

Exploratory spatial data analysis using Stata

Maurizio Pisati

Department of Sociology and Social Research
University of Milano-Bicocca (Italy)
`maurizio.pisati@unimib.it`

2012 German Stata Users Group meeting
WZB Social Science Research Center, Berlin
June 1, 2012

Outline

- 1 Introduction
- 2 Spatial data
- 3 Visualizing spatial data
 - Overview
 - Dot maps
 - Proportional symbol maps
 - Diagram maps
 - Choropleth maps
 - Multivariate maps

Outline

- ④ Exploring spatial point patterns
 - Overview
 - Kernel density estimation

- ⑤ Detecting spatial autocorrelation
 - Overview
 - Spatial weights matrices
 - Measuring spatial autocorrelation
 - Global indices of spatial autocorrelation
 - Local indices of spatial autocorrelation

- ⑥ References

INTRODUCTION

Exploratory spatial data analysis

- **Exploratory spatial data analysis** (ESDA) is the extension of exploratory data analysis (EDA) to the problem of detecting patterns in spatial data (Haining *et al.* 1998: 457)
- ESDA involves seeking good descriptions of spatial data, so as to help the analyst to develop hypotheses and models for such data (Bailey and Gatrell 1995: 23)
- ESDA emphasizes graphical views of the data, designed to highlight meaningful clusters of observations, unusual observations, or relationships between variables. These views often take the form of maps

ESDA in Stata

- Stata users can perform ESDA using a variety of user-written commands published in the *Stata Technical Bulletin*, the *Stata Journal*, or the SSC Archive
- In this talk, I will briefly illustrate the use of six such commands: `spmap`, `spgrid`, `spkde`, `spatwmat`, `spatgsa`, and `spatlsa`

ESDA in Stata

- `spmap` is a general command aimed at visualizing several kinds of spatial data
- `spgrid` generates two-dimensional grids covering rectangular or irregular study areas
- `spkde` implements a variety of nonparametric kernel-based estimators of the probability density function and the intensity function of two-dimensional spatial point patterns
- `spatwmat` imports or generates several kinds of spatial weights matrices
- `spatgsa` computes global indices of spatial autocorrelation
- `spatlisa` computes local indices of spatial autocorrelation

SPATIAL DATA

Spatial data: a discrete view

- For simplicity, let us represent **space** as a plane, i.e., as a flat two-dimensional surface
- In spatial data analysis, we can distinguish two conceptions of space (Bailey and Gatrell 1995: 18):
 - *Entity view*: Space as an area filled with a set of discrete objects
 - *Field view*: Space as an area covered with essentially continuous surfaces
- Here we take the former view and define **spatial data** as information regarding a given set of discrete spatial objects located within a study area \mathcal{A}

Attributes of spatial objects

- Information about spatial objects can be classified into two categories:
 - Spatial attributes
 - Non-spatial attributes
- The **spatial attributes** of a spatial object consist of one or more pairs of coordinates that represent its shape and/or its location within the study area
- The **non-spatial attributes** of a spatial object consist of its additional features that are relevant to the analysis at hand

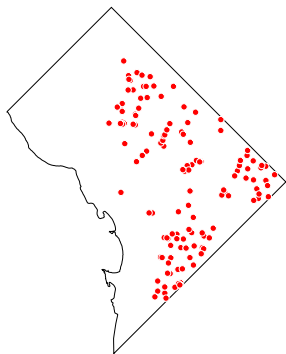
Types of spatial objects

- According to their spatial attributes, **spatial objects** can be classified into several types
- Here, we focus on two basic types:
 - Points (point data)
 - Polygons (area data)

Points

- A point s_i is a zero-dimensional spatial object located within study area \mathcal{A} at coordinates (s_{i1}, s_{i2})
- Points can represent several kinds of real entities, e.g., dwellings, buildings, places where specific events took place, pollution sources, trees

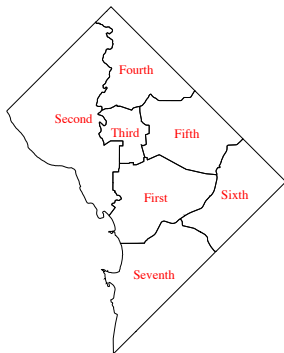
Homicides
Washington D.C. (2009)



Polygons

- A polygon \mathbf{r}_i is a *region* of study area \mathcal{A} bounded by a closed polygonal chain whose $M \geq 4$ vertices are defined by the coordinate set $\{(r_{i1(1)}, r_{i2(1)}), (r_{i1(2)}, r_{i2(2)}), \dots, (r_{i1(m)}, r_{i2(m)}), \dots, (r_{i1(M)}, r_{i2(M)})\}$, where $r_{i1(1)} = r_{i1(M)}$ and $r_{i2(1)} = r_{i2(M)}$
- Polygons can represent several kinds of real entities, e.g., states, provinces, counties, census tracts, electoral districts, parks, lakes

Police Districts
Washington D.C.



Introduction

Spatial data

Visualizing spatial data

Exploring spatial point patterns

Detecting spatial autocorrelation

References

Overview

Dot maps

Proportional symbol maps

Diagram maps

Choropleth maps

Multivariate maps

VISUALIZING SPATIAL DATA

Mapping

- Most exploratory analyses of spatial data have their natural starting point in displaying the information of interest by one or more **maps**
- If properly designed, maps can help the analyst to detect interesting patterns in the data, spatial relationships between two or more phenomena, unusual observations, and so on

Thematic maps

- In this talk I consider only the kind of maps most useful to ESDA: thematic maps
- **Thematic maps** represent the spatial distribution of a phenomenon of interest within a given study area (Slocum *et al.* 2005)

Thematic maps in Stata

- Stata users can generate thematic maps using `spmap`, a user-written command freely available from the SSC Archive (latest version: 1.2.0)
- `spmap` is a very flexible command that allows for creating a large variety of thematic maps, from the simplest to the most complex
- While providing sensible defaults for most options and supoptions, `spmap` gives the user full control over the formatting of almost every map element, thus allowing the production of highly customized maps

Thematic maps in Stata

- In the following, I will show how to use `spmap` for creating the types of thematic maps most commonly used in ESDA:
 - Dot maps
 - Proportional symbol maps
 - Diagram maps
 - Choropleth maps
 - Multivariate maps

Dot maps

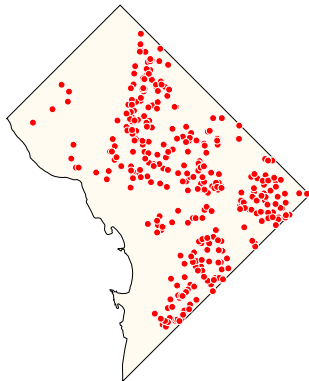
- A **dot map** shows the spatial distribution of a set of point spatial objects $\mathbf{S} \equiv \{\mathbf{s}_i; i = 1, \dots, N\}$, i.e., their location within a given study area \mathcal{A}
- If the point spatial objects have variable attributes, it is possible to represent this information using symbols of different colors and/or of different shape

Dot maps: example 1

Spatial distribution of 359 cases of sex abuse, Washington D.C. (2009)

```
use "Crime2009.dta", clear
generate _ID = _n
spmap using "Boundaries.dta", id(_ID) fcolor(eggshell) ///
point(x(x_coord) y(y_coord) select(keep if offense==6) ///
size(*1.2) fcolor(red) ocolor(white) osize(*0.5) ///
title("Sex abuses", size(*1.2)) ///
subtitle("Washington D.C. (2009)" " ", size(*1.2))
```

Sex abuses
Washington D.C. (2009)

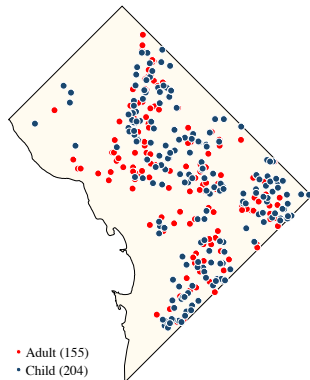


Dot maps: example 2

Spatial distribution of 359 cases of sex abuse, Washington D.C. (2009). Different colors are used to distinguish adult victims from child victims

```
use "Crime2009.dta", clear
generate _ID = _n
generate victim = method
recode victim (4/7-1)(17/18-2)(*=-.)
label define victim 1 "Adult" 2 "Child"
label values victim victim
spmap using "Boundaries.dta", id(_ID) fcolor(eggshell) ///
point(x(x_coord) y(y_coord) select(keep if offense==6) ///
by(victim) size(*1.2) fcolor(red navy) ///
ocolor(white ..) osize(*0.5 ..) legenda(on) ///
legcount) ///
legend(size(*1.8) rowgap(1.2)) ///
title("Sex abuses, by victim age", size(*1.2)) ///
subtitle("Washington D.C. (2009)" " ", size(*1.2))
```

Sex abuses, by victim age
Washington D.C. (2009)

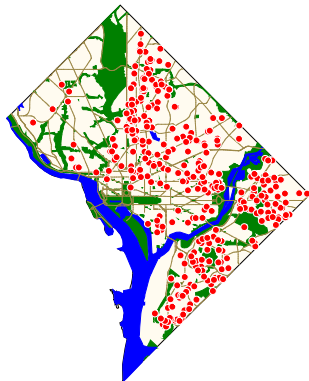


Dot maps: example 3

Spatial distribution of 359 cases of sex abuse, Washington D.C. (2009). Major roads, watercourses and parks are added to the map for reference

```
use "Crime2009.dta", clear
generate _ID = _n
spmap using "Boundaries.dta", id(_ID) fcolor(eggshell)    ///
point(x(x_coord) y(y_coord) select(keep if offense==6)  ///
size(*1.2) fcolor(red) ocolor(white) osize(*0.5))      ///
polygon(data("Water&Parks.dta") by(type)                ///
ocolor(none .) fcolor(green blue))                    ///
line(data("MajorRoads.dta") color(brown))              ///
title("Sex abuses", size(*1.2))                          ///
subtitle("Washington D.C. (2009)" " ", size(*1.2))
```

Sex abuses
Washington D.C. (2009)



Proportional symbol maps

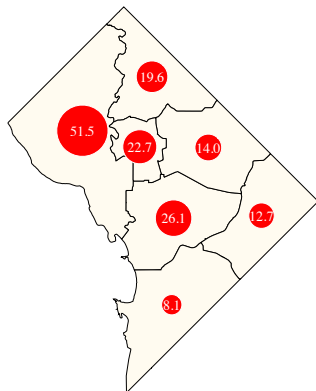
- A **proportional symbol map** represents the values taken by a numeric variable of interest Y on a set of point spatial objects \mathbf{S} located within a given study area \mathcal{A}
- Proportional symbol maps can be used with two types of point data (Slocum *et al.* 2005: 310):
 - **True point data** are measured at actual point locations
 - **Conceptual point data** are collected over a set of regions $\mathbf{R} \equiv \{\mathbf{r}_i; i = 1, \dots, N\}$, but are conceived as being located at representative points within the regions, typically at their centroids
- The area of each point symbol is sized in direct proportion to the corresponding value of Y

Proportional symbol maps: example 1

Mean family income in the seven Police Districts of Washington D.C. (2000)

```
use "PoliceDistricts-Data.dta", clear
generate Y = income_ma/1000
format Y %4.1f
spmap using "PoliceDistricts-Coordinates.dta", id(id) ///
    fcolor(eggshell) ///
    point(x(x_coord) y(y_coord) proportional(Y) fcolor(red) ///
        ocolor(white) size(*3.5)) ///
    label(x(x_coord) ycoord(y_coord) label(Y) color(white) ///
        size(*1.4)) ///
    title("Mean family income (in thousands of US dollars)") ///
    subtitle("Washington D.C. (2000)" " ")
```

Mean family income (in thousands of US dollars)
Washington D.C. (2000)



Proportional symbol maps: example 2

Mean family income in the 188 Census Tracts of Washington D.C. (2000). Solid circles denote *positive* deviations from the overall mean income, hollow circles denote *negative* deviations. Circles are drawn with size proportional to the absolute value of the deviation

```
use "Census2000-Data.dta", clear
generate Y = income_ma/1000
spmap using "Census2000-Coordinates.dta", id(id) ///
    fcolor(eggshell) ///
    point(x(x_coord) y(y_coord) fcolor(red) ocolor(white) ///
        size(*0.6) deviation(Y) refweight(poptot)) ///
    title("Mean family income") ///
    subtitle("Washington D.C. (2000)" " ")
```

Mean family income
Washington D.C. (2000)

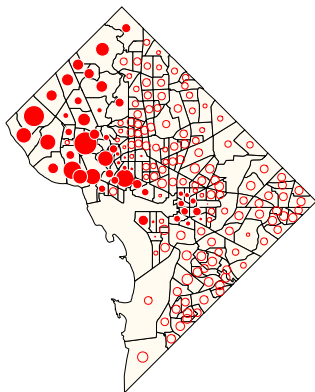


Diagram maps

- A **diagram map** follows the same logic as a proportional symbol map, but represents the values of the variable of interest using bar charts, pie charts, or other types of diagram
- The use of pie charts allows to display the spatial distribution of compositional data, i.e., of two or more numeric variables that represent parts of a whole

Diagram maps: example 1

Mean family income in the seven Police Districts of Washington D.C. (2000).
Data are represented by framed-rectangle charts, with the overall mean income as the reference value

```
use "PoliceDistricts-Data.dta", clear  
spmap using "PoliceDistricts-Coordinates.dta", id(id) //  
    fcolor(eggshell) //  
    diagram(var(income_ma) refweight(poptot) fcolor(green) //  
        x(x_coord) y(y_coord) size(1.3)) //  
    title("Mean family income") //  
    subtitle("Washington D.C. (2000)" " ") //
```

Mean family income
Washington D.C. (2000)

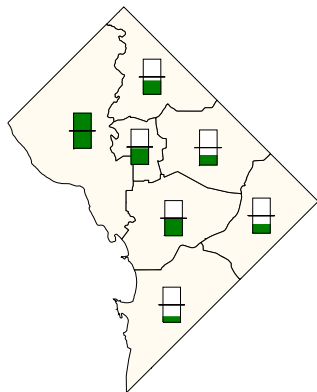
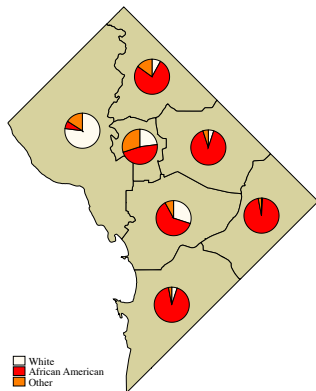


Diagram maps: example 2

Race/Ethnic composition of the population of the seven Police Districts of Washington D.C. (2000). Data are represented by pie charts

```
use "PoliceDistricts-Data.dta", clear
generate white_pct = pop_white/poptot*100
generate afroam_pct = pop_afroam/poptot*100
generate other_pct = pop_other/poptot*100
label variable white_pct "White"
label variable afroam_pct "African American"
label variable other_pct "Other"
spmap using "PoliceDistricts-Coordinates.dta", id(id) ///
    fcolor(stone) ///
    diagram(var(white_pct afroam_pct other_pct) x(x_coord) ///
        y(y_coord) fcolor(eggshell red orange) size(1.3) ///
        legenda(on)) ///
    legend(size(*1.4)) ///
    title("Race/Ethnic composition of the population") ///
    subtitle("Washington D.C. (2000)" " ")
```

Race/Ethnic composition of the population
Washington D.C. (2000)



Choropleth maps

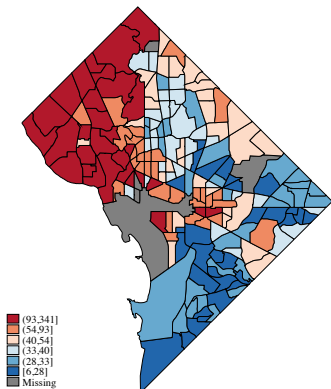
- A **choropleth map** displays the values taken by a variable of interest Y on a set of regions \mathbf{R} within a given study area \mathcal{A}
- When Y is numeric, each region is colored or shaded according to a discrete scale based on its value on Y
- The number of classes k that make up the discrete scale, and the corresponding class breaks, can be based on several different criteria – e.g., quantiles, equal intervals, boxplot, standard deviates

Choropleth maps: example 1

Mean family income in the 188 Census Tracts of Washington D.C. (2000).
Income is divided into six classes based on the *quantiles* method

```
use "Census2000-Data.dta", clear
generate Y = income_ma/1000
format Y %3.0f
spmap Y using "Census2000-Coordinates.dta", id(id) ///
    clnumber(6) clmethod(quantile) fcolor(BuRd) ///
    ndfcolor(gs8) ndlab("Missing") ///
    legend(size(*1.4)) ///
    title("Mean family income (in thousands of US dollars)") ///
    subtitle("Washington D.C. (2000)" " ")
```

Mean family income (in thousands of US dollars)
Washington D.C. (2000)

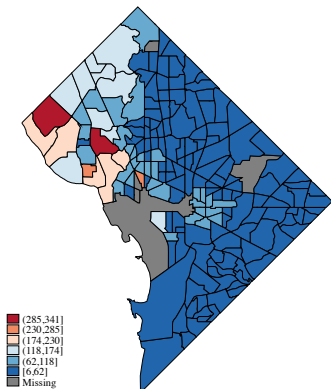


Choropleth maps: example 2

Mean family income in the 188 Census Tracts of Washington D.C. (2000).
Income is divided into six classes based on the *equal intervals* method

```
use "Census2000-Data.dta", clear
generate Y = income_ma/1000
format Y %3.0f
spmap Y using "Census2000-Coordinates.dta", id(id) ///
    clnumber(6) clmethod(eqint) fcolor(BuRd) ///
    ndfcolor(gs8) ndlab("Missing") ///
    legend(size(*1.4)) ///
    title("Mean family income (in thousands of US dollars)") ///
    subtitle("Washington D.C. (2000)" " ")
```

Mean family income (in thousands of US dollars)
Washington D.C. (2000)

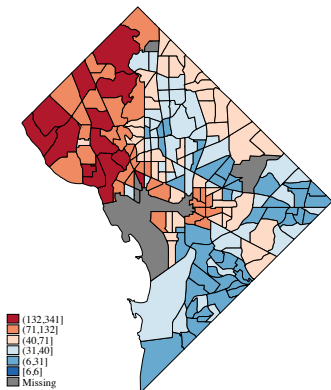


Choropleth maps: example 3

Mean family income in the 188 Census Tracts of Washington D.C. (2000).
Income is divided into six classes based on the *boxplot* method

```
use "Census2000-Data.dta", clear
generate Y = income_ma/1000
format Y %3.0f
spmap Y using "Census2000-Coordinates.dta", id(id) ///
      clnumber(6) clmethod(boxplot) fcolor(BuRd) ///
      ndfcolor(gs8) ndlab("Missing") ///
      legend(size(*1.4)) ///
      title("Mean family income (in thousands of US dollars)") ///
      subtitle("Washington D.C. (2000)" " ")
```

Mean family income (in thousands of US dollars)
Washington D.C. (2000)

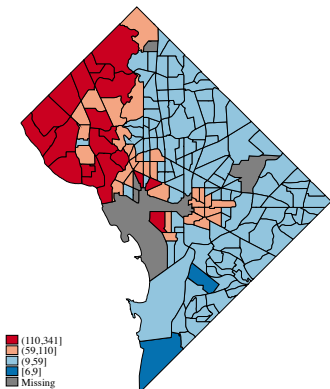


Choropleth maps: example 4

Mean family income in the 188 Census Tracts of Washington D.C. (2000).
Income is divided into four classes based on the *standard deviates* method

```
use "Census2000-Data.dta", clear
generate Y = income_ma/1000
format Y %3.0f
spmap Y using "Census2000-Coordinates.dta", id(id)
cnumber(4) cmethod(stdev) fcolor(BuRd)
ndfcolor(gs8) ndlab("Missing")
legend(size(*1.4))
title("Mean family income (in thousands of US dollars)")
subtitle("Washington D.C. (2000)" " ")
```

Mean family income (in thousands of US dollars)
Washington D.C. (2000)



Multivariate maps

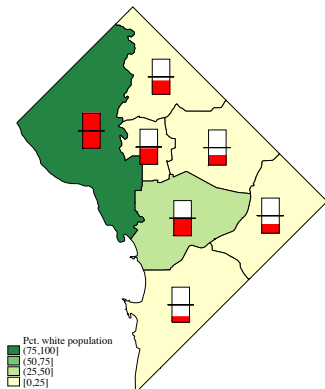
- A **multivariate map** combines several types of thematic mapping to simultaneously display the spatial distribution of multiple phenomena within a given study area \mathcal{A}

Multivariate maps: example 1

The map shows the relationship between mean family income (represented by framed-rectangle charts) and pct. white population (represented by shades of color) across the seven Police Districts of Washington D.C. (2000)

```
use "PoliceDistricts-Data.dta", clear
generate Y = pop_white/poptot*100
format Y %2.0f
spmap Y using "PoliceDistricts-Coordinates.dta", id(id) ///
    clmethod(custom) clbreaks(0 25 50 75 100) fcolor(Y1Gn) ///
    legtit("Pct. white population") ///
    diagram(var(income_ma) refweight(poptot) fcolor(red) ///
        x(x_coord) y(y_coord) size(1.3)) ///
    legend(size(*1.4)) ///
    title("Mean family income and pct. white population") ///
    subtitle("Washington D.C. (2000)" " ")
```

Mean family income and pct. white population
 Washington D.C. (2000)

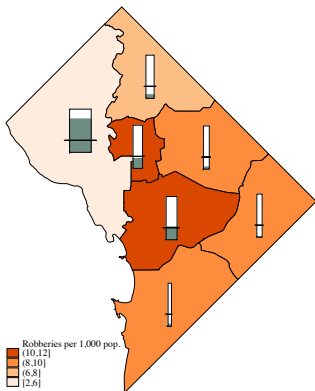


Multivariate maps: example 2

The map shows the relationship between pct. white population (represented by framed-rectangle charts), mean family income (represented by the width of framed-rectangle charts) and robbery rate (represented by shades of color) across the seven Police Districts of Washington D.C. (2000/2009)

```
use "PoliceDistricts-Data.dta", clear
generate Y = pop_white/poptot*100
format Y %2.0f
spmap Y using "PoliceDistricts-Coordinates.dta", id(id) ///
    cmethod(custom) clbreaks(0 25 50 75 100) fcolor(Y1Gn) ///
    legtit("Pct. white population") ///
    diagram(var(income_ma) refweight(poptot) fcolor(red) ///
        x(x_coord) y(y_coord) size(1.3)) ///
    legend(size(*1.4)) ///
    title("Mean family income and pct. white population") ///
    subtitle("Washington D.C. (2000)" " ")
```

Pct. white population, income and robberies
Washington D.C. (2000/2009)



EXPLORING SPATIAL POINT PATTERNS

Two-dimensional spatial point patterns

- A **two-dimensional spatial point pattern** can be defined as a set of N point spatial objects \mathbf{S} located within a given study area \mathcal{A}
- Usually, each point $\mathbf{s}_i \in \mathbf{S}$ represents a real entity of some kind: people, events, sites, buildings, plants, cases of a disease, etc.
- Alternatively, each point \mathbf{s}_i represents the centroid of a region
- Points \mathbf{s}_i are referred to as the *data points*

Two-dimensional spatial point patterns

- In the analysis of spatial point patterns, we are often interested in determining whether the observed data points exhibit some form of *clustering*, as opposed to being distributed uniformly within \mathcal{A}
- To explore the possibility of point clustering, it may be useful to describe the spatial point pattern of interest by means of its probability density function $p(\mathbf{s})$ and/or its intensity function $\lambda(\mathbf{s})$ (Waller and Gotway 2004)

Two-dimensional spatial point patterns

- The **probability density function** $p(\mathbf{s})$ defines the probability of observing an object per unit area at location $\mathbf{s} \in \mathcal{A}$
- The **intensity function** $\lambda(\mathbf{s})$ defines the expected number of objects per unit area at location $\mathbf{s} \in \mathcal{A}$
- The probability density function and the intensity function differ only by a constant of proportionality

Kernel estimators

- Both the probability density function $p(\mathbf{s})$ and the intensity function $\lambda(\mathbf{s})$ of a two-dimensional spatial point pattern can be estimated by means of nonparametric estimators, e.g., kernel estimators (Waller and Gotway 2004)
- **Kernel estimators** are used to generate a spatially smooth estimate of $p(\mathbf{s})$ and/or $\lambda(\mathbf{s})$ at a fine grid of points \mathbf{s}_g ($g = 1, \dots, G$) covering the study area \mathcal{A}
- In the context of spatial data analysis, a **grid** is a regular tessellation of the study area \mathcal{A} that divides it into a set of G contiguous cells whose centers are referred to as the *grid points* and denoted by \mathbf{s}_g

Kernel estimator of $\lambda(\mathbf{s})$

- The intensity $\lambda(\mathbf{s}_g)$ at each grid point \mathbf{s}_g is estimated by:

$$\hat{\lambda}(\mathbf{s}_g) = \frac{c}{A_g} \sum_{i=1}^N k \left(\frac{d(\mathbf{s}_i, \mathbf{s}_g)}{h_i} \right) y_i$$

where $k(\cdot)$ is the *kernel function* – usually a unimodal symmetrical bivariate probability density function; h_i is the *kernel bandwidth*, i.e., the radius of the kernel function; $d(\mathbf{s}_i, \mathbf{s}_g)$ is the Euclidean distance between data point \mathbf{s}_i and grid point \mathbf{s}_g ; y_i is the value taken by an optional variable of interest Y at data point \mathbf{s}_i ; A_g is the area of the region of \mathcal{A} over which the kernel function is evaluated, possibly corrected for *edge effects*; and c is a constant of proportionality

Kernel estimator of $p(\mathbf{s})$

- In turn, the probability density $p(\mathbf{s}_g)$ at each grid point \mathbf{s}_g is estimated by:

$$\hat{p}(\mathbf{s}_g) = \frac{\hat{\lambda}(\mathbf{s}_g)}{\sum_{j=1}^G \hat{\lambda}(\mathbf{s}_j)}$$

Kernel estimation in Stata

- Stata users can generate kernel estimates of the probability density function $p(\mathbf{s})$ and the intensity function $\lambda(\mathbf{s})$ using two user-written commands freely available from the SSC Archive: `spgrid` and `spkde`
- `spgrid` (latest version: 1.0.1) generates several kinds of two-dimensional grids covering rectangular or irregular study areas
- `spkde` (latest version: 1.0.0) implements a variety of kernel estimators of $p(\mathbf{s})$ and $\lambda(\mathbf{s})$
- `spmap` can then be used to visualize the kernel estimates generated by `spgrid` and `spkde`

Kernel estimation: example

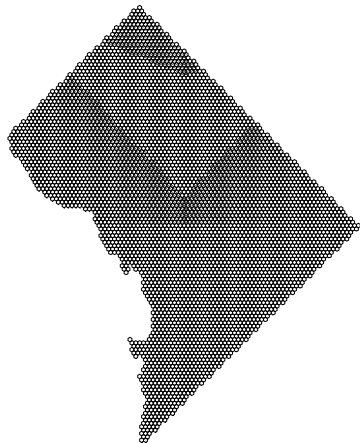
Our purpose is to estimate the probability density function of a set of 139 points representing the **homicides** committed in Washington D.C. in 2009

Kernel estimation: example

Step 1

We use `spgrid` to generate a grid covering the area of Washington D.C. We choose a relatively fine grid resolution (grid cell width = 200 meters). `spmap` is used to display the grid

```
spgrid using "Boundaries.dta", resolution(w200)   ///  
  dots compress unit(meters) cells("ctemp.dta")  ///  
  points("ptemp.dta") replace  
  
use "ptemp.dta", clear  
spmap using "ctemp.dta", id(spgrid_id)
```



Kernel estimation: example

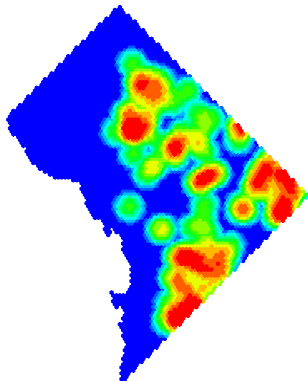
Step 2

We use `spkde` to generate kernel estimates of the probability distribution of homicides in Washington D.C. We choose a quartic kernel function with fixed bandwidth equal to 1,000 meters and edge correction. `spmap` is used to display the results

```
use "Crime2009.dta", clear
keep if offense==4
spkde using "ptemp.dta", x(x_coord) y(y_coord) ///
    kernel(quartic) bandwidth(fbw) fbw(1000) ///
    edgcorrect dots saving("kde.dta", replace)

use "kde.dta", clear
spmap p using "ctemp.dta", id(spgrid_id) clmethod(quantile) ///
    clnumber(20) fcolor(Rainbow) ocolor(none ..) legend(off) ///
    title("Homicides", size(*1.2)) ///
    subtitle("Washington D.C. (2009)" " ", size(*1.2))
```

Homicides
Washington D.C. (2009)



Kernel estimation: example

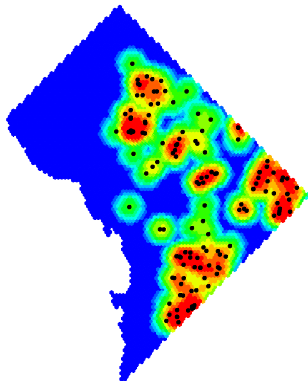
Step 2

We use `spkde` to generate kernel estimates of the probability distribution of homicides in Washington D.C. We choose a quartic kernel function with fixed bandwidth equal to 1,000 meters and edge correction. `spmap` is used to display the results

```
use "Crime2009.dta", clear
keep if offense==4
spkde using "ptemp.dta", x(x_coord) y(y_coord) ///
kernel(quartic) bandwidth(fbw) fbw(1000) ///
edgcorrect dots saving("kde.dta", replace)

use "kde.dta", clear
spmap p using "ctemp.dta", id(spgrid_id) clmethod(quantile) ///
cnumber(20) fcolor(Rainbow) ocolor(none ..) legend(off) ///
title("Homicides", size(*1.2)) ///
subtitle("Washington D.C. (2009)" " ", size(*1.2))
```

Homicides
Washington D.C. (2009)



DETECTING SPATIAL AUTOCORRELATION

Spatial autocorrelation

- Forty years ago, the geographer and statistician Waldo Tobler formulated the *first law of geography*: “Everything is related to everything else, but near things are more related than distant things” (Tobler 1970: 234)
- This “law” defines the statistical concept of (positive) **spatial autocorrelation**, according to which two or more objects that are spatially close tend to be more similar to each other – with respect to a given attribute Y – than are spatially distant objects
- In general, spatial autocorrelation implies **spatial clustering**, i.e., the existence of sub-areas of the study area where the attribute of interest Y takes higher than average values (*hot spots*) or lower than average values (*cold spots*)

Spatial weights matrix

- The analysis of spatial autocorrelation requires the measurement of the degree of **spatial proximity** among the spatial objects of interest
- Typically, the degree of spatial proximity among a given set of N spatial objects is represented by a $N \times N$ matrix called **spatial weights matrix** and denoted by **\mathbf{W}**
- Each element (i, j) of **\mathbf{W}** – which we denote by w_{ij} – expresses the degree of spatial proximity between the pair of objects i and j
- Depending on the application, the N main diagonal elements of **\mathbf{W}** are assigned value $w_{ii} = 0$ or value $w_{ii} > 0$

Spatial weights matrix

- A common variant of \mathbf{W} is the **row-standardized spatial weights matrix** \mathbf{W}_{std} , whose elements are defined as follows:

$$w_{ij}^{std} = \frac{w_{ij}}{\sum_{j=1}^N w_{ij}}$$

Spatial weights matrices in Stata

- Stata users can generate several kinds of spatial weights matrices using `spatwmat`, a user-written command published in the *Stata Technical Bulletin* (Pisati 2001)
- `spatwmat` (latest version: 1.0) imports or generates from scratch the spatial weights matrices required by the commands `spatgsa` and `spatlisa` described below

Indices of spatial autocorrelation

- We consider measures of spatial autocorrelation that apply to *area data*
- Measures of spatial autocorrelation can be classified into two broad categories:
 - Global indices of spatial autocorrelation
 - Local indices of spatial autocorrelation

Global indices of spatial autocorrelation

- A **global index of spatial autocorrelation** expresses the overall degree of similarity between spatially close regions observed in a given study area \mathcal{A} with respect to a numeric variable Y (Pfeiffer *et al.* 2008)
- Since global indices of spatial autocorrelation summarize the phenomenon of interest in a single value, they are intended not so much for identifying specific spatial clusters, as for detecting the presence of a general tendency to clustering within the study area

Global indices of spatial autocorrelation

- In general, the computation of a global index of spatial autocorrelation follows a three-step procedure:
 - First, we compute the degree of similarity ρ_{ij} between every possible pair of regions \mathbf{r}_i and \mathbf{r}_j with respect to the numeric variable of interest Y
 - Second, we weight – i.e., multiply – each value ρ_{ij} by the degree of proximity w_{ij} between regions \mathbf{r}_i and \mathbf{r}_j
 - Finally, we sum up all the products $w_{ij}\rho_{ij}$ and divide the total by a constant of proportionality
- The greater the number of regions that are similar with respect to Y and spatially close, the greater the value taken by the global index of spatial autocorrelation

Moran's I

- **Moran's** global index of spatial autocorrelation I (Moran 1948) defines ρ_{ij} as $(y_i - \bar{y})(y_j - \bar{y})$, where y_i is the value taken by Y in region \mathbf{r}_i , y_j is the value taken by Y in region \mathbf{r}_j , and \bar{y} is the average value of Y :

$$I = \frac{\sum_{i=1}^N \sum_{j=1}^N w_{ij} (y_i - \bar{y})(y_j - \bar{y})}{\frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2 \sum_{i=1}^N \sum_{j=1}^N w_{ij}}$$

where $w_{ii} = 0$

Moran's I

- Under the null hypothesis of no global spatial autocorrelation, the expected value of I is:

$$E(I) = -\frac{1}{N-1}$$

- $I > E(I)$ indicates **positive** spatial autocorrelation – nearby regions tend to exhibit similar values of Y
- $I < E(I)$ indicates **negative** spatial autocorrelation – nearby regions tend to exhibit dissimilar values of Y

Getis and Ord's G

- **Getis and Ord's** global index of spatial autocorrelation G (Getis and Ord 1992) defines ρ_{ij} as $y_i y_j$:

$$G = \frac{\sum_{i=1}^N \sum_{j=1}^N w_{ij} y_i y_j}{\sum_{i=1}^N \sum_{j=1}^N y_i y_j} \quad (i \neq j)$$

where w_{ij} denotes the elements of a non-standardized binary spatial weights matrix, with $w_{ii} = 0$

- G is defined only for $Y > 0$

Getis and Ord's G

- Under the null hypothesis of no global spatial autocorrelation, the expected value of G is:

$$E(G) = \frac{\sum_{i=1}^N \sum_{j=1}^N w_{ij}}{N(N-1)}$$

- $G > E(G)$ indicates (a) positive spatial autocorrelation, and (b) prevalence of spatial clusters with relatively **high** values of Y
- $G < E(G)$ indicates (a) negative spatial autocorrelation, and (b) prevalence of spatial clusters with relatively **low** values of Y

Global indices of spatial autocorrelation in Stata

- Stata users can compute global indices of spatial autocorrelation using **spatgsa**, a user-written command published in the *Stata Technical Bulletin* (Pisati 2001)
- **spatgsa** (latest version: 1.0) computes three global indices of spatial autocorrelation: Moran's I , Getis and Ord's G , and Geary's c . For each index and each numeric variable of interest, **spatgsa** computes and displays in tabular form the value of the index itself, the expected value of the index under the null hypothesis of no global spatial autocorrelation, the standard deviation of the index, the z -value, and the corresponding one- or two-tailed p -value

Global indices of spatial autocorrelation: example

- **Study area:** Ohio
- **Regions:** 88 counties
- **Variables of interest:**
 - Pct. population aged 18+ with poor-to-fair health status (`pct_poorhealth`)
 - Pct. population aged 18+ currently smoking (`pct_currsmoker`)
 - Pct. population aged 18+ ever diagnosed with high blood pressure (`pct_hibloodprs`)
 - Pct. population aged 18+ obese (`pct_obese`)

Global indices of spatial autocorrelation: example

Step 1

We use `spatwmat` to import an existing binary spatial weights matrix – stored in the Stata dataset `Counties-Contiguity.dta` – and convert it into a properly formatted row-standardized spatial weights matrix `Ws`

```
spatwmat using "Counties-Contiguity.dta", ///  
name(Ws) standardize
```

Global indices of spatial autocorrelation: example

The following matrix has been created:

1. Imported binary weights matrix **Ws** (row-standardized)
Dimension: **88x88**

Global indices of spatial autocorrelation: example

Step 2

We use `spatgsa` with the spatial weights matrix `Ws` to compute Moran's I on the variables of interest

```
use "Counties-Data.dta", clear
spatgsa pct_poorhealth pct_currsmoker pct_hibloodprs  ///
pct_obese, w(Ws) moran
```

Global indices of spatial autocorrelation: example

Measures of global spatial autocorrelation

Weights matrix

Name: **Ws**

Type: **Imported (binary)**

Row-standardized: **Yes**

Moran's I

Variables	I	E(I)	sd(I)	z	p-value*
pct_poorhealth	0.399	-0.011	0.065	6.337	0.000
pct_currrsmoker	0.339	-0.011	0.065	5.367	0.000
pct_hibloodprs	0.126	-0.011	0.065	2.119	0.017
pct_obese	0.167	-0.011	0.065	2.730	0.003

*1-tail test

Global indices of spatial autocorrelation: example

Step 3

We use again `spatwmat` to import the binary spatial weights matrix stored in the Stata dataset `Counties-Contiguity.dta`. This time, however, we convert it into a properly formatted non-standardized spatial weights matrix `W`

```
spatwmat using "Counties-Contiguity.dta", ///  
name(W)
```

Global indices of spatial autocorrelation: example

The following matrix has been created:

1. Imported binary weights matrix **W**
Dimension: **88x88**

Global indices of spatial autocorrelation: example

Step 4

We use `spatgsa` with the spatial weights matrix W to compute Getis and Ord's G on the variables of interest

```
use "Counties-Data.dta", clear
spatgsa pct_poorhealth pct_currsmoker pct_hibloodprs  ///
pct_obese, w(W) go
```

Global indices of spatial autocorrelation: example

Measures of global spatial autocorrelation

Weights matrix

Name: **W**

Type: **Imported (binary)**

Row-standardized: **No**

Getis & Ord's G

Variables	G	E(G)	sd(G)	z	p-value*
pct_poorhealth	0.061	0.060	0.001	0.372	0.355
pct_currsmoker	0.060	0.060	0.001	-0.392	0.348
pct_hibloodprs	0.060	0.060	0.000	-1.585	0.056
pct_obese	0.061	0.060	0.000	0.458	0.323

*1-tail test

Local indices of spatial autocorrelation

- A **local index of spatial autocorrelation** expresses, for each region \mathbf{r}_i of a given study area \mathcal{A} , the degree of similarity between that region and its neighboring regions with respect to a numeric variable Y (Pfeiffer *et al.* 2008)
- The local indices of spatial autocorrelation considered here are derived from the two global indices described above – Moran's I and Getis and Ord's G – and, therefore, share their fundamental properties

Moran's I_i

- **Moran's** local index of spatial autocorrelation I_i is defined as follows:

$$I_i = \sum_{j=1}^N w_{ij}^{std} \left(\frac{y_i - \bar{y}}{\sigma_y} \right) \left(\frac{y_j - \bar{y}}{\sigma_y} \right)$$

where σ_y denotes the standard deviation of Y ; and w_{ij}^{std} denotes the elements of a row-standardized spatial weights matrix, with $w_{ii}^{std} = 0$

Moran's I_i

- $I_i > E(I_i)$ indicates that region \mathbf{r}_i is surrounded by regions that, on average, are similar to \mathbf{r}_i with respect to Y (*positive* spatial autocorrelation). Moreover:
 - If $(y_i - \bar{y}) > 0$, then \mathbf{r}_i represents a *hot spot*
 - If $(y_i - \bar{y}) < 0$, then \mathbf{r}_i represents a *cold spot*
- $I_i < E(I_i)$ indicates that region \mathbf{r}_i is surrounded by regions that, on average, are different from \mathbf{r}_i with respect to Y (*negative* spatial autocorrelation)

Getis and Ord's G_i

- **Getis and Ord's** local index of spatial autocorrelation G_i is defined as follows:

$$G_i = \frac{\sum_{j=1}^N w_{ij} y_j}{\sum_{j=1}^N y_j} \quad (i \neq j)$$

where w_{ij} denotes the elements of a non-standardized binary spatial weights matrix, with $w_{ii} = 0$

Getis and Ord's G_i^*

- **Getis and Ord's** local index of spatial autocorrelation G_i^* is defined as follows:

$$G_i^* = \frac{\sum_{j=1}^N w_{ij} y_j}{\sum_{j=1}^N y_j}$$

where w_{ij} denotes the elements of a non-standardized binary spatial weights matrix, with $w_{ii} = 1$

Getis and Ord's G_i and G_i^*

- Like G , G_i and G_i^* can identify only regions characterized by *positive* spatial autocorrelation
- If $G_i > E(G_i)$ or $G_i^* > E(G_i^*)$, then region \mathbf{r}_i represents a *hot spot*
- If $G_i < E(G_i)$ or $G_i^* < E(G_i^*)$, then region \mathbf{r}_i represents a *cold spot*

Local indices of spatial autocorrelation in Stata

- Stata users can compute local indices of spatial autocorrelation using `spatlisa`, a user-written command published in the *Stata Technical Bulletin* (Pisati 2001)
- `spatlisa` (latest version: 1.0) computes four indices of spatial autocorrelation: Moran's I_i , Getis and Ord's G_i and G_i^* , and Geary's c_i . For each index and each region in the analysis, `spatlisa` computes and displays in tabular form the value of the index itself, the expected value of the index under the null hypothesis of no local spatial autocorrelation, the standard deviation of the index, the z -value, and the corresponding one- or two-tailed p -value

Local indices of spatial autocorrelation: example

- **Study area:** Ohio
- **Regions:** 88 counties
- **Variable of interest:** Pct. population aged 18+ with poor-to-fair health status (`pct_poorhealth`)

Local indices of spatial autocorrelation: example

Step 1

We use `spatlsa` with the standardized spatial weights matrix `Ws` – previously generated by `spatwmat` – to compute Moran's I_i on the variable of interest. In the output, counties are sorted by z -value

```
spatwmat using "Counties-Contiguity.dta", name(Ws) standardize  
use "Counties-Data.dta", clear  
spatlsa pct_poorhealth, w(Ws) moran id(name) sort
```

Local indices of spatial autocorrelation: example

Measures of local spatial autocorrelation

(Output omitted)

Moran's Ii (Poor-to-fair health status (pct. pop 18+))

name	Ii	E(Ii)	sd(Ii)	z	p-value*
Knox	-0.816	-0.011	0.358	-2.246	0.012
Hardin	-0.760	-0.011	0.358	-2.090	0.018
Paulding	-1.089	-0.011	0.560	-1.924	0.027
Licking	-0.457	-0.011	0.358	-1.244	0.107

(Output omitted)

Hancock	0.545	-0.011	0.358	1.555	0.060
Williams	0.927	-0.011	0.560	1.675	0.047
Delaware	0.677	-0.011	0.389	1.769	0.038
Mercer	0.908	-0.011	0.482	1.906	0.028
Putnam	0.949	-0.011	0.358	2.682	0.004
Henry	1.087	-0.011	0.358	3.067	0.001
Vinton	1.246	-0.011	0.389	3.230	0.001
Gallia	2.433	-0.011	0.482	5.069	0.000
Pike	2.197	-0.011	0.429	5.150	0.000
Adams	3.578	-0.011	0.482	7.442	0.000
Lawrence	5.503	-0.011	0.560	9.844	0.000
Jackson	3.911	-0.011	0.389	10.077	0.000
Scioto	5.400	-0.011	0.482	11.220	0.000

*1-tail test

Local indices of spatial autocorrelation: example

Step 2

We use Stata command `genmsp` – a simple wrapper to `spatlsa` (source code in Appendix) – to generate the variables required for drawing the Moran scatterplot and the corresponding cluster map (Anselin 1995)

```
genmsp pct_poorhealth, w(Ws)
```

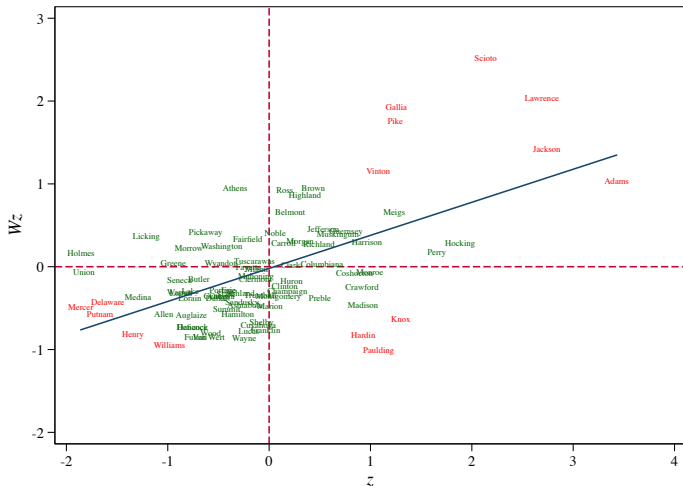
Local indices of spatial autocorrelation: example

Step 3

We use `graph twoway` to draw the Moran scatterplot

```
graph twoway                                     ///
  (scatter Wstd_pct_poorhealth std_pct_poorhealth  ///
    if pval_pct_poorhealth >= 0.05,             ///
    msymbol(i) mlabel(name) mlabsz(*0.6) mlabpos(c)  ///
  (scatter Wstd_pct_poorhealth std_pct_poorhealth  ///
    if pval_pct_poorhealth < 0.05,             ///
    msymbol(i) mlabel(name) mlabsz(*0.6) mlabpos(c)  ///
    mlabcol(red))                               ///
  (lfit Wstd_pct_poorhealth std_pct_poorhealth),   ///
  yline(0, lpattern(--)) xline(0, lpattern(--))  ///
  xlabel(-2(1)4, labsz(*0.8)) xtitle("{it:z}")   ///
  ylabel(-2(1)3, angle(0) labsz(*0.8))          ///
  ytitle("{it:Wz}") legend(off) scheme(s1color)
```

Local indices of spatial autocorrelation: example



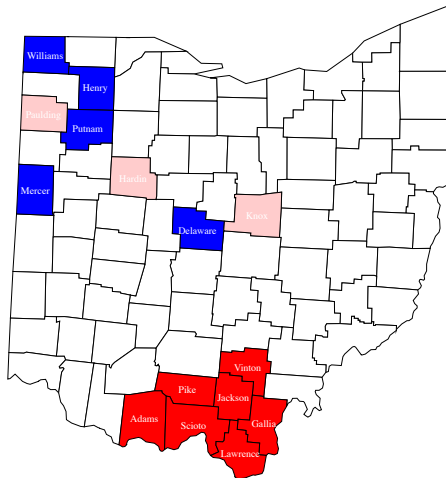
Local indices of spatial autocorrelation: example

Step 4

We use `spmap` to draw the cluster map corresponding to the Moran scatterplot

```
spmap msp_pct_poorhealth using "Counties-Coordinates.dta", ///  
  id(id) clmethod(unique) fcolor(blue red*0.2 red) ///  
  label(x(x_coord) y(y_coord) label(name) color(white) ///  
    size(*0.8) select(keep if msp_pct_poorhealth!=.)) ///  
  legend(off)
```

Local indices of spatial autocorrelation: example



REFERENCES

References I

- Anselin, L. 1995. Local indicators of spatial association – LISA. *Geographical Analysis* 27: 93–115.
- Bailey, T.C. and A.C. Gatrell. 1995. *Interactive Spatial Data Analysis*. Harlow: Longman.
- Getis, A. and J.K. Ord. 1992. The analysis of spatial association by use of distance statistics. *Geographical Analysis* 24: 189–206.
- Haining, R., Wise, S. and J. Ma. 1998. Exploratory spatial data analysis in a geographic information system environment. *The Statistician* 47: 457–469.
- Moran, P. 1948. The interpretation of statistical maps. *Journal of the Royal Statistical Society, Series B* 10: 243–251.
- Pfeiffer, D., Robinson, T., Stevenson, M., Stevens, K., Rogers, D. and A. Clements. 2008. *Spatial Analysis in Epidemiology*. Oxford: Oxford University Press.
- Pisati, M. 2001. sg162: Tools for spatial data analysis. *Stata Technical Bulletin* 60: 21–37. In *Stata Technical Bulletin Reprints*, vol. 10, 277–298. College Station, TX: Stata Press.

References II

- Slocum, T.A., McMaster, R.B., Kessler, F.C. and H.H. Howard. 2005. *Thematic Cartography and Geographic Visualization*. 2nd ed. Upper Saddle River, NJ: Pearson Prentice Hall.
- Tobler, W.R. 1970. A computer movie simulating urban growth in the Detroit region. *Economic Geography* 46: 234–240.
- Waller, L.A. and C.A. Gotway. 2004. *Applied Spatial Statistics for Public Health Data*. Hoboken, NJ: Wiley.

APPENDIX

Source code for genmsp I

```
program genmsp, sortpreserve
version 12.1

syntax varname, Weights(name) [Pvalue(real 0.05)]

unab Y : 'varlist'
tempname W
matrix 'W' = 'weights'
tempvar Z
qui summarize 'Y'
qui generate 'Z' = ('Y' - r(mean)) / sqrt( r(Var) * ( r(N)-1) / r(N) )
qui cap drop std_ 'Y'
qui generate std_ 'Y' = 'Z'
tempname z Wz
qui mkmat 'Z', matrix('z')
matrix 'Wz' = 'W'*'z'
matrix colnames 'Wz' = Wstd_ 'Y'
qui cap drop Wstd_ 'Y'
qui svmat 'Wz', names(col)
qui spatlsa 'Y', w('W') moran
tempname M
matrix 'M' = r(Moran)
matrix colnames 'M' = __c1 __c2 __c3 zval_ 'Y' pval_ 'Y'
```

Source code for genmsp II

```
qui cap drop __c1 __c2 __c3
qui cap drop zval_`Y'
qui cap drop pval_`Y'
qui svmat `M', names(col)
qui cap drop __c1 __c2 __c3
qui cap drop msp_`Y'
qui generate msp_`Y' = .
qui replace msp_`Y' = 1 if std_`Y'<0 & Wstd_`Y'<0 & pval_`Y'<'pvalue'
qui replace msp_`Y' = 2 if std_`Y'<0 & Wstd_`Y'>0 & pval_`Y'<'pvalue'
qui replace msp_`Y' = 3 if std_`Y'>0 & Wstd_`Y'<0 & pval_`Y'<'pvalue'
qui replace msp_`Y' = 4 if std_`Y'>0 & Wstd_`Y'>0 & pval_`Y'<'pvalue'
lab def __msp 1 "Low-Low" 2 "Low-High" 3 "High-Low" 4 "High-High", modify
lab val msp_`Y' __msp

end
exit
```