbiased regression estimates using a Maximum Likelihood approach (ML Spatial Lag or Spatial Error models).

3. Spatial regression in GeoDa

Start a project and create weights matrix

a. Start GeoDa by,

b. Go to **File > New Project**, click on the **folder** icon choose **ESRI shapefile**, **navigate** to **newyork** and click **Connect**.



c. Create a weights matrix. Go to **Tools > Weights > Create** to open the

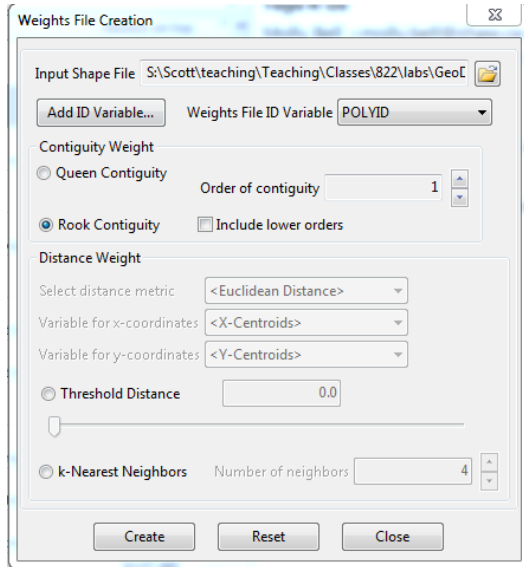
Creating Weights dialogue box.

In the **Creating weights dialogue box**:

d. Select **POLYID** as the “Weights File ID Variable” (by scrolling and clicking on variable name). The ID Variable must have a unique value for each observation (i.e., in this case Census tract). The unique value is used to implement the link between maps and statistical graphs.

Select **newyork.shp** as the input, type “**rook**” and click **Create** (**in the save window** use a name that means something to you, **NewYorkROOK**, and use the default extension, **.gal**).

Click **Close**.

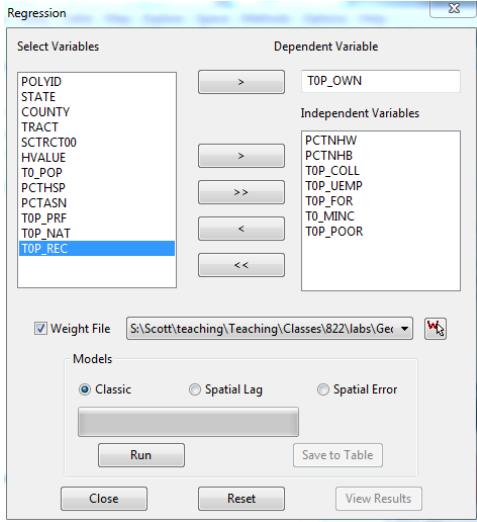


3.1. Classical OLS regression with diagnostics.

Now, we are ready to perform an OLS regression and evaluate the spatial dependence in this regression.

a. On the menu bar, choose **Methods**, then **Regression**. A dialogue box will appear:

b. A variable selection box will appear:



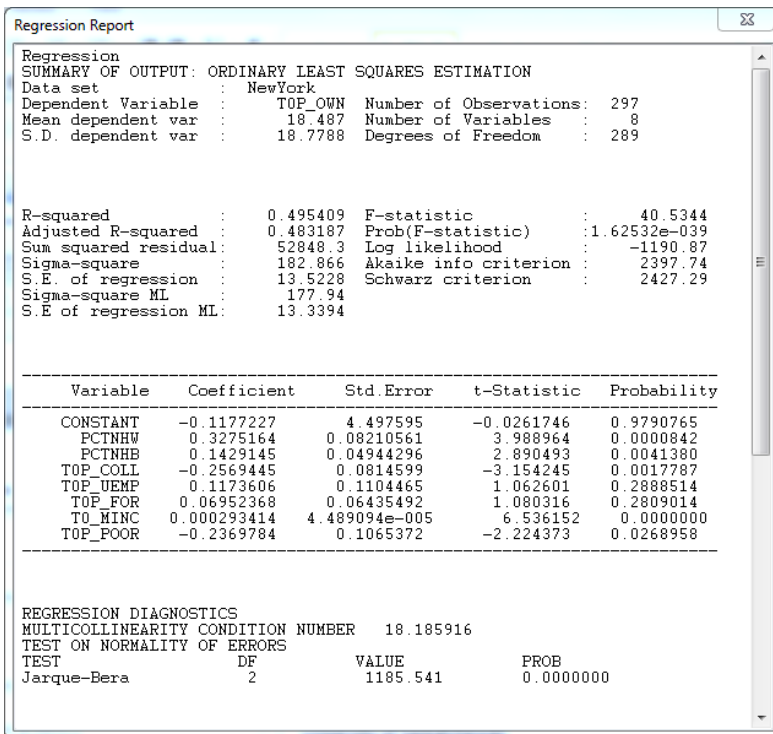
In this example, we will predict neighborhood homeownership with several indicators, including % non-Hispanic white, % non-Hispanic black, % college+ educated, unemployment rate, % foreign born, median household income, and % person below poverty.

Browse to locate the weights matrix file you just created and check the **Weight Files** icon to the right of the file name window.

c. Check **Classic**- this will run classical OLS regression with spatial dependence diagnostics, click **Run**. When done, click **Run**.

d. A new window of regression output will appear, and it has several sections.

1) The first section is the summary output of OLS regression:



It first shows general information of the run, including the mean and standard deviation of the dependent variable, the model coefficient of determination, F-test probability, and Log likelihood. Then, the coefficients, standard errors, and significance are shown. We can see that among the seven indicators, % non-Hispanic white, % non-Hispanic black, % college+ educated, and median household income are positively related to homeownership; while poverty rate is negatively related to homeownership; and unemployment rate and % foreign born have insignificant effects.

2) The next section deals with regression diagnostics:

```
REGRESSION DIAGNOSTICS
MULTICOLLINEARITY CONDITION NUMBER    18.185916
TEST ON NORMALITY OF ERRORS
TEST      DF      VALUE      PROB
Jarque-Bera      2      1185.541    0.0000000
```

GeoDa tests multicollinearity of the model- one should be alarmed when **MULTICOLLINEARITY CONDITION NUMBER** is greater than 20. The Jarque-Bera test is used to examine the normality of the distribution of the errors. This test is a test of the combined effects of both skewness and Kurtosis. The low probability of the test score indicates non-normal distribution of the error term. Since the following tests of variance and spatial dependence are conditioned upon normal distribution, in real research, one should be very cautious to interpret the test results. Here, we simply give the illustration of how to interpret the diagnostics when non-normality is not encountered.

3) **DIAGNOSTICS FOR HETEROSKEDASTICITY-** a test of the variance of the error term as the BLUE requires constant error variance.

```
DIAGNOSTICS FOR HETEROSKEDASTICITY
RANDOM COEFFICIENTS
TEST      DF      VALUE      PROB
Breusch-Pagan test      7      102.959    0.0000000
Koenker-Bassett test     7      18.45908   0.0100618
SPECIFICATION ROBUST TEST
TEST      DF      VALUE      PROB
White      35      185.7326   0.0000000
```

The low probabilities of the three tests point to existence of heteroskedasticity. This is not necessarily a surprise because the error variance could well be affected by the spatial dependence in the data.

4) **DIAGNOSTICS FOR SPATIAL DEPENDENCE**

```
DIAGNOSTICS FOR SPATIAL DEPENDENCE
FOR WEIGHT MATRIX : NewYorkROOK.gal
(row-standardized weights)
TEST      MI/DF      VALUE      PROB
Moran's I (error)      0.196095    5.7432856   0.0000000
Lagrange Multiplier (lag)      1      21.5140684   0.0000035
Robust LM (lag)      1      0.1141221    0.7354991
Lagrange Multiplier (error)      1      27.8417603   0.0000001
Robust LM (error)      1      6.4418140    0.0111465
Lagrange Multiplier (SARMA)      2      27.9558824   0.0000009
===== END OF REPORT =====
```

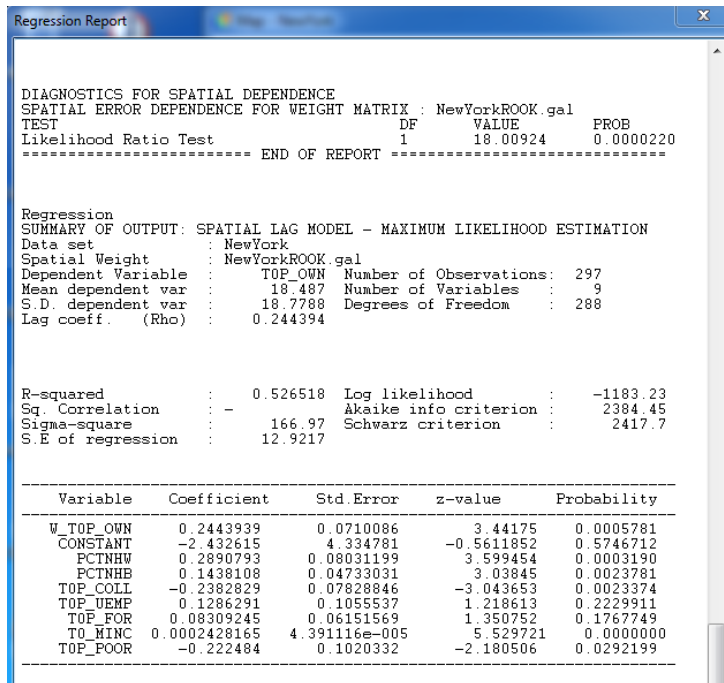
There are six tests performed to assess the spatial dependence of the model. First, Moran's I score of 0.196 is highly significant, indicating strong spatial autocorrelation of the residuals. In addition, the function reports the estimates of tests chosen among five statistics for testing for spatial dependence in linear models. The statistics are the simple LM test for a missing spatially lagged dependent variable (Lagrange Multiplier (lag)), the simple LM test for error dependence (Lagrange Multiplier (error)), variants of these robust to the presence of the other (Robust LM (lag) and Robust LM (error)- which tests for error dependence in the possible presence of a missing lagged dependent variable, Robust LM (lag) is the other way round), and a portmanteau test (SARMA, in fact Lagrange Multiplier (error) + Robust LM (lag)).

We can see both simple tests of the lag and error are significant, indicating presence of spatial dependence. The robust tests help us understand what type of spatial dependence may be at work. The robust measure for error is still significant, but the robust lag test becomes insignificant, which means that when a lag dependent variable is present the error dependence disappears.

3.2. Spatial lag model vs. spatial error model

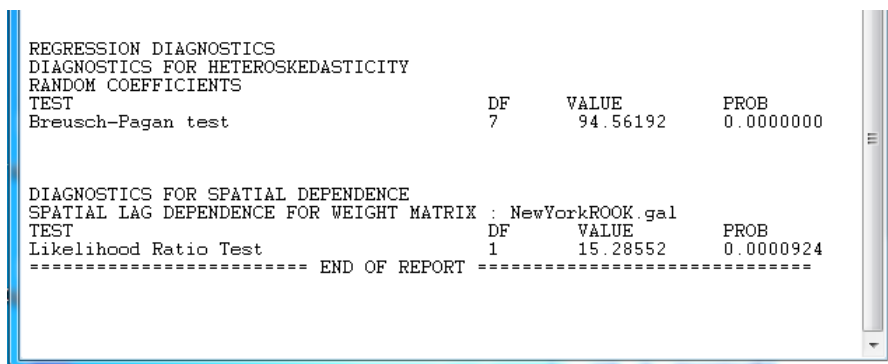
After identifying the presence of spatial dependence, we will use GeoDa to re-estimate the model with maximum likelihood approach while controlling for the spatial dependence.

- a. On the menu bar, again choose **Regress**. Check **Moran's I z-value** in the output box, and click **OK**.
 - b. In the variable selection box, set dependent and independent variables as before. Check **Weight Files**, and browse to locate the weights matrix file.
 - c. Check **Spatial Lag** in the **Models** selection, and click **Run**.
 - d. Check **Spatial Error** and click **Run**. Then click **View Results**.
 - e. A new window of regression output will appear presenting results from each of the three models (Classic, Lag, and Error)
- 1) First, we see the results from the already run OLS model. This is followed by the Lag and Error models; each model consists of multiple sections and begins with general information and regression coefficients with significance tests.



Notice, besides the information that appeared in previous OLS regression output spatial weights file is specified below the data file: rook.GAL. For the lag model we also have a new variable for the spatial lag term of homeownership W_TOP_OWN. (Its coefficient parameter (Rho) reflects the spatial dependence inherent in our sample data, measuring the average influence on observations by their neighboring observations.) It has a positive effect and it is highly significant. As a result, the general model fit improved, as indicated in higher values of R-squared, Log likelihood, AIC (Akaike info Criterion). You should note the absence of adjusted R-squared (which takes into account the degrees of freedom in the model); generally when comparing models of the same set of observations (the same census tracts, in this case) it is advised to use AIC as a measure of model fit, smaller AIC indicates a better fit. The effects of other independent variables remain virtually the same.

2) The following sections show tests of Heteroskedasticity and Spatial dependence.



We can see that the low probability in the Breusch-Pagan test suggests that there is still Heteroskedasticity in the model after introducing the spatial lag term. And in the Likelihood Ratio Test

of Spatial Lag Dependence, the result is still significant. Therefore, we conclude that although the introduction of spatial lag term improved the model fit, it didn't make the spatial effects go away.

3) Now let's review results for the Spatial Error model.

Regression Report

Regression
SUMMARY OF OUTPUT: SPATIAL ERROR MODEL - MAXIMUM LIKELIHOOD ESTIMATION

Data set : NewYork
Spatial Weight : NewYorkROOK.gal
Dependent Variable : TOP_OWN Number of Observations: 297
Mean dependent var : 18.487003 Number of Variables : 8
S.D. dependent var : 18.778784 Degrees of Freedom : 289
Lag coeff. (Lambda) : 0.299370

R-squared : 0.533820 R-squared (BUSE) : -
Sq. Correlation : - Log likelihood : -1181.865225
Sigma-square : 164.395 Akaike info criterion : 2379.73
S.E of regression : 12.8217 Schwarz criterion : 2409.28

Variable	Coefficient	Std. Error	z-value	Probability
CONSTANT	-0.5624482	4.470884	-0.1258025	0.8998882
PCTNHV	0.3534593	0.07917175	4.464462	0.0000080
PCTNHB	0.1685745	0.05518363	3.054791	0.0022523
TOP_COLL	-0.2548426	0.08248224	-3.089666	0.0020040
TOP_UEMP	0.1141583	0.1054082	1.083012	0.2788034
TOP_FOR	0.08322937	0.06332577	1.314305	0.1887438
TO_MINC	0.0002531697	4.59578e-005	5.508744	0.0000000
TOP_POOR	-0.2535516	0.1014969	-2.49812	0.0124854
LAMBDA	0.2993704	0.08044892	3.721248	0.0001983

In comparison with the Spatial Lag model output, we also have a designated spatial weight file: rook.GAL. And a coefficient on the spatially correlated errors (LAMBDA) is added as an additional indicator. It has a positive effect and it is highly significant. As a result, the general model fit improved, as indicated in higher values of R-squared, Log likelihood, AIC (Akaike info Criterion). Like the lag model, the effects of other independent variables remain virtually the same.

4) Diagnostics of Heteroskedasticity and spatial independence

REGRESSION DIAGNOSTICS
DIAGNOSTICS FOR HETEROSKEDASTICITY
RANDOM COEFFICIENTS
TEST

TEST	DF	VALUE	PROB
Breusch-Pagan test	7	102.8063	0.0000000

DIAGNOSTICS FOR SPATIAL DEPENDENCE
SPATIAL ERROR DEPENDENCE FOR WEIGHT MATRIX : NewYorkROOK.gal
TEST

TEST	DF	VALUE	PROB
Likelihood Ratio Test	1	18.00924	0.0000220

===== END OF REPORT =====

Similar to the lag model, the Heteroskedasticity test remains significant. Also, the Likelihood Ratio Test of Spatial Error Dependence has a significant result. Therefore, we conclude that although allowing the error terms to be spatially correlated improved the model fit, it didn't make the spatial effects go away.

5) What should we conclude?

Comparing the spatial lag and spatial error models, we can see both alternative models yield improvement to the original OLS model. Therefore, we should conclude that controlling spatial dependence will improve our model performance. Now the question is which of the two models is better? To some extent, this is an open question. The general advice is first to look for a theoretical basis to inform your choice. If there are strong substantive grounds for one model instead of the other, you should adopt it. When it is not so clear theoretically, you can compare the model performance parameters: R-squared, Log likelihood, AIC (Akaike info Criterion). In our case, the spatial error model has greater R-squared and Log likelihood values and a smaller AIC. That provides a statistical basis to adopt this solution.

4. Another example

Now, do another exercise with % person below poverty as dependent variable, and % non-Hispanic white, % non-Hispanic black, % college+ educated, unemployment rate, % foreign born, and median household income as independent variables. Run OLS regression, interpret the regression and spatial dependence diagnostics; if spatial dependence exists, try to re-estimate the models with Spatial lag model and Spatial Error model, then compare the output and make your choice.

In this session, we learned to use GeoDa to examine and control spatial dependence in the regression models. By now, you should know the general strategy applied in GeoDa to deal with the spatial dependence issue. You should be able to understand the various diagnostics and regression tests, and able to interpret the results and make decisions accordingly.