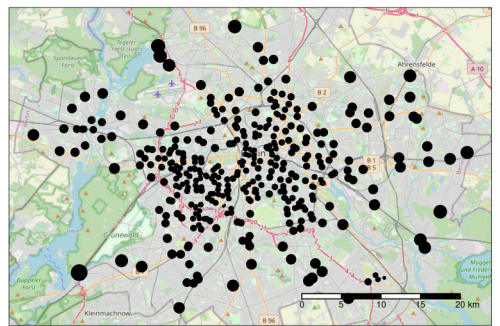
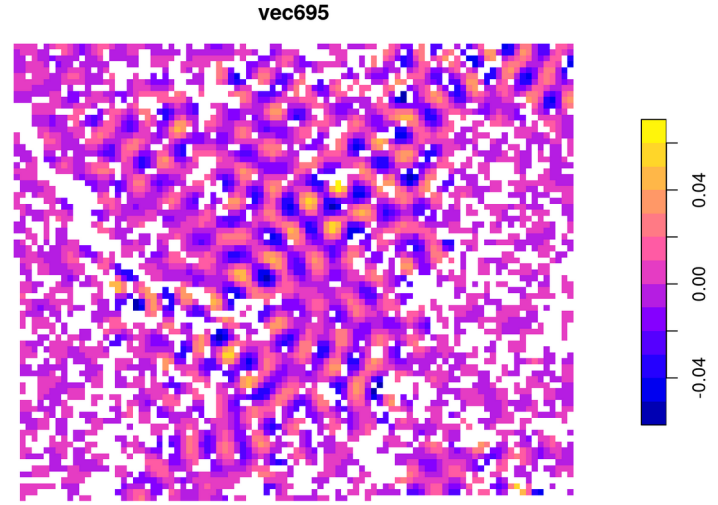


HEIDELBERG INSTITUTE  
FOR GEOINFORMATION  
TECHNOLOGY

# Spatial analysis



Sven Lautenbach  
Urban Data Science Workshop  
19.-23.10.2020



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386

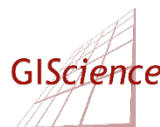


Klaus Tschira Stiftung  
gemeinnützige GmbH



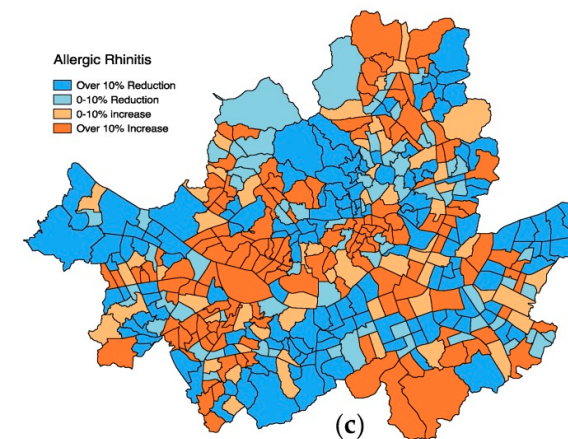
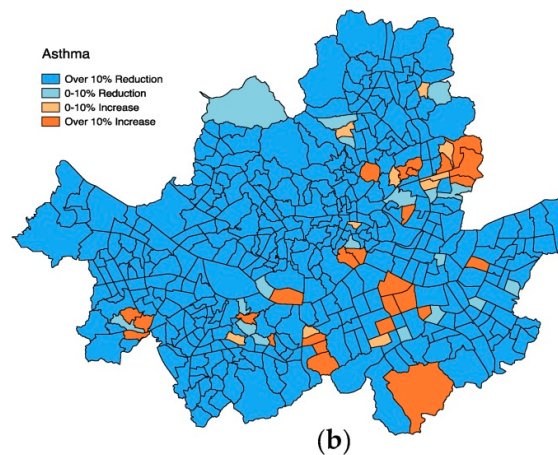
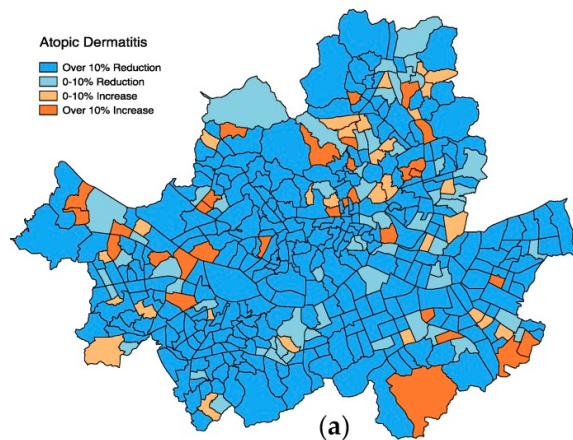
# Overview

- Intro
- Spatial explorative analysis
- Defining neighborhoods and assigning weights
- Detecting spatial autocorrelation
- Spatial clusters
- Dealing with it in regression context
  - Scoping the problem
  - Spatial Eigenvector Mapping
  - Spatial lag and spatial error model
  - Other approaches



# What is special about space?

- Tobler's first law of geography "everything is related to everything else, but near things are more related than distant things."



Kim et al. A Closer Look at the Bivariate Association between Ambient Air Pollution and Allergic Diseases: The Role of Spatial Analysis. *Int J Environ Res Public Health*. 2018 Aug; 15(8): 1625

# Important tasks in spatial statistics

- Point pattern analysis
  - Analysis of spatial configuration of a population (not a sample)
  - Requires point data
- Interpolation
  - Estimation of surfaces from a sample of observation
- **Spatial explorative analysis**
  - Analysis of **spatial autocorrelation**, hotspots and coldspots
  - **Spatial Cluster analysis**
    - Geographically weighted regression
- **Regression analysis under incorporation of spatial autocorrelation**



# Short history of spatial statistics

- Development independent in different disciplines:
  - Geosciences especially for mining → Geostatistics (Kriging), e.g. Cressie
  - Ecology, e.g. Legendre, Levin
  - Spatial econometrics, e.g. Anselin
  - Geography, e.g. Griffith, Haining
- Therefore related concepts with varying terminology



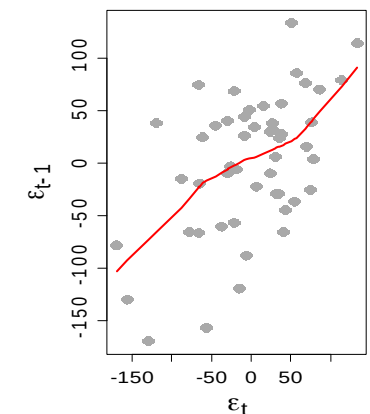
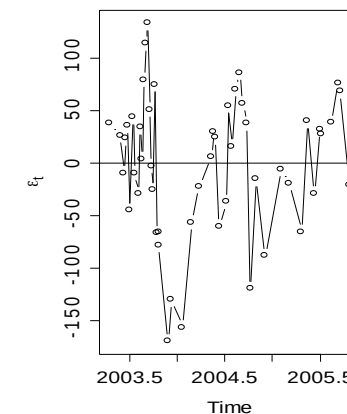
# Spatial explorative analysis

- Visual interpretation of spatial data is e.g. complicated due to the necessary classification of continuous data
- It is not always obvious if a spatial pattern occurs in the data
  - Clustered
  - Regular
  - Random
- Spatial autocorrelation is a way to quantify this
- A number of tools exist that quantified the relationship
- In difference to point pattern attribute values have a higher importance
  - Not only geometry but distribution of z-values in space
  - Data might be points, polylines or polygons



# Autocorrelation

- Much classic statistical theory assumes that errors are independent and identically distributed (iid), often conforming to a Gaussian distribution
  - Simplifies equations since covariation terms can be set to zero
- However, errors might be structured, violating the simplifying assumptions
- Hierarchical structure, e.g. by unaccounted effects of groups
- Temporal autocorrelation
- Spatial autocorrelation



# Effects of autocorrelation

- Violating assumptions
- Increasing variance for positive spatial autocorrelation, decreasing variance for negative spatial autocorrelation
  - Estimated standard errors are too small (or too big for negative sac) which effects p-values
  - Mixing up model selection
- Increase or decrease correlation coefficients
- Sample sizes goes down (for positive sac)
  - Effecting standard errors and p-values and model selection procedures

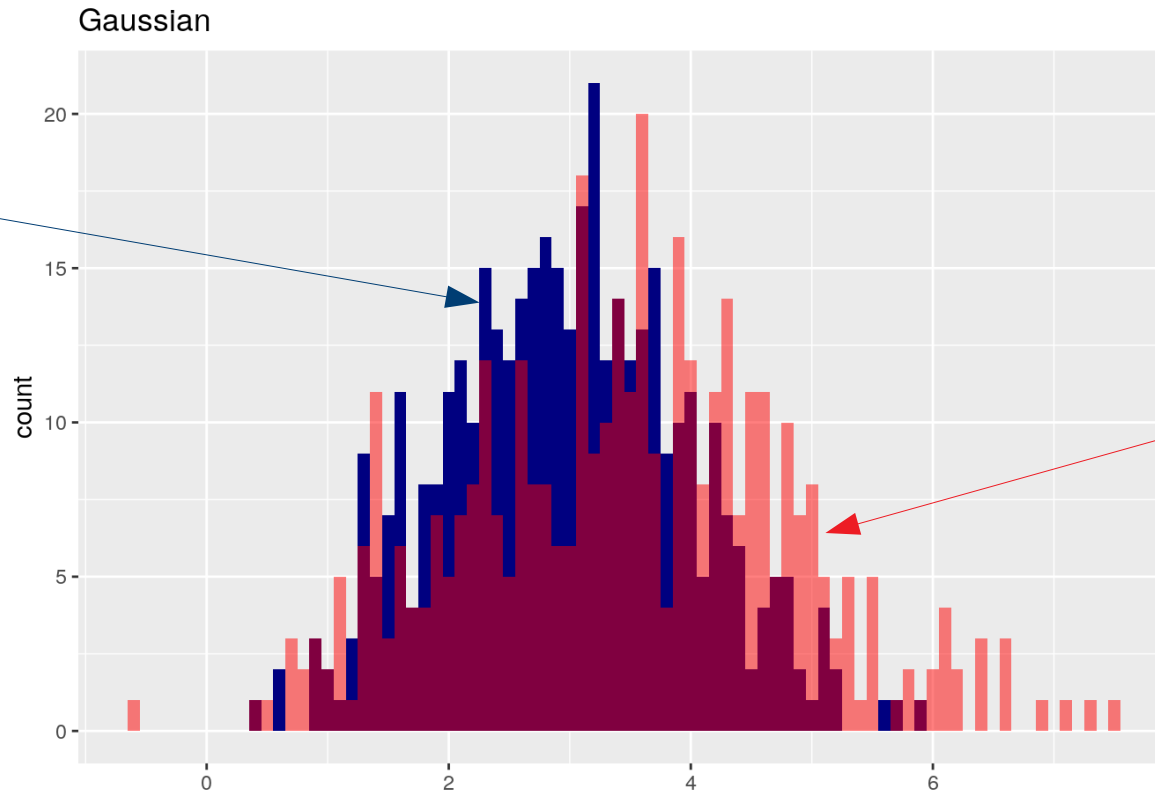




# Effects of positive autocorrelation

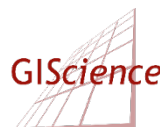
- 

No sac



SAC

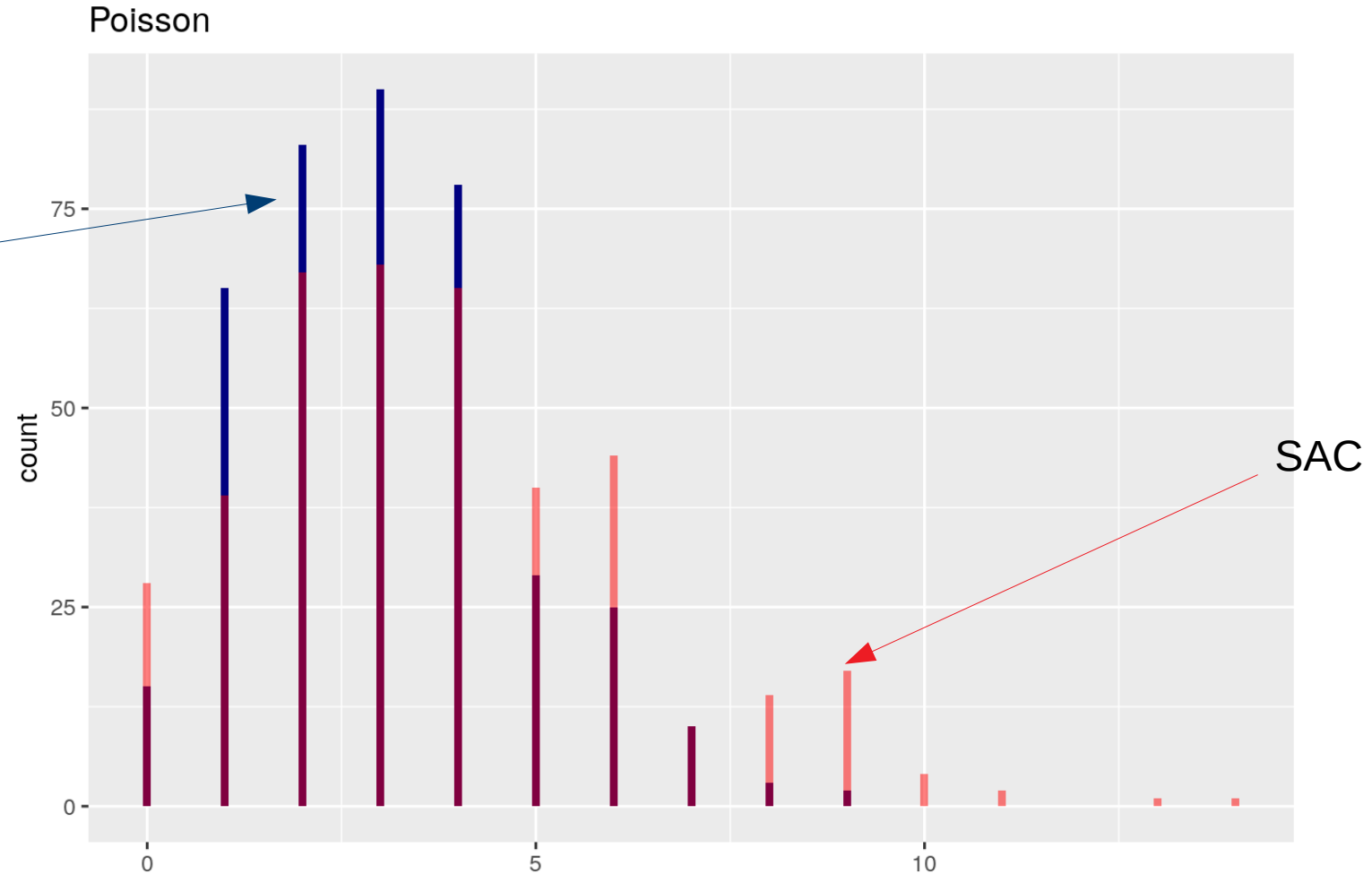
Increasing variance  
Tails become heavier  
Preservation of the mean



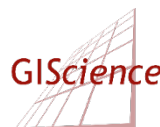
# Effects of positive autocorrelation

- 

No sac



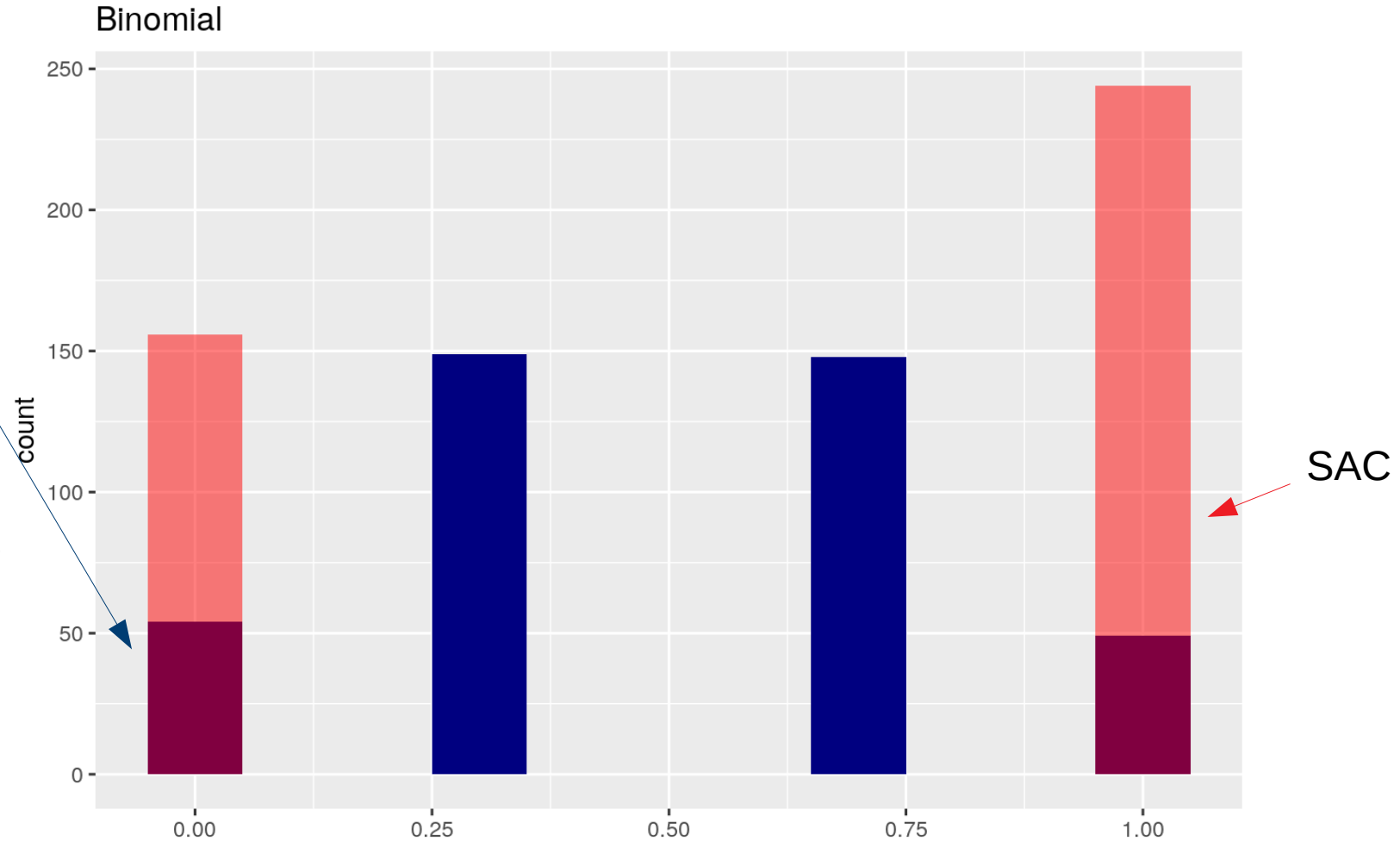
Increasing variance  
More zeros and outliers  
Preservation of the mean



# Effects of positive autocorrelation

No sac

Distribution moving to the All-or nothing situation  
Preservation of the mean



# Mechanisms behind spatial autocorrelation

- **Induced spatial dependence**  
function dependency of the response on spatially structured predictors  
e.g. ecological niche for mosquitos of Aedes type, distribution of charging stations follows population density
- **True autocorrelation**  
functional dependency between the response and adjacent response values, e.g. infectious diseases, distribution of charging stations depends on charging stations in adjacent units
- **Historical dynamics**  
Past events have led to a spatial structure that influences the response  
e.g. infrastructure (railroads, highways) led to spatially structured path dependencies



# Reasons for spatial autocorrelation

- The response is autocorrelated
  - High densities of charging stations in one neighborhood might reduce number in adjacent neighborhoods (negative autocorrelation, spill over effect)
- A predictor with a spatial structure is missing which might lead to autocorrelation of the residuals
- We cannot distinguish between the two cases from the residuals
- From a theoretical perspective we might be able to develop a hypothesis on the reason for the presence of spatial autocorrelation



# Spatial autocorrelation

Complete independence

$$Y_i = \varepsilon_i, \varepsilon_i \approx N(0, \sigma_\varepsilon^2)$$

Spatial independence

$$y_i = \beta z_i + \varepsilon_i$$

$$z_i = \xi_i$$

$$\xi_i \approx N(0, \sigma_\xi^2)$$

Inherent autoregressive

$$y_i = \rho y_{i-1} + \varepsilon_i$$

$$-1 \leq \rho \leq 1$$

Induced autoregressive

$$y_i = \beta z_i + \varepsilon_i$$

$$z_i = \rho_z z_i + \xi_i$$

Doubly autoregressive

$$y_i = \beta z_i + \rho_y y_{i-1} + \varepsilon_i$$

$$z_i = \rho_z z_i + \xi_i$$

$$\xi_i \approx N(0, \sigma_\xi^2)$$



# Representation of autocorrelation

- Autocorrelation can be e.g. represented in the variance-covariance matrix of the error term

- 

$$Y = X \beta + \varepsilon$$

$$\varepsilon \approx N(0, I \sigma^2)$$

- 

- 

- Independent errors

$$Y = X \beta + \varepsilon$$

$$\varepsilon \approx N(0, \sum \sigma^2)$$

- Structured errors

- The variance-covariance matrix of the error term is assumed (and tested) to follow a specific structure:

- Correlation in groups: covariance only for members of the same group
- Temporal auto-correlation: covariance depends on temporal lag
- Spatial auto-correlation: covariance depends on spatial lag / neighborhood definition



# Defining spatial relationship Neighbors and spatial weights





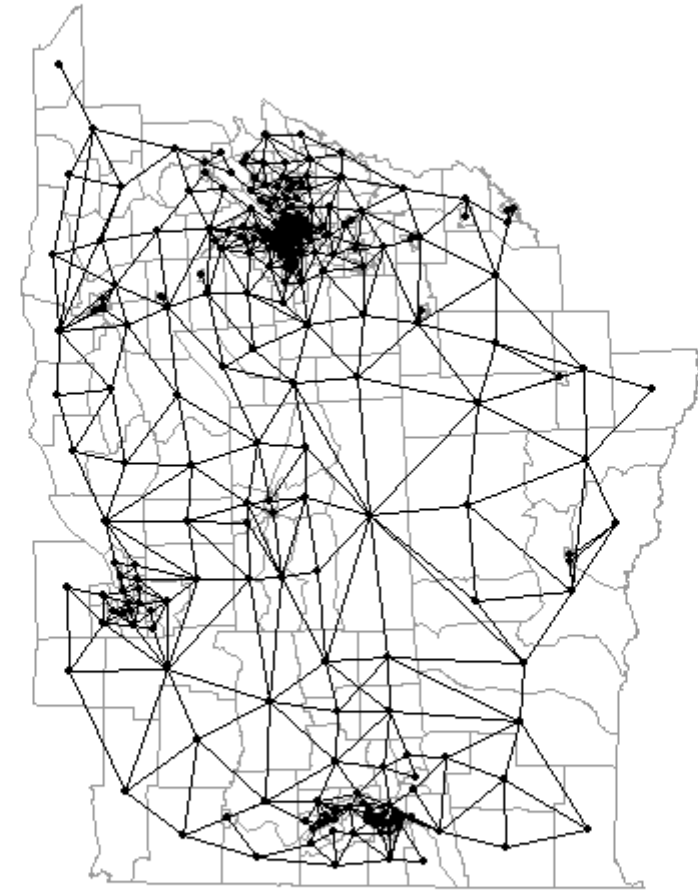
# Defining neighborhoods

- A neighborhood (or contiguity) matrix **C** represents if pairs of spatial features are to be considered as neighbors or not
- A spatial weight **W** matrix is a weighted form of such a neighborhood matrix
- **W** represents the possible spatial interactions for the selected neighborhood+ weighting approach

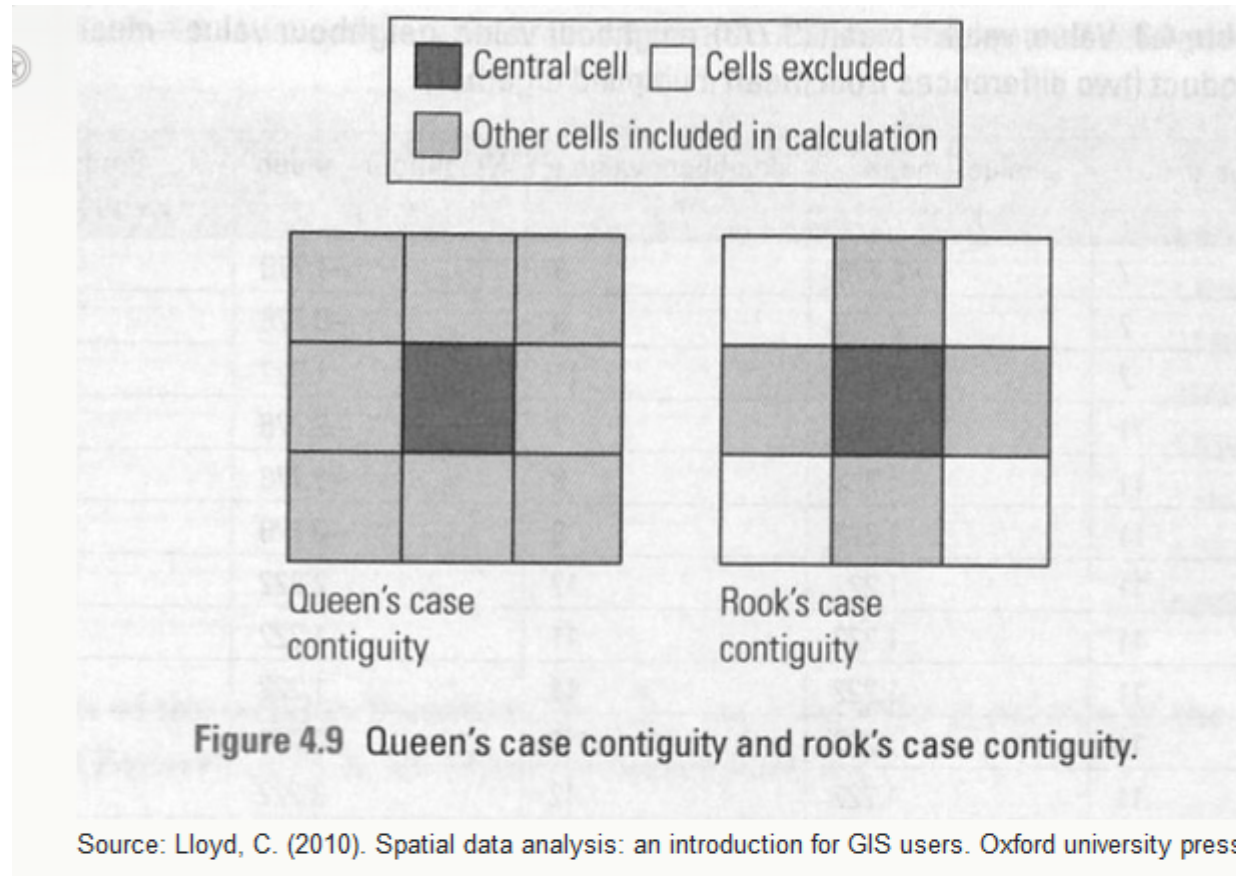


# Neighborhood

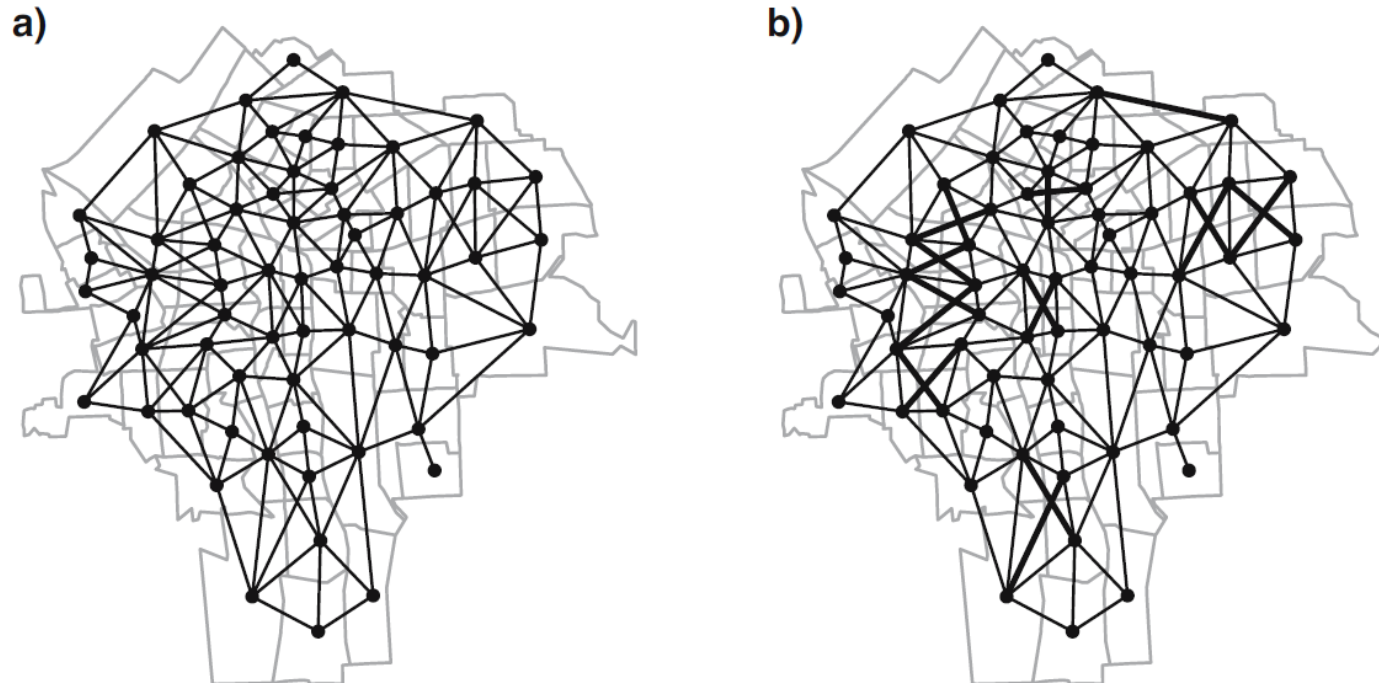
- Spatial autocorrelation depends on a defined neighborhood
- What is a neighbor?
- Several approaches possible
- Polygon or raster data: contiguity relationship
  - Polygons that share a border are neighbors
  - Rook or queen



# Queen and Rook neighborhood



# Queen and Rook neighborhood



**Fig. 9.3.** (a) Queen-style census tract contiguities, Syracuse; (b) Rook-style contiguity differences shown as thicker lines

Bivand, R.S., Pebesma, E.J., Gómez-rubio, V., 2008. Applied Spatial Data Analysis with R. Springer, New York, NY

# Other types of relationships

- Distance based
  - Based on Euclidean distance
  - K-nearest neighbors
  - Every neighbor inside of search distance
- Based on graph measures
  - Based on topological position based on centroids (or points on surface)
  - Delaunay Triangulation, Sphere of Influence, Gabriel graph, relative graph neighbors, minimal spanning tree
- Higher order neighborhood definitions possible
  - Neighbors of neighbors



# K-nearest neighbors

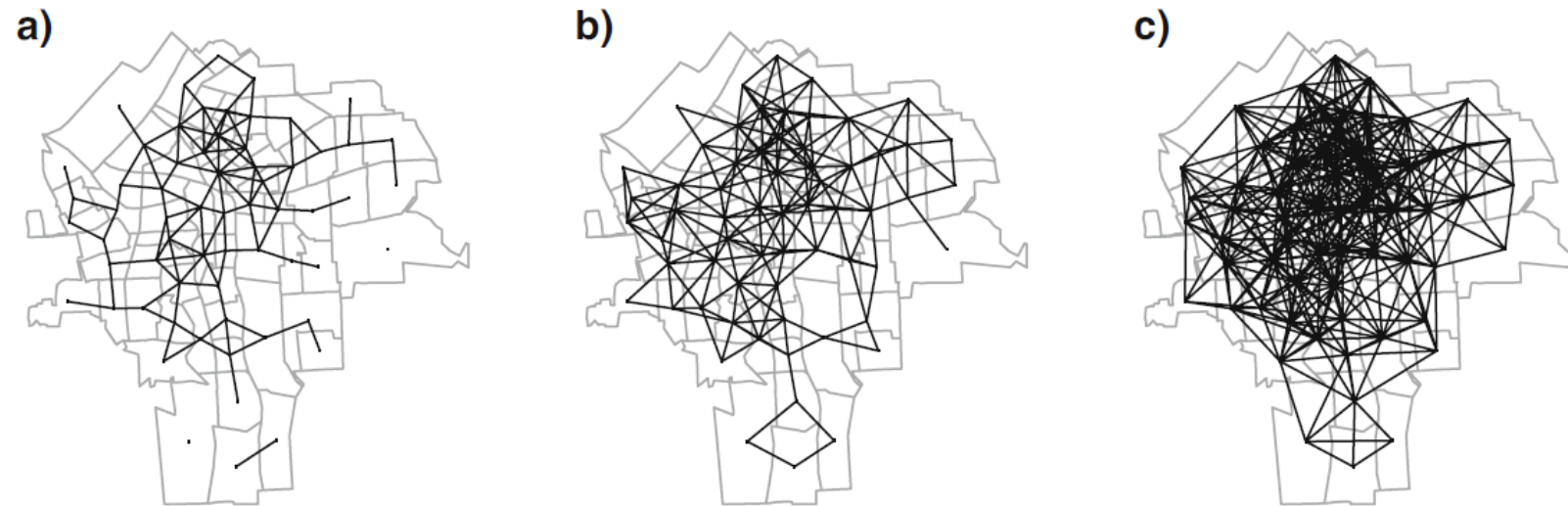


**Fig. 9.5.** (a)  $k = 1$  neighbours; (b)  $k = 2$  neighbours; (c)  $k = 4$  neighbours

Bivand, R.S., Pebesma, E.J., Gómez-rubio, V., 2008. Applied Spatial Data Analysis with R. Springer, New York, NY



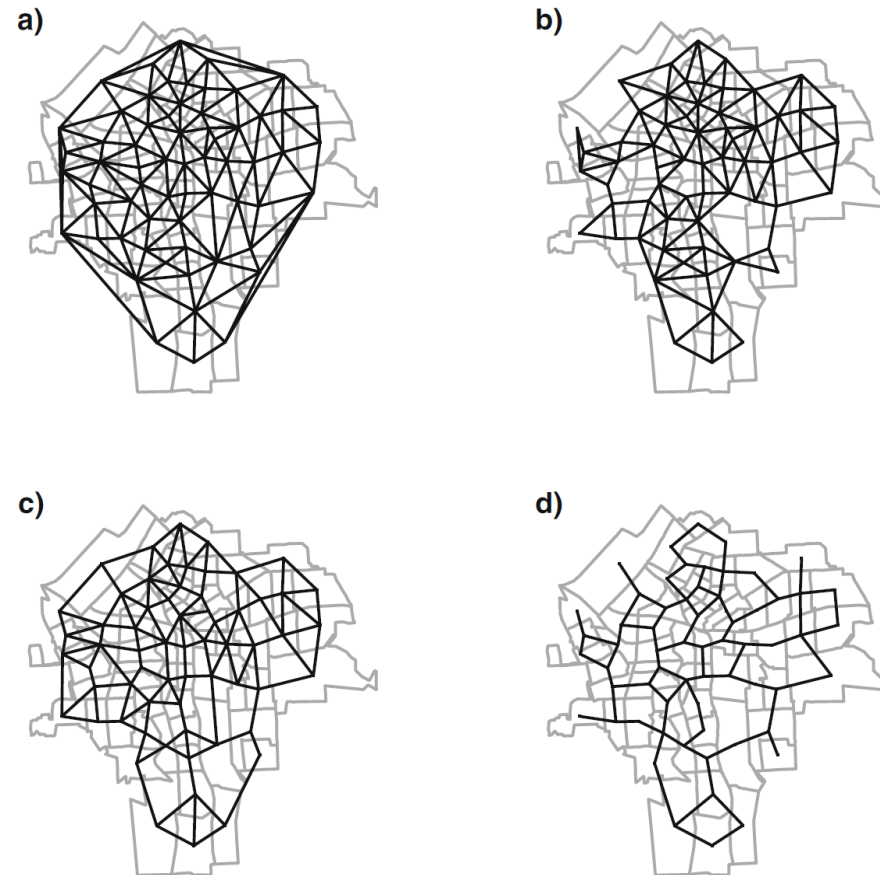
# Distance threshold based neighbors



**Fig. 9.6.** (a) Neighbours within 1,158 m; (b) neighbours within 1,545 m; (c) neighbours within 2,317 m



# Topological neighborhood definitions



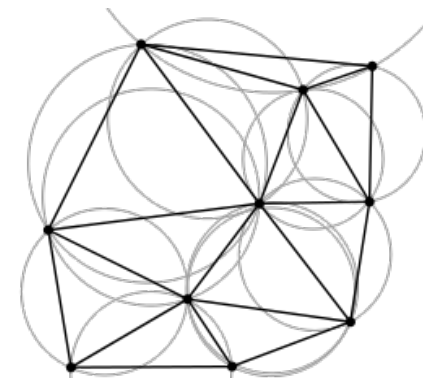
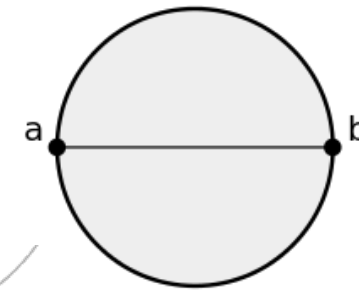
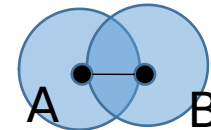
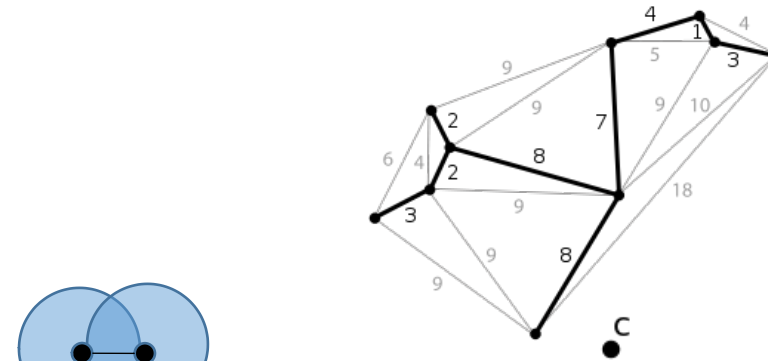
**Fig. 9.4.** (a) Delauney triangulation neighbours; (b) Sphere of influence neighbours; (c) Gabriel graph neighbours; (d) Relative graph neighbours





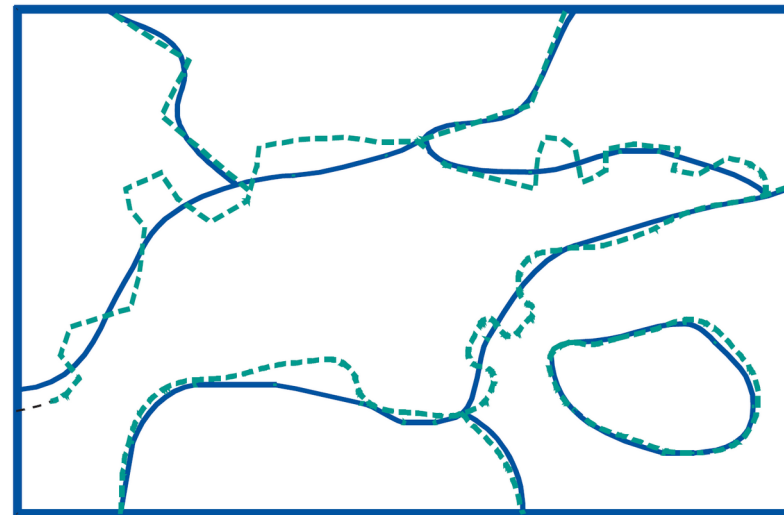
# Topological neighborhood definitions

- Minimum spanning tree connects all nodes together while minimizing total edge length
- Relative neighborhood graph all nodes connected for those the lens formed by the radii of their circles contains no other points
- Gabriel graph all nodes connected if there is no other node inside a circle with their distance
- Delaunay triangulation all nodes connected for which the circumcircle around ABC contains no other nodes



# Selecting a neighborhood – things to consider

- GIS precision might lead to artifacts in neighborhoods for contiguity definition
  - Sliver polygons etc.
  - Fix it beforehand



# Selecting a neighborhood – things to consider

- Is Euclidean distance important for the assumed process?
  - If you don't know you could investigate using different approaches
  - If distance is important we might assign weights based on distance. Use different approaches to investigate importance of distance



## Selecting a neighborhood – things to consider

- Does the definition lead to objects without neighbors?
  - Artifacts or true islands?
  - Problematic for analysis
  - Exclude or set zero.policy
  - If zero policy is set to TRUE, weights vectors of zero length are inserted for regions without neighbor in the neighbors list
  - These will in turn generate lag values of zero
  - The spatially lagged value of  $x$  for the zero-neighbor region will then be zero, which may (or may not) be a sensible choice



# Assigning weights

- How to deal with uneven neighborhood distributions?
- Binary coding  $B$  – neighbor (1) or not (0)
  - Objects with many neighbors get more weight
- **Row standardized coding**  $W$  : weight (0 or 1) divided by the number of neighbors, i.e. weights sum to unity for each object (sums over all links to  $n$ )
  - Mostly used in practice, assumed for some approaches, the effect of the neighbors is expressed as the weighted sum
- Globally standardized coding  $C$ : weight (1 or 0) divided by sum of all weights, i.e. sum to  $n$  for all objects
- $U$ : equal to  $C$  divided by the number of neighbors (sums over all links to unity)
- $S$ : variance-stabilizing coding scheme proposed by Tiefelsdorf et al. 1999, p. 167-168 (sums over all links to  $n$ ). Between  $W$  and  $C$



## Selecting a neighborhood – things to consider

- Neighborhood frequency distribution even or skewed?
  - Contiguity -> typically relatively equal distribution
  - Distance based -> skewed
  - K-nearest -> equally distributed
- K-nearest does **not** lead to a symmetric relationship!
- Some analysis require symmetric relationships
  - It is possible to make a nb relationship symmetric by adding neighbors (*make.sym.nb*)
- It is also possible to use set operations on nbs (intersect, union, difference, complement) or to manually modify nbs



## Selecting a weighting scheme - things to consider

- Beside coding style  $U$  all coding styles sum to  $n$  over all links  $\rightarrow$  estimated spatial auto-correlation should be comparable
- In addition, weights can be inversely weighted by distance prior to standardization
- $S$  and  $W$  might lead to non-symmetric weight matrix



# Spatial exploratory analysis



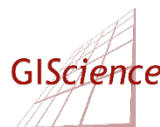


# What are we looking for?

- DATA = SMOOTH + ROUGH (Tukey, 1975)
- SPATIAL DATA = SPATIAL SMOOTH + SPATIAL ROUGH
-

# What are we looking for?

- Spatial smooth
  - Presence of spatial trend?
  - Spatial heterogeneity – is variation in data values as smooth as implied by some trend?
  - Global spatial dependence – are high/low values close to other high/low values, anywhere on the map?
  - Spatial heterogeneity – localized patterns of dependence? Hot-/coldspots?



# What are we looking for?

- Spatial rough
  - Presence of outliers?
    - In addition to distributional outliers: are there spatial outliers (which might not be distributional outliers)
    - i.e. are there some data points special with respect to their neighborhood?



# Moran's I

- Measures global spatial auto-correlation (for all observations)
- $I > 0$  positive auto-correlation, clustered
- $I < 0$  negative auto-correlation, dispersed
- Expected value under absence of spatial auto-correlation  $E(I) = -n(n-1)^{-1}$

$$I = \frac{1}{\sum_{i=1}^n \sum_{j=1, j \neq i}^n w_{i,j}} \frac{\sum_{i=1}^n \sum_{j=1, j \neq i}^n w_{i,j} (x_i - \bar{x})(x_j - \bar{x})}{\frac{1}{n} \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

- $w_{i,j}$  – weights for observation  $i$  and  $j$  from the weight matrix  $W$
- $x_i$  – value of observation  $i$

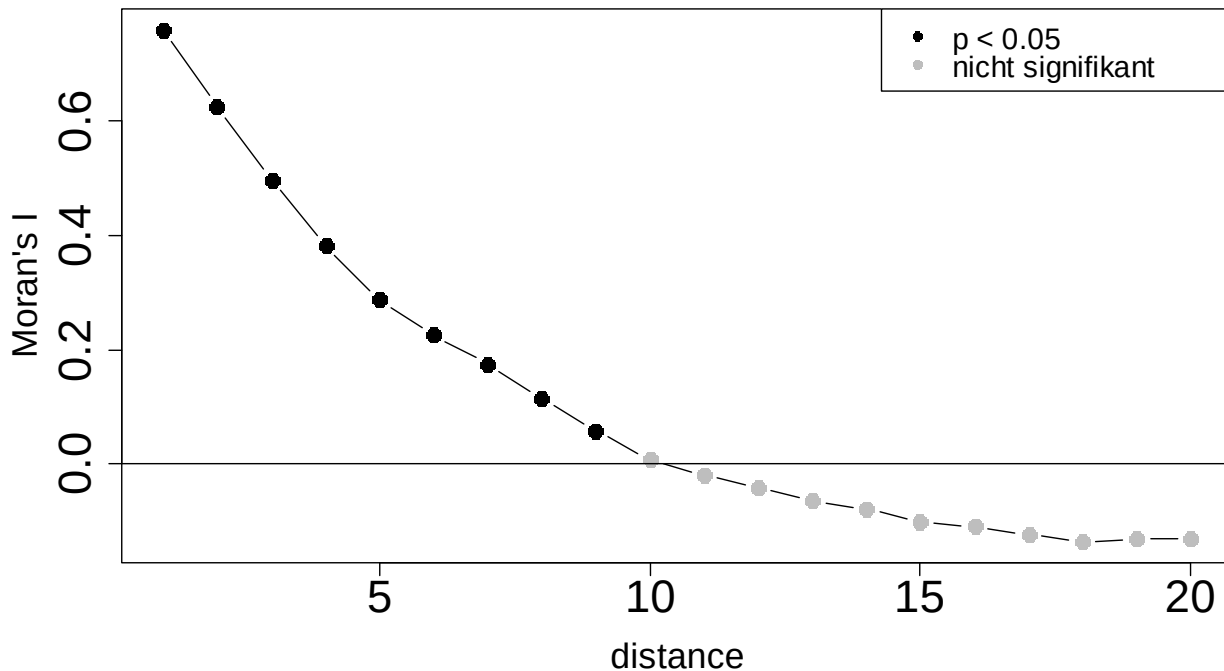
# Moran's I

- Not bound to the interval  $[-1,1]$ 
  - Might shrink but typically expand
  - Interval defined by the largest and second largest eigenvalue of the weight Matrix  $W$
  - Often the interval is more in the range  $-0.5$  to  $1.15$
  - The expected value for no autocorrelation is not zero but  $E(I) = -1/(n-1)$
- For regression residuals a modified test statistic is used
- For rates the empirical Bayes index modification is used



# Correlogram

- Plot of Moran's I against distance classes
  - Significance testing via permutation under  $H_0$

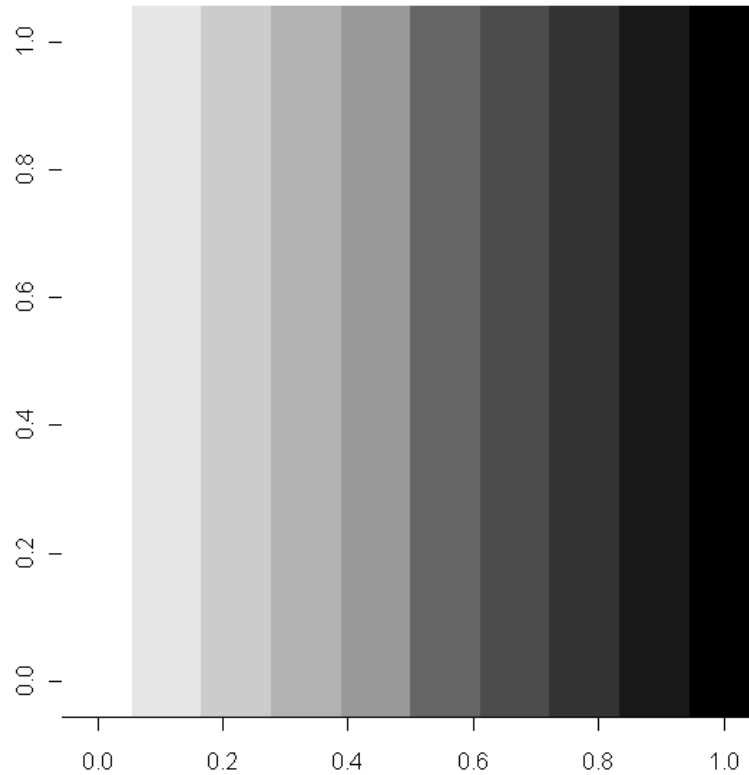


$$I(d) = \left( \frac{1}{W(d)} \right) \frac{\sum_{\substack{i=1 \\ i \neq j}}^n \sum_{\substack{j=1 \\ j \neq i}}^n w_{i,j}(d) (x_i - \bar{x})(x_j - \bar{x})}{\frac{1}{n} \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

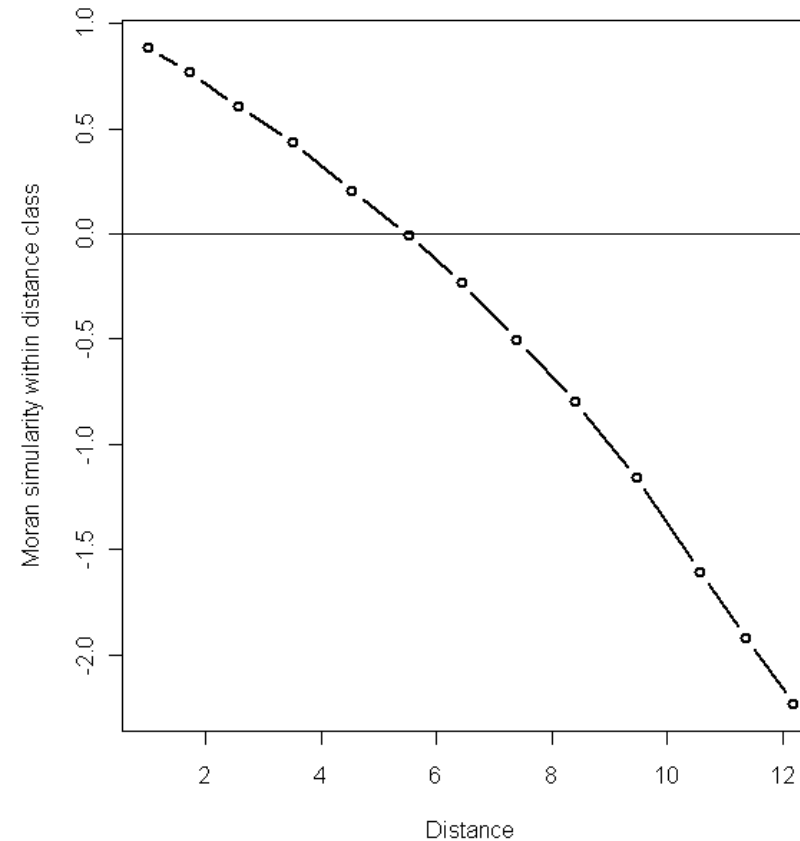
D – Distance class  
 $w_{i,j}$  – neighbor or not  
 $W(d)$  – sum over all  $w_{i,j}$

# Examples - Gradient

Value raster

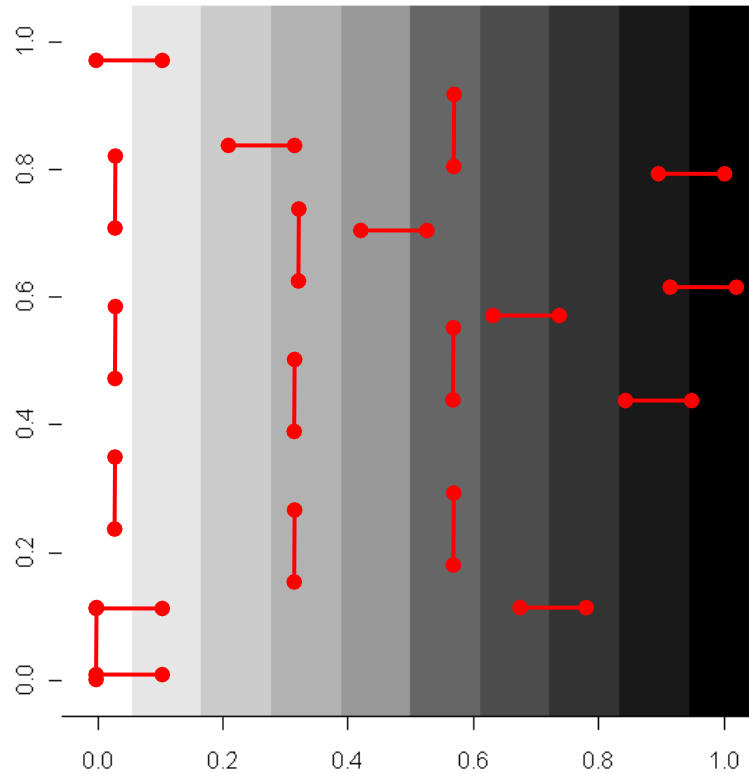


Correlogram that shows how spatial autocorrelation changes with distance

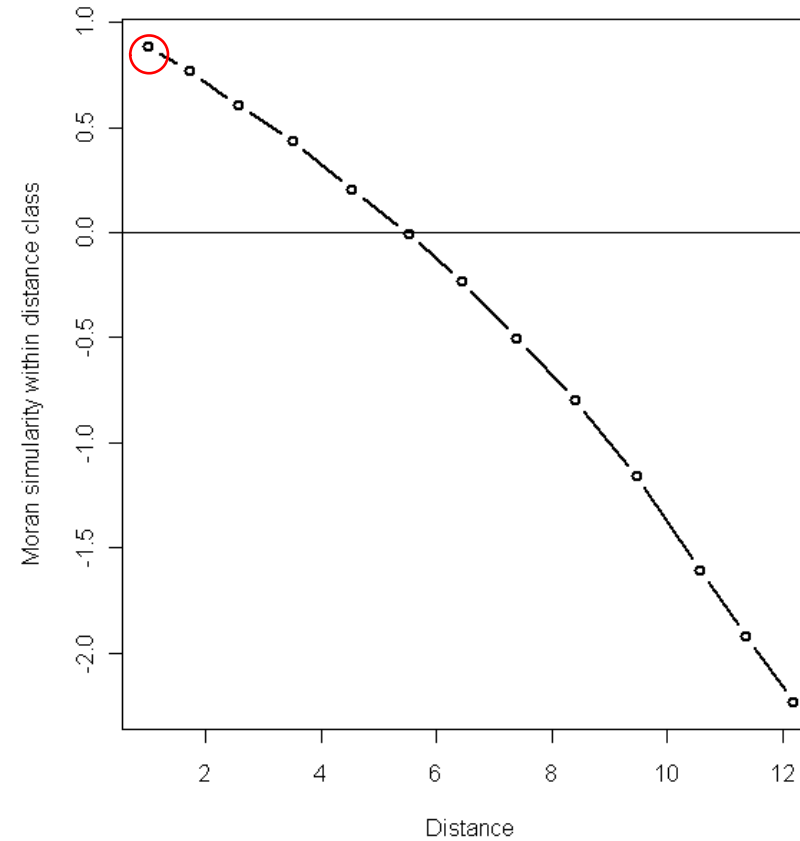


# Example - Gradient

Some pairs of distance class 1



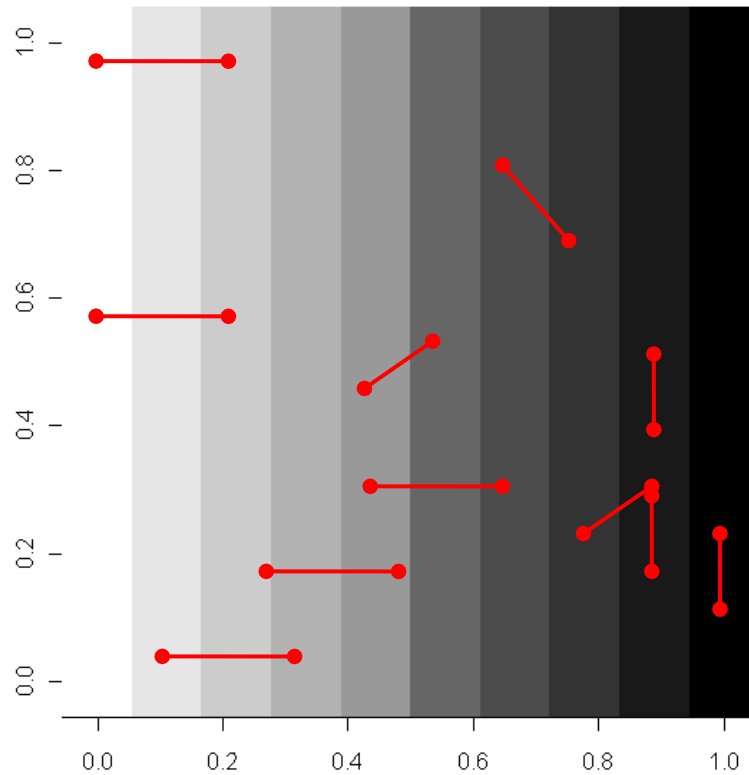
Moran's I for all pairs of distance class 1



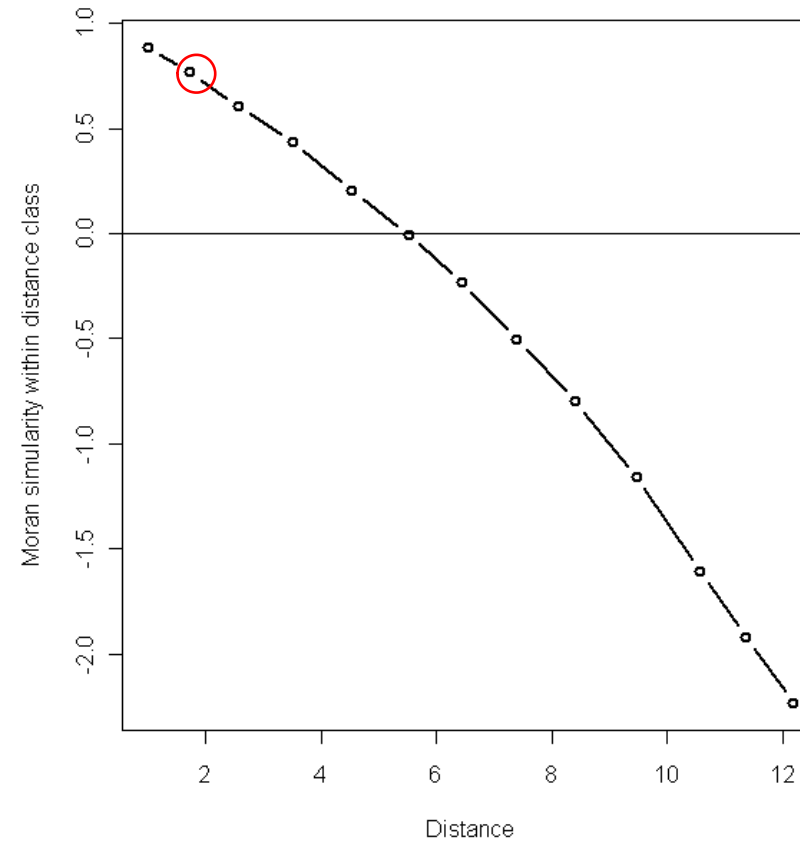


# Example - Gradient

Some pairs of distance class 2

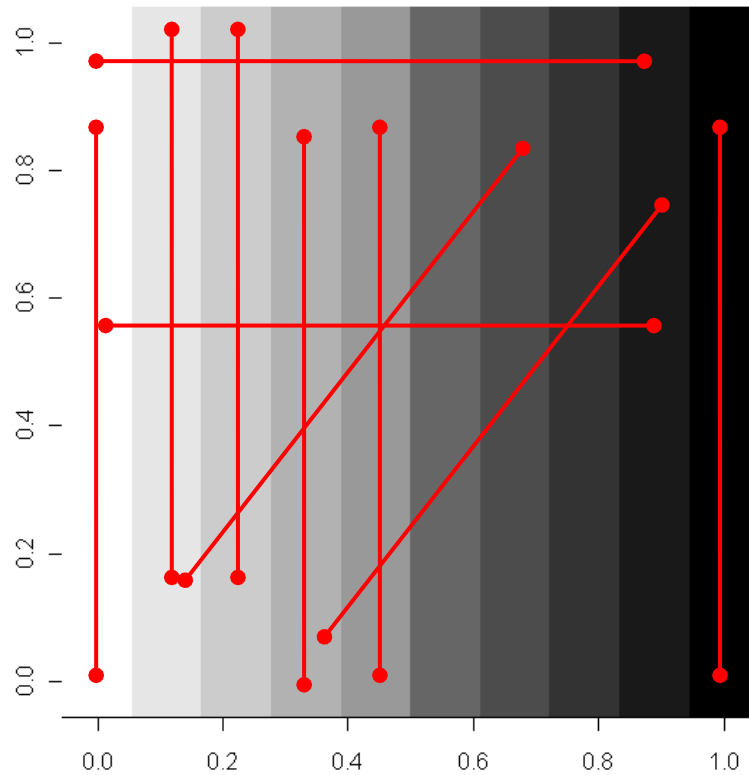


Moran's I for all pairs of distance class 2

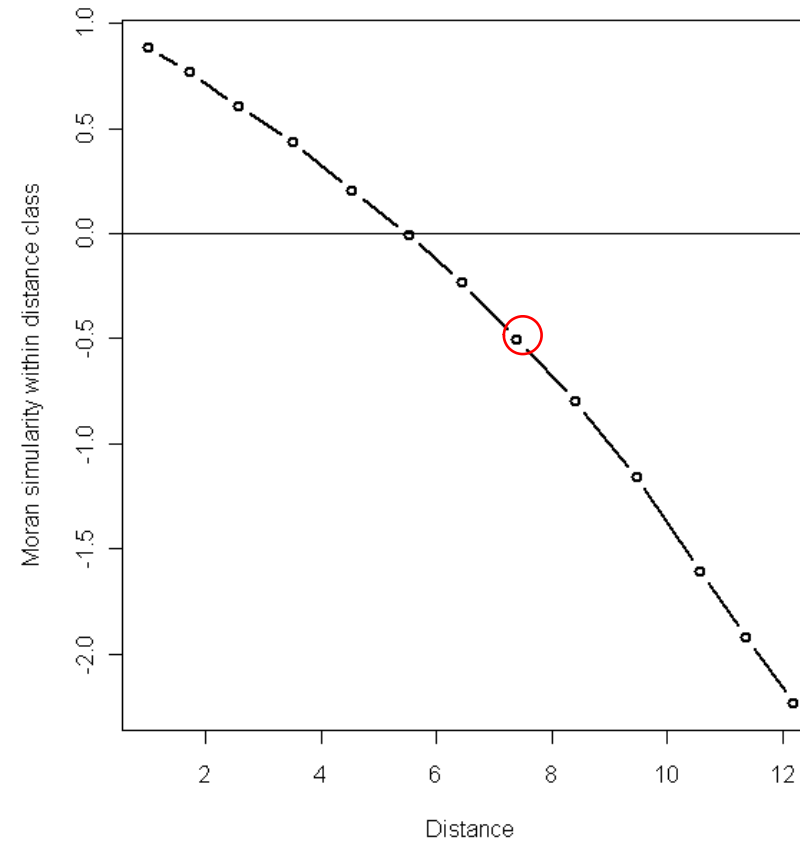


# Example - Gradient

Some pairs of distance class 8

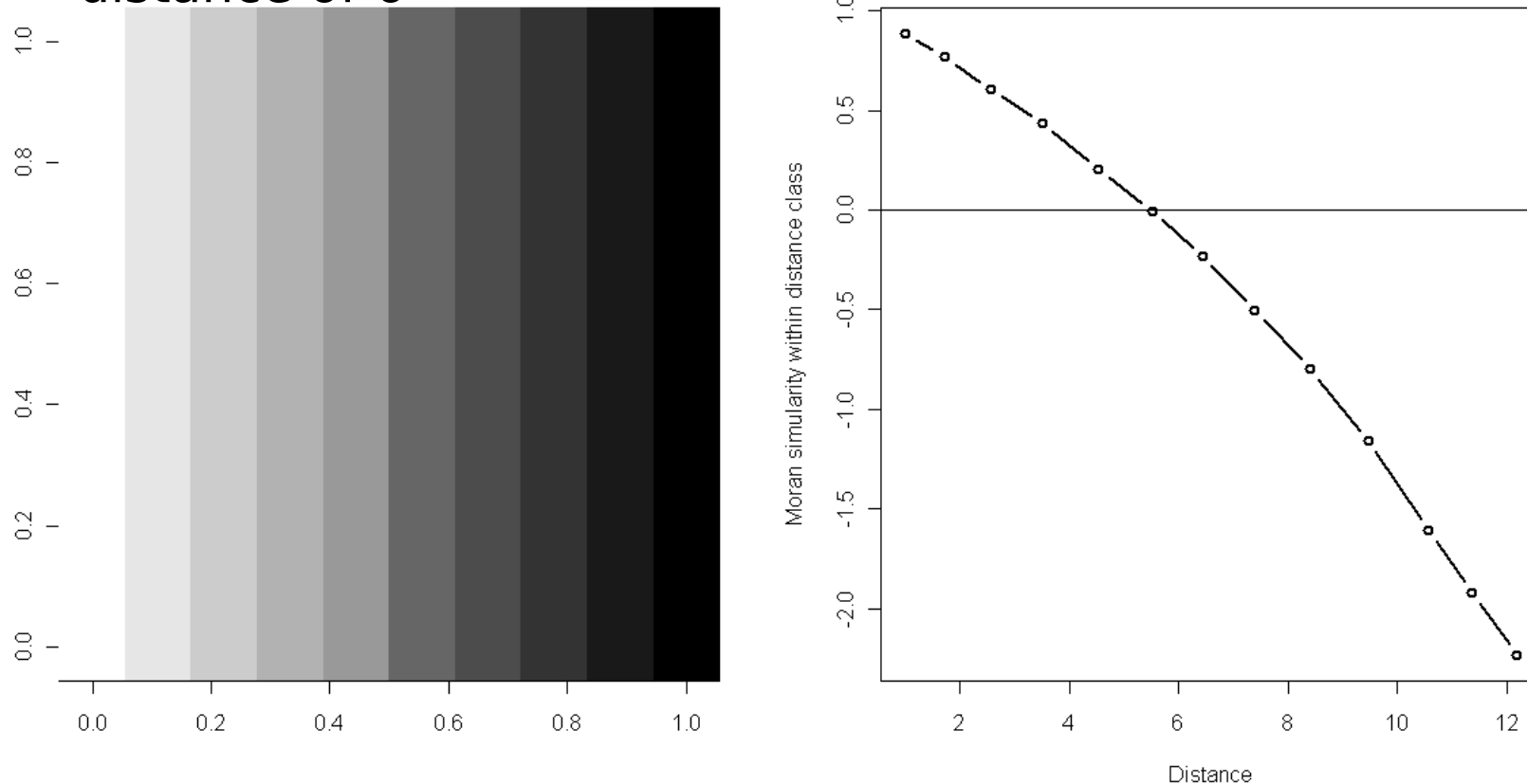


Moran's I for all pairs of distance class 8



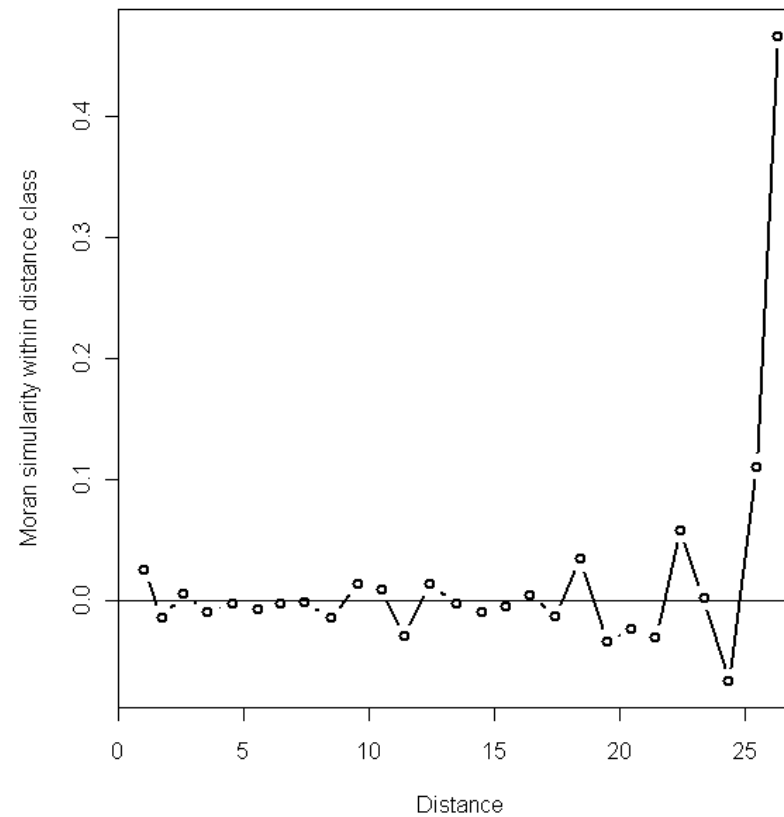
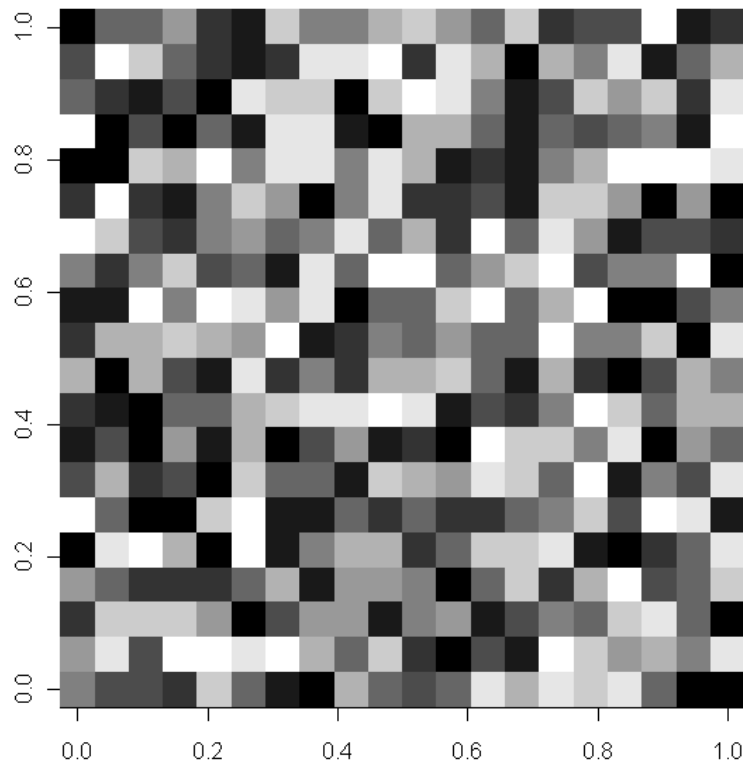
# Example - Gradient

Close observations are very similar. Effect diminishes with increasing distance and turns negative at a distance of 6



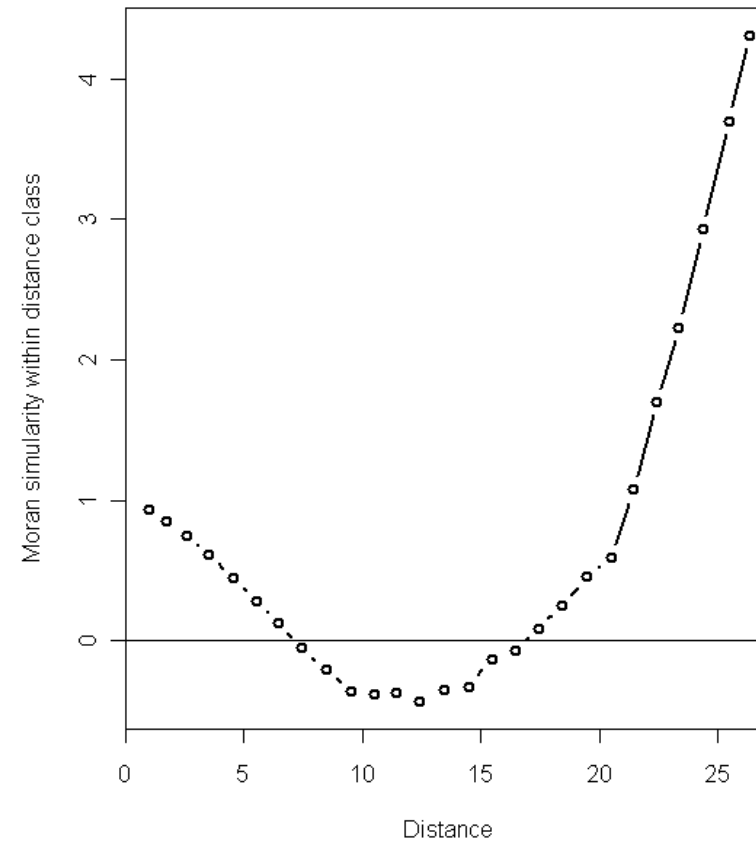
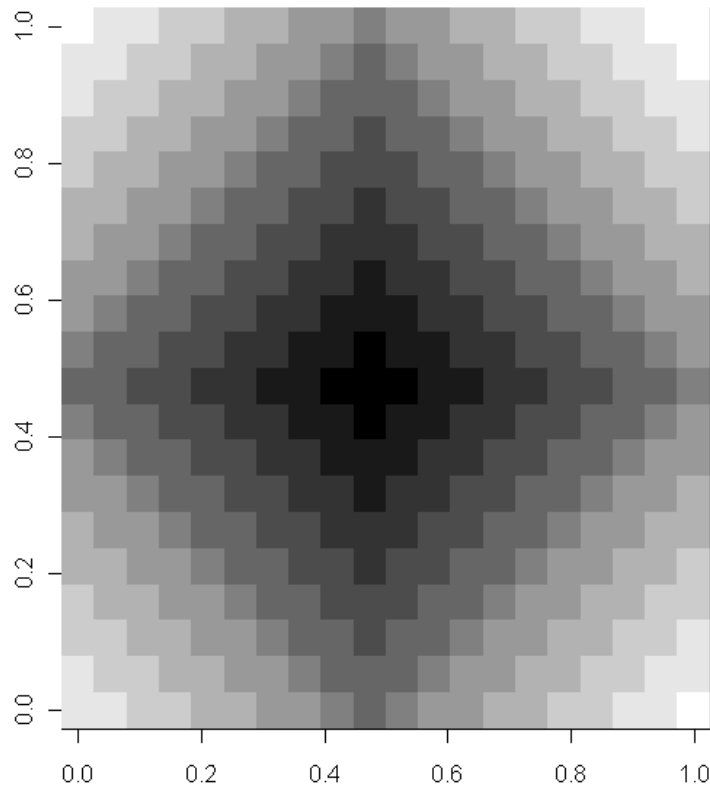
# Random pattern

Peak at distance class 26 is an artefact – based on only a few observation. Therefore, large distance classes should not be plotted!

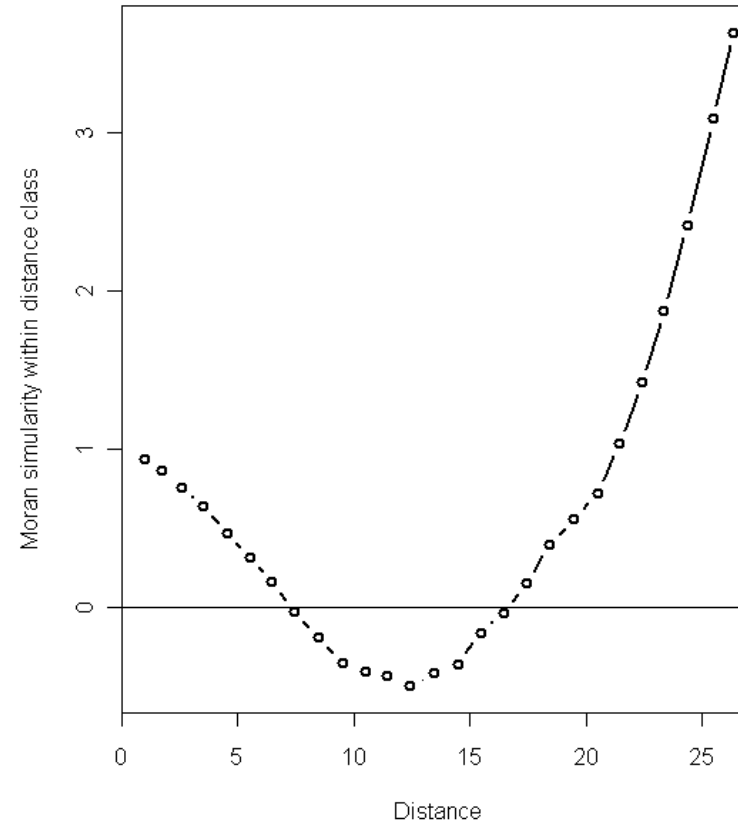
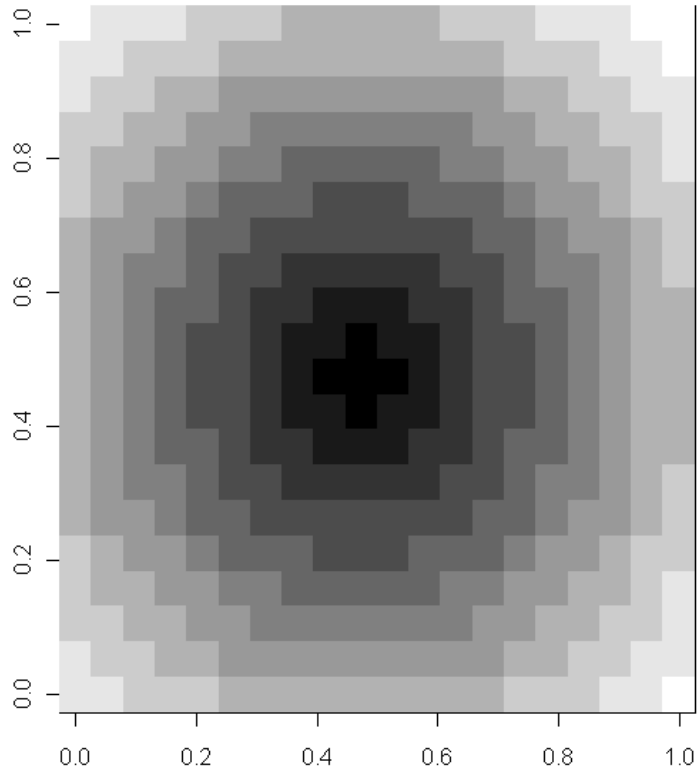


# Concentric pattern

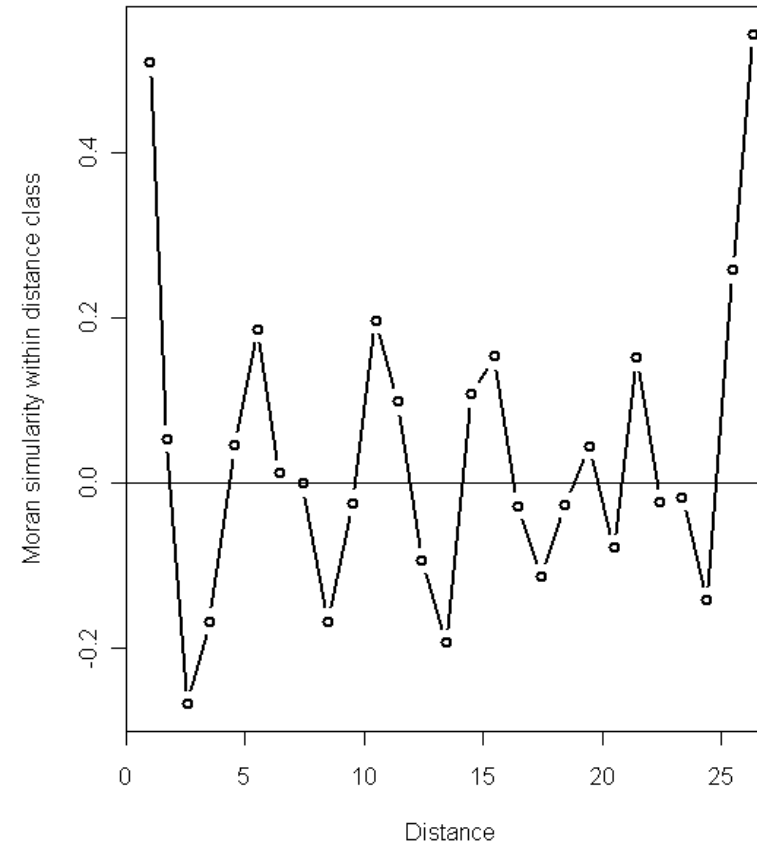
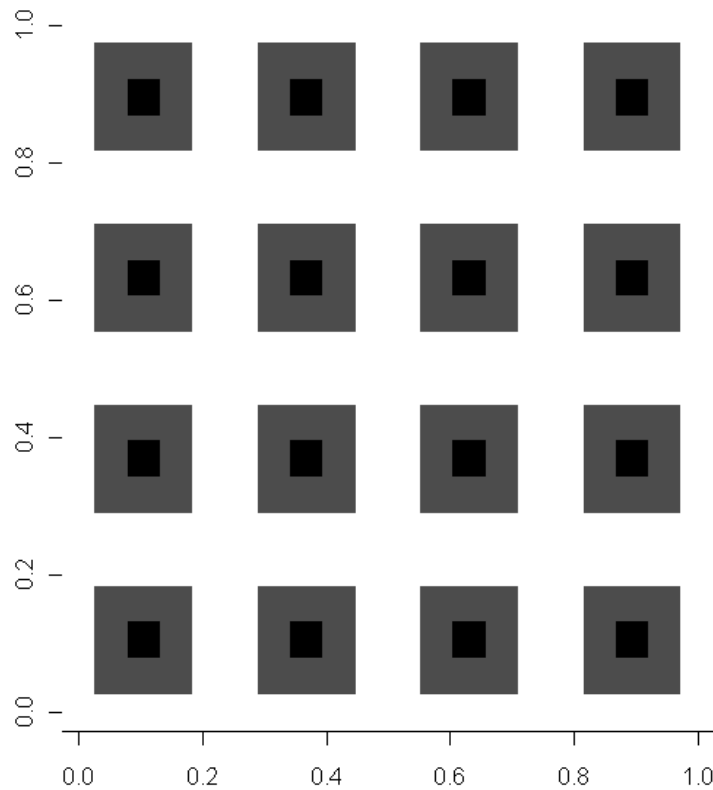
Pattern for distance classes larger than 15 are not reliable here



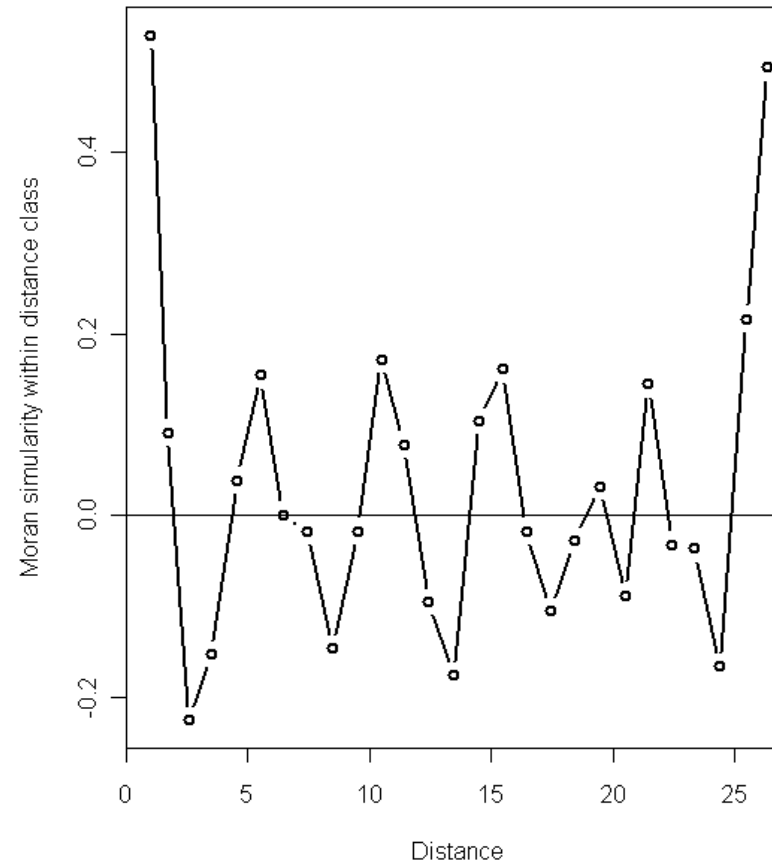
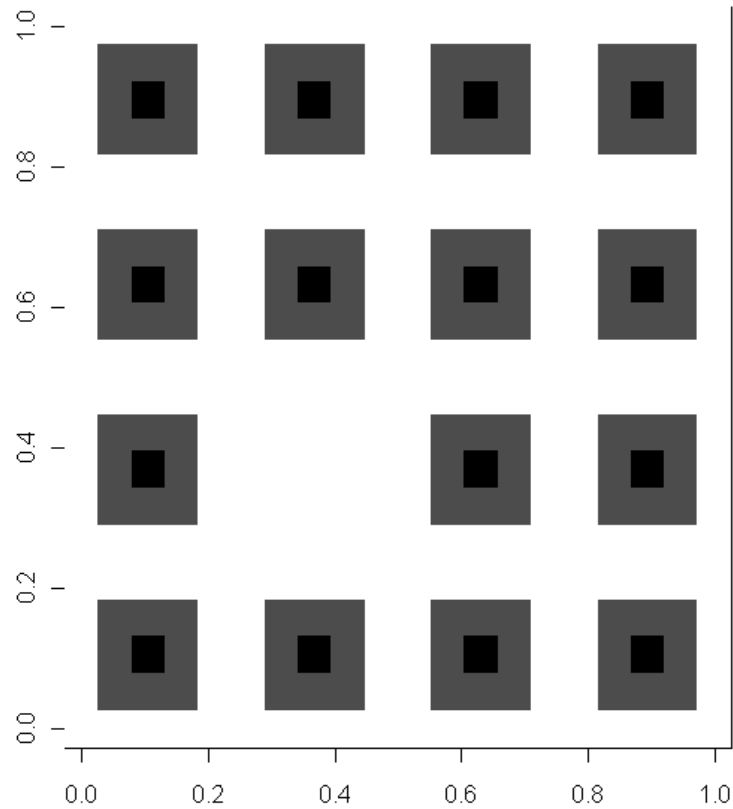
# Concentric pattern



# Patchy pattern

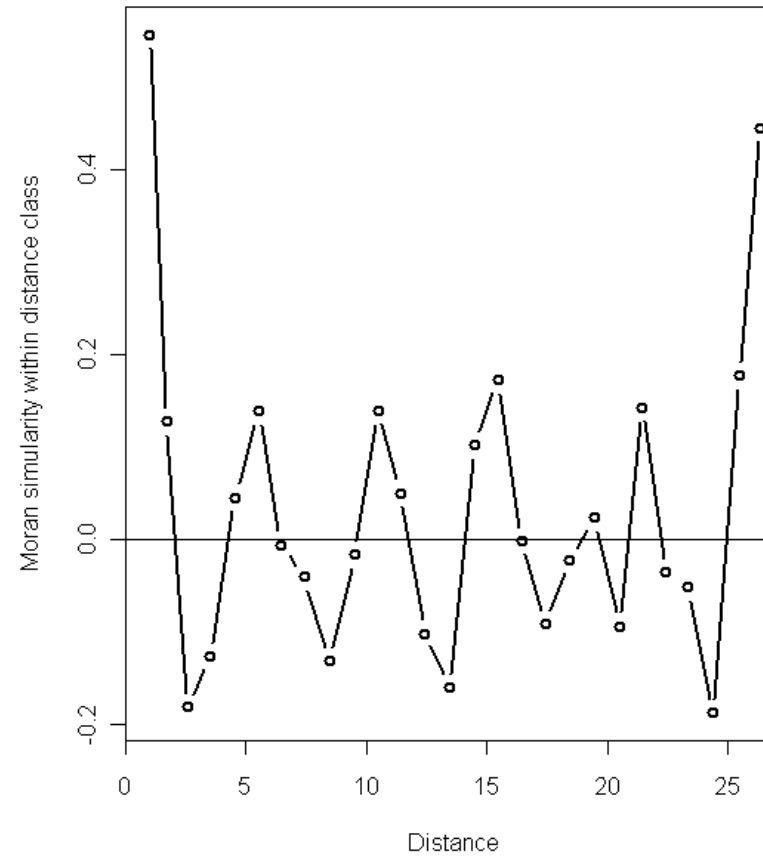
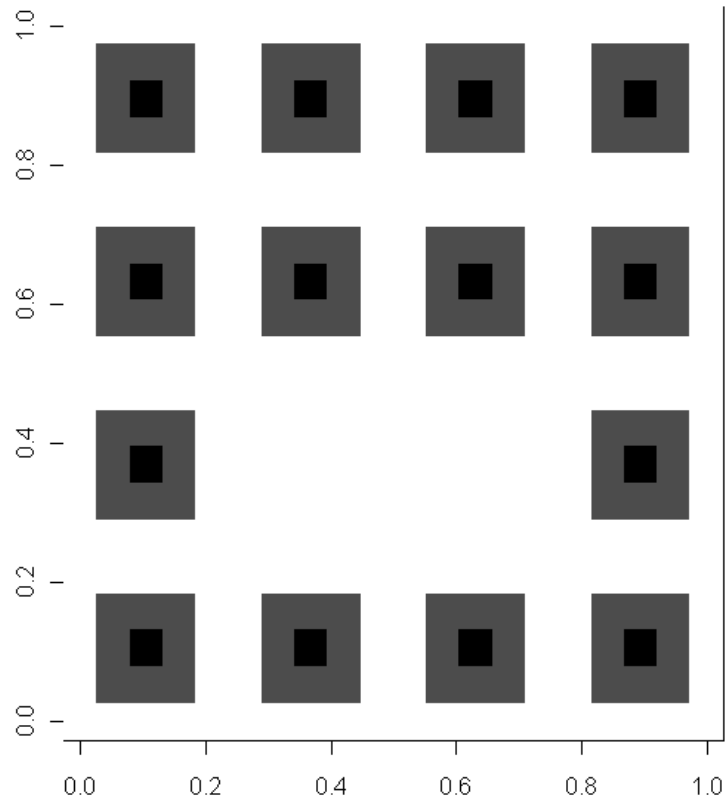


# Patchy pattern with holes

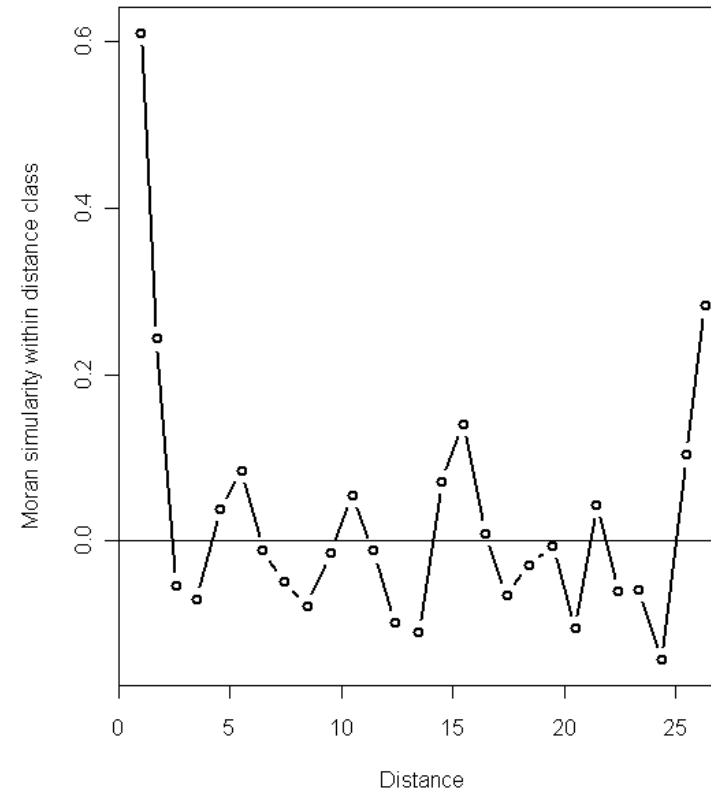
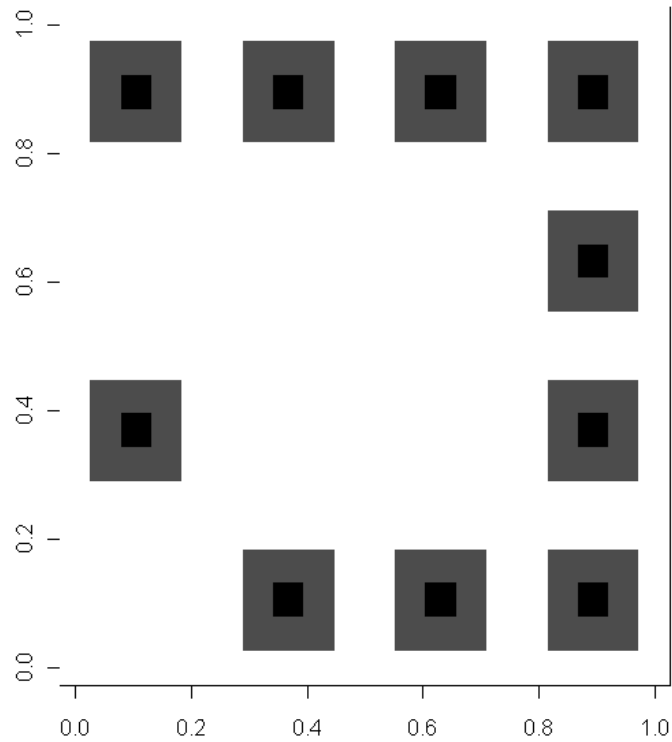




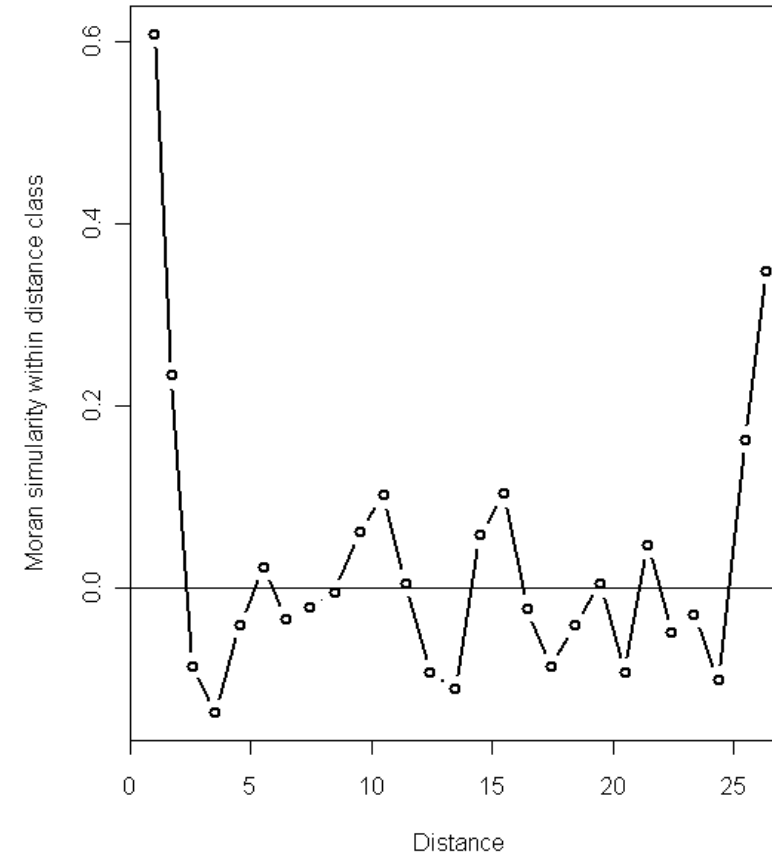
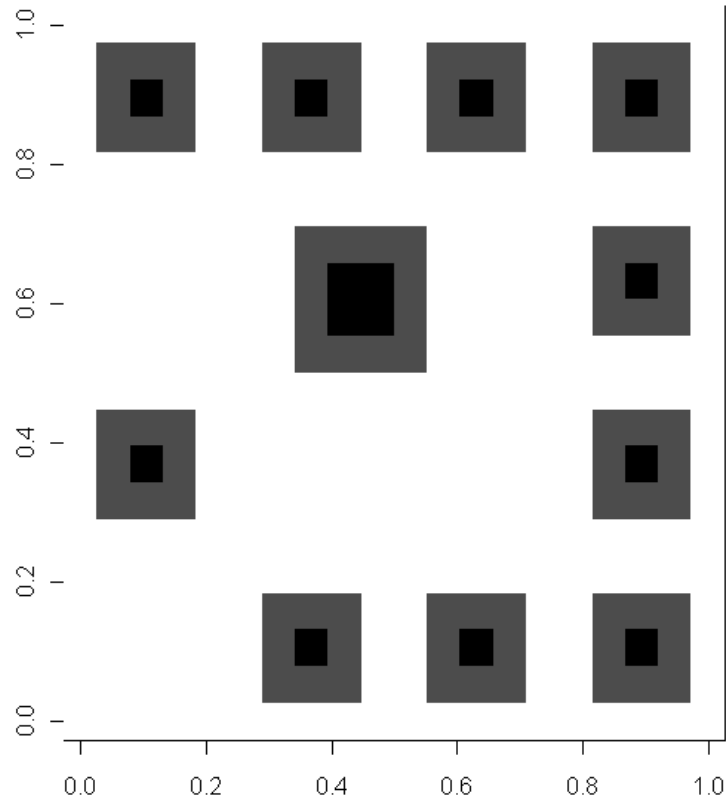
# Patchy pattern with holes



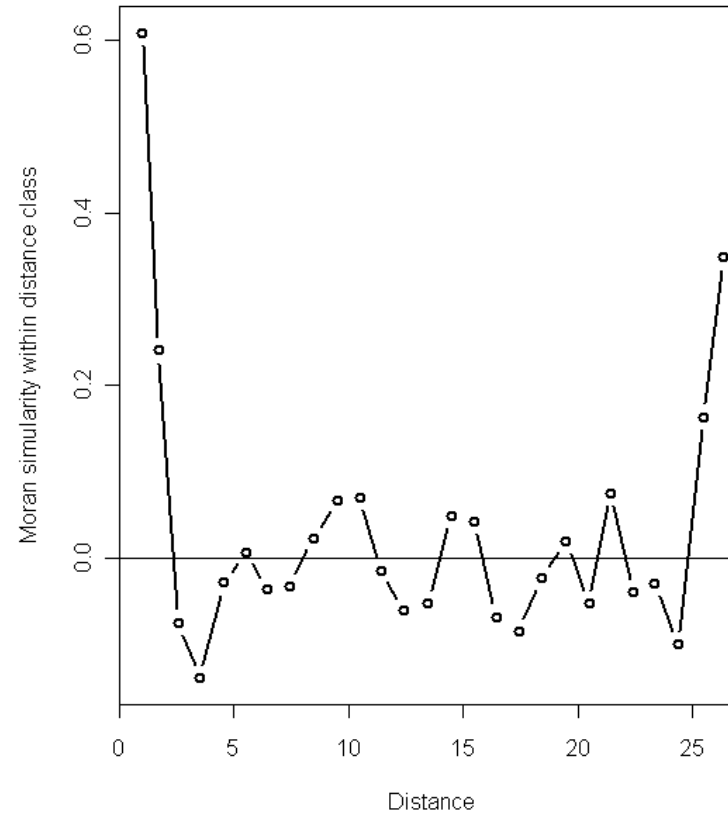
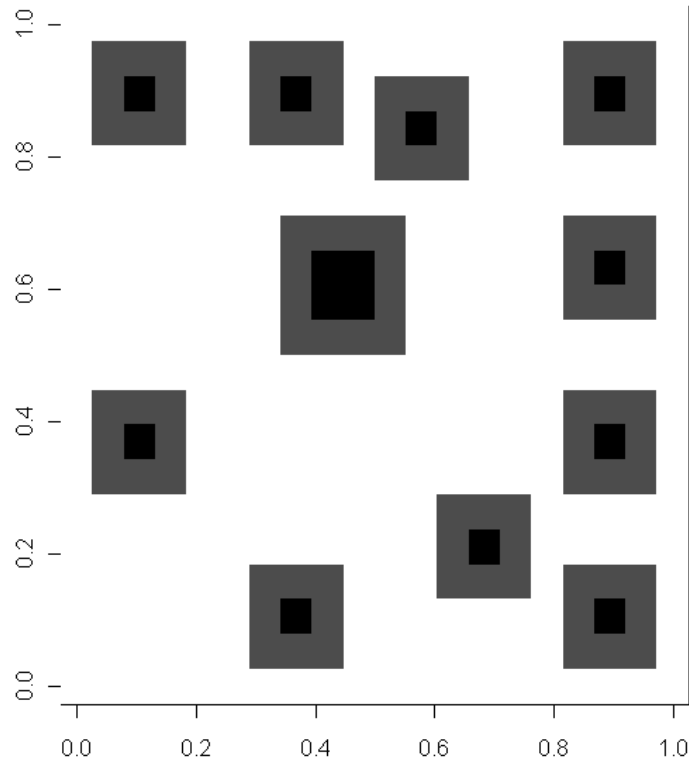
# Patchy pattern with holes



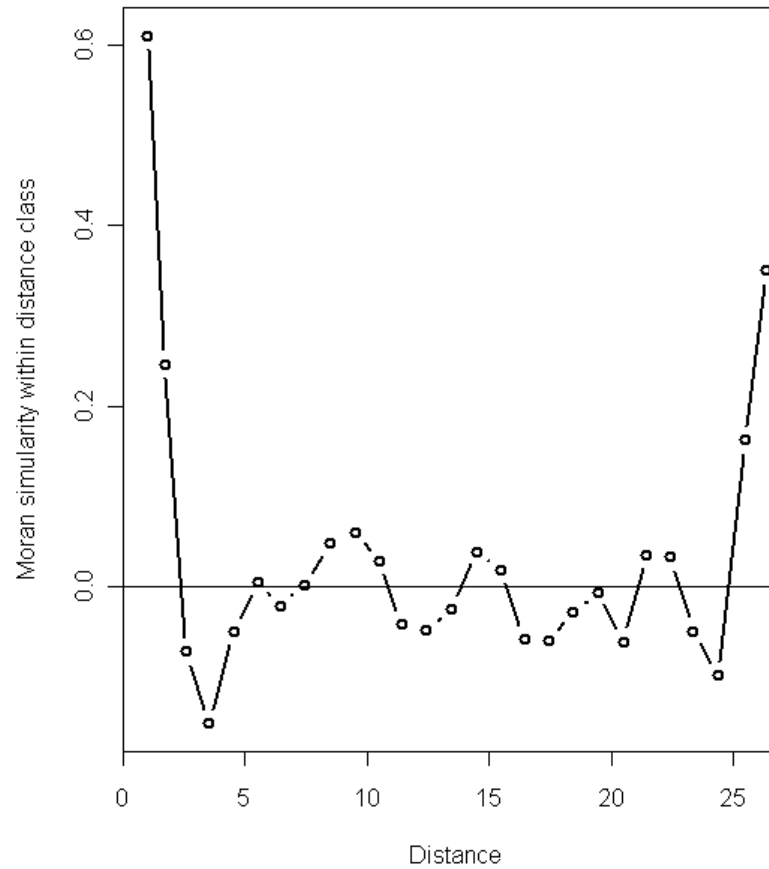
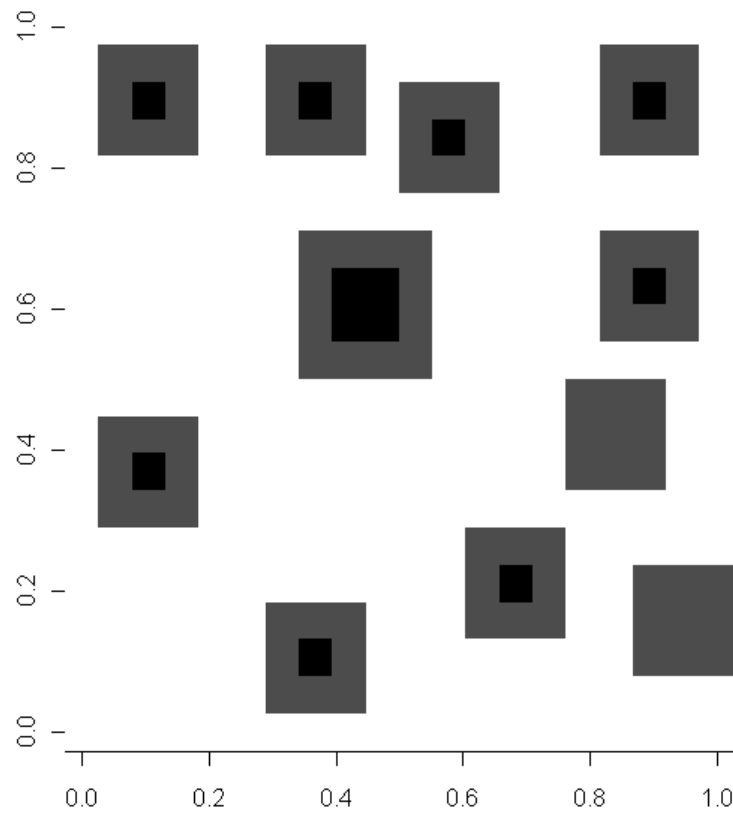
# Patchy pattern with holes, different size of patches



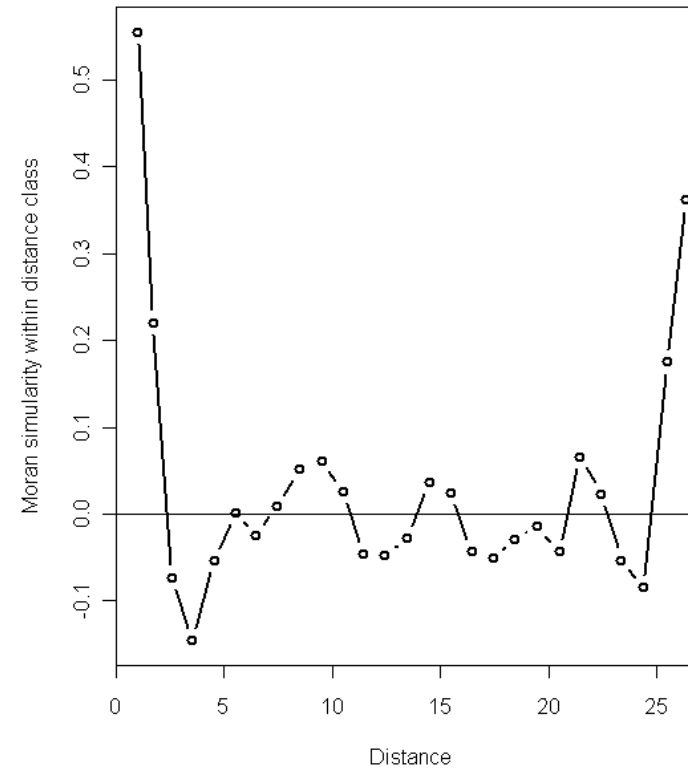
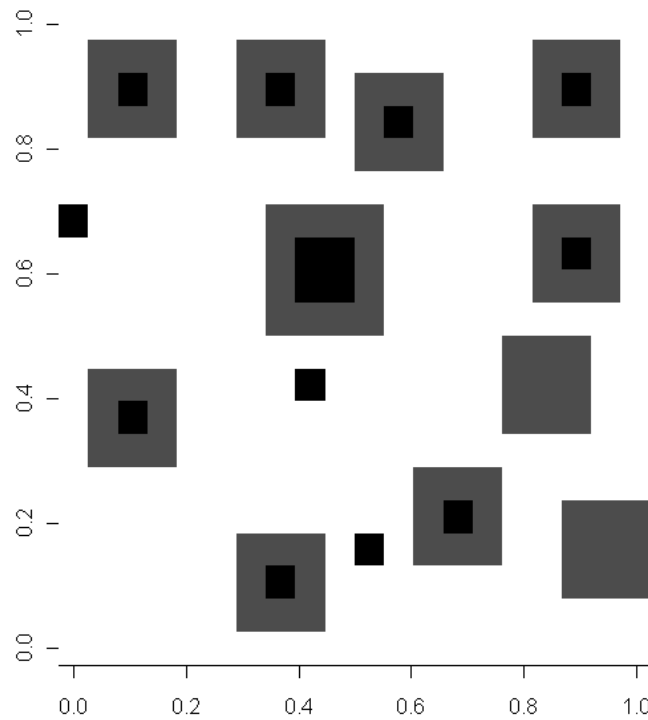
# Patchy pattern with holes, different size of patches



# Patchy pattern with holes, different size of patches



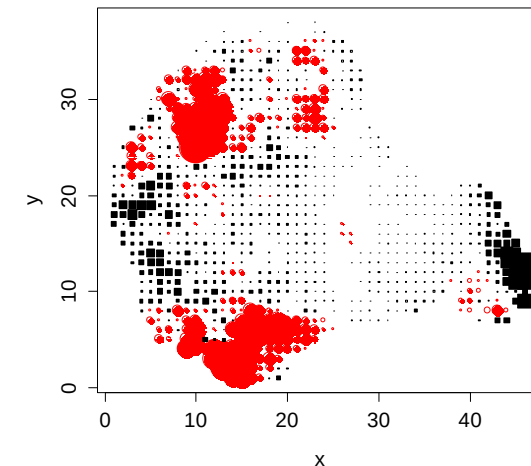
# Patchy pattern with holes, different size of patches



# LISA - local indicator of spatial association

- E.g. local Moran's I
- Spatial auto-correlation based on a focal function
- Takes in stationarity into account

$$I_i = n(x_i - \bar{x}) \frac{\sum_{i \neq j}^n w_{i,j} (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})}$$



# Local Moran's I

- Calculates Moran's for the neighborhood of a point and relates it to the mean over all data

Mean for all observations

$$I_i = n(x_i - \bar{x}) \frac{\sum_{i \neq j}^n w_{i,j} (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})}$$

$$I = \frac{1}{n} \sum_{i=1}^n I_i$$

No double sum as in global Moran's I  
Only sum over the observations in the neighborhood



# Critic local Moran's I

- Direction of change is not available
  - High values at location  $i$ , surrounded by high values in neighborhood -> positive value
  - Low values at location  $i$ , surrounded by low values in neighborhood -> positive value
  - High values at location  $i$ , surrounded by low values in neighborhood -> negative value
  - Low values at location  $i$ , surrounded by high values in neighborhood -> negative value



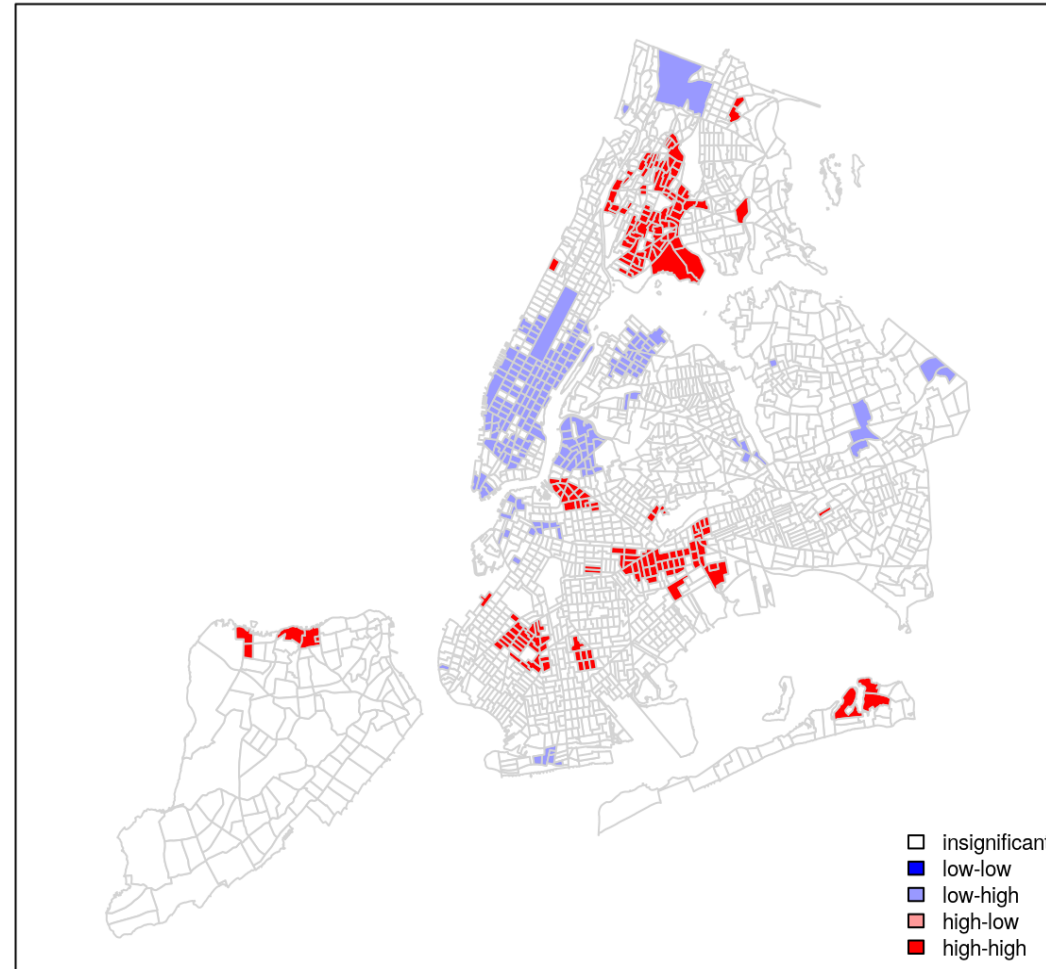
# Local Moran's I

- Therefore often the value of the observation is reported together with the average value in the neighborhood
- High - high - cluster of high values
- Low - low - low values (low relative to the global mean) in a neighborhood of high values
- High -low - high value (compared to the mean) in a neighborhood of negative spatial autocorrelation, i.e. high local outlier
- Low - high - low value (compared to the mean) in a neighborhood of positive spatial autocorrelation, i.e. a cluster of low values



# Local Moran's I

LISA on percent inhabitants younger than 21, SOI nb

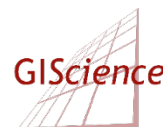


# Getis G and G\*

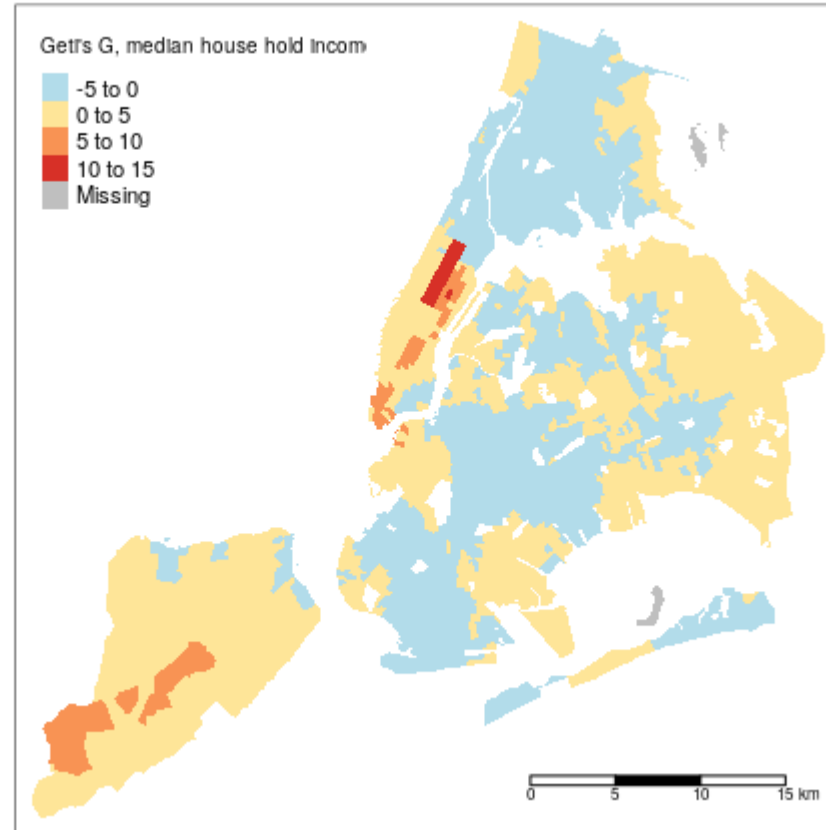
- Relation between the local mean at location  $i$  to the global mean (local moving average, hotspot/coldspot)
- Sign indicates direction of deviation
- $G_i^*$  and  $G_i$  differ by that  $G_i^*$  does include the point itself

$$G_i = \frac{\sum_{j=1, j \neq i}^n w_{ij} X_i X_j}{\sum_{j=1, j \neq i}^n X_i X_j}$$

$$G_i^{star} = \frac{\sum_{j=1}^n w_{ij} X_i X_j}{\sum_{j=1}^n X_i X_j}$$



# Getis G and G\*

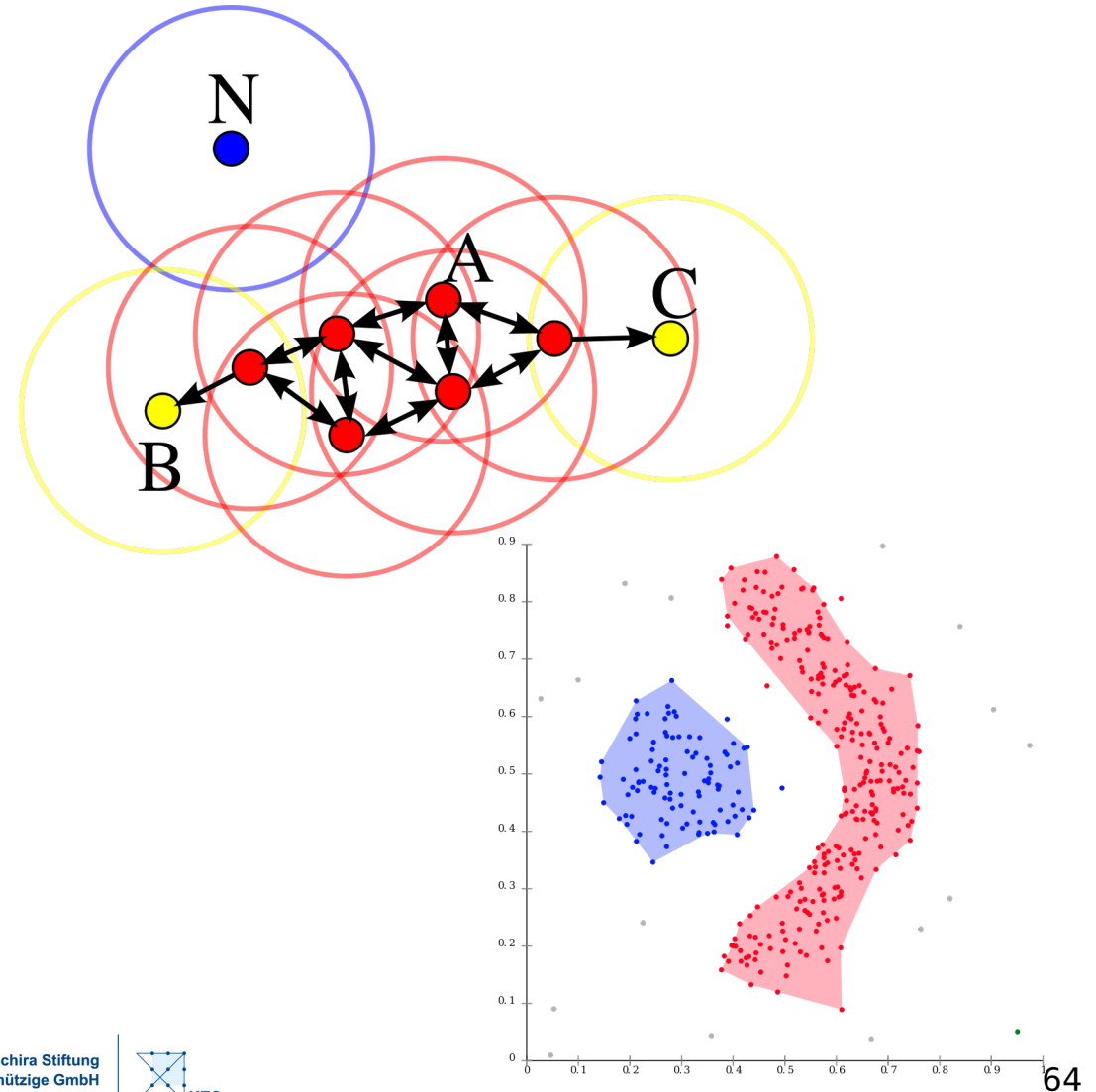


# Spatial cluster identification



# dbScan

- Aims to identify spatially clustered points
- A cluster is defined as the subset of points that can be reached from all points in the cluster by a distance less than a threshold
  - In addition to core points there are points at the edges that can be reached from the cluster
  - Core points: at least minP points in distance
- Point outside clusters are considered noise points



# Adjusting for spatial autocorrelation in regression models





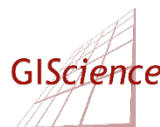
# How to deal with it?

- Incorporate it in additional covariates
  - Capture the spatial configuration in additional covariates and add them to the model
  - **Spatial eigenvector mapping**
  - Auto-covariate Regression
  - Wavelet analysis
- Adjust the error term
  - Fit a variance-covariance matrix based on the non-independence of spatial observations
  - GLS and GLMM – error structure needs to be assumed
  - Simultaneous autoregressive error models (SAR) and conditional autoregressive models (CAR)
  - Generalized estimating equations (GEE) split the data into smaller clusters before modelling the variance-covariance relationship



## How to deal with it? (2)

- Adjustments of test statistics
  - Dutilleul's modified t-test, {SpatialPack}
  - CRH-correction for correlations {SpatialPack}
- Lagged response models
  - The response depends on the response of the neighboring units
- Lagged predictor effects (SLX model)
  - Spill-over effect from in neighboring units
  - Artifacts due to the artificial spatial discretization
    - Are administrative units well suited for health data?
    - MAUP



# Properties of some approaches

method	residuals	computational intensity
GAM	normal, Poisson, binomial	low
autogressive models (SAR/CAR)	normal	medium-high
GLS	normal	medium-high
GEE	normal, Poisson, binomial	low
autocovariate regression	normal, Poisson, binomial	low
spatial GLMM	normal, Poisson, binomial	very high
Spatial Eigenvector Mapping	normal, Poisson, binomial	very high

- Wavelet and GEE: flexible with distributional assumption, no categorical variables possible

# Spatial autocorrelation in residuals

## Spatial error model

- Incorporates spatial effects through error term

- Where:

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

$$\varepsilon = \lambda W\varepsilon + \xi$$

$\varepsilon$  : vector of error terms, spatially weighted using weights matrix  $W$

$\lambda$  : spatial error coefficient

$\xi$  : vector of uncorrelated error terms

- If there is no spatial correlation between the errors, then  $\lambda = 0$



# Spatial autocorrelation in response

## Spatial lag model

- Incorporates spatial effects by including a spatially lagged dependent variable as an additional predictor

$$y = \rho W y + x \beta + \varepsilon$$

- Where:
  - $W y$  : the spatially lagged response for weights matrix  $W$
  - $x$  : matrix of observations on the explanatory variables
  - $\varepsilon$  : vector of error terms
  - $\rho$  : spatial coefficient
- If there is no spatial dependence, and  $y$  does not depend on neighboring  $y$  values,  $\rho = 0$



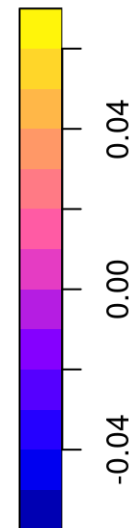
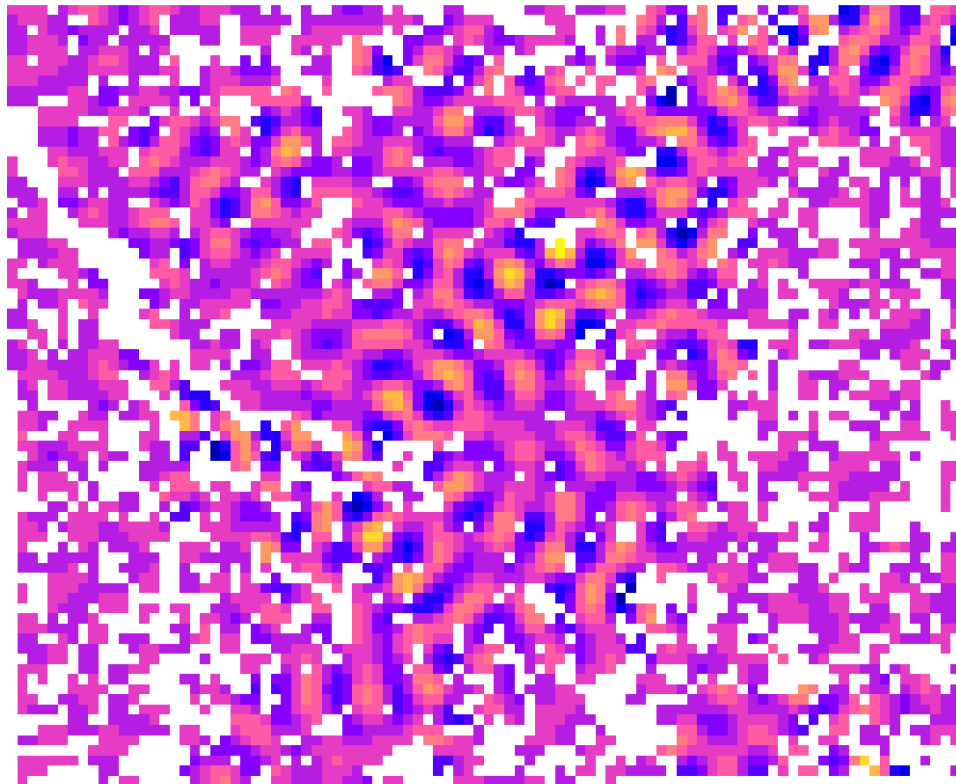
# Spatial eigenvector mapping

- Based on the eigenfunction decomposition of spatial connectivity matrices
- Eigenvectors from these matrices represent the decompositions the spatial weight Matrix into all mutually orthogonal  $m$
- Eigenvectors with positive eigenvalues represent positive autocorrelation, whereas eigenvectors with negative eigenvalues represent negative autocorrelation
- Only eigenvectors with positive eigenvalues are used



# Spatial eigenvectos

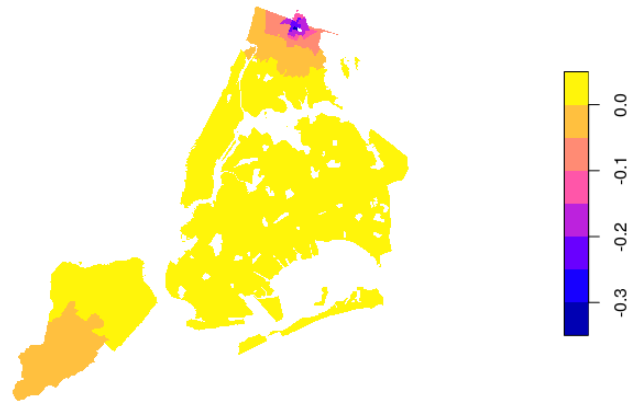
vec695



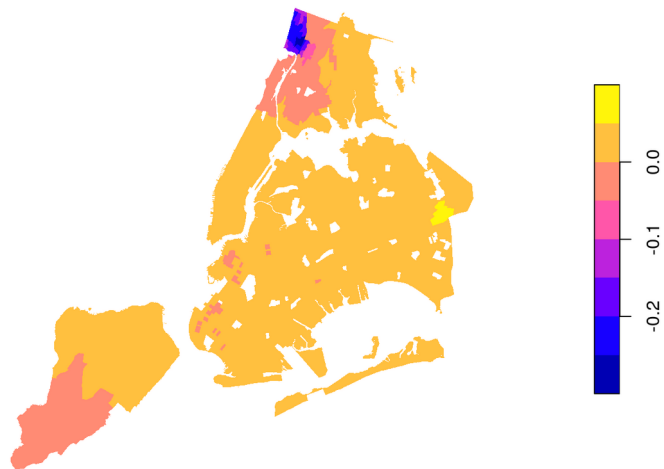
- Spatial pattern for independent dimensions of spatial structure
- Could be used to build hypothesis about missing covariates etc.

# Spatial eigenvectos

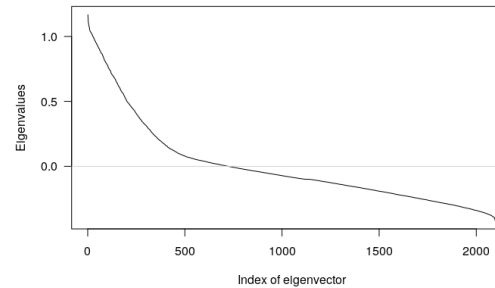
10nn, W: EV2



10nn, W: EV3



- Correlation length decreases with decreasing eigenvalue of the eigenvectors



...

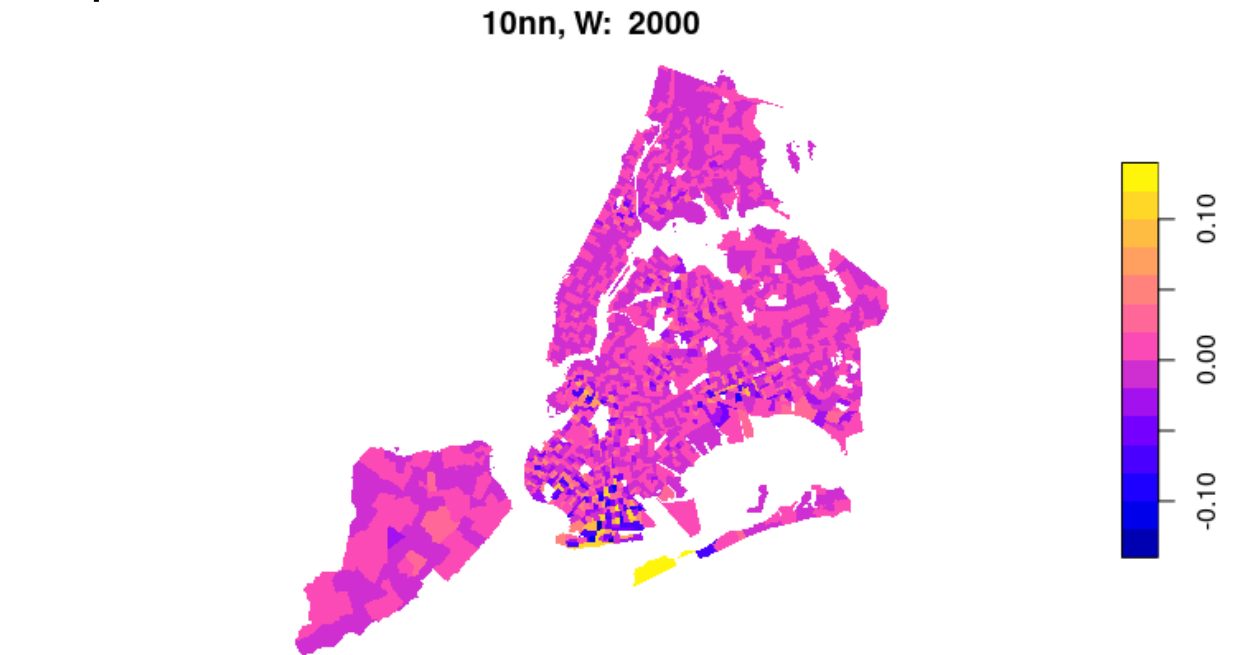
10nn, W: EV311





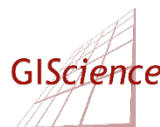
# Spatial eigenvectos

- Eigenvectors with negative eigenvalue represent negative spatial autocorrelation



# Spatial eigenvector mapping

- Compute connectivity matrix
  - Needs to be symmetric!
- Compute eigenvectors of the centered similarity matrix
- Select eigenvectors to be included in the GLM
  - eigenvectors are added to a model until the spatial autocorrelation in the residuals, measured by Moran's I, is non-significant



# Literature for further studies

- Anselin, L., Rey, S.J., 2014. Modern Spatial Econometrics in Practice. GeoDa Press LLC, Chicago, IL.
  - Spatial lag and error model
  - Not with R but GeoDa, GeoDaSpatial and Python
  - Packed with Theory and Examples
- Bivand, R.S., Pebesma, E.J., Gómez-rubio, V., 2008. Applied Spatial Data Analysis with R. Springer, New York, NY.
  - Overview about spatial tasks in R
  - A bit geeky, not much theory
- Fortin, M.-J., Dale, M., 2005. Spatial analysis: a guide for ecologists. Cambridge University Press, Cambridge (UK).
  - Excellent overview about spatial statistics
  - No code



# Literature for further studies

- Dormann et al. 2007. Methods to account for spatial autocorrelation in the analysis of species distributional data: a review. *Ecography* 30:609–628.
  - Excellent overview about methods
  - R code
  - Strengths and weaknesses of the different methods
- Carl, G., C. F. Dormann, and I. Kühn. 2008. A wavelet-based method to remove spatial autocorrelation in the analysis of species distributional data. *Web Ecology*:22–29
  - Add on to Dormann et al (2007)
  - R code
- Griffith, D. A. 2006. Spatial Modelling in Ecology: The Flexibility of Eigenfunction Spatial Analysis. *Ecology* 87:2603–2613.
  - Spatial eigenvector mapping and filtering



# Literature for further studies

- Bivand, R., Piras, G., 2015. Comparing Implementations of Estimation Methods for Spatial Econometrics. J. Stat. Softw. 63, 1–36. doi:10.18637/jss.v063.i18
  - Comparison of the different implementations of the spatial lag and error model in different software packages
  - Good theoretical overview
- Vignettes for the different spatial R packages
- Spatial task view on CRAN



# Literature for further studies

- Chun, Y. & Griffith, D. A. (2013): Spatial Statistics & Geostatistics, SAGE
- Griffith, D.A., Chun, Y. and Li, B. (2019): Spatial Regression Analysis Using Eigenvector Spatial Filtering, Elsevier Academic Press
- Haining & Li (2020) Modelling Spatial and Spatio-Temporal Data – A Bayesian Approach, CRC Press

