

From dots to the results: spatial methods for geo-located data

Katarzyna Kopczewska

Faculty of Economic Sciences

University of Warsaw, Poland

kkopczewska@wne.uw.edu.pl



UNIVERSITY
OF WARSAW



FACULTY OF
ECONOMIC SCIENCES

We have dots...

*/geo-referenced observations
possibly with values attached/*

What can we do with this information?

Statistics

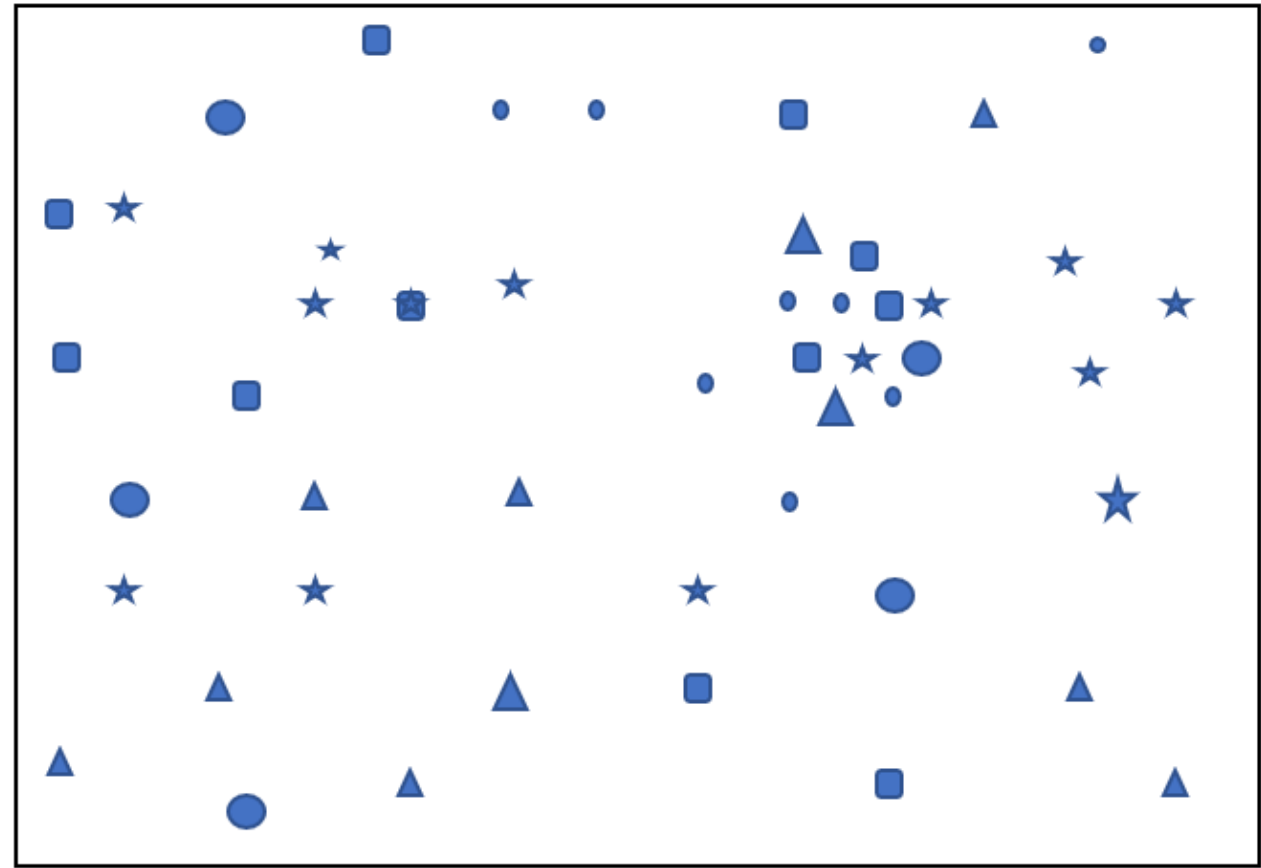
Econometrics

Machine Learning

Spatial Statistics

Spatial Econometrics

Spatial Machine Learning



Neighbourhood, distance, relative and absolute location matter

This will be the story what to do with data that have values, location and time label.

→ Spatial statistics

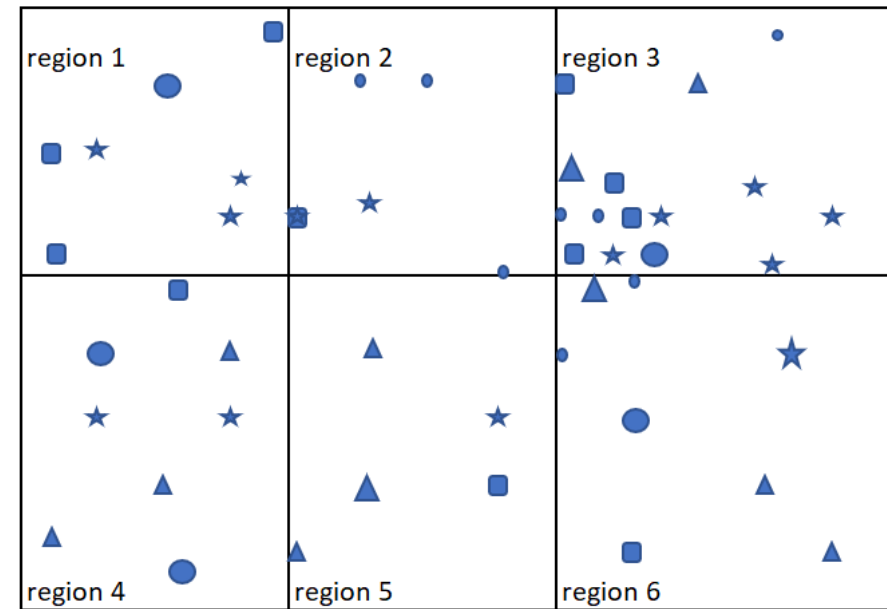
- Most „intuitive” is to **aggregate** dots to boxes

Great paper

Briant, A., Combes, P. P., & Lafourcade, M. (2010). Dots to boxes: Do the size and shape of spatial units jeopardize economic geography estimations?. *Journal of Urban Economics*, 67(3), 287-302.

- So, we get cross-table and we can apply **all concentration /clustered-based/ measures**

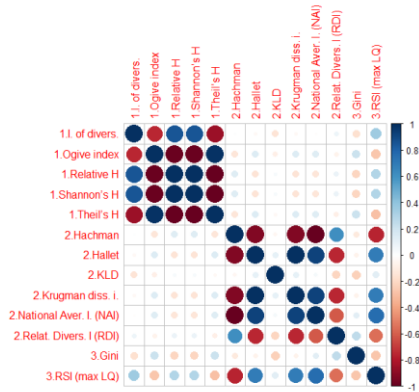
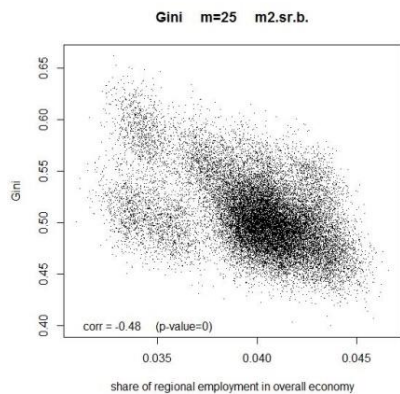
Gini, Krugman dissimilarity, entropy, LQ &
 Relative Specialisation, Hachman, Ogive, Refined Diversification,
 Relative Diversity, Hallet, National Averages, Lillen,
 Kullback-Leibler divergence, Bruelhart-Traeger, Ellison-Glaeser,
 Maurell-Sedillot, Geographic concentration and a few more



sectors	region 1	region 2	region 3	region 4	region 5	region 6	Total
▲	0	0	2	3	3	3	11
●	1	3	4	2	0	3	13
★	3	1	5	2	1	1	13
■	3	1	4	1	1	1	11
Total	7	5	15	8	5	8	48

➔ Spatial statistics

- First, let's **identify** all cluster-based measures and see how to get them
- Second, **any problems** in calculations?
 - MAUP / aggregation level
 - Can the same table bring different results?



Aggregation (x) matters for the results (y)

Measures depending on the design can bring the same or new info




Book | © 2017

Measuring Regional Specialisation A New Approach


Home > Book

Authors: [Katarzyna Kopczewska](#), [Paweł Churski](#), [Artur Ochojski](#), [Adam Polko](#)



Spatial Statistics

Volume 27, October 2018, Pages 31-57



Cluster-based measures of regional concentration. Critical overview

[Katarzyna Kopczewska](#) ✉

Show more ▾

+ Add to Mendeley 🔗 Share 🗣️ Cite

<https://doi.org/10.1016/j.spasta.2018.07.008> ↗

[Get rights and content](#) ↗

<https://www.sciencedirect.com/science/article/abs/pii/S2211675317302956>
<https://link.springer.com/book/10.1007/978-3-319-51505-2>

→ Spatial statistics

- Clustered-based measures **do not look inside** the table cells – spatial distribution can be any
- How to express **spatial agglomeration** with one number with regard to sector, size and location?

FULL ARTICLE

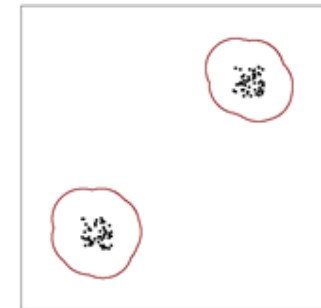
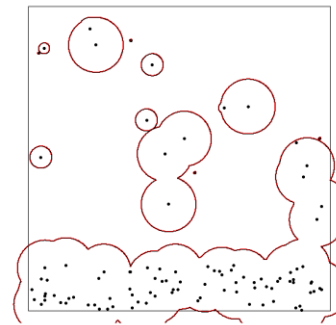
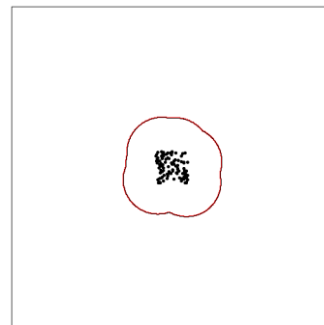
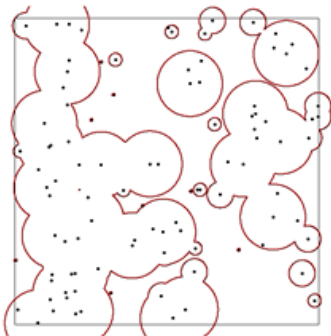
SPAG: Index of spatial agglomeration

Katarzyna Kopczewska ✉, Paweł Churski, Artur Ochojski, Adam Polko

First published: 12 July 2019 | <https://doi.org/10.1111/pirs.12470> | Citations: 3

<https://rsaiconnect.onlinelibrary.wiley.com/doi/abs/10.1111/pirs.12470>

Random pattern	High agglomeraton	Costal & interrior locations	Two agglomerations
SPAG = 0.72	SPAG = 0.008	SPAG=0.41	SPAG=0.1



➔ Spatial statistics

- Alternative to SPAG measure of spatial agglomeration: using **entropy for tessellated** point pattern (Voronoi polygons)

EE Elgaronline Get Access or Sign In

Browse Librarian Services Help Products Subjects Journals

< Previous Chapter Next Chapter >

Chapter 6: Entropy as a measure of agglomeration 🔒 Restricted access

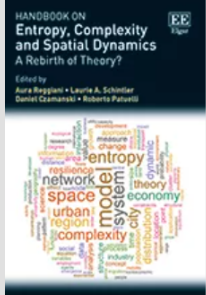
Katarzyna Kopczewska

Category: Handbook Chapter Collection: Economics 2021

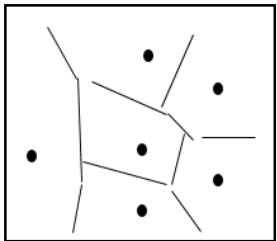
Published: 17 Dec 2021 DOI: <https://doi.org/10.4337/9781839100598.00013>

Page Range: 97–117

Abstract



Handbook on Entropy, Complexity and Spatial Dynamics



Tessellated regular points



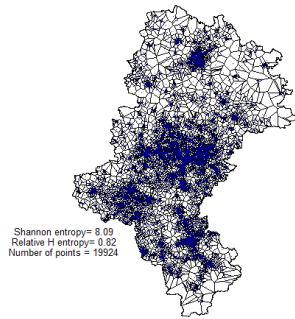
Shannon entropy= 5.93
Relative H entropy= 1

relH=1



Shannon entropy: $H = - \sum_{n=1}^N s \cdot \log s$
 Max entropy (uniform distr): $H_{max} = \log n$
 Relative entropy: $RelH = \frac{H_{empir}}{H_{max}} = \frac{H_{empir}}{\ln n}$

Tessellated business locations in Silesia region

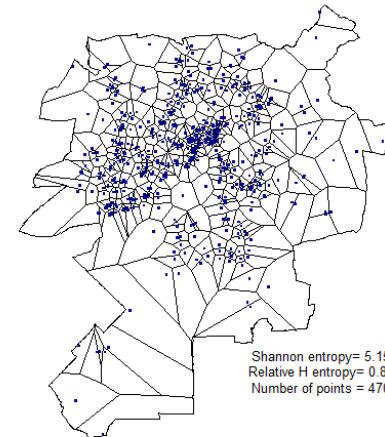


Shannon entropy= 6.09
Relative H entropy= 0.82
Number of points = 19924

relH=0.82



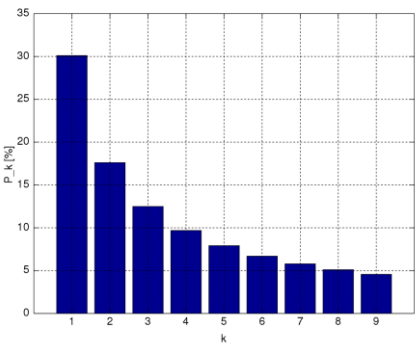
Tessellated business locations in Lublin city



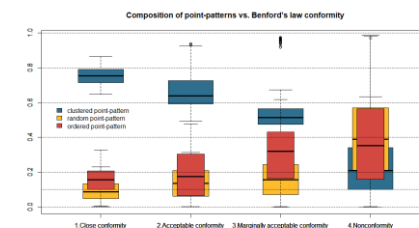
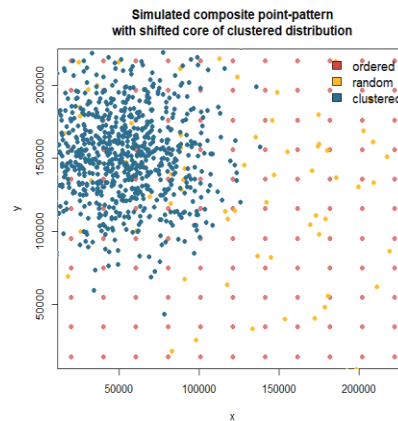
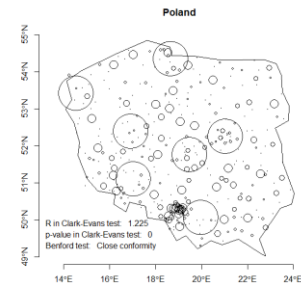
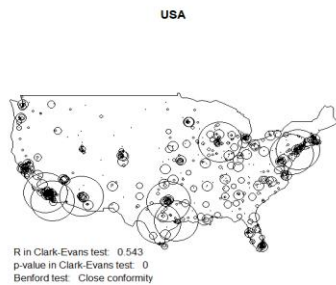
Shannon entropy= 5.15
Relative H entropy= 0.84
Number of points = 470

→ Spatial statistics

- We know the **uniform, random, agglomerated** point patterns.
- Can point patterns be **natural and follow Benford's law**? Can we generate it?
- If yes, are **cities with their location and population** (3D) following Benford's law? What about Zipf law?



Benford's law

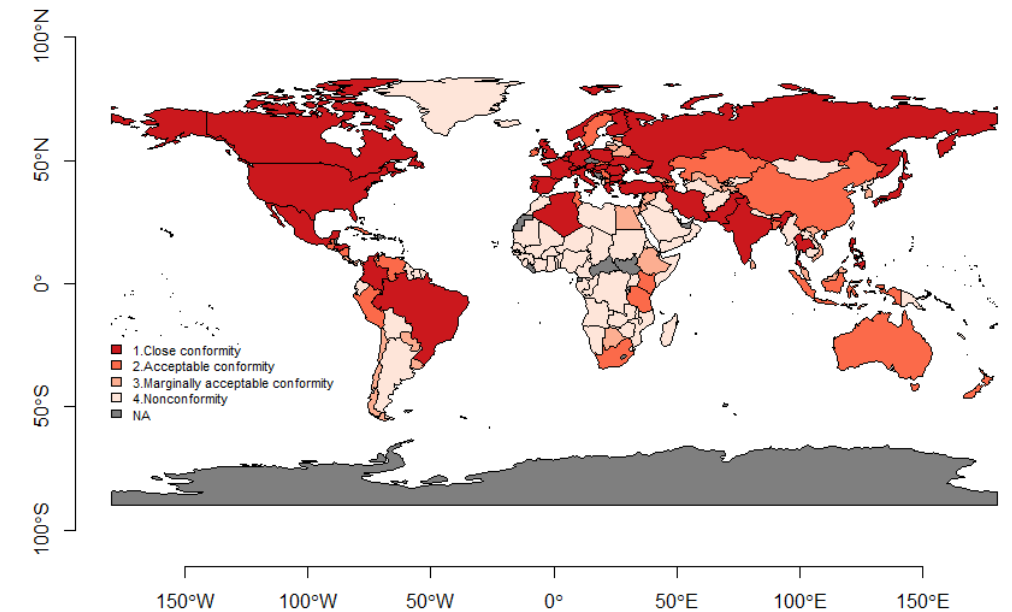


Natural spatial pattern—When mutual socio-geo distances between cities follow Benford's law

Katarzyna Kopczewska, Tomasz Kopczewski

Published: October 20, 2022 • <https://doi.org/10.1371/journal.pone.0276450>

Conformity to Benford distribution
Analysis of 3D socio-geo distances between cities within countries



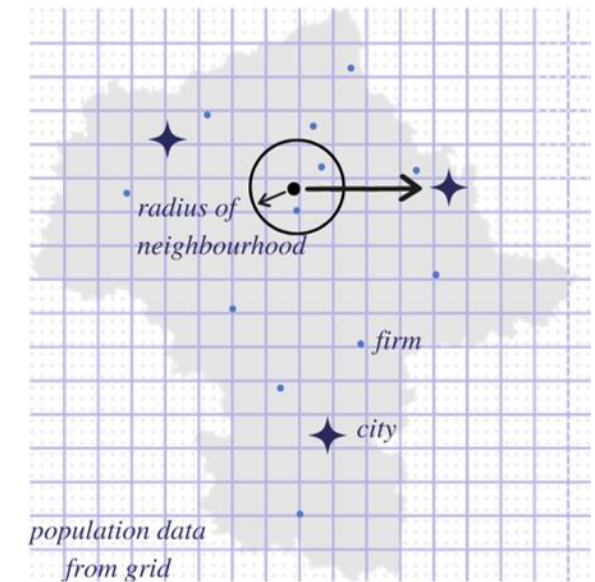
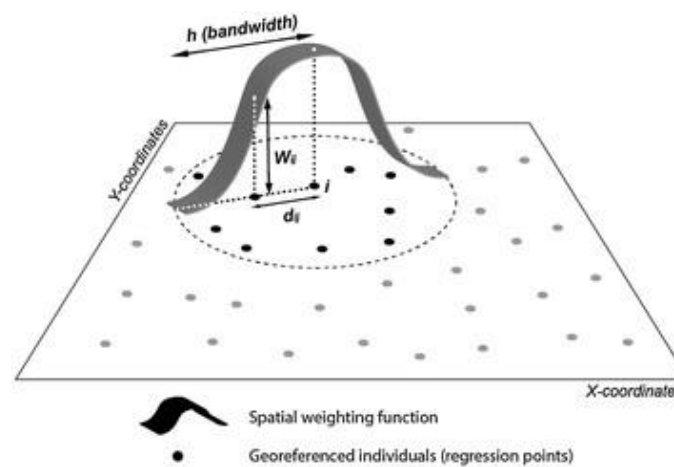
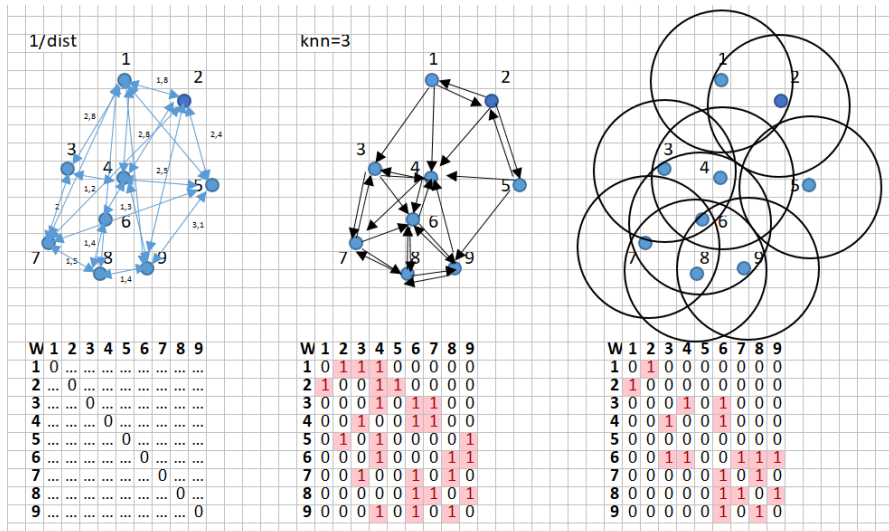
Criteria used: MAD Conformity - Nigrini (2012)

$$dist_{ij,3D} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2}$$

→ → Spatial econometrics

- How to incorporate info from a neighbourhood in econometric model on points:

Spatial econometrics with spatial weights matrix W	Geographically weighted regression GWR	Microgeography and agglomerations
Spatial lag for Y , X or e is the average value from the neighbourhood	Local regression on neighbours only; each observation gets its own coefficients	Making a radius around the point and counting what's inside



What can go wrong? 😊

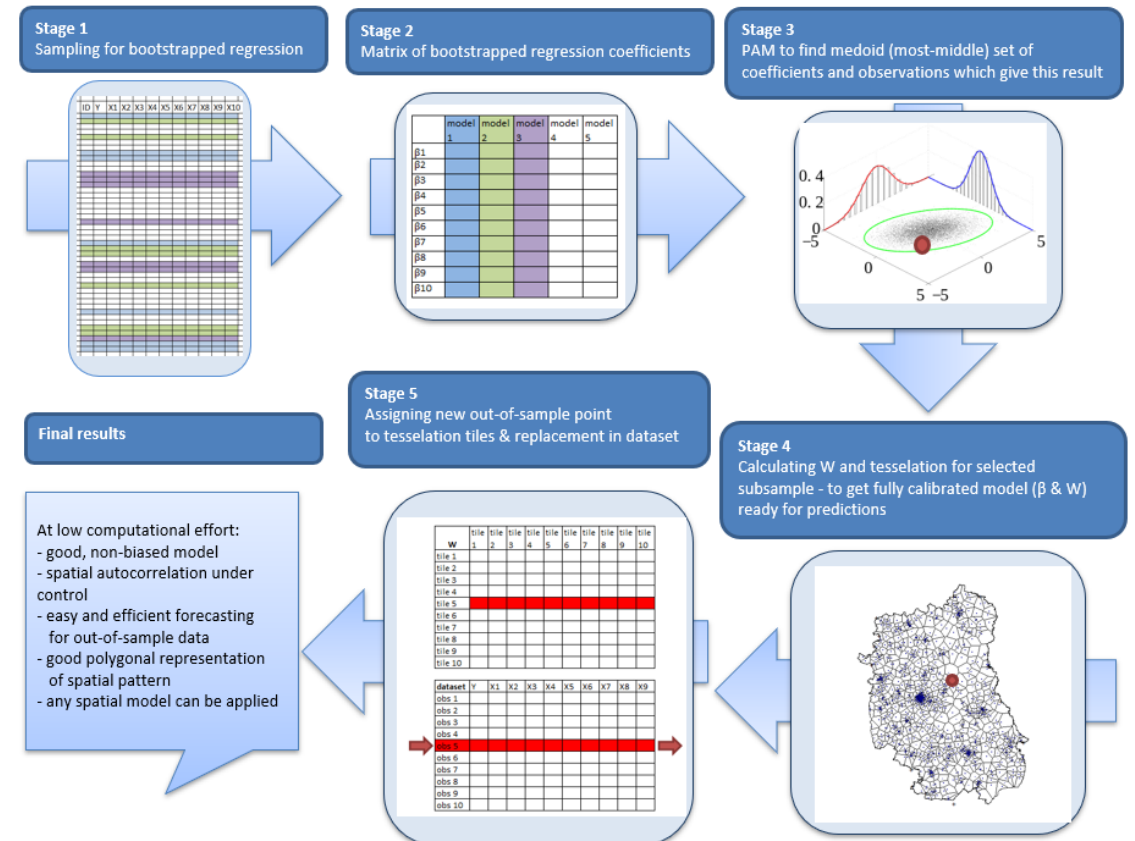
→ → Spatial econometrics

- How to deal with big spatial data (for which W is unavailable)? → like 1 mln points (spatial econometrics works up to 100K points)
 - How to predict out-of-sample for point data using spatial econometrics? → W keeps old locations, no possibility to input new locations
- ↓
- We can bootstrap data and choose the best model from the candidates – estimate small model quickly
 - We know well the parameters (mean is mean, sd is a function of n & sd)
 - We can approximate exchange old point into new point – we use tessellation tiles to link them in pairs

Spatial bootstrapped microeconometrics: forecasting for out-of-sample geo-locations in big data

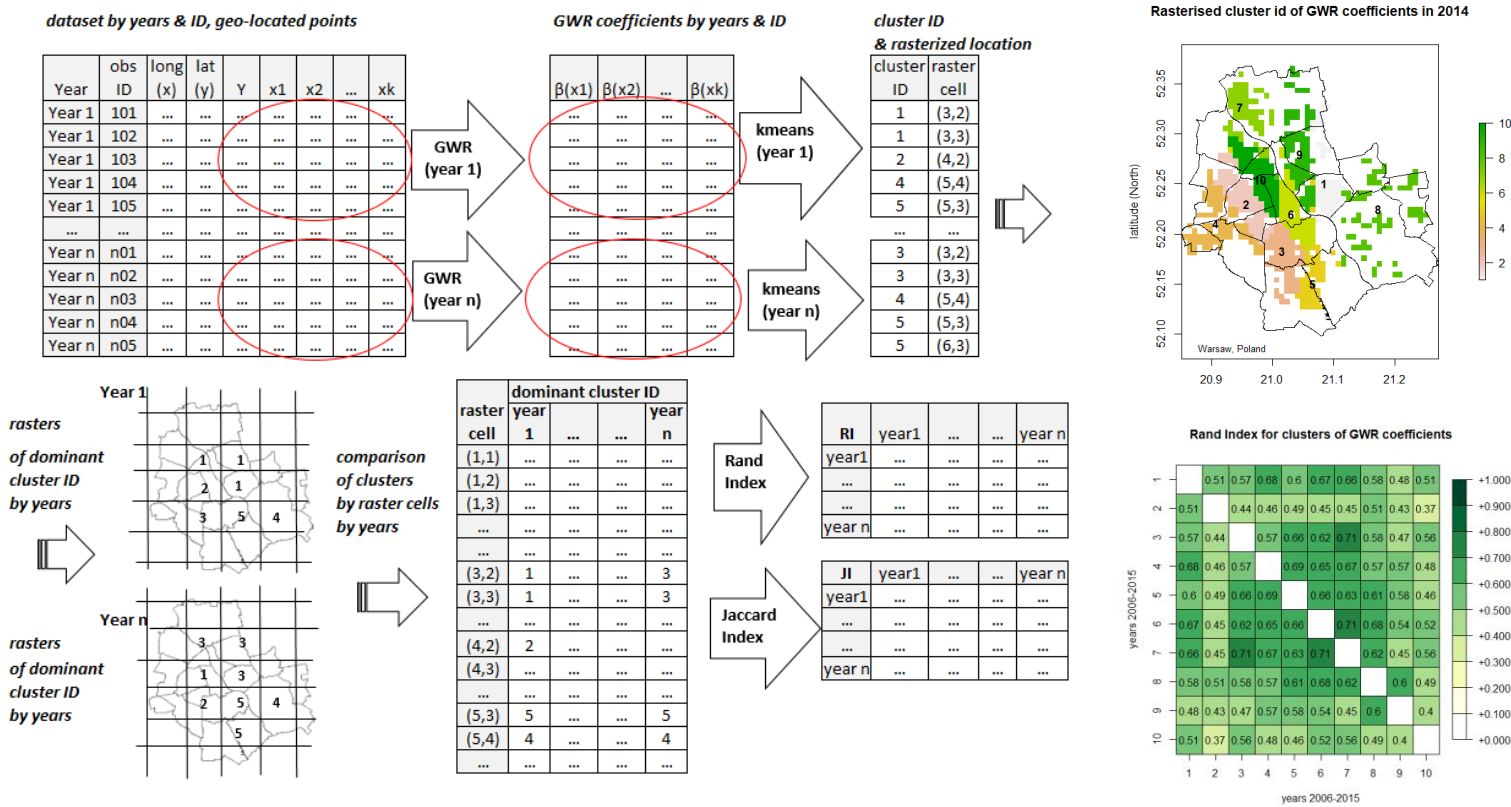
Katarzyna Kopczewska 

First published: 09 March 2023 | <https://doi.org/10.1111/sjos.12636>



→ → Spatial econometrics

- For spatio-temporal data (e.g. housing transactions every year) we can get annual GWR hedonic models – are the valuations stable over time and space?
- We can cluster GWR coefficients to get submarkets – are they the same place every year? How long is model valid?



Spatio-temporal stability of housing submarkets. Tracking spatial location of clusters of geographically weighted regression estimates of price determinants ☆

Katarzyna Kopczewska, Piotr Cwiakowski

Show more

Add to Mendeley Share Cite

<https://doi.org/10.1016/j.landusepol.2021.105292>

Get rights and content

Under a Creative Commons license

open access

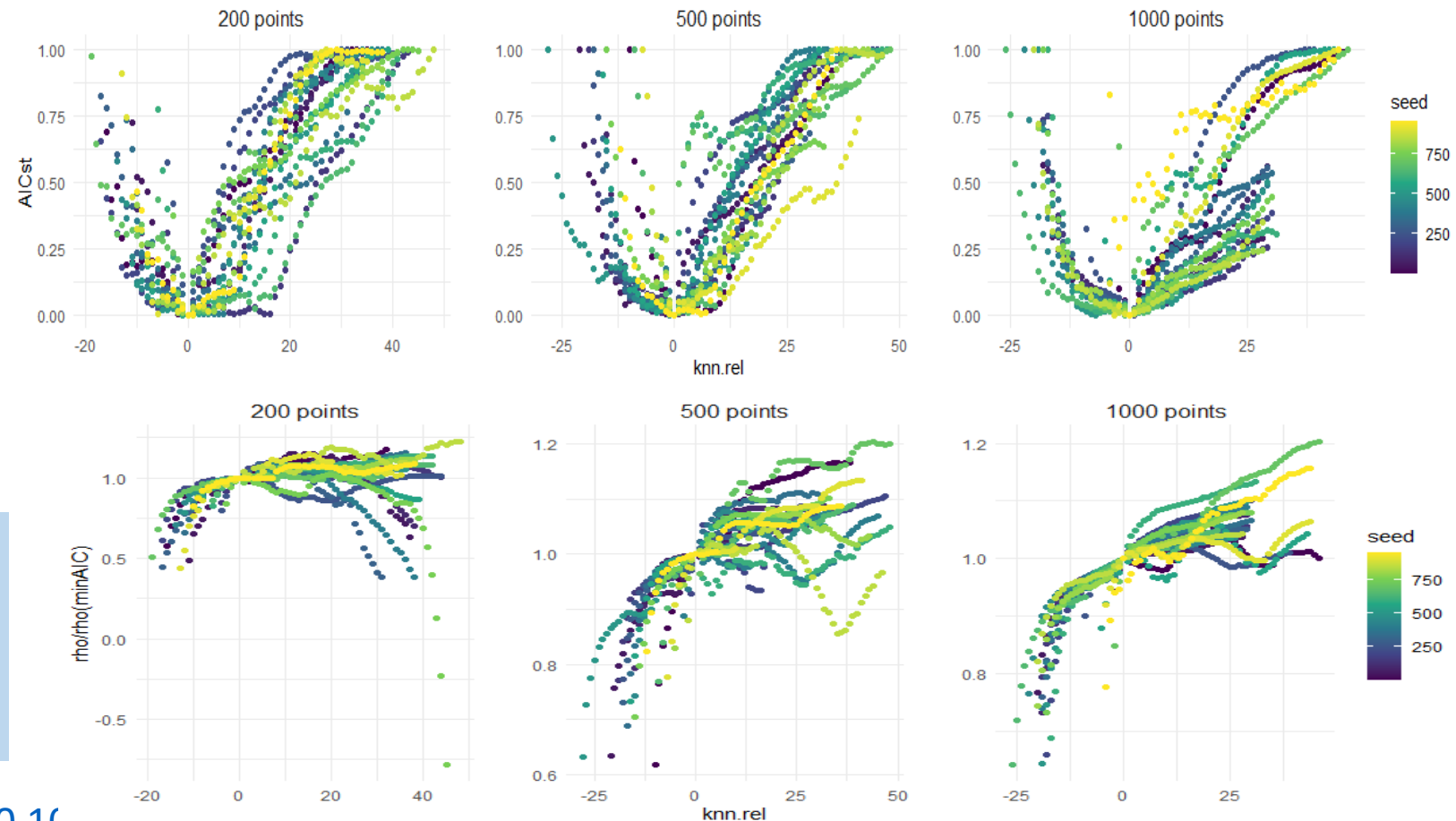
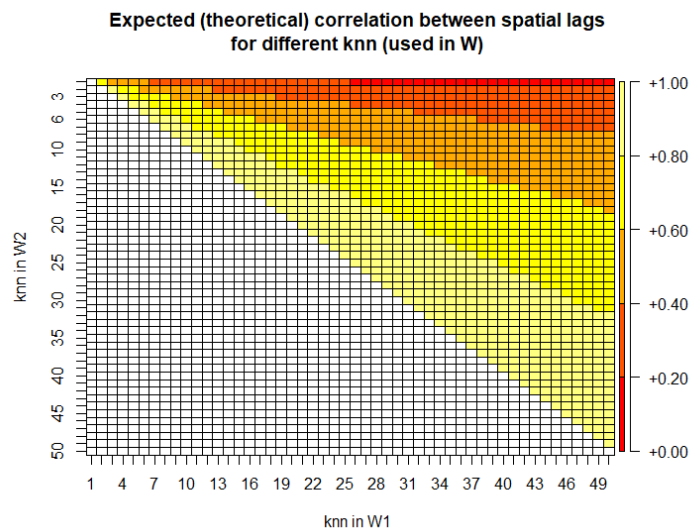
- Cluster annual GWR coefficients and grid it – then take cluster ID as representation for grid
- Compare grids if they belong to the same clusters over time – use **Rand index** (comparing pairs of pairs of points) – output is a stability index

→ → Spatial econometrics

- How many k nearest neighbours (knn) should be included in W ?
- Does knn matter for the result?

Akaike information criterion in choosing the optimal k -nearest neighbours of the spatial weight matrix

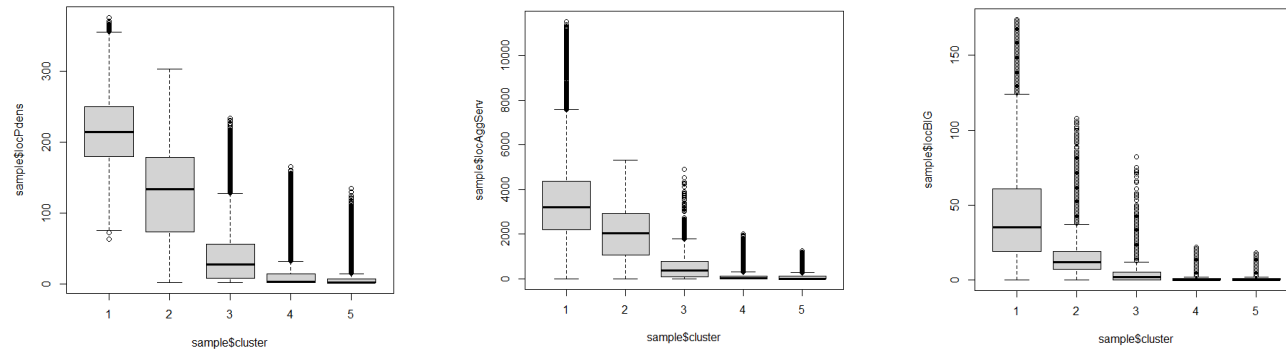
Maria Kubara  and Katarzyna Kopczewska 



- **Knn matters** – wrongly set generates significant **bias of coefficients**
- **Use AIC** – check candidate knn and choose the model with lowest AIC

→ → Spatial econometrics

- Does local density of points matters for their nature and relations they generate?
- How to deal with spatial heterogeneity of these data?



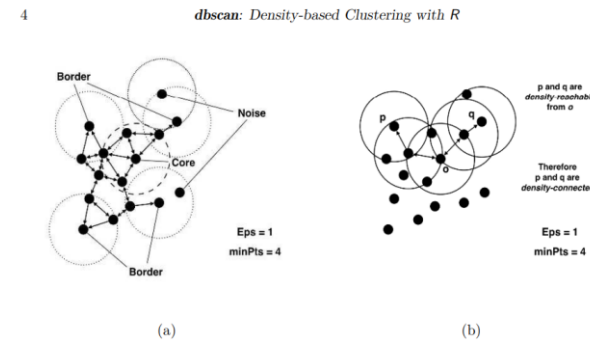
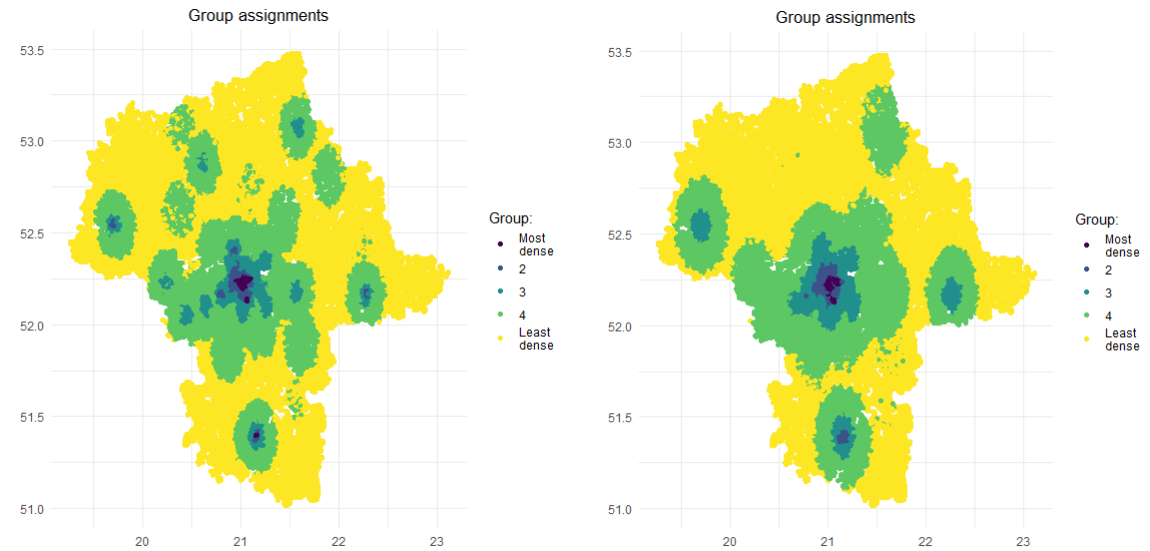
Close neighbourhood of firms differs among density groups
 → Economic relations differ too - in direction and strength!



- Divide data into DBSCAN density clusters (ML algorithm)
- Estimate models in subgroups and compare them
- Information from neighbourhood calculated in radius from point (e.g number of firms from given sector, number of inhabitants)

Spatial switching regimes regression: density-based cluster approach to deal with spatial heterogeneity

Katarzyna Kopczewska, Ewa Dobrowolska, Anil Bera – will be ready soon



DBSCAN checks if number of points within the radius exceeds a threshold – if yes, it is high density cluster

Figure 1: Concepts used the DBSCAN family of algorithms. (a) shows examples for the three point classes, core, border, and noise points, (b) illustrates the concept of density-reachability and density-connectivity.

→ → → Spatial machine learning

- Novel solution with huge potential
- Challenge: how to include spatial component?
- What is done and what not?
- New perspectives for regional and urban studies

Why use ML?

- Why not? 😊
- Dealing with non-linearity, partial impact
- Finding new patterns, also spatio-temporal
- Higher computational efficiency
- New types of data possible to include (photos, pixels, etc)
- Better possibilities to forecast

[Home](#) > [The Annals of Regional Science](#) > [Article](#)

Original Paper | [Open Access](#) | [Published: 24 December 2021](#)

Spatial machine learning: new opportunities for regional science

[Katarzyna Kopczewska](#) ✉

[The Annals of Regional Science](#) **68**, 713–755 (2022) | [Cite this article](#)

10k Accesses | **18** Citations | **32** Altmetric | [Metrics](#)

Unsupervised
learning

Clustering:

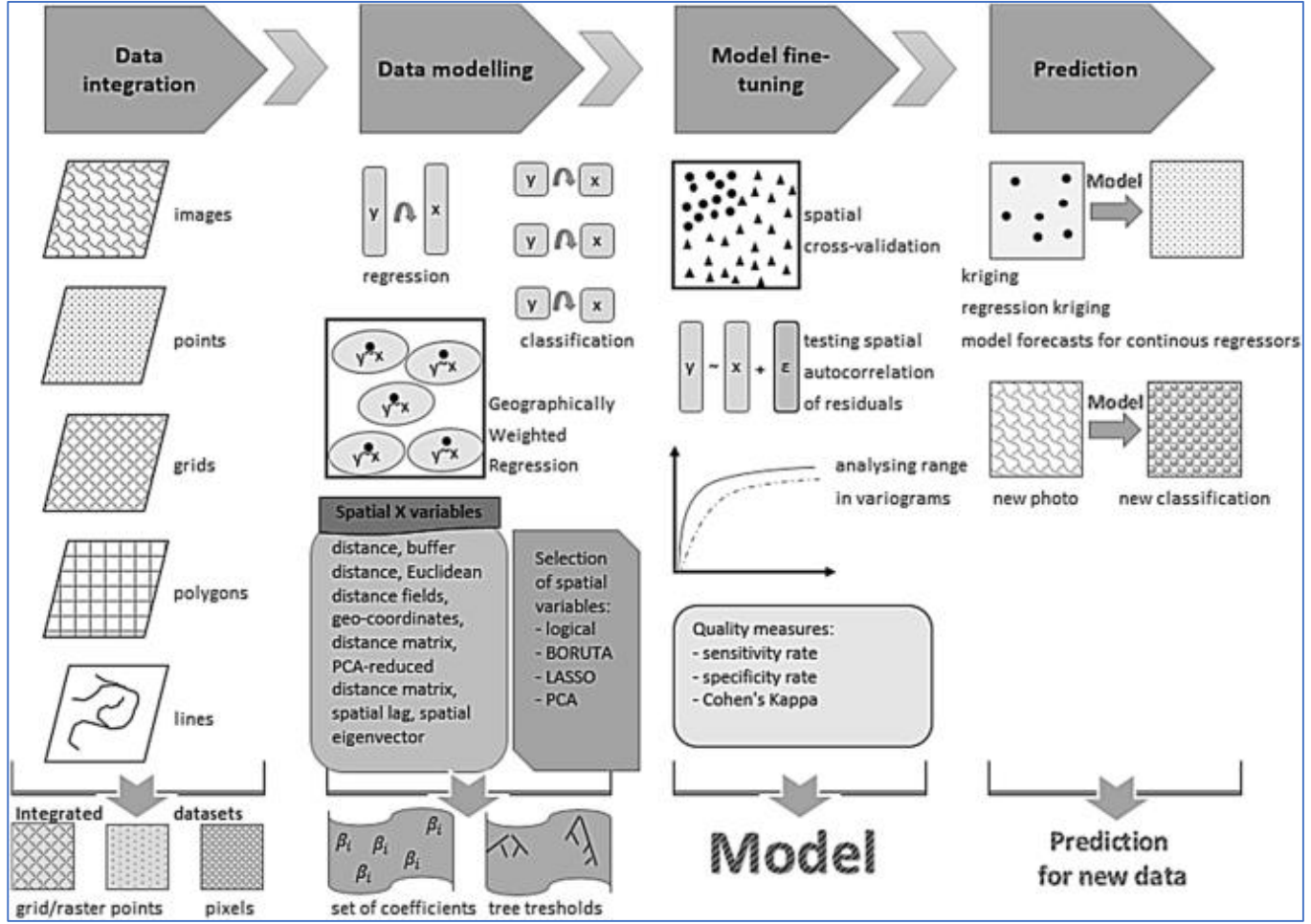
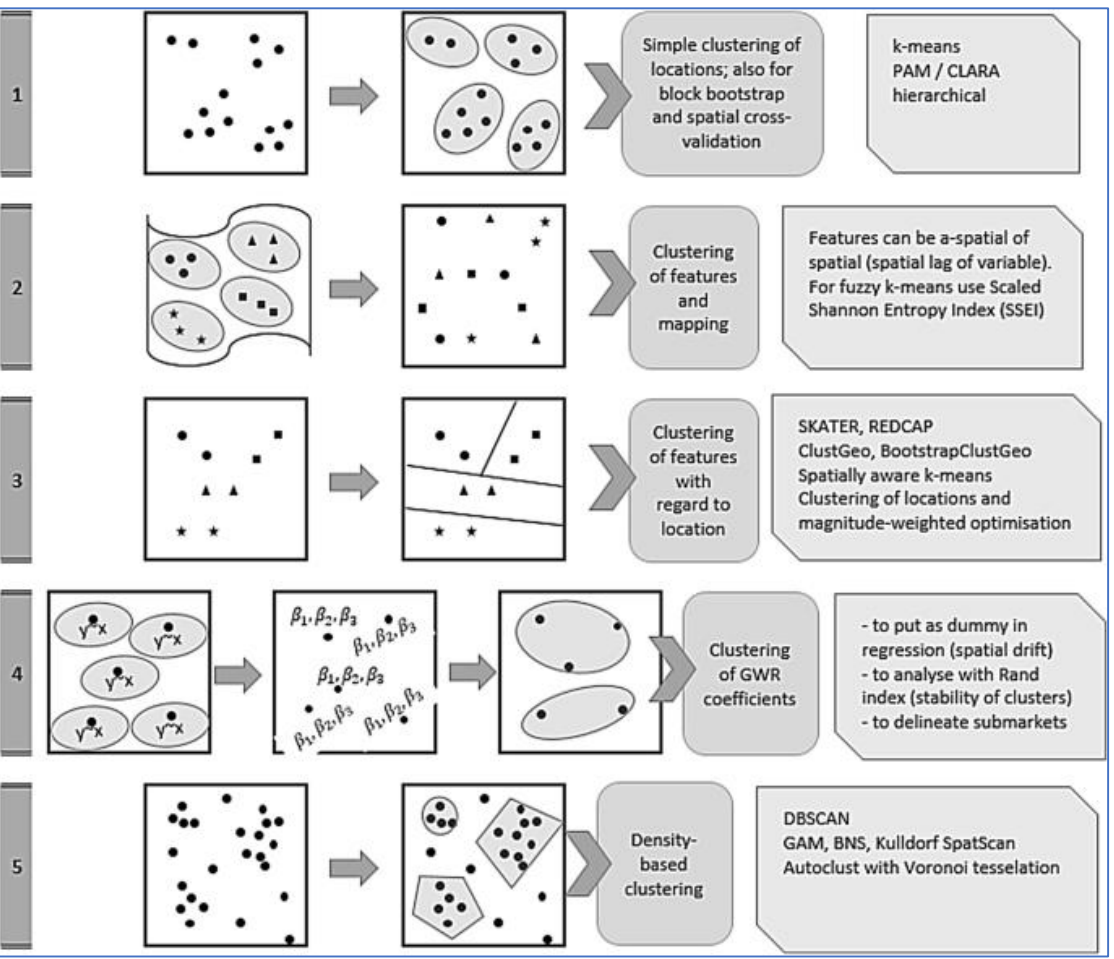
- locations
- features and mapping
- features in space
- GWR coefficients
- density-based

Supervised
learning

Regression models:

- Simple regressions
- Spatial cross-validation
- Image recognition in spatial classification
- Mixtures of GWR and ML
- Spatial variables

Unsupervised and supervised spatial machine learning



→ → → Spatial Machine Learning

Perspectives:

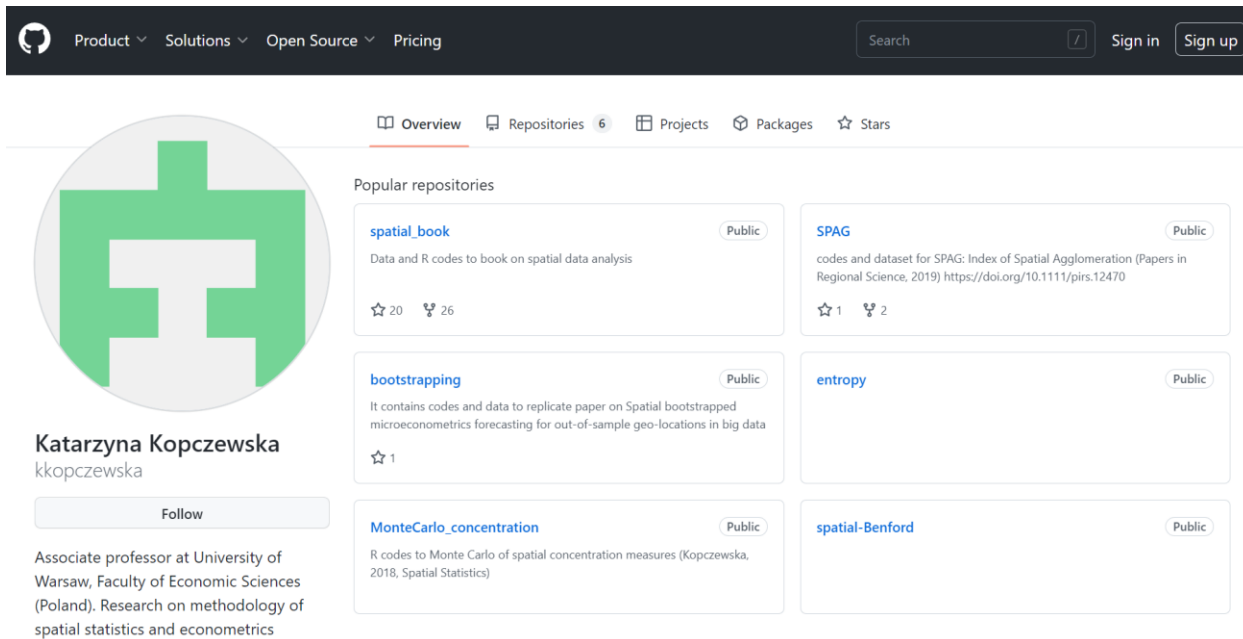
- Big data analytics (to help spatial econometrics)
- Using new sources of data – nightlights, day photo of landscape, spectral data, vegetation indicators
- Dealing with spatial heterogeneity and isotropy
- Spatio-temporal modeling with many layers (also of different granulation)
- Better forecasting (new approach - „econometric” models not only to explain)

What will come soon?

- Spatial 3D solutions
- Like in life sciences: standards of reporting, computational reproducibility, workflow managers (ready-to-use environments)
- Communication with those who believe in p-value 😊

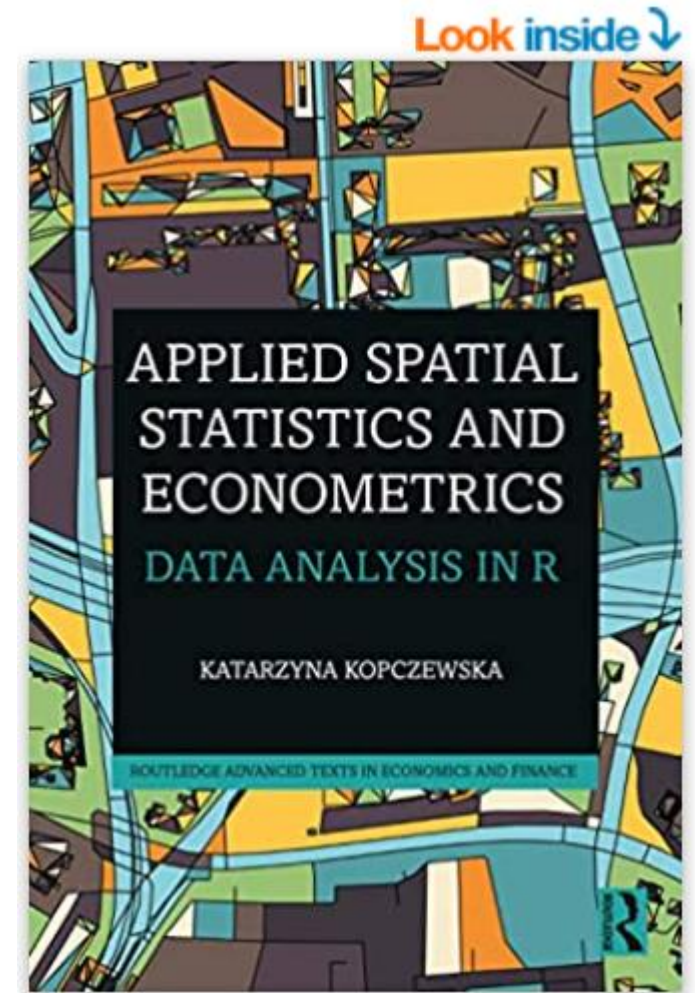
→ → → → R implementations


- Toolbox to run the analysis is finally the core of our research
- R still wins with Python – especially since C++ implementations are available
- Making codes available at Github makes the science really open



The screenshot shows the GitHub profile of Katarzyna Kopczewska. The profile includes a green and white geometric logo, her name, and a bio: "Associate professor at University of Warsaw, Faculty of Economic Sciences (Poland). Research on methodology of spatial statistics and econometrics". Below the profile, there is a section for "Popular repositories" with six entries:

Repository Name	Description	Stars	Forks
spatial_book	Data and R codes to book on spatial data analysis	20	26
SPAG	codes and dataset for SPAG: Index of Spatial Agglomeration (Papers in Regional Science, 2019) https://doi.org/10.1111/pirs.12470	1	2
bootstrapping	It contains codes and data to replicate paper on Spatial bootstrapped microeconometrics forecasting for out-of-sample geo-locations in big data	1	0
entropy		0	0
MonteCarlo_concentration	R codes to Monte Carlo of spatial concentration measures (Kopczewska, 2018, Spatial Statistics)	0	0
spatial-Benford		0	0



Applied Spatial Statistics and Econometrics: Data Analysis in R (Routledge Advanced Texts in Economics and Finance) 1st Edition 

<https://www.routledge.com/Applied-Spatial-Statistics-and-Econometrics-Data-Analysis-in-R/Kopczewska/p/book/9780367470760>
<https://github.com/kkopczewska>

Thank you!

Katarzyna Kopczewska

Faculty of Economic Sciences
University of Warsaw, Poland

kkopczewska@wne.uw.edu.pl



@KathyKopczewska
@SpatialWarsaw