

# **Geodata & Spatial Regression**

Tobias Rüttenauer

2024-06-29



# Table of contents

<b>Introduction</b>	<b>1</b>
Schedule . . . . .	2
Some useful packages . . . . .	3
Further Readings . . . . .	4
Course materials . . . . .	5
<b>1 Refresher</b>	<b>7</b>
Required packages . . . . .	7
Session info . . . . .	7
1.1 Packages . . . . .	9
1.2 Coordinates . . . . .	9
1.2.1 Coordinate reference system (CRS) . . . . .	10
1.2.2 Projected CRS . . . . .	14
1.2.3 Why different projections? . . . . .	14
1.3 Importing some real world data . . . . .	17
1.3.1 London shapefile (polygon) . . . . .	18
1.3.2 Census API (admin units) . . . . .	23
1.3.3 Gridded data . . . . .	28
1.3.4 OpenStreetMap (points) . . . . .	30
1.3.5 Save . . . . .	33
<b>2 Data Manipulation &amp; Visualization</b>	<b>35</b>
Required packages . . . . .	35
Session info . . . . .	35
Reload data from previous session . . . . .	37

*Table of contents*

2.1	Manipulation and linkage . . . . .	37
2.1.1	Subsetting . . . . .	39
2.1.2	Point in polygon . . . . .	43
2.1.3	Distance measures . . . . .	44
2.1.4	Intersections + Buffers . . . . .	45
2.1.5	and more . . . . .	49
2.1.6	Air pollution and ethnic minorities . . . . .	50
2.1.7	Save spatial data . . . . .	53
2.2	Visualization . . . . .	53
2.2.1	Tmaps . . . . .	54
2.2.2	ggplot . . . . .	63
2.3	Exercises . . . . .	66
<b>3</b>	<b>Spatial Relationships <math>\mathbf{W}</math></b>	<b>73</b>
	Required packages . . . . .	73
	Session info . . . . .	73
	Reload data from previous session . . . . .	75
3.1	Spatial interdependence . . . . .	75
3.2	$\mathbf{W}$ : Connectivity between units . . . . .	76
3.2.1	Contiguity weights . . . . .	78
3.2.2	Distance based weights . . . . .	82
3.3	Normalization of $\mathbf{W}$ . . . . .	85
3.3.1	Row-normalization . . . . .	85
3.3.2	Maximum eigenvalues normalization . . . . .	92
3.4	Islands / missings . . . . .	94
<b>4</b>	<b>Exercises I</b>	<b>95</b>
	Required packages . . . . .	95
	Session info . . . . .	95
	Reload data from previous session . . . . .	97

## *Table of contents*

4.1 General Exercises . . . . .	98
4.1.1 1) Can you import the spatial administrative units of Germany (“Kreisgrenzen_2020_mit_Einwohn- erzahl” in _data folder) and make a simple plot of the boundaries? {unnumbered} . . . . .	98
2) What is the Coordinate reference system of this German shape file? . . . . .	98
3) Please use the msoa.spdf and calculate a neighbours weights matrix of the nearest 10 neighbours (see spdep::knearneigh()), and create a listw object using row normalization. . . . .	98
4) OPTIONAL: Can you create a map containing the City of London (MSOA11CD = “E02000001”) and its ten nearest neighbours? . . . . .	98
5) Please use the msoa.spdf and calculate a neighbours weights matrix of the nearest 10 neighbours (see spdep::knearneigh()), and create a listw object using row normalization. . . . .	98
6) Please calculate the queens neighbours and make a listw object that includes the second order neighbours (see nblag()). . . . .	98
7) Generate a matrix from the listw object . . . . .	98
8) What do you get when you multiply a variable (data column) such as the home owner rate with your weights matrix? . . . . .	98
<b>5 Detecting Spatial Dependence</b>	<b>99</b>
Required packages . . . . .	99
Session info . . . . .	99
Reload data from previous session . . . . .	101
5.1 Global Autocorrelation . . . . .	101
5.1.1 Visualization . . . . .	101
5.1.2 Moran’s I . . . . .	104
5.1.3 Residual-based Moran’s I . . . . .	107

*Table of contents*

5.1.4	Semivariogram . . . . .	109
5.1.5	Example . . . . .	114
5.2	Local Autocorrelation . . . . .	114
5.3	Local Moran's I . . . . .	116
5.4	Example . . . . .	121
	Tate.2021 . . . . .	121
5.5	Exercises . . . . .	124
<b>6</b>	<b>Spatial Regression Models</b>	<b>125</b>
	Required packages . . . . .	125
	Session info . . . . .	125
	Reload data from previous session . . . . .	127
6.1	Why do we need spatial regression models . . . . .	127
6.1.1	Non-spatial OLS . . . . .	127
6.1.2	Problem of ignoring spatial dependence . . . . .	129
6.2	Spatial Regression Models . . . . .	131
6.2.1	Spatial Error Model (SEM) . . . . .	132
6.2.2	Spatial Autoregressive Model (SAR) . . . . .	132
6.2.3	Spatially lagged X Model (SLX) . . . . .	133
6.2.4	Spatial Durbin Model (SDM) . . . . .	133
6.2.5	Spatial Durbin Error Model (SDEM) . . . . .	134
6.2.6	Combined Spatial Autocorrelation Model (SAC) . .	134
6.2.7	General Nesting Spatial Model (GNS) . . . . .	134
6.2.8	A note on missings . . . . .	135
6.3	Mini Example . . . . .	136
6.4	Real Example . . . . .	142
6.4.1	SAR . . . . .	142
6.4.2	SEM . . . . .	144
6.4.3	SLX . . . . .	146
6.4.4	SDEM . . . . .	151
6.4.5	SDM . . . . .	152
<b>7</b>	<b>Spatial Regression Models: Estimation</b>	<b>155</b>
	Required packages . . . . .	155

*Table of contents*

Session info . . . . .	155
Reload data from pervious session . . . . .	157
7.1 Simulataneity bias . . . . .	158
7.2 Instrumental variable . . . . .	159
7.3 Generalized Method of Moments . . . . .	162
7.4 Maximum likelihood estimation . . . . .	164
7.4.1 ML SAR . . . . .	164
7.4.2 ML SEM . . . . .	166
<b>8 Exercises II</b>	<b>171</b>
Required packages . . . . .	171
Session info . . . . .	171
Reload data from pervious session . . . . .	173
8.1 Environmental inequality . . . . .	173
1) Define a neigbours weights object of your choice . . . . .	173
2) Estimate the extent of spatial auto-correlation in air pollution . . . . .	174
3) Estimate a Spatial SAR regression model . . . . .	174
4) Estimate a Spatial SEM regression model . . . . .	174
5) Estimate a Spatial SLX regression model . . . . .	174
6) Estimate a Spatial Durbin regression model . . . . .	174
7) Estimate a Spatial Durbin Error regression model . . . . .	174
8.1.1 8) Sneak preview on tomorrow: Which of the spatial model specifications about would you choose / prefer in a real world example? . . . . .	174
8.1.2 9) Please calculate the spatially lagged value of the median house price. . . . .	174
8.1.3 10) Can you use the results of the previous task to run a non-linear SLX model, where you predict if an MSOA is within the ulez zone based on the house prices? Can you make sense of the result? . . . . .	174
<b>9 Spatial Impacts</b>	<b>175</b>
Required packages . . . . .	175

*Table of contents*

Session info . . . . .	175
Reload data from previous session . . . . .	177
9.1 Coefficient estimates $\neq$ ‘marginal’ effects . . . . .	177
9.2 Global and local spillovers . . . . .	182
9.2.1 Local spillovers . . . . .	182
9.2.2 Global spillovers . . . . .	184
9.3 Summary impact measures . . . . .	186
9.4 Examples . . . . .	192
Boillat, Ceddia, and Bottazzi (2022) . . . . .	192
Fischer et al. (2009) . . . . .	193
Rüttenauer (2018) . . . . .	193
<b>10 Comparing and Selecting Models</b>	<b>195</b>
Required packages . . . . .	195
Session info . . . . .	195
Reload data from previous session . . . . .	197
10.1 Specific-to-general . . . . .	197
10.1.1 Lagrange Multiplier Test . . . . .	198
10.1.2 Problem . . . . .	203
10.2 General-to-specific approach . . . . .	203
10.3 General advice? . . . . .	205
10.4 Design and Theory . . . . .	206
10.5 Monte Carlo simulation . . . . .	206
10.5.1 Without omitted variable bias . . . . .	209
10.5.2 With omitted variable bias . . . . .	210
10.5.3 Indirect impacts if DGP = GNS . . . . .	210
10.6 Example: House prices in London . . . . .	213
<b>11 Exercises III</b>	<b>231</b>
Required packages . . . . .	231
Session info . . . . .	231
Reload data from previous session . . . . .	233

*Table of contents*

11.1 Environmental inequality (continued) . . . . .	233
1) Please calculate the true multiplier matrix of this SAR model. . . . .	236
2) Create an N x N effects matrix for the effect of the non- EU citizens. What is the effect of unit 6 on unit 10? Why is this larger than the effect of unit 5 on unit 8? . . . . .	237
3) Calculate and interpret the summary impact measures of the SAR model. . . . .	237
4) Is SAR the right model choice or would you rather esti- mate a different model? Please run a Durbin model and caculate its impact summary measures . . . . .	237
5) Please repeat with a Durbin Error model. Why are the impacts here idenptical to the coefficients? . . . . .	237
11.2 Inkar data: the effect of regional characteristics on life expectancy . . . . .	237
11.3 County shapes . . . . .	238
1) Please map the life expectancy across Germany . . . . .	239
2) Chose some variables that could predict life expectancy. See for instance the following paper. . . . .	239
3) Generate a neighbours object (e.g. the 10 nearest neigh- bours). . . . .	239
4) Estimate a cross-sectional spatial model for the year 2020 and calculate the impacts. . . . .	239
5) Calculate the spatial lagged variables for your covariates (e.g. use <code>create_WX()</code> , which needs a non-spatial df as input) . . . . .	239
6) Can you run a spatial machine learning model? (for instance, using <code>randomForest</code> )? . . . . .	239
11.4 Esimate an FE model with SLX specification . . . . .	239
<b>12 Spatio-temporal models</b>	<b>241</b>
Required packages . . . . .	241
Session info . . . . .	241
12.1 Static panel data models . . . . .	243

*Table of contents*

12.2	Dynamic panel data models . . . . .	246
12.2.1	Impacts in spatial panel models . . . . .	246
12.3	Example: Local employment impacts of immigration . . . . .	247
12.4	Estimation in R . . . . .	248
12.5	Example: Industrial facilities and municipal income . . . . .	259
<b>13</b>	<b>Other Models</b>	<b>261</b>
13.0.1	Required packages . . . . .	261
13.0.2	Session info . . . . .	261
13.0.3	Reload data from previous session . . . . .	263
13.1	Geographically weighted regression . . . . .	263
13.2	Non-Linear Models . . . . .	270
13.2.1	Problem with non-linear models . . . . .	270
13.2.2	Estimation . . . . .	271
13.2.3	Suggestion . . . . .	272
	<b>References</b>	<b>273</b>

# Introduction

This course material is designed for the 3-days GESIS workshop on geodata and spatial regression analysis. Rüttenauer (2024) provides a handbook chapter accompanying these workshop materials.

In recent years, more and more spatial data has become available, providing the possibility to combine otherwise unrelated data, such as social, economic, and environmental data. This also opens up the possibility of analysing spatial patterns and processes (e.g., spillover effects or diffusion).

Many social science research questions are spatially dependent such as voting outcomes, housing prices, labour markets, protest behaviour, or migration decisions. Observing an event in one region or neighbourhood increases the likelihood that we observe similar processes in proximate areas. As Tobler's first law of geography puts it: "Everything is related to everything else, but near things are more related than distant things". This dependence can stem from spatial contagion, spatial spillovers, or common confounders. Therefore, basic assumptions of standard regression models are violated when analysing spatial data. However, more importantly, spatial processes are interesting for their own sake. Spatial regression models can detect spatial dependence and explicitly model spatial relations, identifying spatial clustering, spillovers or diffusion processes.

The main objective of the course is the theoretical understanding and practical application of spatial regression models. This course will first give an overview on how to perform common spatial operations using spatial information, such as aggregating spatial units, calculating distances, merging spatial data as well as visualizing them. The course will further focus

## *Introduction*

on the analysis of geographic data and the application of spatial regression techniques to model and analyse spatial processes, and furthermore, the course addresses several methods for defining spatial relationships, detecting and diagnosing spatial dependence and autocorrelation. Finally, we will discuss various spatial regression techniques to model processes, clarify the assumptions of these models, and show how they differ in their applications and interpretations.

The field has developed very quickly over the past few years, and *R* now provides a rich set of packages for various spatial data operations. For a more in-depth introduction into spatial data analysis in *R*, have a look into the materials references below.

The material introduces the use of geographical information to connect and analyze different spatial data sources very briefly. This introduction is limited to the fundamentals of using geographical information in *R*. Stefan Jünger & Anne-Kathrin Stroppe have provided a comprehensive GESIS workshop on geospatial techniques in *R*. The focus of this workshop will be on techniques for spatial data analysis, such as spatial regression models.

## **Schedule**

Day 1	Working with Spatial Data
10:00 - 11:30	Refresher on R as GIS Coffee break
11:45 - 13:00	Spatial Data Manipulation & Visualization Lunch break
14:00 - 15:30	Defining Spatial Relationships (W) Coffee break
15:45 - 17:15	Lab Exercises in R

## *Some useful packages*

Day 2      Spatial Regression Models I	
10:00 - 11:30	Detecting Spatial Dependence Coffee break
11:45 - 13:00	Spatial Regression Models: Theory Lunch break
14:00 - 15:30	Estimating Spatial Regression Models Coffee break
15:45 - 17:15	Lab Exercises in R

Day 3      Spatial Regression Models II	
10:00 - 11:30	Interpreting Results: Spatial Impacts Coffee break
11:45 - 13:00	Comparing and Selecting Models Lunch break
14:00 - 15:30	Lab Exercises in R Coffee break
15:45 - 17:15	Other Models

## **Some useful packages**

By now, *R* provides a lot of functionalities for GIS applications and spatial econometrics, and further extensions. There are lots of packages providing a huge variety of spatial functionalities and methods (see e.g. R. Bivand, Millo, and Piras 2021). Important packages for fundamental spatial operations are:

- Spatial data workhorses: sf (Pebesma 2018) and terra
- Visualization: mapview (Appelhans et al. 2021) and tmap (Tennekes 2018)

## *Introduction*

- Spatial weights and other relations: spdep (R. S. Bivand and Rudec 2018)
- Spatial interpolation and kriging: gstat (Gräler, Pebesma, and Heuvelink 2016)
- Spatial regression models: spatialreg and spphet (R. Bivand and Piras 2015)
- The packages have constantly developed over the past years, and older packages such as rgdal, rgeos, and sp are currently retiring (Blog post)

## **Further Readings**

- Great up-to-date introduction to spatial R: Lovelace, Nowosad, and Muenchow (2019), updated version available online
- Great open-science book on Spatial Data Science Pebesma and Bivand (2023)
- Comprehensive introduction to spatial econometrics: LeSage and Pace (2009)
- Relative intuitive introduction to spatial econometrics: Ward and Gleditsch (2008)
- Article-length introductions to spatial econometrics: Elhorst (2012), Halleck Vega and Elhorst (2015), LeSage (2014a), Rüttenauer (2024), and Rüttenauer (2022)

*Course materials*

## **Course materials**

- I highly recommend the great Introduction to Geospatial Techniques for Social Scientists in R including, see Stefan Jünger & Anne-Kathrin Stroppe's GESIS workshop materials. Nice materials on GIS, spatial operations and spatial data visualisation!
- For those looking for a more in-depth introduction, I highly recommend Roger Bivand's course on Spatial Data Analysis: Youtube recordings, Course Materials
- I've learned most of what I know about spatial econometrics from Scott J. Cook and his workshop on Spatial Econometrics at the Essex Summer school.



# 1 Refresher

## Required packages

```
pkgs <- c("sf", "gstat", "mapview", "rngeo", "rnatural-earth", "dplyr",
         "nomisr", "osmdata", "tidyverse", "texreg")
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()

R version 4.4.1 (2024-06-14 ucrt)
Platform: x86_64-w64-mingw32/x64
Running under: Windows 11 x64 (build 22631)

Matrix products: default

locale:
[1] LC_COLLATE=English_United Kingdom.utf8
[2] LC_CTYPE=English_United Kingdom.utf8
[3] LC_MONETARY=English_United Kingdom.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United Kingdom.utf8
```

## 1 Refresher

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets   methods    base

other attached packages:
[1] texreg_1.39.3     tidyverse_1.3.1    osmdata_0.2.5
[4] nomisr_0.4.7      dplyr_1.1.4       rnaturalearth_1.0.1
[7] nngeo_0.4.8       mapview_2.11.2    gstat_2.1-1
[10] sf_1.0-16

loaded via a namespace (and not attached):
 [1] xfun_0.45          raster_3.6-26      htmlwidgets_1.6.4  lattice_0.22-
6
 [5] vctrs_0.6.5        tools_4.4.1       crosstalk_1.2.1   generics_0.1.3
 [9] stats4_4.4.1       tibble_3.2.1      proxy_0.4-27     spacetime_1.3-
1
[13] fansi_1.0.6        xts_0.14.0       pkgconfig_2.0.3   KernSmooth_2.2-
24
[17] satellite_1.0.5    data.table_1.15.4 leaflet_2.2.2    lifecycle_1.0.4
[21] compiler_4.4.1     FNN_1.1.4        rsdmx_0.6-3     munsell_0.5.1
[25] terra_1.7-78      codetools_0.2-20  snakecase_0.11.1 htmltools_0.5.8
[29] class_7.3-22      pillar_1.9.0     classInt_0.4-
10  tidyselect_1.2.1
[33] digest_0.6.35      purrr_1.0.2      fastmap_1.2.0   grid_4.4.1
[37] colorspace_2.1-0   cli_3.6.2       magrittr_2.0.3   base64enc_0.1-
3
[41] XML_3.99-0.16.1   utf8_1.2.4      leafem_0.2.3   e1071_1.7-
14
[45] scales_1.3.0       sp_2.1-4       rmarkdown_2.27   httr_1.4.7
[49] zoo_1.8-12         png_0.1-8       evaluate_0.24.0 knitr_1.47
[53] rlang_1.1.4        Rcpp_1.0.12    glue_1.7.0     DBI_1.2.3
```

## 1.1 Packages

```
[57] rstudioapi_0.16.0  jsonlite_1.8.8      R6_2.5.1          plyr_1.8.9
[61] intervals_0.15.4   units_0.8-5
```

### 1.1 Packages

*Please make sure that you have installed the following packages:*

```
pks <- c("dplyr",
"gstat",
"mapview",
"nngeo",
"nomisr",
"osmdata",
"rnaturalearth",
"sf",
"spatialreg",
"spdep",
"texreg",
"tidyR",
"tmap",
"viridisLite")
```

The most important package is sf: Simple Features for R. users are strongly encouraged to install the sf binary packages from CRAN. If that does not work, please have a look at the installation instructions. It requires software packages GEOS, GDAL and PROJ.

### 1.2 Coordinates

In general, spatial data is structured like conventional/tidy data (e.g. data.frames, matrices), but has one additional dimension: every

## 1 Refresher

observation is linked to some sort of geo-spatial information. Most common types of spatial information are:

- Points (one coordinate pair)
- Lines (two coordinate pairs)
- Polygons (at least three coordinate pairs)
- Regular grids (one coordinate pair for centroid + raster / grid size)

### 1.2.1 Coordinate reference system (CRS)

In its raw form, a pair of coordinates consists of two numerical values. For instance, the pair `c(51.752595, -1.262801)` describes the location of Nuffield College in Oxford (one point). The first number represents the latitude (north-south direction), the second number is the longitude (west-east direction), both are in decimal degrees.

However, we need to specify a reference point for latitudes and longitudes (in the Figure above: equator and Greenwich). For instance, the pair of coordinates above comes from Google Maps which returns GPS coordinates in ‘WGS 84’ (EPSG:4326).

```
# Coordinate pairs of two locations
coords1 <- c(51.752595, -1.262801)
coords2 <- c(51.753237, -1.253904)
coords <- rbind(coords1, coords2)

# Conventional data frame
nuffield.df <- data.frame(name = c("Nuffield College", "Radcliffe Camera"),
                            address = c("New Road", "Radcliffe Sq"),
                            lat = coords[,1], lon = coords[,2])

head(nuffield.df)
```

## 1.2 Coordinates

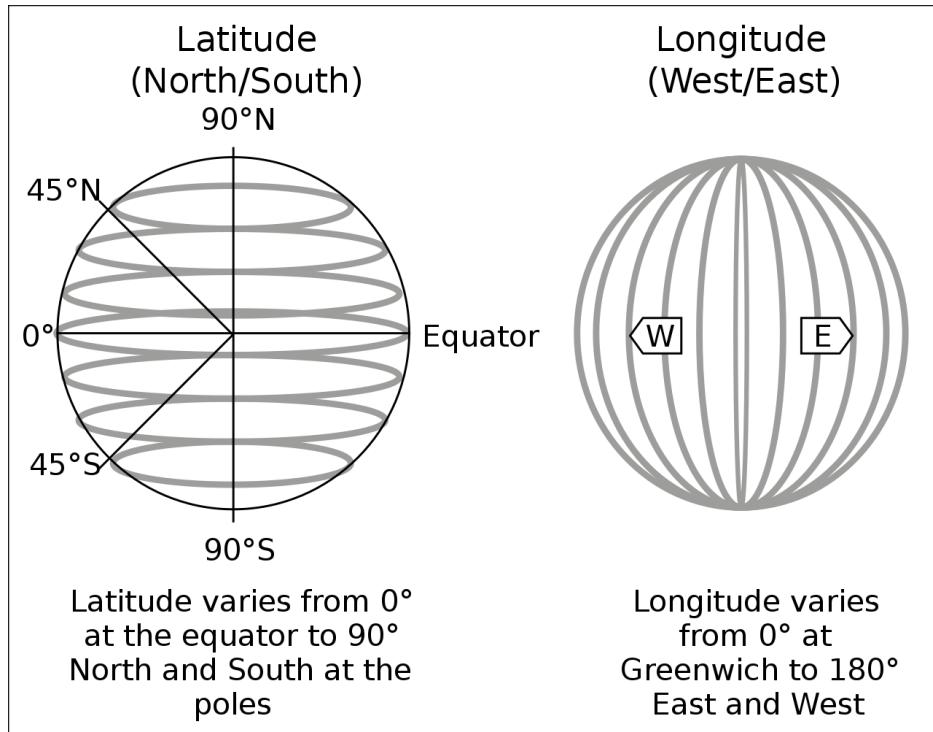


Figure 1.1: Figure: Latitude and longitude, Source: Wikipedia

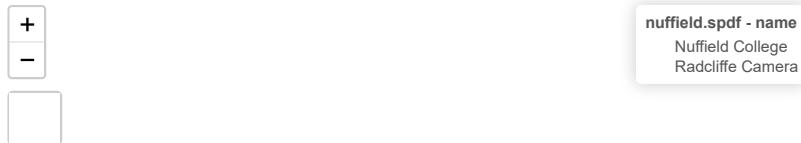
## 1 Refresher

```
          name      address      lat      lon
coords1 Nuffield College    New Road 51.75259 -1.262801
coords2 Radcliffe Camera Radcliffe Sq 51.75324 -1.253904
```

```
# Combine to spatial data frame
nuffield.spdf <- st_as_sf(nuffield.df,
                           coords = c("lon", "lat"), # Order is important
                           crs = 4326) # EPSG number of CRS

# Map
mapview(nuffield.spdf, zcol = "name")
```

## 1.2 Coordinates



### 1.2.2 Projected CRS

However, different data providers use different CRS. For instance, spatial data in the UK usually uses ‘OSGB 1936 / British National Grid’ (EPSG:27700). Here, coordinates are in meters, and projected onto a planar 2D space.

There are a lot of different CRS projections, and different national statistics offices provide data in different projections. Data providers usually specify which reference system they use. This is important as using the correct reference system and projection is crucial for plotting and manipulating spatial data.

If you do not know the correct CRS, try starting with a standards CRS like EPSG:4326 if you have decimal degree like coordinates. If it looks like projected coordinates, try searching for the country or region in CRS libraries like <https://epsg.io/>. However, you must check if the projected coordinates match their real location, e.g. using `mapview()`.

### 1.2.3 Why different projections?

By now, (most) people agree that the earth is not flat. So, to plot data on a 2D planar surface and to perform certain operations on a planar world, we need to make some re-projections. Depending on where we are, different re-projections of our data (globe in this case) might work better than others.

```
world <- ne_countries(scale = "medium", returnclass = "sf")
class(world)

[1] "sf"           "data.frame"
```

## 1.2 Coordinates

```
st_crs(world)
```

Coordinate Reference System:

User input: WGS 84

wkt:

```
GEOGCRS["WGS 84",
    DATUM["World Geodetic System 1984",
        ELLIPSOID["WGS 84",6378137,298.257223563,
            LENGTHUNIT["metre",1]]],
    PRIMEM["Greenwich",0,
        ANGLEUNIT["degree",0.0174532925199433]],
    CS[ellipsoidal,2],
        AXIS["latitude",north,
            ORDER[1],
            ANGLEUNIT["degree",0.0174532925199433]],
        AXIS["longitude",east,
            ORDER[2],
            ANGLEUNIT["degree",0.0174532925199433]],
    ID["EPSG",4326]]
```

```
# Extract a country and plot in current CRS (WGS84)
ger.spdf <- world[world$name == "Germany", ]
plot(st_geometry(ger.spdf))
```

## 1 Refresher



```
# Now, let's transform Germany into a CRS optimized for Iceland
ger_rep.spdf <- st_transform(ger.spdf, crs = 5325)
plot(st_geometry(ger_rep.spdf))
```

### *1.3 Importing some real world data*



Depending on the angle, a 2D projection of the earth looks different. It is important to choose a suitable projection for the available spatial data. For more information on CRS and re-projection, see e.g. Lovelace, Nowosad, and Muenchow (2019) or Stefan Jünger & Anne-Kathrin Stroppe's GESIS workshop materials.

## **1.3 Importing some real world data**

`sf` imports many of the most common spatial data files, like geojson, gpkg, or shp.

## 1 Refresher

### 1.3.1 London shapefile (polygon)

Let's get some administrative boundaries for London from the London Datastore. We use the `sf` package and its function `st_read()` to import the data.

```
# Create subdir (all data will be stored in "_data")
dn <- "_data"
ifelse(dir.exists(dn), "Exists", dir.create(dn))

# Download zip file and unzip
tmpf <- tempfile()
boundary.link <- "https://data.london.gov.uk/download/statistical-gis-boundaries-london/MSOA2011_London_gen_MHW.zip"
download.file(boundary.link, tmpf)
unzip(zipfile = tmpf, exdir = paste0(dn))
unlink(tmpf)

dn <- "_data"
# This is a shapefile
# We only need the MSOA layer for now
msoa.spdf <- st_read(dsn = paste0(dn, "/statistical-gis-boundaries-london/MSOA2011_London_gen_MHW")) # Note: no file extension

Reading layer `MSOA_2011_London_gen_MHW' from data source
`C:\work\Lehre\Geodata_Spatial_Regression\_data\statistical-gis-boundaries-london\ESRI'
using driver `ESRI Shapefile'
Simple feature collection with 983 features and 12 fields
Geometry type: MULTIPOLYGON
Dimension:      XY
Bounding box:  xmin: 503574.2 ymin: 155850.8 xmax: 561956.7 ymax: 200933.6
Projected CRS: OSGB36 / British National Grid
```

### 1.3 Importing some real world data

The object `msoa.spdf` is our spatial data.frame. It looks essentially like a conventional data.frame, but has some additional attributes and geographical information stored with it. Most importantly, notice the column `geometry`, which contains a list of polygons. In most cases, we have one polygon for each line / observation.

```
head(msoa.spdf)
```

	MSOA11CD	MSOA11NM	LAD11CD	LAD11NM	RGN11CD		
1	E02000001	City of London	001 E09000001	City of London	E12000007		
2	E02000002	Barking and Dagenham	001 E09000002	Barking and Dagenham	E12000007		
3	E02000003	Barking and Dagenham	002 E09000002	Barking and Dagenham	E12000007		
4	E02000004	Barking and Dagenham	003 E09000002	Barking and Dagenham	E12000007		
5	E02000005	Barking and Dagenham	004 E09000002	Barking and Dagenham	E12000007		
6	E02000007	Barking and Dagenham	006 E09000002	Barking and Dagenham	E12000007		
	RGN11NM	USUALRES	HHOLDRES	COMESTRES	POPDEN	HHOLDS	AVHHOLDSZ
1	London	7375	7187	188	25.5	4385	1.6
2	London	6775	6724	51	31.3	2713	2.5
3	London	10045	10033	12	46.9	3834	2.6
4	London	6182	5937	245	24.8	2318	2.6
5	London	8562	8562	0	72.1	3183	2.7
6	London	8791	8672	119	50.6	3441	2.5
1	531667.6, 531647.2, 531626.7, 531667.6, 180535.0, 180532.9, 180539.0, 180535.0, 532135.1,						
2							
3				549102.4, 548954.5, 548889.7, 548874.2, 548952.5,			
4							
5							
6							

Shapefiles are still among the most common formats to store and transmit spatial data, despite them being inefficient (file size and file number).

## 1 Refresher

However, `sf` reads everything spatial, such as `geo.json`, which usually is more efficient, but less common (but we're getting there).

```
# Download file
ulez.link <- "https://data.london.gov.uk/download/ultra_low_emissions_zone"
download.file(ulez.link, paste0(dn, "/ulez.json"))

# Read geo.json
st_layers(paste0(dn, "/ulez.json"))

Driver: GeoJSON
Available layers:
      layer_name geometry_type features fields
1 CentralUltraLowEmissionZone Multi Polygon      1      4
      crs_name
1 OSGB36 / British National Grid

ulez.spdf <- st_read(dsn = paste0(dn, "/ulez.json")) # here dsn is simply

Reading layer `CentralUltraLowEmissionZone' from data source
`C:\work\Lehre\Geodata_Spatial_Regression\_data\ulez.json'
using driver `GeoJSON'
Simple feature collection with 1 feature and 4 fields
Geometry type: MULTIPOLYGON
Dimension:     XY
Bounding box:  xmin: 527271.5 ymin: 178041.5 xmax: 533866.3 ymax: 183133.4
Projected CRS: OSGB36 / British National Grid

head(ulez.spdf)
```

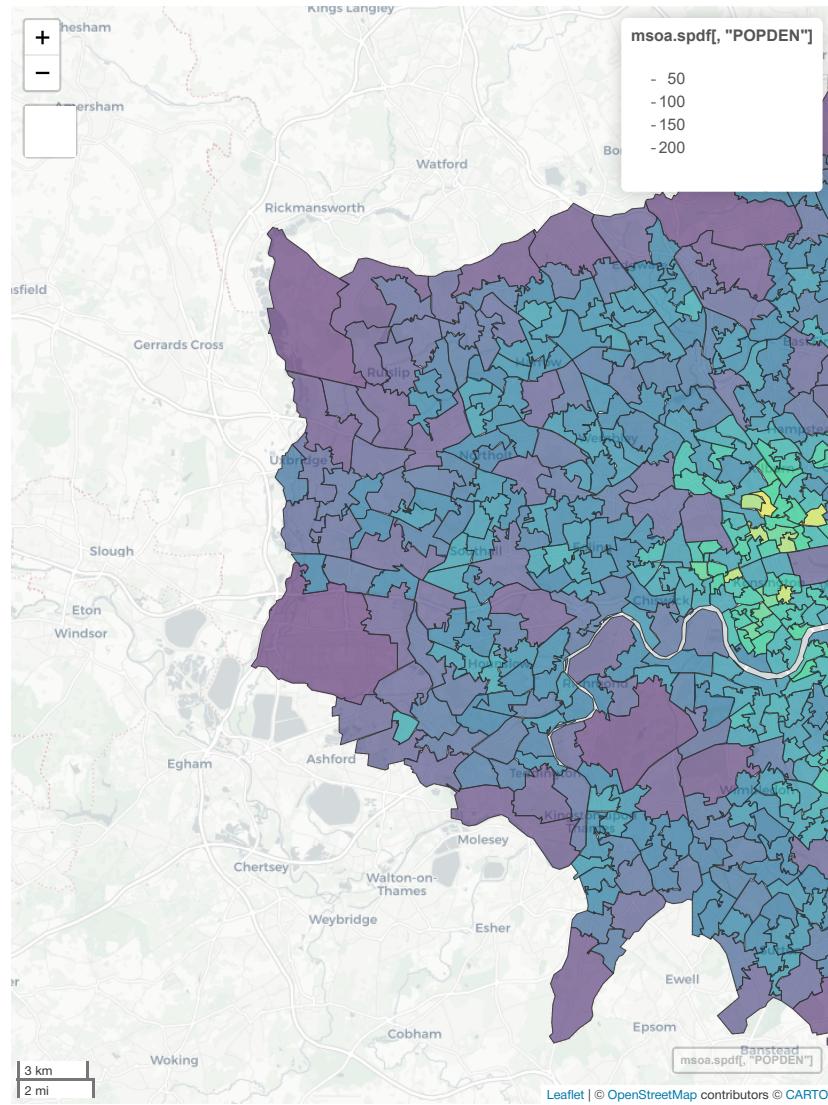
### 1.3 Importing some real world data

```
Simple feature collection with 1 feature and 4 fields
Geometry type: MULTIPOLYGON
Dimension:      XY
Bounding box:   xmin: 527271.5 ymin: 178041.5 xmax: 533866.3 ymax: 183133.4
Projected CRS:  OSGB36 / British National Grid
  fid OBJECTID BOUNDARY Shape_Area           geometry
1 1          1 CSS Area  21.37557 MULTIPOLYGON (((531562.7 18...
```

Again, this looks like a conventional `data.frame` but has the additional column `geometry` containing the coordinates of each observation. `st_geometry()` returns only the geographic object and `st_drop_geometry()` only the `data.frame` without the coordinates. We can plot the object using `mapview()`.

```
mapview(msoa.spdf[, "POPDEN"])
```

## 1 Refresher



### 1.3 Importing some real world data

#### 1.3.2 Census API (admin units)

Now that we have some boundaries and shapes of spatial units in London, we can start looking for different data sources to populate the geometries.

A good source for demographic data is for instance the 2011 census. Below we use the nomis API to retrieve population data for London, See the Vignette for more information (Guest users are limited to 25,000 rows per query). Below is a wrapper to avoid some errors with sex and urban-rural cross-tabulation in some of the data.

```
### For larger request, register and set key
# Sys.setenv(NOMIS_API_KEY = "XXX")
# nomis_api_key(check_env = TRUE)

x <- nomis_data_info()

# Get London ids
london_ids <- msoa.spdf$MSOA11CD

### Get key statistics ids
# select requires tables (https://www.nomisweb.co.uk/sources/census\_2011\_ks)
# Let's get KS201EW (ethnic group), KS205EW (passport held), and KS402EW (housing tenure)

# Get internal ids
stats <- c("KS201EW", "KS402EW", "KS205EW")
oo <- which(grepl(paste(stats, collapse = "|"), x$name.value))
ksids <- x$id[oo]
ksids # This are the internal ids

[1] "NM_608_1" "NM_612_1" "NM_619_1"
```

## 1 Refresher

```
### look at meta information
q <- nomis_overview(ksids[1])
head(q)

# A tibble: 6 x 2
  name          value
  <chr>        <list>
1 analyses     <named list [1]>
2 analysisname <chr [1]>
3 analysisnumber <int [1]>
4 contact      <named list [4]>
5 contenttypes <named list [1]>
6 coverage      <chr [1]>

a <- nomis_get_metadata(id = ksids[1], concept = "GEOGRAPHY", type = "type")
a # TYPE297 is MSOA level

# A tibble: 24 x 3
  id      label.en      description.en
  <chr>   <chr>        <chr>
1 TYPE265 NHS area teams NHS area teams
2 TYPE266 clinical commissioning groups clinical commissioning groups
3 TYPE267 built-up areas including subdivisions built-up areas including subdivisions
  up areas in~ built-up areas in~
4 TYPE269 built-up areas built-up areas
  up areas built-up areas
5 TYPE273 national assembly for wales electoral regions 2010 national assembly for wales electoral regions 2010
6 TYPE274 postcode areas postcode areas
7 TYPE275 postcode districts postcode districts
8 TYPE276 postcode sectors postcode sectors
9 TYPE277 national assembly for wales constituencies 2010 national assembly for wales constituencies 2010
```

### 1.3 Importing some real world data

```
10 TYPE279 parishes 2011                                parishes 2011
# i 14 more rows

  b <- nomis_get_metadata(id = ksids[1], concept = "MEASURES", type = "TYPE297")
  b # 20100 is the measure of absolute numbers

# A tibble: 2 x 3
  id    label.en description.en
  <chr> <chr>     <chr>
1 20100 value     value
2 20301 percent   percent

#### Query data in loop over the required statistics
for(i in ksids){

  # Determin if data is divided by sex or urban-rural
  nd <- nomis_get_metadata(id = i)
  if("RURAL_URBAN" %in% nd$conceptref){
    UR <- TRUE
  }else{
    UR <- FALSE
  }
  if("C_SEX" %in% nd$conceptref){
    SEX <- TRUE
  }else{
    SEX <- FALSE
  }

  # make data request
  if(UR == TRUE){
    if(SEX == TRUE){
      tmp_en <- nomis_get_data(id = i, time = "2011",
```

## 1 Refresher

```
geography = london_ids, # replace with "TYP"
measures = 20100, RURAL_URBAN = 0, C_SEX =
}else{
  tmp_en <- nomis_get_data(id = i, time = "2011",
                            geography = london_ids, # replace with "TYP"
                            measures = 20100, RURAL_URBAN = 0)
}
}else{
  if(SEX == TRUE){
    tmp_en <- nomis_get_data(id = i, time = "2011",
                              geography = london_ids, # replace with "TYP"
                              measures = 20100, C_SEX = 0)
  }else{
    tmp_en <- nomis_get_data(id = i, time = "2011",
                              geography = london_ids, # replace with "TYP"
                              measures = 20100)
  }
}

# Append (in case of different regions)
ks_tmp <- tmp_en

# Make lower case names
names(ks_tmp) <- tolower(names(ks_tmp))
names(ks_tmp)[names(ks_tmp) == "geography_code"] <- "msoa11"
names(ks_tmp)[names(ks_tmp) == "geography_name"] <- "name"

# replace weird cell codes
onlynum <- which(grepl("^[[[:digit:]]]+$", ks_tmp$cell_code))
if(length(onlynum) != 0){
  code <- substr(ks_tmp$cell_code[-onlynum][1], 1, 7)
```

### 1.3 Importing some real world data

```
if(is.na(code)){
  code <- i
}
ks_tmp$cell_code[onlynum] <- paste0(code, "_", ks_tmp$cell_code[onlynum])
}

# save codebook
ks_cb <- unique(ks_tmp[, c("date", "cell_type", "cell", "cell_code", "cell_name")])

#### Reshape
ks_res <- tidyr::pivot_wider(ks_tmp, id_cols = c("msoa11", "name"),
                             names_from = "cell_code",
                             values_from = "obs_value")

#### Merge
if(i == ksids[1]){
  census_keystat.df <- ks_res
  census_keystat_cb.df <- ks_cb
}else{
  census_keystat.df <- merge(census_keystat.df, ks_res, by = c("msoa11", "name"), all =
  census_keystat_cb.df <- rbind(census_keystat_cb.df, ks_cb)
}

}

# Descriptions are saved in the codebook
save(census_keystat.df, file = "_data/Census_ckeystat.RData")
save(census_keystat_cb.df, file = "_data/Census_codebook.RData")
```

Now, we have one file containing the geometries of MSOAs and one file with the census information on ethnic groups. Obviously, we can easily merge them together using the MSOA identifiers.

## 1 Refresher

```
load("_data/Census_ckeystat.RData")
msoa.spdf <- merge(msoa.spdf, census_keystat.df,
                     by.x = "MSOA11CD", by.y = "msoa11", all.x = TRUE)
```

And we can, for instance, plot the spatial distribution of ethnic groups.

```
msoa.spdf$per_white <- msoa.spdf$KS201EW_100 / msoa.spdf$KS201EW0001 * 100
msoa.spdf$per_mixed <- msoa.spdf$KS201EW_200 / msoa.spdf$KS201EW0001 * 100
msoa.spdf$per_asian <- msoa.spdf$KS201EW_300 / msoa.spdf$KS201EW0001 * 100
msoa.spdf$per_black <- msoa.spdf$KS201EW_400 / msoa.spdf$KS201EW0001 * 100
msoa.spdf$per_other <- msoa.spdf$KS201EW_500 / msoa.spdf$KS201EW0001 * 100

mapview(msoa.spdf[, "per_white"])
```

If you're interested in more data sources, see for instance APIs for social scientists: A collaborative review by Paul C. Bauer, Camille Landesvatter, Lion Behrens. It's a collection of several APIs for social sciences.

### 1.3.3 Gridded data

So far, we have queried data on administrative units. However, often data comes on other spatial scales. For instance, we might be interested in the amount of air pollution, which is provided on a regular grid across the UK from Defra.

```
# Download
pol.link <- "https://uk-air.defra.gov.uk/datastore/pcm/mapno22011.csv"
download.file(pol.link, paste0(dn, "/mapno22011.csv"))

pol.df <- read.csv(paste0(dn, "/mapno22011.csv"), skip = 5, header = T, sep = ",",
                     stringsAsFactors = F, na.strings = "MISSING")
```

### 1.3 Importing some real world data

```
head(pol.df)
```

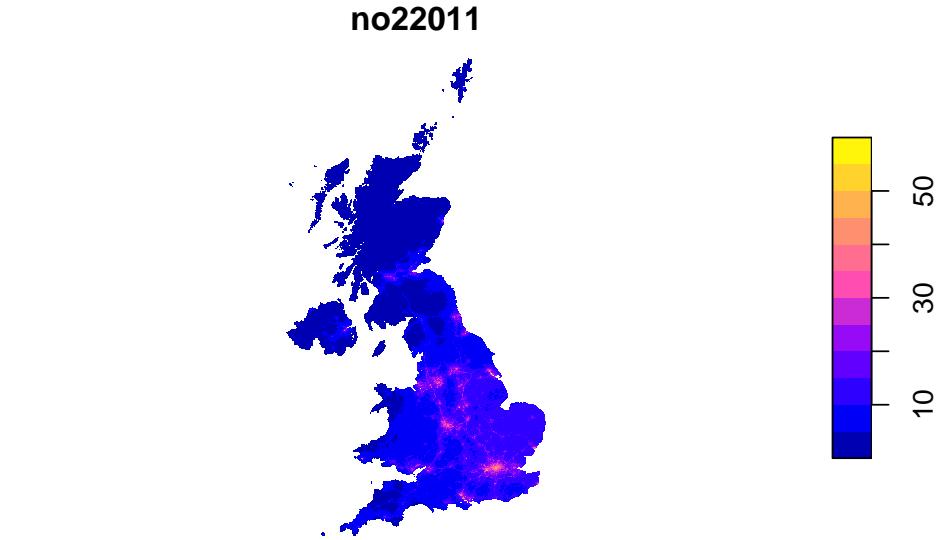
```
ukgridcode      x      y no22011
1      54291 460500 1221500     NA
2      54292 461500 1221500     NA
3      54294 463500 1221500     NA
4      54979 458500 1220500     NA
5      54980 459500 1220500     NA
6      54981 460500 1220500     NA
```

The data comes as point data with x and y as coordinates. We have to transform this into spatial data first. We first setup a spatial points object with `st_as_sf`. Subsequently, we transform the point coordinates into a regular grid. We use a buffer method `st_buffer` with “diameter”, and only one segment per quadrant (`nQuadSegs`). This gives us a 1x1km regular grid.

```
# Build spatial object
pol.spdf <- st_as_sf(pol.df, coords = c("x", "y"),
                      crs = 27700)

# we transform the point coordinates into a regular grid with "diameter" 500m
pol.spdf <- st_buffer(pol.spdf, dist = 500, nQuadSegs = 1,
                        endCapStyle = 'SQUARE')

# Plot NO2
plot(pol.spdf[, "no22011"], border = NA)
```



### 1.3.4 OpenStreetMap (points)

Another interesting data source is the OpenStreetMap API, which provides information about the geographical location of a series of different indicators. Robin Lovelace provides a nice introduction to the osmdata API. Available features can be found on OSM wiki.

First we create a bounding box of where we want to query data. `st_bbox()` can be used to get bounding boxes of an existing spatial object (needs CRS = 4326). An alternative would be to use `opq(bbox = 'greater london uk')`.

```
# bounding box of where we want to query data
q <- opq(bbox = st_bbox(st_transform(msoa.spdf, 4326)))
```

### 1.3 Importing some real world data

And we want to get data for all pubs and bars which are within this bounding box.

```
# First build the query of location of pubs in London  
osmq <- add_osm_feature(q, key = "amenity", value = "pub")  
  
# And then query the data  
pubs.osm <- osmdata_sf(osmq)
```

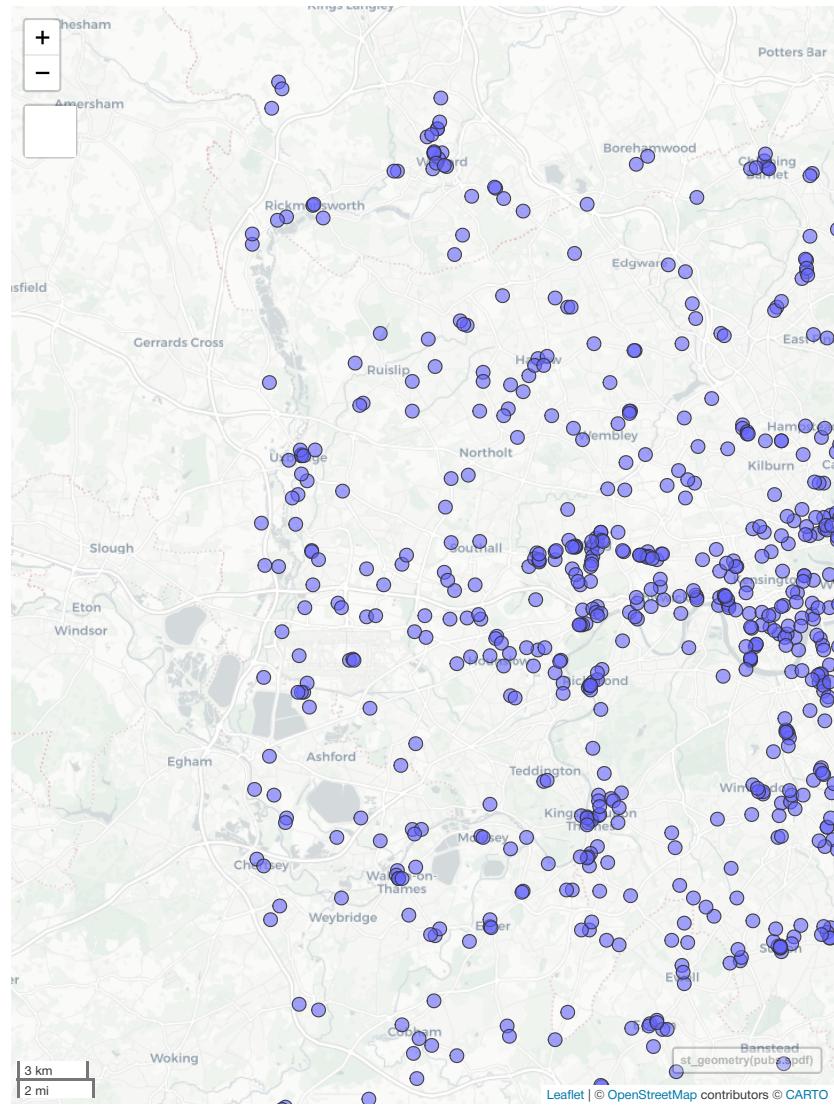
Right now there are some results in polygons, some in points, and they overlap. Often, data from OSM needs some manual cleaning. Sometimes the same features are represented by different spatial objects (e.g. points + polygons).

```
# Make unique points / polygons  
pubs.osm <- unique_osmdata(pubs.osm)  
  
# Get points and polygons (there are barely any pubs as polygons, so we ignore them)  
pubs.points <- pubs.osm$osm_points  
pubs.polys <- pubs.osm$osm_multipolygons  
  
# # Drop OSM file  
# rm(pubs.osm); gc()  
  
# Reduce to point object only  
pubs.spdf <- pubs.points  
  
# Reduce to a few variables  
pubs.spdf <- pubs.spdf[, c("osm_id", "name", "addr:postcode", "diet:vegan")]
```

Again, we can inspect the results with `mapview`.

```
mapview(st_geometry(pubs.spdf))
```

## 1 Refresher



### 1.3 Importing some real world data

Note that OSM is solely based on contribution by users, and the **quality of OSM data varies**. Usually data quality is better in larger cities, and better for more stable features (such as hospitals, train stations, highways) rather than pubs or restaurants which regularly appear and disappear. However, data from London Datastore would indicate more pubs than what we find with OSM.

#### 1.3.5 Save

We will store the created data to use them again in the next session.

```
save(msoa.spdf, file = "_data/msoa_spatial.RData")
save(ulez.spdf, file = "_data/ulez_spatial.RData")
save(pol.spdf, file = "_data/pollution_spatial.RData")
save(pubs.spdf, file = "_data/pubs_spatial.RData")
```



## 2 Data Manipulation & Visualization

### Required packages

```
pkgs <- c("sf", "gstat", "mapview", "nngeo", "rnatural-earth", "dplyr",
         "nomisr", "osmdata", "OpenStreetMap", "tidyverse", "texreg", "downlit", "xml2")
lapply(pkgs, require, character.only = TRUE)
```

For mapping

```
pkgs <- c("tmap", "tmaptools", "viridisLite",
         "ggplot2", "ggthemes", "rmapshaper", "cowplot")
lapply(pkgs, require, character.only = TRUE)
```

### Session info

```
sessionInfo()
```

```
R version 4.4.1 (2024-06-14 ucrt)
Platform: x86_64-w64-mingw32/x64
Running under: Windows 11 x64 (build 22631)
```

```
Matrix products: default
```

## 2 Data Manipulation & Visualization

```
locale:  
[1] LC_COLLATE=English_United Kingdom.utf8  
[2] LC_CTYPE=English_United Kingdom.utf8  
[3] LC_MONETARY=English_United Kingdom.utf8  
[4] LC_NUMERIC=C  
[5] LC_TIME=English_United Kingdom.utf8  
  
time zone: Europe/Berlin  
tzcode source: internal  
  
attached base packages:  
[1] stats      graphics   grDevices utils      datasets  methods    base  
  
other attached packages:  
[1] cowplot_1.1.3      rmapshaper_0.5.0      ggthemes_5.1.0  
[4] ggplot2_3.5.1      viridisLite_0.4.2      tmapproj_3.1-1  
[7] tmap_3.3-4         xml2_1.3.6          downlit_0.4.4  
[10] texreg_1.39.3     tidyverse_1.3.1      OpenStreetMap_0.4.0  
[13] osmdata_0.2.5     nomisr_0.4.7        dplyr_1.1.4  
[16] rnaturalearth_1.0.1 nngeo_0.4.8       mapview_2.11.2  
[19] gstat_2.1-1       sf_1.0-16  
  
loaded via a namespace (and not attached):  
[1] tidyselect_1.2.1      fastmap_1.2.0      leaflet_2.2.2      XML_3.99-  
0.16.1  
[5] digest_0.6.35        lifecycle_1.0.4      terra_1.7-78      magrittr_2.0.3  
[9] compiler_4.4.1       rlang_1.1.4        tools_4.4.1       utf8_1.2.4  
[13] rsdmx_0.6-3        data.table_1.15.4  knitr_1.47       FNN_1.1.4  
[17] htmlwidgets_1.6.4    curl_5.2.1        sp_2.1-4        classInt_0.4-  
10  
[21] plyr_1.8.9         RColorBrewer_1.1-3  abind_1.4-5      KernSmooth_2.2-  
24  
[25] withr_3.0.0        purrrr_1.0.2      leafsync_0.1.0     grid_4.4.1  
[29] stats4_4.4.1       fansi_1.0.6       xts_0.14.0       e1071_1.7-
```

## 2.1 Manipulation and linkage

```
14
[33] leafem_0.2.3      colorspace_2.1-0    scales_1.3.0      dichromat_2.0-
0.1
[37] cli_3.6.2        rmarkdown_2.27     intervals_0.15.4  generics_0.1.3
[41] rstudioapi_0.16.0 httr_1.4.7       DBI_1.2.3        cachem_1.1.0
[45] proxy_0.4-27     stringr_1.5.1     stars_0.6-5      parallel_4.4.1
[49] base64enc_0.1-3   vctrs_0.6.5      V8_4.4.2         jsonlite_1.8.8
[53] crosstalk_1.2.1  units_0.8-5      glue_1.7.0       lwgeom_0.2-
14
[57] codetools_0.2-20  stringi_1.8.4    rJava_1.0-11     gtable_0.3.5
[61] raster_3.6-26    munsell_0.5.1    tibble_3.2.1     pillar_1.9.0
[65] htmltools_0.5.8.1 satellite_1.0.5  R6_2.5.1        evaluate_0.24.0
[69] lattice_0.22-6   png_0.1-8       memoise_2.0.1    snakecase_0.11.1
[73] class_7.3-22    Rcpp_1.0.12      spacetime_1.3-
1 xfun_0.45
[77] zoo_1.8-12       pkgconfig_2.0.3
```

### Reload data from previous session

```
load("_data/msoa_spatial.RData")
load("_data/ulez_spatial.RData")
load("_data/pollution_spatial.RData")
load("_data/pubs_spatial.RData")
```

## 2.1 Manipulation and linkage

Having data with geo-spatial information allows to perform a variety of methods to manipulate and link different data sources. Commonly used methods include 1) subsetting, 2) point-in-polygon operations, 3) distance measures, 4) intersections or buffer methods.

## 2 Data Manipulation & Visualization

The online Vignettes of the sf package provide a comprehensive overview of the multiple ways of spatial manipulations.

### 2.1.0.1 Check if data is on common projection

```
st_crs(msoa.spdf) == st_crs(pol.spdf)
```

```
[1] FALSE
```

```
st_crs(msoa.spdf) == st_crs(pubs.spdf)
```

```
[1] FALSE
```

```
st_crs(msoa.spdf) == st_crs(ulez.spdf)
```

```
[1] FALSE
```

The spatial data files are on different projections. Before we can do any spatial operations with them, we have to transform them into a common projection.

```
# MSOA in different crs --> transform  
pol.spdf <- st_transform(pol.spdf, crs = st_crs(msoa.spdf))  
pubs.spdf <- st_transform(pubs.spdf, crs = st_crs(msoa.spdf))  
ulez.spdf <- st_transform(ulez.spdf, crs = st_crs(msoa.spdf))
```

```
# Check if all geometries are valid, and make valid if needed
```

## 2.1 Manipulation and linkage

```
msoa.spdf <- st_make_valid(msoa.spdf)
```

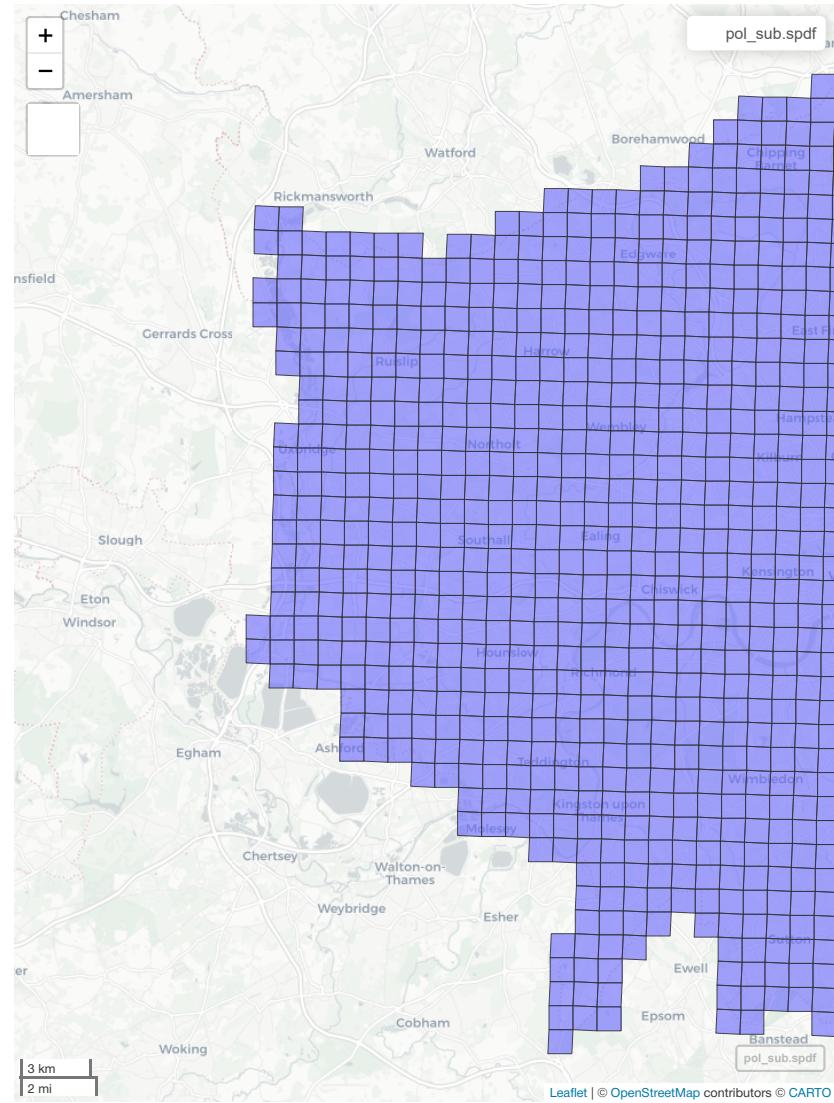
The `st_make_valid()` function can help if the spatial geometries have some problems such as holes or points that don't match exactly.

### 2.1.1 Subsetting

We can subset spatial data in a similar way as we subset conventional data.frames or matrices. For instance, below we simply reduce the pollution grid across the UK to observations in London only.

```
# Subset to pollution estimates in London
pol_sub.spdf <- pol.spdf[msoa.spdf, ] # or:
pol_sub.spdf <- st_filter(pol.spdf, msoa.spdf)
mapview(pol_sub.spdf)
```

## 2 Data Manipulation & Visualization

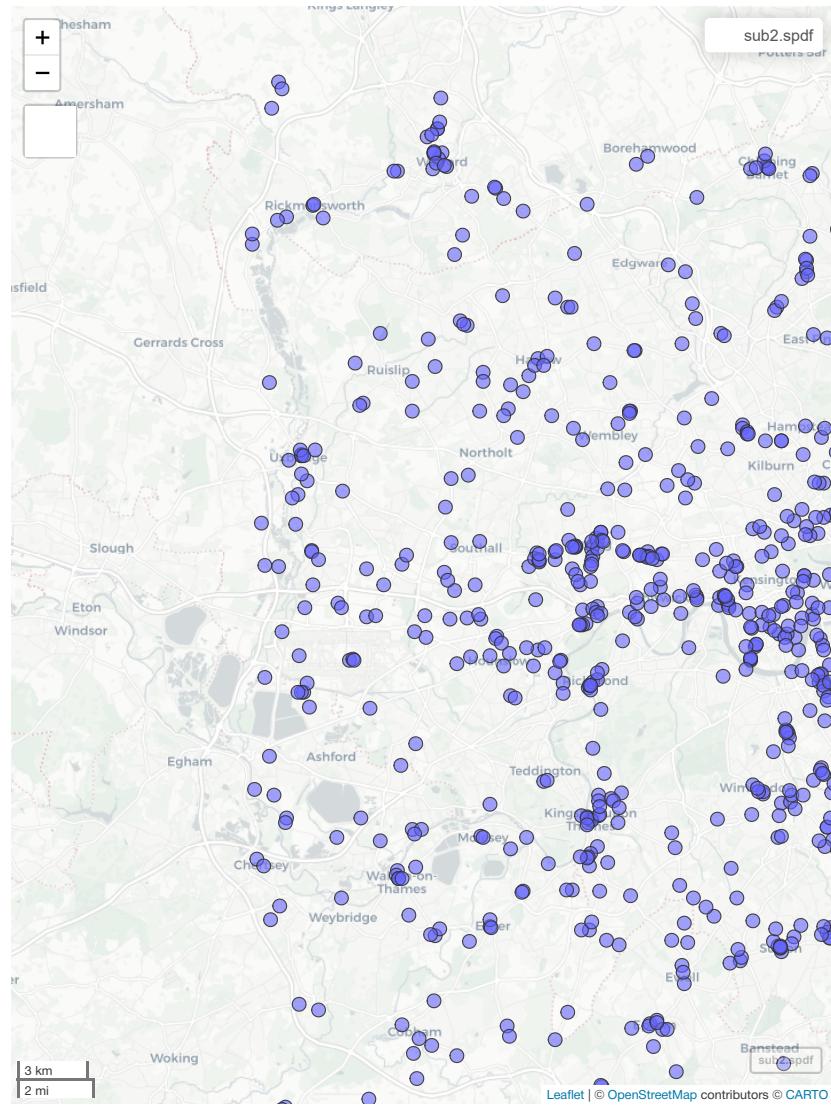


## 2.1 Manipulation and linkage

Or we can reverse the above and exclude all intersecting units by specifying `st_disjoint` as alternative spatial operation using the `op =` option (note the empty space for column selection). `st_filter()` with the `.predicate` option does the same job. See the sf Vignette for more operations.

```
# Subset pubs to pubs not in the ulez area
sub2.spdf <- pubs.spdf[ulez.spdf, , op = st_disjoint] # or:
sub2.spdf <- st_filter(pubs.spdf, ulez.spdf, .predicate = st_disjoint)
mapview(sub2.spdf)
```

## 2 Data Manipulation & Visualization



## 2.1 Manipulation and linkage

We can easily create indicators of whether an MSOA is within ulez or not.

```
msoa.spdf$ulez <- 0

# intersecting lsoas
within <- msoa.spdf[ulez.spdf,]

# use their ids to create binary indicator
msoa.spdf$ulez[which(msoa.spdf$MSOA11CD %in% within$MSOA11CD)] <- 1
table(msoa.spdf$ulez)
```

```
0    1
955  28
```

### 2.1.2 Point in polygon

We are interested in the number of pubs in each MSOA. So, we count the number of points in each polygon.

```
# Assign MSOA to each point
pubs_msoa.join <- st_join(pubs.spdf, msoa.spdf, join = st_within)

# Count N by MSOA code (drop geometry to speed up)
pubs_msoa.join <- dplyr::count(st_drop_geometry(pub_msoa.join),
                                MSOA11CD = pub_msoa.join$MSOA11CD,
                                name = "pubs_count")
sum(pub_msoa.join$pubs_count)
```

```
[1] 1601
```

## 2 Data Manipulation & Visualization

```
# Merge and replace NAs with zero (no matches, no pubs)
msoa.spdf <- merge(msoa.spdf, pubs_msoa.join,
                     by = "MSOA11CD", all.x = TRUE)
msoa.spdf$pubs_count[is.na(msoa.spdf$pubs_count)] <- 0
```

### 2.1.3 Distance measures

We might be interested in the distance to the nearest pub. Here, we use the package `nngeo` to find k nearest neighbours with the respective distance.

```
# Use geometric centroid of each MSOA
cent.sp <- st_centroid(msoa.spdf[, "MSOA11CD"])
```

Warning: `st_centroid` assumes attributes are constant over geometries

```
# Get K nearest neighbour with distance
knb.dist <- st_nn(cent.sp,
                  pubs.spdf,
                  k = 1, # number of nearest neighbours
                  returnDist = TRUE, # we also want the distance
                  progress = FALSE)
```

projected points

```
msoa.spdf$dist_pubs <- unlist(knb.dist$dist)
summary(msoa.spdf$dist_pubs)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
9.079	305.149	565.018	701.961	948.047	3735.478

## 2.1 Manipulation and linkage

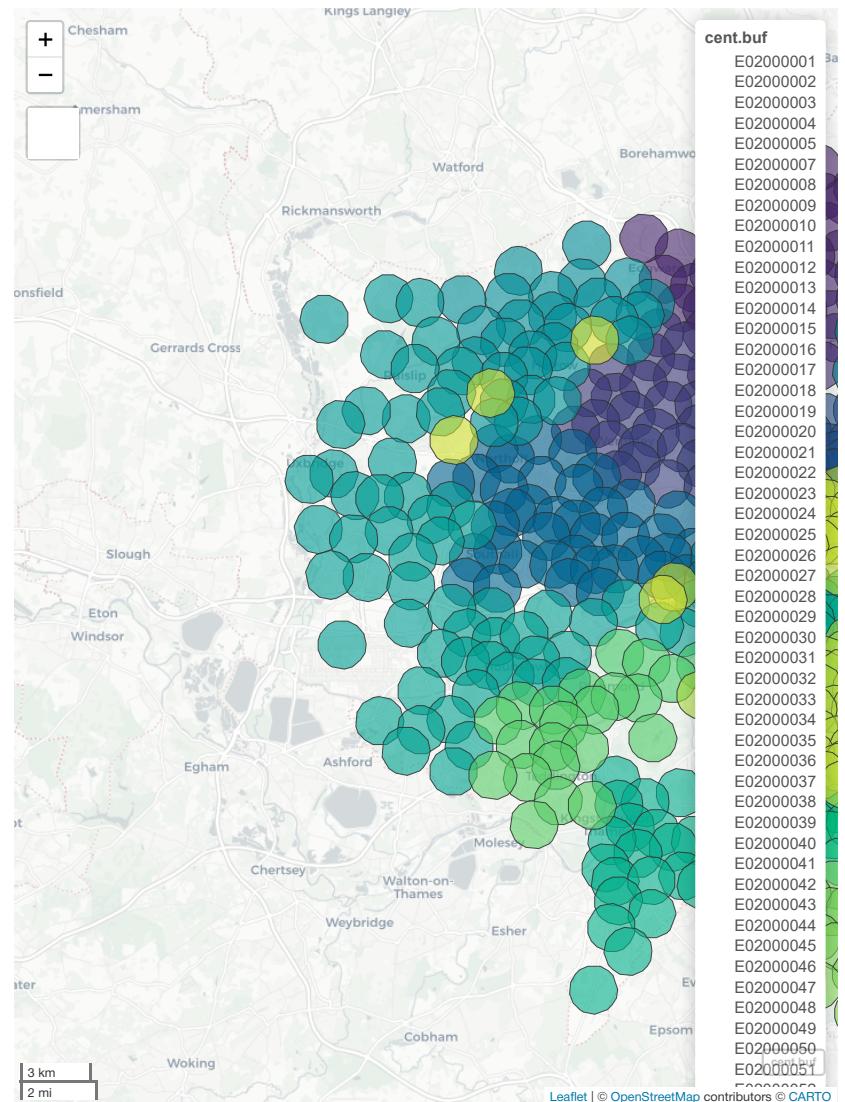
### 2.1.4 Intersections + Buffers

We may also want the average pollution within 1 km radius around each MSOA centroid. Note that it is usually better to use a ego-centric method where you calculate the average within a distance rather than using the characteristic of the intersecting cells only (B. A. Lee et al. 2008; Mohai and Saha 2007).

Therefore, we first create a buffer with `st_buffer()` around each midpoint and subsequently use `st_intersection()` to calculate the overlap.

```
# Create buffer (1km radius)
cent.buf <- st_buffer(cent.sp,
                      dist = 1000) # dist in meters
mapview(cent.buf)
```

## 2 Data Manipulation & Visualization



## 2.1 Manipulation and linkage

```
# Add area of each buffer (in this constant)
cent.buf$area <- as.numeric(st_area(cent.buf))

# Calculate intersection of pollution grid and buffer
int.df <- st_intersection(cent.buf, pol.spdf)
```

Warning: attribute variables are assumed to be spatially constant throughout all geometries

```
int.df$int_area <- as.numeric(st_area(int.df)) # area of intersection

# Area of intersection as share of buffer
int.df$area_per <- int.df$int_area / int.df$area
```

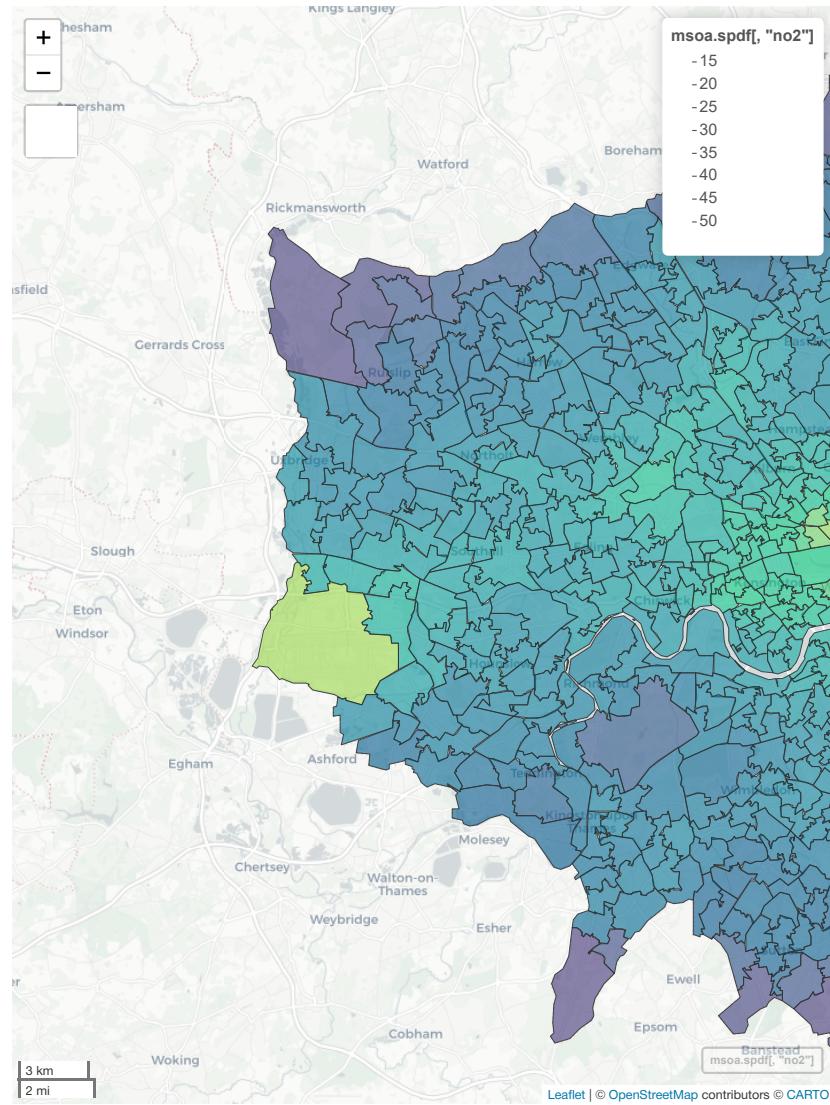
And we use the percent overlap areas as the weights to calculate a weighted mean.

```
# Aggregate as weighted mean
int.df <- st_drop_geometry(int.df)
int.df$no2_weighted <- int.df$no22011 * int.df$area_per
int.df <- aggregate(list(no2 = int.df[, "no2_weighted"]),
                      by = list(MSOA11CD = int.df$MSOA11CD),
                      sum)

# Merge back to spatial data.frame
msoa.spdf <- merge(msoa.spdf, int.df, by = "MSOA11CD", all.x = TRUE)

mapview(msoa.spdf[, "no2"])
```

## 2 Data Manipulation & Visualization



## 2.1 Manipulation and linkage

Note: for buffer related methods, it often makes sense to use population weighted centroids instead of geographic centroids (see here for MSOA population weighted centroids). However, often this information is not available.

### 2.1.5 and more

There are more spatial operation possible using sf. Have a look at the sf Cheatsheet.

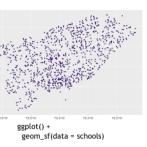
## Spatial manipulation with sf: : CHEAT SHEET

The sf package provides a set of tools for working with geospatial vectors, i.e. points, lines, polygons, etc.



**Geometric confirmation**

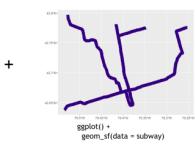
- st\_contains(x, y, ...) Identifies if y is within x (i.e. point within polygon)
- st\_covered\_by(x, y, ...) Identifies if x is completely within y (i.e. polygon completely within polygon)
- st\_covers(x, y, ...) Identifies if any point from x is outside of y (i.e. polygon outside polygon)
- st\_crosses(x, y, ...) Identifies if any geometry of x have commonalities with y
- st\_disjoint(x, y, ...) Identifies when geometries from x do not share space with y
- st\_equals(x, y, ...) Identifies if x and y share the same geometry
- st\_intersects(x, y, ...) Identifies if x and y geometry share any space
- st\_overlaps(x, y, ...) Identifies if geometries of x and y share space, are of the same dimension, but are not completely contained by each other
- st\_touches(x, y, ...) Identifies if geometries of x and y share a common point but their interiors do not intersect
- st\_within(x, y, ...) Identifies if x is in a specified distance to y



```
ggplot() + geom_sf(data = schools)
```

**Geometric operations**

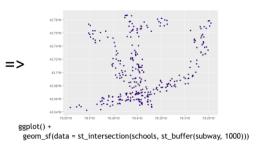
- st\_boundary(x) Creates a polygon that encompasses the full extent of the geometry
- st\_buffer(x, dist, nQuadSegs) Creates a polygon covering all points of the geometry within a given distance
- st\_centroid(x, ..., of\_largest\_polygon) Creates a point at the geometric centre of the geometry
- st\_convex\_hull(x) Creates geometry that represents the minimum convex geometry of x
- st\_line\_merge(x) Creates linestring geometry from sewing multi linestring geometry together
- st\_node(x) Creates nodes on overlapping geometry where nodes do not exist
- st\_point\_on\_surface(x) Creates a point that is guaranteed to fall on the surface of the geometry
- st\_polyonize(x) Creates polygon geometry from linestring geometry
- st\_segmentize(x, dfMaxLength, ...) Creates linestring geometry from x based on a specified length
- st\_simplify(x, preserveTopology, dTolerance) Creates a simplified version of the geometry based on a specified tolerance



```
ggplot() + geom_sf(data = subway)
```

**Geometry creation**

- st\_triangulate(x, dTolerance, bOnlyEdges) Creates polygon geometry as triangles from point geometry
- st\_voronoi(x, envelope, dTolerance, bOnlyEdges) Creates polygon geometry covering the envelope of x, with x at the centre of the geometry
- st\_point(x, cnumeric vector), dim = "XYZ") Creating point geometry from numeric values
- st\_multipoint(x = matrix<numeric values in rows>, dim = "XYZ") Creating multi point geometry from numeric values
- st\_linestring(x = matrix<numeric values in rows>, dim = "XYZ") Creating linestring geometry from numeric values
- st\_multilinestring(x = list<numeric matrices in rows>, dim = "XYZ") Creating multi linestring geometry from numeric values
- st\_polygon(x = list<numeric matrices in rows>, dim = "XYZ") Creating polygon geometry from numeric values
- st\_multipolygon(x = list<numeric matrices in rows>, dim = "XYZ") Creating multi polygon geometry from numeric values



```
ggplot() + geom_sf(data = st_intersection(schools, st_buffer(subway, 1000)))
```

This cheatsheet presents the sf package [Edzer Pebesma 2018] in version 0.6.3. See <https://github.com/r-spatial/sf> for more details.

CC BY Ryan Garnett <http://github.com/yngarnett>

<https://creativecommons.org/licenses/by/4.0/>

### 2.1.6 Air pollution and ethnic minorities

With a few lines of code, we have compiled an original dataset containing demographic information, air pollution, and some infrastructural information.

Let's see what we can do with it.

```
# Define ethnic group shares
msoa.spdf$per_mixed <- msoa.spdf$KS201EW_200 / msoa.spdf$KS201EW0001 * 100
msoa.spdf$per_asian <- msoa.spdf$KS201EW_300 / msoa.spdf$KS201EW0001 * 100
msoa.spdf$per_black <- msoa.spdf$KS201EW_400 / msoa.spdf$KS201EW0001 * 100
msoa.spdf$per_other <- msoa.spdf$KS201EW_500 / msoa.spdf$KS201EW0001 * 100

# Define tenure
msoa.spdf$per_owner <- msoa.spdf$KS402EW_100 / msoa.spdf$KS402EW0001 * 100
msoa.spdf$per_social <- msoa.spdf$KS402EW_200 / msoa.spdf$KS402EW0001 * 100

# Non British passport
msoa.spdf$per_nonUK <- (msoa.spdf$KS205EW0001 - msoa.spdf$KS205EW0003) / msoa.spdf$KS205EW0001
msoa.spdf$per_nonEU <- (msoa.spdf$KS205EW0001 - msoa.spdf$KS205EW0003 -
                           msoa.spdf$KS205EW0004 - msoa.spdf$KS205EW0005 -
                           msoa.spdf$KS205EW0006) / msoa.spdf$KS205EW0001 *
msoa.spdf$per_nonUK_EU <- (msoa.spdf$KS205EW0005 + msoa.spdf$KS205EW0006)

# Run regression
mod1.lm <- lm(no2 ~ per_mixed + per_asian + per_black + per_other +
                per_owner + per_social + pubs_count + POPDEN + ulez,
                data = msoa.spdf)

# summary
screenreg(list(mod1.lm), digits = 3)
```

## 2.1 Manipulation and linkage

```
=====
      Model 1
-----
(Intercept) 37.112 ***
                  (1.308)
per_mixed    -0.090
                  (0.099)
per_asian     0.018 *
                  (0.007)
per_black     -0.085 ***
                  (0.016)
per_other     0.462 ***
                  (0.047)
per_owner     -0.207 ***
                  (0.013)
per_social    -0.058 ***
                  (0.013)
pubs_count    0.218 ***
                  (0.040)
POPDEN        0.037 ***
                  (0.003)
ulez          9.556 ***
                  (0.686)
-----
R^2            0.774
Adj. R^2       0.772
Num. obs.     983
=====

*** p < 0.001; ** p < 0.01; * p < 0.05
```

For some examples later, we also add data on house prices. We use the median house prices in 2017 from the London Datastore.

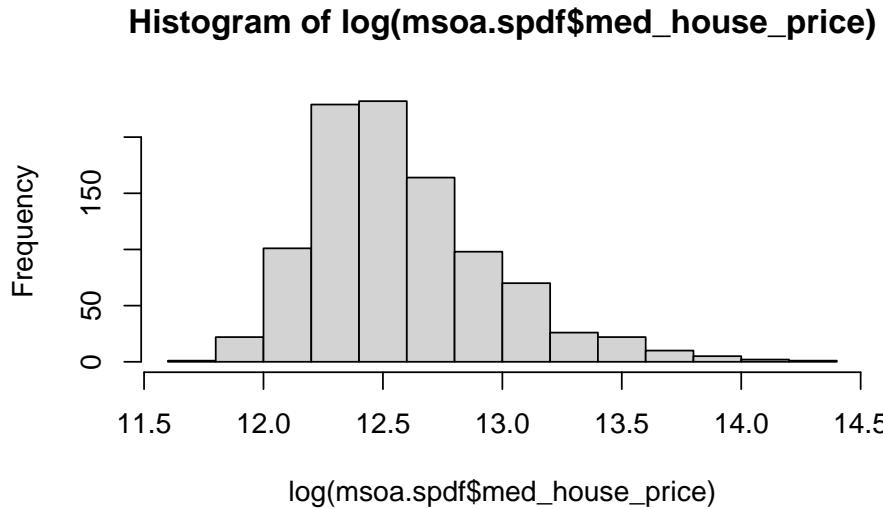
## 2 Data Manipulation & Visualization

```
# Download
hp.link <- "https://data.london.gov.uk/download/average-house-prices/bdf8e"
hp.df <- read.csv(hp.link)
hp.df <- hp.df[which(hp.df$Measure == "Median" &
                      grep("2011", hp.df$Year)), ]
table(hp.df$Year)

Year ending Dec 2011 Year ending Jun 2011 Year ending Mar 2011
983                     983                     983
Year ending Sep 2011
983

# Aggregate across 2011 values
hp.df$med_house_price <- as.numeric(hp.df$value)
hp.df <- aggregate(hp.df[, "med_house_price", drop = FALSE],
                    by = list(MSOA11CD = hp.df$Code),
                    FUN = function(x) mean(x, na.rm = TRUE))

# Merge spdf and housing prices
msoa.spdf <- merge(msoa.spdf, hp.df,
                    by = "MSOA11CD",
                    all.x = TRUE, all.y = FALSE)
hist(log(msoa.spdf$med_house_price))
```



### 2.1.7 Save spatial data

```
# Save
save(msoa.spdf, file = "_data/msoa2_spatial.RData")
```

## 2.2 Visualization

A large advantage of spatial data is that different data sources can be connected and combined. Another nice advantage is: you can create very nice maps. And it's quite easy to do! Stefan Jünger & Anne-Kathrin Stroppe provide more comprehensive materials on mapping in their GESIS workshop on geospatial techniques in R.

## 2 Data Manipulation & Visualization

Many packages and functions can be used to plot maps of spatial data. For instance, `ggplot` as a function to plot spatial data using `geom_sf()`. I am personally a fan of `tmap`, which makes many steps easier (but sometimes is less flexible).

A great tool for choosing colour is for instance Colorbrewer. `viridisLite` provides another great resource to chose colours.

### 2.2.1 Tmaps

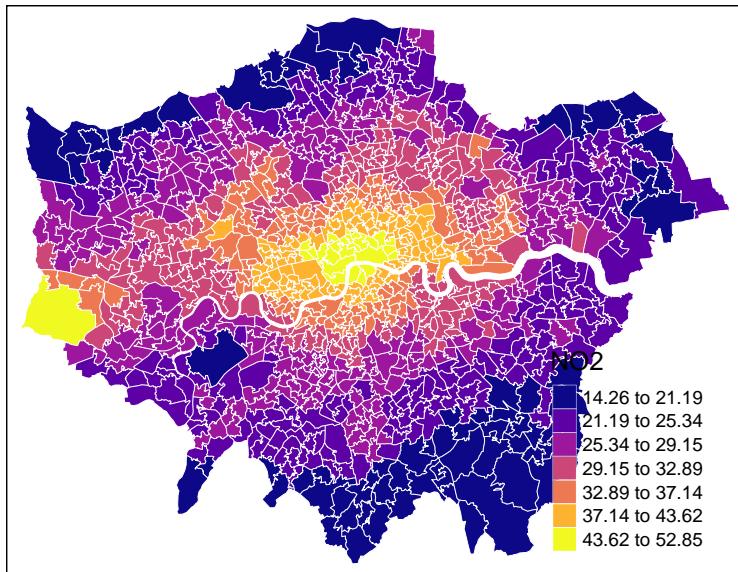
For instance, lets plot the NO2 estimates using `tmap + tm_fill()` (there are lots of alternatives like `tm_shape`, `tm_points()`, `tm_dots()`).

```
# Define colours
cols <- viridis(n = 7, direction = 1, option = "C")

mp1 <- tm_shape(msoa.spdf) +
  tm_fill(col = "no2",
          style = "fisher", # algorithm to def cut points
          n = 7, # Number of requested cut points
          palette = cols, # colours
          alpha = 1, # transparency
          title = "NO2",
          legend.hist = FALSE # histogram next to map?
  ) +
  tm_borders(col = "white", lwd = 0.5, alpha = 0.5)

mp1
```

## 2.2 Visualization



Tmap allows to easily combine different objects by defining a new object via `tm_shape()`.

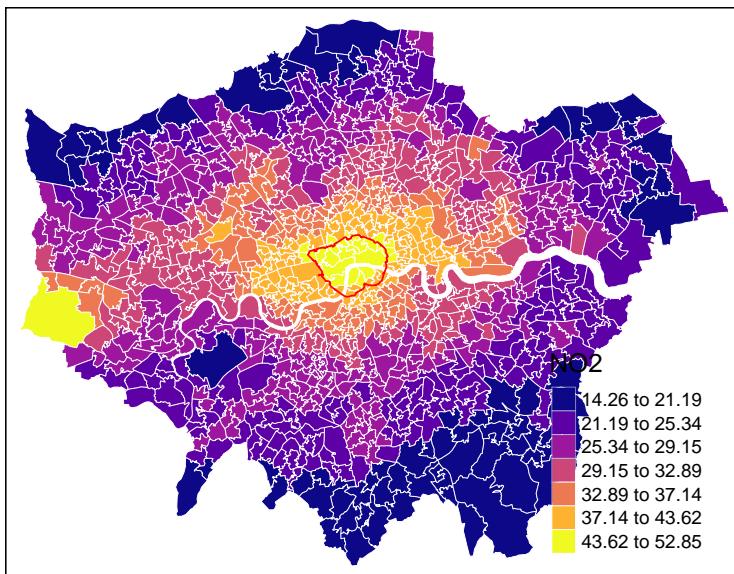
```
# Define colours
cols <- viridis(n = 7, direction = 1, option = "C")

mp1 <- tm_shape(msoa.spdf) +
  tm_fill(col = "no2",
          style = "fisher", # algorithm to def cut points
          n = 7, # Number of requested cut points
          palette = cols, # colours
          alpha = 1, # transparency
          title = "NO2",
          legend.hist = FALSE # histogram next to map?
  ) +
```

## 2 Data Manipulation & Visualization

```
tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +  
tm_shape(ulez.spdf) +  
tm_borders(col = "red", lwd = 1, alpha = 1)
```

```
mp1
```



And it is easy to change the layout.

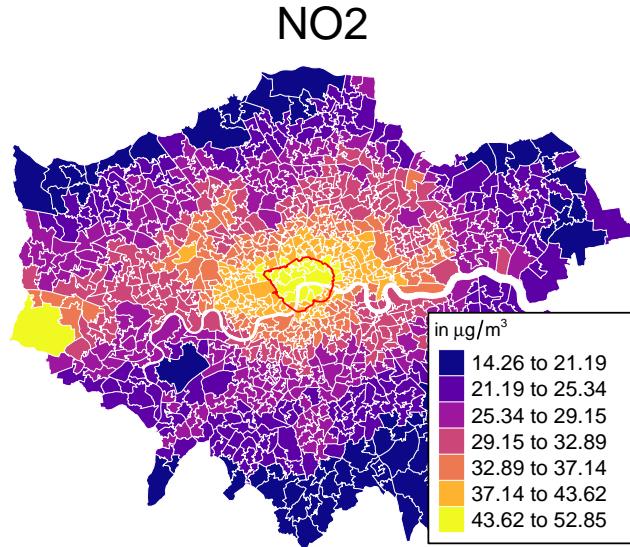
```
# Define colours  
cols <- viridis(n = 7, direction = 1, option = "C")  
  
mp1 <- tm_shape(msoa.spdf) +  
  tm_fill(col = "no2",  
          style = "fisher", # algorithm to def cut points
```

## 2.2 Visualization

```
n = 7, # Number of requested cut points
palette = cols, # colours
alpha = 1, # transparency
title = expression('in'~mu*'g'/m^{3}),
legend.hist = FALSE # histogram next to map?
) +
tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +
tm_shape(ulez.spdf) +
tm_borders(col = "red", lwd = 1, alpha = 1) +
tm_layout(frame = FALSE,
legend.frame = TRUE, legend.bg.color = TRUE,
legend.position = c("right", "bottom"),
legend.outside = FALSE,
main.title = "NO2",
main.title.position = "center",
main.title.size = 1.6,
legend.title.size = 0.8,
legend.text.size = 0.8)
```

mp1

## 2 Data Manipulation & Visualization



We can also add some map information from OSM. However, it's sometimes a bit tricky with the projection. That's why we switch into the OSM projection here. Note that this osm query is build on retiring packages.

```
# Save old projection
crs_orig <- st_crs(msoa.spdf)

# Change projection
ulez.spdf <- st_transform(ulez.spdf, 4326)
msoa.spdf <- st_transform(msoa.spdf, 4326)

# Get OSM data for background
osm_tmp <- read_osm(st_bbox(msoa.spdf), ext = 1.1, type = "osm-german")

# Define colours
```

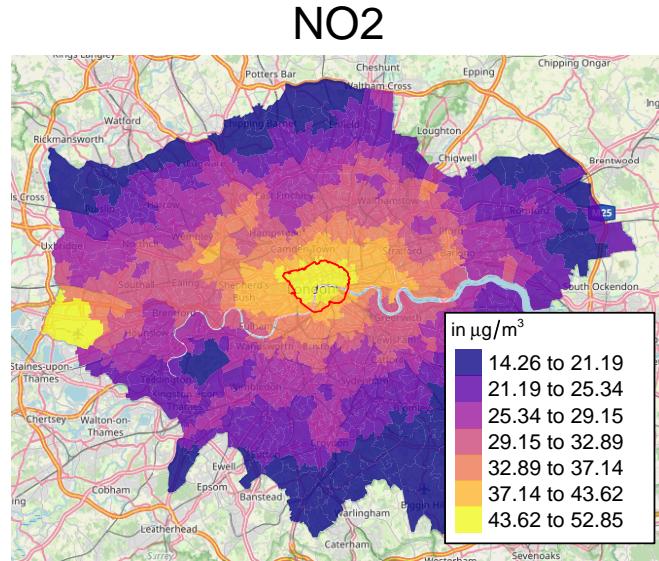
## 2.2 Visualization

```
cols <- viridis(n = 7, direction = 1, option = "C")

mp1 <- tm_shape(osm_tmp) + tm_rgb() +
  tm_shape(msoa.spdf) +
  tm_fill(col = "no2",
    style = "fisher", # algorithm to def cut points
    n = 7, # Number of requested cut points
    palette = cols, # colours
    alpha = 0.8, # transparency
    title = expression('in'~mu*'g'/m^{3}), # histogram next to map?
    legend.hist = FALSE # histogram next to map?
  ) +
  #tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +
  tm_shape(ulez.spdf) +
  tm_borders(col = "red", lwd = 1, alpha = 1) +
  tm_layout(frame = FALSE,
    legend.frame = TRUE, legend.bg.color = TRUE,
    legend.position = c("right", "bottom"),
    legend.outside = FALSE,
    main.title = "NO2",
    main.title.position = "center",
    main.title.size = 1.6,
    legend.title.size = 0.8,
    legend.text.size = 0.8)

mp1
```

## 2 Data Manipulation & Visualization



Tmap also makes it easy to combine single maps

```
# Define colours
cols1 <- viridis(n = 7, direction = 1, option = "C")

# Define colours
cols2 <- viridis(n = 7, direction = 1, option = "D")

mp1 <- tm_shape(osm_tmp) + tm_rgb() +
  tm_shape(msoa.spdf) +
  tm_fill(col = "no2",
          style = "fisher", # algorithm to def cut points
          n = 7, # Number of requested cut points
          palette = cols1, # colours
          alpha = 0.8, # transparency
```

## 2.2 Visualization

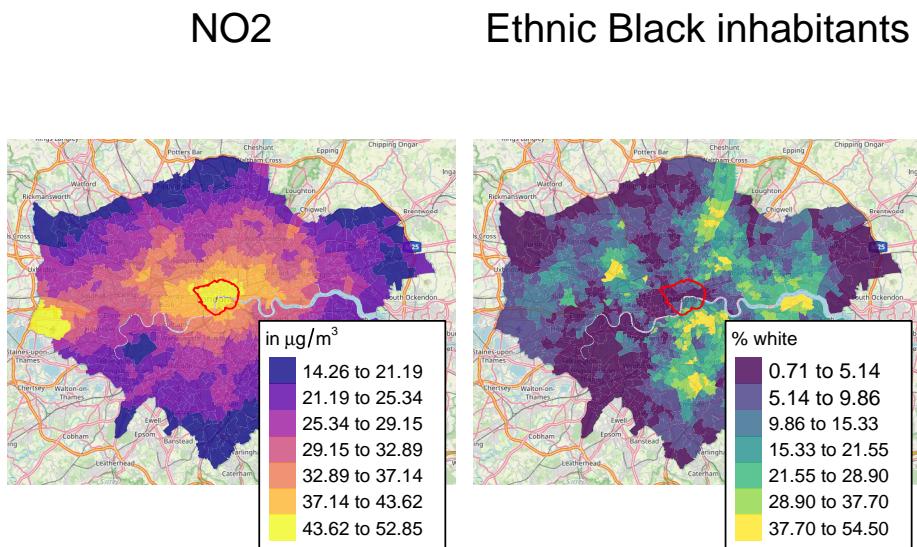
```
title = expression('in'~mu*'g'/m^{3}),  
legend.hist = FALSE # histogram next to map?  
) +  
#tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +  
tm_shape(ulez.spdf) +  
tm_borders(col = "red", lwd = 1, alpha = 1) +  
tm_layout(frame = FALSE,  
          legend.frame = TRUE, legend.bg.color = TRUE,  
          legend.position = c("right", "bottom"),  
          legend.outside = FALSE,  
          main.title = "NO2",  
          main.title.position = "center",  
          main.title.size = 1.4,  
          legend.title.size = 0.8,  
          legend.text.size = 0.8)  
  
mp2 <- tm_shape(osm_tmp) + tm_rgb() +  
  tm_shape(msoa.spdf) +  
  tm_fill(col = "per_black",  
          style = "fisher", # algorithm to def cut points  
          n = 7, # Number of requested cut points  
          palette = cols2, # colours  
          alpha = 0.8, # transparency  
          title = "% white",  
          legend.hist = FALSE # histogram next to map?  
) +  
#tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +  
tm_shape(ulez.spdf) +  
tm_borders(col = "red", lwd = 1, alpha = 1) +  
tm_layout(frame = FALSE,  
          legend.frame = TRUE, legend.bg.color = TRUE,  
          legend.position = c("right", "bottom"),
```

## 2 Data Manipulation & Visualization

```
    legend.outside = FALSE,  
    main.title = "Ethnic Black inhabitants",  
    main.title.position = "center",  
    main.title.size = 1.4,  
    legend.title.size = 0.8,  
    legend.text.size = 0.8)  
  
tmap_arrange(mp1, mp2, ncol = 2, nrow = 1)
```

Legend labels were too wide. The labels have been resized to 0.63, 0.63, 0.63

Some legend labels were too wide. These labels have been resized to 0.68, 0.68



And you can easily export those to png or pdf

## 2.2 Visualization

```
png(file = paste("London.png", sep = ""), width = 14, height = 7, units = "in",
     res = 100, bg = "white")
par(mar=c(0,0,3,0))
par(mfrow=c(1,1),oma=c(0,0,0,0))
tmap_arrange(mp1, mp2, ncol = 2, nrow = 1)
dev.off()
```

```
pdf  
2
```

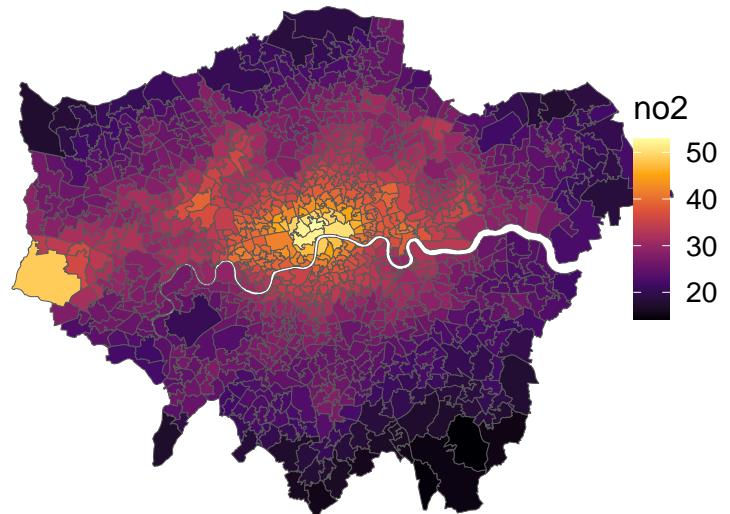
### 2.2.2 ggplot

```
gp <- ggplot(msoa.spdf)+  
  geom_sf(aes(fill = no2))+  
  scale_fill_viridis_c(option = "B") +  
  coord_sf(datum = NA) +  
  theme_map() +  
  theme(legend.position = c(.9, .6))
```

```
Warning: A numeric `legend.position` argument in `theme()` was deprecated in ggplot2  
3.5.0.  
i Please use the `legend.position.inside` argument of `theme()` instead.
```

```
gp
```

## 2 Data Manipulation & Visualization



```
# Get some larger scale boundaries
borough.spdf <- st_read(dsn = paste0("_data", "/statistical-gis-boundaries",
                                         layer = "London_Borough_Excluding_MHW" # Note: no fil
                                         ))
```

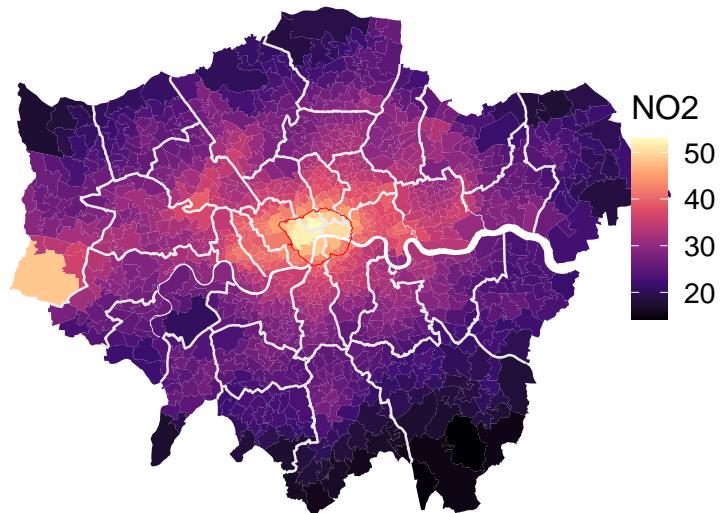
```
Reading layer `London_Borough_Excluding_MHW' from data source
  `C:\work\Lehre\Geodata_Spatial_Regression\_data\statistical-
gis-boundaries-london\ESRI'
  using driver `ESRI Shapefile'
Simple feature collection with 33 features and 7 fields
Geometry type: MULTIPOLYGON
Dimension:     XY
Bounding box:  xmin: 503568.2 ymin: 155850.8 xmax: 561957.5 ymax: 200933.9
Projected CRS: OSGB36 / British National Grid
```

## 2.2 Visualization

```
# transform to only inner lines
borough_inner <- ms_innerlines(borough.spdf)

# Plot with inner lines
gp <- ggplot(msoa.spdf) +
  geom_sf(aes(fill = no2), color = NA) +
  scale_fill_viridis_c(option = "A") +
  geom_sf(data = borough_inner, color = "gray92") +
  geom_sf(data = ulez.spdf, color = "red", fill = NA) +
  coord_sf(datum = NA) +
  theme_map() +
  labs(fill = "NO2") +
  theme(legend.position = c(.9, .6))

gp
```



## 2.3 Exercises

- 1) What is the difference between a spatial “sf” object and a conventional “data.frame”? What’s the purpose of the function `st_drop_geometry()`?

It’s the same. A spatial “sf” object just has an additional column containing the spatial coordinates.

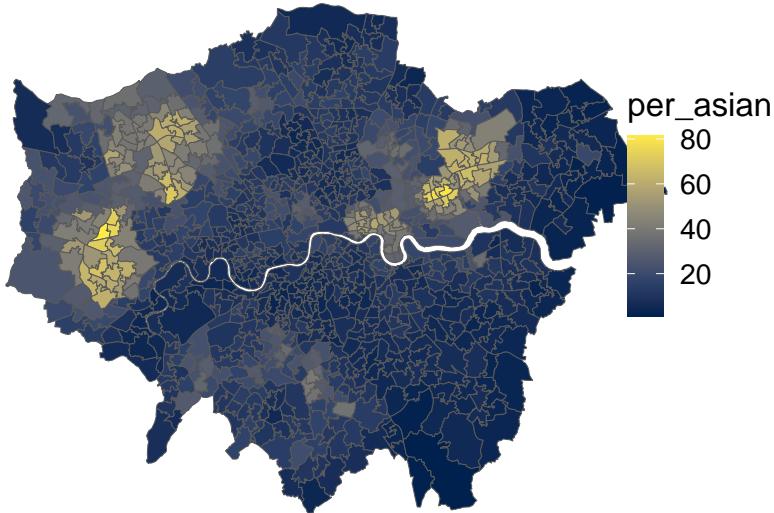
- 2) Using `msoa.spdf`, please create a spatial data frame that contains only the MSOA areas that are within the ulez zone.

```
sub4.spdf <- msoa.spdf[ulez.spdf, ]
```

- 3) Please create a map for London (or only the msoa-ulez subset) which shows the share of Asian residents (or any other ethnic group).

```
gp <- ggplot(msoa.spdf)+  
  geom_sf(aes(fill = per_asian))+  
  scale_fill_viridis_c(option = "E") +  
  coord_sf(datum = NA) +  
  theme_map() +  
  theme(legend.position = c(.9, .6))  
gp
```

### 2.3 Exercises



- 4) Please calculate the distance of each MSOA to the London city centre
- use google maps to get lon and lat,
  - use `st_as_sf()` to create the spatial point
  - use `st_distance()` to calculate the distance

```
### Distance to city center
# Define centre
centre <- st_as_sf(data.frame(lon = -0.128120855701165,
                                lat = 51.50725909644806),
                      coords = c("lon", "lat"),
                      crs = 4326)

# Reproject
centre <- st_transform(centre, crs = st_crs(msoa.spdf))
# Calculate distance
msoa.spdf$dist_centre <- as.numeric(st_distance(msoa.spdf, centre)) / 1000
```

## 2 Data Manipulation & Visualization

```
# hist(msoa.spdf$dist_centre)
```

- 5) Can you create a plot with the distance to the city centre and pub counts next to each other?

```
# Define colours
cols <- viridis(n = 10, direction = 1, option = "B")
cols2 <- viridis(n = 10, direction = 1, option = "E")

mp1 <- tm_shape(msoa.spdf) +
  tm_fill(col = "dist_centre",
          style = "fisher", # algorithm to def cut points
          n = 10, # Number of requested cut points
          palette = cols, # colours
          alpha = 1, # transparency
          title = "Distance",
          legend.hist = FALSE # histogram next to map?
  ) +
  tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +
  tm_layout(frame = FALSE,
            legend.frame = TRUE, legend.bg.color = TRUE,
            legend.position = c("right", "bottom"),
            legend.outside = FALSE,
            main.title = "Dist centre",
            main.title.position = "center",
            main.title.size = 1.6,
            legend.title.size = 0.8,
            legend.text.size = 0.8)

mp2 <- tm_shape(msoa.spdf) +
```

## 2.3 Exercises

```
tm_fill(col = "dist_centre",
        style = "quantile", # algorithm to def cut points
        n = 10, # Number of requested cut points
        palette = cols, # colours
        alpha = 1, # transparency
        title = "Distance",
        legend.hist = FALSE # histogram next to map?
      ) +
tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +
tm_layout(frame = FALSE,
          legend.frame = TRUE, legend.bg.color = TRUE,
          legend.position = c("right", "bottom"),
          legend.outside = FALSE,
          main.title = "Dist centre",
          main.title.position = "center",
          main.title.size = 1.6,
          legend.title.size = 0.8,
          legend.text.size = 0.8)

mp3 <- tm_shape(msoa.spdf) +
  tm_fill(col = "pubs_count",
          style = "fisher", # algorithm to def cut points
          n = 10, # Number of requested cut points
          palette = cols, # colours
          alpha = 1, # transparency
          title = "Count",
          legend.hist = FALSE # histogram next to map?
      ) +
  tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +
  tm_layout(frame = FALSE,
            legend.frame = TRUE, legend.bg.color = TRUE,
```

## 2 Data Manipulation & Visualization

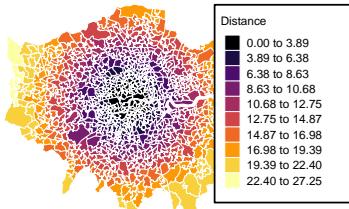
```
    legend.position = c("right", "bottom"),
    legend.outside = FALSE,
    main.title = "Pubs",
    main.title.position = "center",
    main.title.size = 1.6,
    legend.title.size = 0.8,
    legend.text.size = 0.8)

mp4 <- tm_shape(msoa.spdf) +
  tm_fill(col = "pubs_count",
          style = "quantile", # algorithm to def cut points
          n = 10, # Number of requested cut points
          palette = cols, # colours
          alpha = 1, # transparency
          title = "Count",
          legend.hist = FALSE # histogram next to map?
  ) +
  tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +
  tm_layout(frame = FALSE,
            legend.frame = TRUE, legend.bg.color = TRUE,
            legend.position = c("right", "bottom"),
            legend.outside = FALSE,
            main.title = "Pubs",
            main.title.position = "center",
            main.title.size = 1.6,
            legend.title.size = 0.8,
            legend.text.size = 0.8)

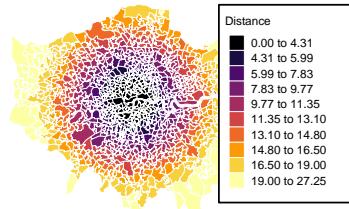
tmap_arrange(mp1, mp2, mp3, mp4, ncol = 2, nrow = 2)
```

### 2.3 Exercises

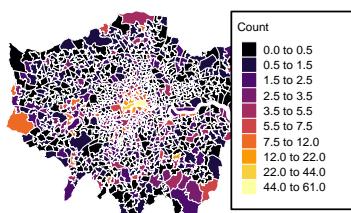
Dist centre



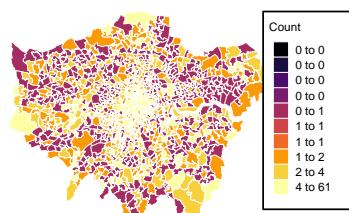
Dist centre



Pubs



Pubs





# 3 Spatial Relationships W

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite") # note: load spd  
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()
```

```
R version 4.4.1 (2024-06-14 ucrt)  
Platform: x86_64-w64-mingw32/x64  
Running under: Windows 11 x64 (build 22631)
```

```
Matrix products: default
```

```
locale:  
[1] LC_COLLATE=English_United Kingdom.utf8  
[2] LC_CTYPE=English_United Kingdom.utf8  
[3] LC_MONETARY=English_United Kingdom.utf8  
[4] LC_NUMERIC=C  
[5] LC_TIME=English_United Kingdom.utf8
```

### 3 Spatial Relationships W

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets   methods    base

other attached packages:
[1] viridisLite_0.4.2 tmap_3.3-4           spatialreg_1.3-
4 Matrix_1.7-0
[5] spdep_1.3-5       spData_2.3.1        mapview_2.11.2     sf_1.0-
16

loaded via a namespace (and not attached):
[1] xfun_0.45          raster_3.6-26        htmlwidgets_1.6.4  lattice_0.22-
6
[5] tools_4.4.1         crosstalk_1.2.1      LearnBayes_2.15.1 parallel_4.4.1
[9] stats4_4.4.1        sandwich_3.1-0       proxy_0.4-27      KernSmooth_2.2
24
[13] satellite_1.0.5    RColorBrewer_1.1-3  leaflet_2.2.2    lifecycle_1.0.4
[17] compiler_4.4.1     deldir_2.0-4       munsell_0.5.1    terra_1.7-
78
[21] codetools_0.2-20   leafsync_0.1.0      stars_0.6-5      htmltools_0.5.8
[25] class_7.3-22       MASS_7.3-60.2      classInt_0.4-
10 lwgeom_0.2-14
[29] wk_0.9.1          abind_1.4-5       boot_1.3-30     multcomp_1.4-
25
[33] nlme_3.1-164       digest_0.6.35     mvtnorm_1.2-
5 splines_4.4.1
[37] fastmap_1.2.0      grid_4.4.1       colorspace_2.1-
0 cli_3.6.2
[41] magrittr_2.0.3     base64enc_0.1-3   dichromat_2.0-
0.1 XML_3.99-0.16.1
[45] survival_3.6-4     leafem_0.2.3      TH.data_1.1-
2 e1071_1.7-14
```

### 3.1 Spatial interdependence

```
[49] scales_1.3.0      sp_2.1-4        rmarkdown_2.27    zoo_1.8-  
12  
[53] png_0.1-8       coda_0.19-4.1   evaluate_0.24.0  knitr_1.47  
[57] tmaptools_3.1-1 s2_1.1.6       rlang_1.1.4     Rcpp_1.0.12  
[61] glue_1.7.0      DBI_1.2.3      rstudioapi_0.16.0 jsonlite_1.8.8  
[65] R6_2.5.1       units_0.8-5
```

#### Reload data from previous session

```
load("_data/msoa2_spatial.RData")
```

## 3.1 Spatial interdependence

We can not only use coordinates and geo-spatial information to connect different data sources, we can also explicitly model spatial (inter)dependence in the analysis of our data. In many instances, accounting for spatial dependence might even be necessary to avoid biased point estimates and standard errors, as observations are often not independent and identically distributed.

Tobler's first law of geography has been used extensively (11,584 citations in 2023-06) to describe spatial dependence: 'Everything is related to everything else, but near things are more related than distant things' (Tobler 1970).

#### Note

Tobler's first law is a bit of a story  
And it has been labeled as an excuse to not think too much about the reasons for spatial dependence or auto-correlation. For instance, measurement error, omitted variables, or inappropriate levels of aggregation.

### 3 Spatial Relationships $W$

gation are among reasons for auto-correlation (Pebesma and Bivand 2023).

We will come back to the reasons of spatial dependence. However, for now, we are interested in some tools to detect and analyse spatial relations.

To analyse spatial relations, we first need to define some sort of connectivity between units (e.g. similar to network analysis). There are some obvious candidates that can be used to define these relations here: adjacency and proximity.

## 3.2 $W$ : Connectivity between units

The connectivity between units is usually represented in a matrix  $\mathbf{W}$ . There is an ongoing debate about the importance of spatial weights for spatial econometrics and about the right way to specify weights matrices (LeSage and Pace 2014; Neumayer and Plümper 2016). The following graph shows some possible options in how to define connectivity between units.

In spatial econometrics, the spatial connectivity (as shown above) is usually represented by a spatial weights matrix  $\mathbf{W}$ :

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \dots & w_{nn} \end{bmatrix}$$

The spatial weights matrix  $\mathbf{W}$  is an  $N \times N$  dimensional matrix with elements  $w_{ij}$  specifying the relation or connectivity between each pair of units  $i$  and  $j$ .

### 3.2 $\mathbf{W}$ : Connectivity between units



Figure 3.1: Figure: Different measures of connectivity, Source: R. S. Bivand and Rudel (2018)

### 3 Spatial Relationships $W$

Note: The diagonal elements  $w_{i,i} = w_{1,1}, w_{2,2}, \dots, w_{n,n}$  of  $\mathbf{W}$  are always zero. No unit is a neighbour of itself. This is not true for spatial multiplier matrices (as we will see later).

#### 3.2.1 Contiguity weights

A very common type of spatial weights. Binary specification, taking the value 1 for neighbouring units (queens: sharing a common edge; rook: sharing a common border), and 0 otherwise.

Contiguity weights  $w_{i,j}$ , where

$$w_{i,j} = \begin{cases} 1 & \text{if } i \text{ and } j \text{ neighbours} \\ 0 & \text{otherwise} \end{cases}$$

A contiguity weights matrix with three units, where unit 1 and unit 3 are neighbours, while unit 2 has no neighbours would look like this:

$$\mathbf{W} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

- Sparse matrices
- Problem of ‘island’: units without neighbours (if I calculate an average of their neighbours, would that be zero, or NA, or a mean?)

Lets create a contiguity weights matrix (Queens neighbours) for the London MSOAs: we create a neighbours list (`nb`) using `poly2nb()`, which is an efficient way of storing  $\mathbf{W}$ . A `snap` of 1 meter accounts for potential lacks of accuracy between lines and points.

### 3.2 W: Connectivity between units

```
# Contiguity (Queens) neighbours weights
queens.nb <- poly2nb(msoa.spdf,
                      queen = TRUE, # a single shared boundary point meets the contiguity
                      snap = 1) # we consider points in 1m distance as 'touching'
summary(queens.nb)
```

Neighbour list object:

Number of regions: 983

Number of nonzero links: 5648

Percentage nonzero weights: 0.5845042

Average number of links: 5.745677

Link number distribution:

2	3	4	5	6	7	8	9	10	11	12	13
9	39	130	264	273	169	66	19	5	6	2	1

9 least connected regions:

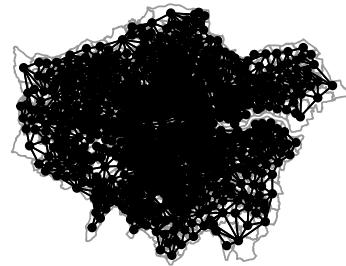
160 270 475 490 597 729 755 778 861 with 2 links

1 most connected region:

946 with 13 links

```
# Lets plot that
plot(st_geometry(msoa.spdf), border = "grey60")
plot(queens.nb, st_centroid(st_geometry(msoa.spdf)),
     add = TRUE, pch = 19, cex = 0.6)
```

### 3 Spatial Relationships W



```
# We can also transform this into a matrix W
W <- nb2mat(queens.nb, style = "B")
print(W[1:10, 1:10])
```

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]	[,10]
1	0	0	0	0	0	0	0	0	0	0
2	0	0	1	0	0	0	0	0	0	0
3	0	1	0	0	1	0	0	0	0	0
4	0	0	0	0	0	1	0	0	0	1
5	0	0	1	0	0	1	1	0	0	0
6	0	0	0	1	1	0	1	0	1	1
7	0	0	0	0	1	1	0	1	1	0
8	0	0	0	0	0	0	1	0	0	0
9	0	0	0	0	0	1	1	0	0	1
10	0	0	0	1	0	1	0	0	1	0

### 3.2 W: Connectivity between units

#### 💡 Question

Among those first 10 units that you see above, which are the neighbours of unit number 6?

Why is the diagonal of this matrix all zero?

Overall, the matrix W has dimensions  $N \times N$ , a row and a column for each observation. The value in a cell shows how units  $i$  (row number) and  $j$  (column number) are related to each other.

```
dim(W)
```

```
[1] 983 983
```

The row and column sums indicate the number of neighbours of each observation.

```
rowSums(W)[1:10]
```

```
1 2 3 4 5 6 7 8 9 10  
11 6 7 5 5 6 6 6 6 5
```

```
colSums(W)[1:10]
```

```
[1] 11 6 7 5 5 6 6 6 6 5
```

Adjacency or graph-based neighbour's weights matrices are usually symmetric. If unit 1 is a neighbour of unit 55, then unit 55 is also a neighbour of unit 1.

### 3 Spatial Relationships W

#### 💡 Higher Order Neighbours

Your neighbours have neighbours too, and they are called higher (second) order neighbours. The neighbours of your neighbour's neighbours are third order neighbours.

You can use `nblag()` to calculate higher order neighbour relations.

#### 3.2.2 Distance based weights

Another common type uses the distance  $d_{ij}$  between each unit  $i$  and  $j$ .

- Inverse distance weights  $w_{i,j} = \frac{1}{d_{ij}^\alpha}$ , where  $\alpha$  define the strength of the spatial decay.

$$\mathbf{W} = \begin{bmatrix} 0 & \frac{1}{d_{ij}^\alpha} & \frac{1}{d_{ij}^\alpha} \\ \frac{1}{d_{ij}^\alpha} & 0 & \frac{1}{d_{ij}^\alpha} \\ \frac{1}{d_{ij}^\alpha} & \frac{1}{d_{ij}^\alpha} & 0 \end{bmatrix}$$

- Dense matrices
- Specifying thresholds may be useful (to get rid of very small non-zero weights)

For now, we will just specify a neighbours list with a distance threshold of 3km using `dnearneigh()`. An alternative would be k nearest neighbours using `knearneigh()`. We will do the inverse weighting later.

```
# Create centroids
coords <- st_geometry(st_centroid(msoa.spdf))
```

Warning: `st_centroid` assumes attributes are constant over geometries

### 3.2 W: Connectivity between units

```
# Neighbours within 3km distance
dist_3.nb <- dnearneigh(coords, d1 = 0, d2 = 3000)
summary(dist_3.nb)
```

Neighbour list object:

Number of regions: 983

Number of nonzero links: 22086

Percentage nonzero weights: 2.285652

Average number of links: 22.46796

2 disjoint connected subgraphs

Link number distribution:

```
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26
4 3 7 13 11 14 14 17 26 22 26 30 33 34 46 34 59 43 38 30 25 19 22 15 21 14
27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47
23 17 17 23 28 19 26 24 29 24 27 25 22 18 8 10 12 5 3 2 1
```

4 least connected regions:

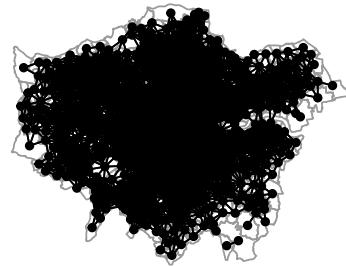
158 160 463 959 with 1 link

1 most connected region:

545 with 47 links

```
# Lets plot that
plot(st_geometry(msoa.spdf), border = "grey60")
plot(dist_3.nb, coords,
      add = TRUE, pch = 19, cex = 0.6)
```

### 3 Spatial Relationships W



And you can see that the matrix is not so sparse anymore:

```
W2 <- nb2mat(dist_3.nb, style = "B")
W2[1:10, 1:10]
```

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]	[,10]
1	0	0	0	0	0	0	0	0	0	0
2	0	0	1	0	1	0	0	0	0	0
3	0	1	0	0	1	1	1	0	0	0
4	0	0	0	0	1	1	1	0	1	1
5	0	1	1	1	0	1	1	1	1	1
6	0	0	1	1	1	0	1	1	1	1
7	0	0	1	1	1	1	0	1	1	1
8	0	0	0	0	1	1	1	0	1	0
9	0	0	0	1	1	1	1	1	0	1

### 3.3 Normalization of $\mathbf{W}$

10    0    0    0    1    1    1    1    0    1    0

## 3.3 Normalization of $\mathbf{W}$

Normalizing ensures that the parameter space of the spatial multiplier is restricted to  $-1 < \rho > 1$ , and the multiplier matrix is non-singular (don't worry, more on this later).

The main message: Normalizing your weights matrix is always a good idea. Otherwise, the spatial parameters might blow up – if you can estimate the model at all. It also ensure easy interpretation of spillover effects.

Again, how to normalize a weights matrix is subject of debate (LeSage and Pace 2014; Neumayer and Plümper 2016).

### 3.3.1 Row-normalization

Row-normalization divides each non-zero weight by the sum of all weights of unit  $i$ , which is the sum of the row.

$$\frac{w_{ij}}{\sum_j^n w_{ij}}$$

- With contiguity weights, spatially lagged variables contain mean of this variable among the neighbours of  $i$
- Proportions between units such as distances get lost (can be bad!)
- Can induce asymmetries:  $w_{ij} \neq w_{ji}$

For instance, we can use row-normalization for the Queens neighbours created above, and create a neighbours list with spatial weights.

### 3 Spatial Relationships W

```
queens.lw <- nb2listw(queens.nb,
                      style = "W") # W ist row-normalization
summary(queens.lw)
```

Characteristics of weights list object:

Neighbour list object:

Number of regions: 983

Number of nonzero links: 5648

Percentage nonzero weights: 0.5845042

Average number of links: 5.745677

Link number distribution:

2	3	4	5	6	7	8	9	10	11	12	13
9	39	130	264	273	169	66	19	5	6	2	1

9 least connected regions:

160 270 475 490 597 729 755 778 861 with 2 links

1 most connected region:

946 with 13 links

Weights style: W

Weights constants summary:

n	nn	S0	S1	S2	
W	983	966289	983	355.1333	4017.47

To see what happened, let's look at our example in matrix format again.

```
# transform into matrix with row-normalization
W_norm <- nb2mat(queens.nb, style = "W")
print(W_norm[1:10, 1:10])
```

[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]
------	------	------	------	------	------	------	------

### 3.3 Normalization of $\mathbf{W}$

```
1 0 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000  
2 0 0.0000000 0.1666667 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000  
3 0 0.1428571 0.0000000 0.0000000 0.1428571 0.0000000 0.0000000 0.0000000 0.0000000  
4 0 0.0000000 0.0000000 0.0000000 0.0000000 0.2000000 0.0000000 0.0000000 0.0000000  
5 0 0.0000000 0.2000000 0.0000000 0.0000000 0.2000000 0.2000000 0.0000000 0.0000000  
6 0 0.0000000 0.0000000 0.1666667 0.1666667 0.0000000 0.1666667 0.0000000 0.0000000  
7 0 0.0000000 0.0000000 0.0000000 0.1666667 0.1666667 0.0000000 0.1666667 0.0000000  
8 0 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000 0.1666667 0.0000000 0.0000000  
9 0 0.0000000 0.0000000 0.0000000 0.0000000 0.1666667 0.1666667 0.0000000 0.0000000  
10 0 0.0000000 0.0000000 0.2000000 0.0000000 0.2000000 0.0000000 0.0000000 0.0000000  
[,9] [,10]  
1 0.0000000 0.0000000  
2 0.0000000 0.0000000  
3 0.0000000 0.0000000  
4 0.0000000 0.2000000  
5 0.0000000 0.0000000  
6 0.1666667 0.1666667  
7 0.1666667 0.0000000  
8 0.0000000 0.0000000  
9 0.0000000 0.1666667  
10 0.2000000 0.0000000
```

#### 💡 Question

Overall, how many neighbours does unit 9 have (including all columns)? How do you know?

```
rowSums(W) [9]
```

We can also use the `nb` object to see which ones the neighbours are. Here, for instance, neighbours of unit 6:

### 3 Spatial Relationships W

```
queens.nb[6]
```

```
[[1]]  
[1] 4 5 7 9 10 462
```

This fits to what we see in the matrix above.

#### ⚠ Warning

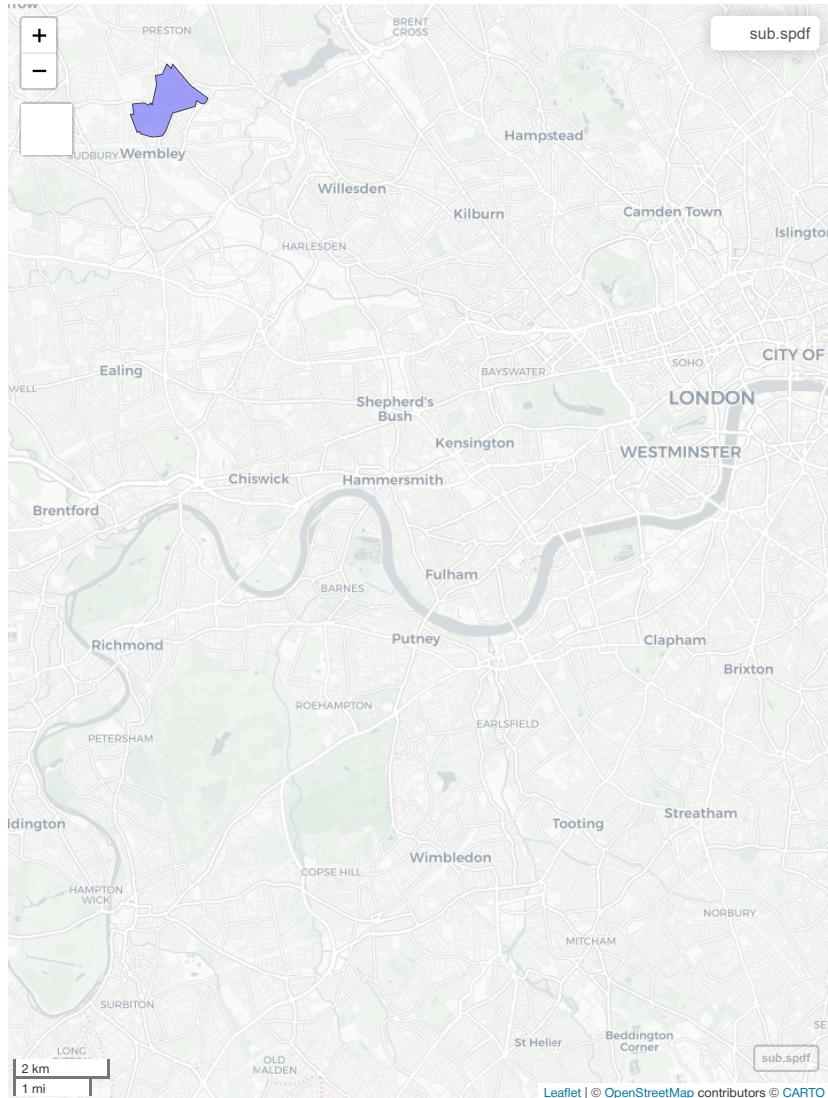
Note that row-normalization has some undesirable properties when we use some non-contigutiy based neighbour relations, such as distance based neighbours.

The problem: It obscures the proportion due to dividing by a row-specific value.

Let's construct a hypothetical example

```
# Subset of 5 units  
sub.spdf <- msoa.spdf[c(4, 5, 6, 102, 150), ]  
mapview(sub.spdf)
```

### 3.3 Normalization of $\mathbf{W}$



### 3 Spatial Relationships $W$

We construct the **inverse-distance weighted 2 nearest neighbors**.

```
# 2 closest neighbours  
sub.coords <- st_geometry(st_centroid(sub.spdf))
```

```
Warning: st_centroid assumes attributes are constant over geometries
```

```
knn.nb <- knearneigh(sub.coords,  
                      k = 2) # number of nearest neighbours
```

```
Warning in knearneigh(sub.coords, k = 2): k greater than one-  
third of the  
number of data points
```

```
knn.nb <- knn2nb(knn.nb)  
summary(knn.nb)
```

```
Neighbour list object:  
Number of regions: 5  
Number of nonzero links: 10  
Percentage nonzero weights: 40  
Average number of links: 2  
Non-symmetric neighbours list  
Link number distribution:  
  
2  
5  
5 least connected regions:  
1 2 3 4 5 with 2 links  
5 most connected regions:  
1 2 3 4 5 with 2 links
```

### 3.3 Normalization of $\mathbf{W}$

```
# listw with inverse-distance based weights
sub.lw <- nb2listwdist(knn.nb,
                        x = sub.coords, # needed for idw
                        type = "idw", # inverse distance weighting
                        alpha = 1, # the decay parameter for distance weighting
                        style = "raw") # without normalization
W_sub <- listw2mat(sub.lw)
formatC(W_sub, format = "f", digits = 6)
```

```
[,1]      [,2]      [,3]      [,4]      [,5]
1 "0.000000" "0.000414" "0.000723" "0.000000" "0.000000"
2 "0.000414" "0.000000" "0.000962" "0.000000" "0.000000"
3 "0.000723" "0.000962" "0.000000" "0.000000" "0.000000"
4 "0.000000" "0.000033" "0.000032" "0.000000" "0.000000"
5 "0.000049" "0.000000" "0.000049" "0.000000" "0.000000"
```

As you can see, units 1, 2, 3 have relatively proximate neighbours (.e.g inverse distance 0.000962: 3 zeros). Units 4 and 5, in contrast, have only very distant neighbours (e.g. inverse distance 0.000049: 4 zeros).

Now, see what happens when we use row-normalization.

```
sub.lw <- nb2listwdist(knn.nb,
                        x = sub.coords, # needed for idw
                        type = "idw", # inverse distance weighting
                        alpha = 1, # the decay parameter for distance weighting
                        style = "W") # for row normalization
W_sub <- listw2mat(sub.lw)
formatC(W_sub, format = "f", digits = 6)
```

```
[,1]      [,2]      [,3]      [,4]      [,5]
1 "0.000000" "0.364083" "0.635917" "0.000000" "0.000000"
```

### 3 Spatial Relationships $\mathbf{W}$

```
2 "0.300879" "0.000000" "0.699121" "0.000000" "0.000000"
3 "0.429123" "0.570877" "0.000000" "0.000000" "0.000000"
4 "0.000000" "0.507955" "0.492045" "0.000000" "0.000000"
5 "0.499360" "0.000000" "0.500640" "0.000000" "0.000000"
```

All rows sum up to 1, but the strength of the relation is now similar for the distant units 4 and 5, and the proximate units 1, 2, 3.

#### 3.3.2 Maximum eigenvalues normalization

Maximum eigenvalues normalization divides each non-zero weight by the overall maximum eigenvalue  $\lambda_{max}$ . Each element of  $\mathbf{W}$  is divided by the same scalar parameter, which preserves the relations.

$$\frac{\mathbf{W}}{\lambda_{max}}$$

- Interpretation may become more complicated
- Keeps proportions of connectivity strengths across units (relevant esp. for distance based  $\mathbf{W}$ )

We use eigenvalue normalization for the inverse distance neighbours. We use `nb2listwdist()` to create weight inverse distance based weights and normalize in one step.

```
coords <- st_geometry(st_centroid(msoa.spdf))
```

Warning: `st_centroid` assumes attributes are constant over geometries

```
idw.lw <- nb2listwdist(dist_3.nb,
                         x = coords, # needed for idw
```

### 3.3 Normalization of $\mathbf{W}$

```
type = "idw", # inverse distance weighting
alpha = 1, # the decay parameter for distance weighting
style = "minmax") # for eigenvalue normalization
summary(idw.lw)
```

Characteristics of weights list object:

Neighbour list object:

Number of regions: 983

Number of nonzero links: 22086

Percentage nonzero weights: 2.285652

Average number of links: 22.46796

2 disjoint connected subgraphs

Link number distribution:

```
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26
4 3 7 13 11 14 14 17 26 22 26 30 33 34 46 34 59 43 38 30 25 19 22 15 21 14
27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47
```

```
23 17 17 23 28 19 26 24 29 24 27 25 22 18 8 10 12 5 3 2 1
```

4 least connected regions:

```
158 160 463 959 with 1 link
```

1 most connected region:

```
545 with 47 links
```

Weights style: minmax

Weights constants summary:

n	nn	S0	S1	S2
minmax	983	966289	463.6269	23.92505
				1117.636

Examples from above: See how this keeps the proportions in our example.

Instead of transforming values to sum up to 1 in each row, we now have much smaller values for 4 and 5 then we have for the proximate units 1, 2, 3.

### 3 Spatial Relationships W

```
sub.lw <- nb2listwdist(knn.nb,
                        x = sub.coords, # needed for idw
                        type = "idw", # inverse distance weighting
                        alpha = 1, # the decay parameter for distance weight
                        style = "minmax") # for eigenvalue normalization
W_sub <- listw2mat(sub.lw)
formatC(W_sub, format = "f", digits = 6)

[,1]      [,2]      [,3]      [,4]      [,5]
1 "0.000000" "0.245687" "0.429123" "0.000000" "0.000000"
2 "0.245687" "0.000000" "0.570877" "0.000000" "0.000000"
3 "0.429123" "0.570877" "0.000000" "0.000000" "0.000000"
4 "0.000000" "0.019663" "0.019047" "0.000000" "0.000000"
5 "0.029099" "0.000000" "0.029174" "0.000000" "0.000000"
```

### 3.4 Islands / missings

In practice, we often have a problem with islands. If we use contiguity based or distance based neighbour definitions, some units may end up with empty neighbours sets: they just do not touch any other unit and do not have a neighbour within a specific distance. This however creates a problem: what is the value in the neighbouring units?

The `zero.policy` option in `spdep` allows to proceed with empty neighbours sets. However, many further functions may run into problems and return errors. It often makes sense to either drop islands, to choose weights which always have neighbours (e.g. k nearest), or impute empty neighbours sets by using the nearest neighbours.

# 4 Exercises I

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite",
         "ggplot2", "ggthemes") # note: load spdep first, then spatialreg
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()
```

```
R version 4.4.1 (2024-06-14 ucrt)
Platform: x86_64-w64-mingw32/x64
Running under: Windows 11 x64 (build 22631)
```

```
Matrix products: default
```

```
locale:
[1] LC_COLLATE=English_United Kingdom.utf8
[2] LC_CTYPE=English_United Kingdom.utf8
[3] LC_MONETARY=English_United Kingdom.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United Kingdom.utf8
```

## 4 Exercises I

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets   methods    base

other attached packages:
[1] ggthemes_5.1.0    ggplot2_3.5.1    viridisLite_0.4.2 tmap_3.3-
4
[5] spatialreg_1.3-4 Matrix_1.7-0     spdep_1.3-5      spData_2.3.1
[9] mapview_2.11.2   sf_1.0-16

loaded via a namespace (and not attached):
[1] tidyselect_1.2.1    dplyr_1.1.4       fastmap_1.2.0    leaflet_2.2.2
[5] TH.data_1.1-2      XML_3.99-0.16.1   digest_0.6.35   lifecycle_1.0.4
[9] LearnBayes_2.15.1  survival_3.6-4    terra_1.7-78    magrittr_2.0.3
[13] compiler_4.4.1    rlang_1.1.4       tools_4.4.1     utf8_1.2.4
[17] knitr_1.47        htmlwidgets_1.6.4 sp_2.1-4      classInt_0.4-
10
[21] RColorBrewer_1.1-3 multcomp_1.4-25   abind_1.4-5     KernSmooth_2.2-
24
[25] purrr_1.0.2        withr_3.0.0       leafsync_0.1.0   grid_4.4.1
[29] stats4_4.4.1       fansi_1.0.6      e1071_1.7-14   leafem_0.2.3
[33] colorspace_2.1-0   scales_1.3.0     MASS_7.3-60.2  dichromat_2.0-
0.1
[37] cli_3.6.2          mvtnorm_1.2-5    rmarkdown_2.27  generics_0.1.3
[41] rstudioapi_0.16.0  tmaptools_3.1-1   DBI_1.2.3     proxy_0.4-
27
[45] stringr_1.5.1      splines_4.4.1    stars_0.6-5    parallel_4.4.1
[49] s2_1.1.6            base64enc_0.1-3   vctrs_0.6.5    boot_1.3-
30
[53] sandwich_3.1-0     jsonlite_1.8.8    crosstalk_1.2.1 units_0.8-
5
```

```
[57] glue_1.7.0          lwgeom_0.2-14      codetools_0.2-
20   stringi_1.8.4      deldir_2.0-4       raster_3.6-
[61] gtable_0.3.5        munsell_0.5.1      htmltools_0.5.8.1  satellite_1.0.5
26   tibble_3.2.1        pillar_1.9.0       evaluate_0.24.0    lattice_0.22-
[65] R6_2.5.1            wk_0.9.1          Rcpp_1.0.12        coda_0.19-
[69] png_0.1-8           class_7.3-22      zoo_1.8-12         pkgconfig_2.0.3
4.1
[73] nlme_3.1-164         xfun_0.45
```

## Reload data from previous session

97

```
load("_data/msoa2_spatial.RData")
```

## 4.1 General Exercises

- 4.1.1 1) Can you import the spatial administrative units of Germany (“Kreisgrenzen\_2020\_mit\_Einwohnerzahl” in \_data folder) and make a simple plot of the boundaries?  
{.unnumbered}
- 2) What is the Coordinate reference system of this German shape file?
- 3) Please use the msoa.spdf and calculate a neighbours weights matrix of the nearest 10 neighbours (see `spdep::knearneigh()`), and create a listw object using row normalization.
- 4) OPTIONAL: Can you create a map containing the City of London (MSOA11CD = “E02000001”) and its ten nearest neighbours?
- 5) Please use the msoa.spdf and calculate a neighbours weights matrix of the nearest 10 neighbours (see `spdep::knearneigh()`), and create a listw object using row normalization.
- 6) Please calculate the queens neighbours and make a listw object that includes the second order neighbours (see `nblag()`).
- 7) Generate a matrix from the listw object
- 8) What do you get when you multiply a variable (data column) such as the home owner rate with your weights matrix?

# 5 Detecting Spatial Dependence

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite", "gstat") # note:  
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()
```

```
R version 4.4.1 (2024-06-14 ucrt)  
Platform: x86_64-w64-mingw32/x64  
Running under: Windows 11 x64 (build 22631)
```

```
Matrix products: default
```

```
locale:  
[1] LC_COLLATE=English_United Kingdom.utf8  
[2] LC_CTYPE=English_United Kingdom.utf8  
[3] LC_MONETARY=English_United Kingdom.utf8  
[4] LC_NUMERIC=C  
[5] LC_TIME=English_United Kingdom.utf8
```

## 5 Detecting Spatial Dependence

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets   methods    base

other attached packages:
[1] gstat_2.1-1       viridisLite_0.4.2 tmap_3.3-4        spatialreg_1.3-
4
[5] Matrix_1.7-0      spdep_1.3-5      spData_2.3.1      mapview_2.11.2
[9] sf_1.0-16

loaded via a namespace (and not attached):
[1] xfun_0.45          raster_3.6-26     htmlwidgets_1.6.4  lattice_0.22-
6
[5] tools_4.4.1         crosstalk_1.2.1   LearnBayes_2.15.1  parallel_4.4.1
[9] stats4_4.4.1       sandwich_3.1-0    spacetime_1.3-
1
[13] proxy_0.4-27      xts_0.14.0       KernSmooth_2.23-24 satellite_1.0.5   RColorBrewer_1
3
[17] leaflet_2.2.2     lifecycle_1.0.4   FNN_1.1.4        compiler_4.4.1
[21] deldir_2.0-4      munsell_0.5.1    terra_1.7-78    codetools_0.2-
20
[25] leafsync_0.1.0    stars_0.6-5      htmltools_0.5.8.1 class_7.3-
22
[29] MASS_7.3-60.2     classInt_0.4-10  lwgeom_0.2-
14
[33] abind_1.4-5       boot_1.3-30     multcomp_1.4-
25
[37] nlme_3.1-164      digest_0.6.35   mvtnorm_1.2-5    fastmap_1.2.0
[41] grid_4.4.1        colorspace_2.1-0  cli_3.6.2       magrittr_2.0.3
[45] base64enc_0.1-3   dichromat_2.0-0.1 XML_3.99-0.16.1 survival_3.6-
4
[49] leafem_0.2.3     TH.data_1.1-2    e1071_1.7-14    scales_1.3.0
```

## 5.1 Global Autocorrelation

```
[53] sp_2.1-4           rmarkdown_2.27    zoo_1.8-12      png_0.1-  
8  
[57] coda_0.19-4.1     evaluate_0.24.0   knitr_1.47     tmaptools_3.1-  
1  
[61] s2_1.1.6          rlang_1.1.4       Rcpp_1.0.12    glue_1.7.0  
[65] DBI_1.2.3         rstudioapi_0.16.0 jsonlite_1.8.8 R6_2.5.1  
[69] intervals_0.15.4  units_0.8-5
```

### Reload data from previous session

```
load("_data/msoa2_spatial.RData")
```

## 5.1 Global Autocorrelation

If spatially close observations are more likely to exhibit similar values, we cannot handle observations as if they were independent.

$$E(\varepsilon_i \varepsilon_j) \neq E(\varepsilon_i)E(\varepsilon_j) = 0$$

This violates a basic assumption of the conventional OLS model. We will talk more about whether that is good or bad (any guess?).

### 5.1.1 Visualization

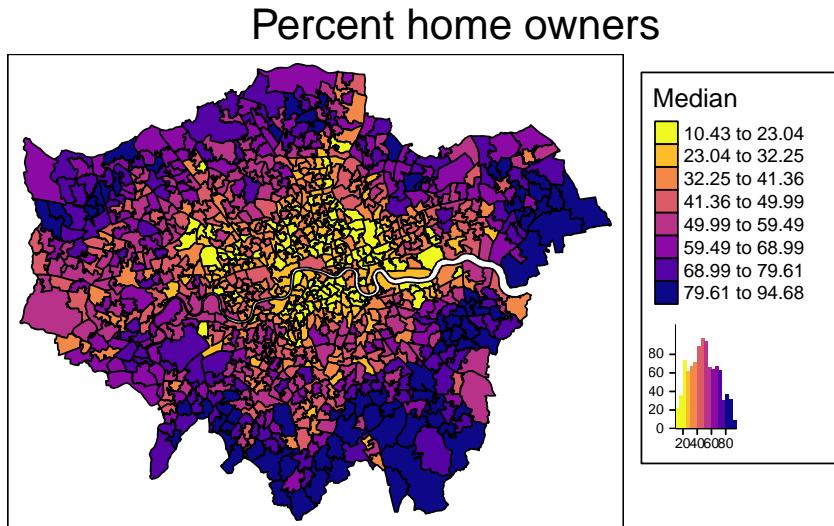
There is one very easy and intuitive way of detecting spatial autocorrelation: Just look at the map. We do so by using `tmap` for plotting the share of home owners.

## 5 Detecting Spatial Dependence

```
mp1 <- tm_shape(msoa.spdf) +
  tm_fill(col = "per_owner",
    #style = "cont",
    style = "fisher", n = 8,
    title = "Median",
    palette = viridis(n = 8, direction = -1, option = "C"),
    legend.hist = TRUE) +
  tm_borders(col = "black", lwd = 1) +
  tm_layout(legend.frame = TRUE, legend.bg.color = TRUE,
    #legend.position = c("right", "bottom"),
    legend.outside = TRUE,
    main.title = "Percent home owners",
    main.title.position = "center",
    title.snap.to.legend = TRUE)
```

mp1

## 5.1 Global Autocorrelation



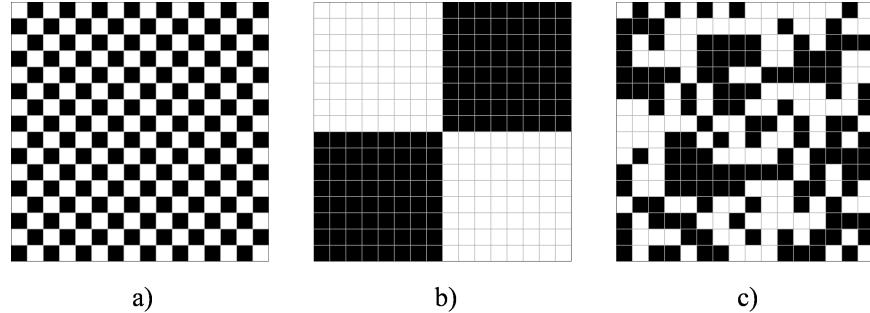
We definitely see some clusters with spatial units having a low share of home owner (e.g. in the city center), and other clusters where home ownership is high (e.g. suburbs in the south and east, such as Bromley or Havering).

However, this is (to some degree) dependent on how we define cutoffs and coloring of the map: the Modifiable Areal Unit Problem (Wong 2009).

### 💡 Question

Which of the following three checkerboards has no (or the lowest) autocorrelation?

## 5 Detecting Spatial Dependence



Would your answer be the same if we would aggregate the data to four larger areas / districts using the average within each of the four districts?

### 5.1.2 Moran's I

The most common and well known statistic for spatial dependence or autocorrelation is Moran's I, which goes back to Moran (1950) and Cliff and Ord (1972). For more extensive materials on Moran's I see for instance Kelejian and Piras (2017), Chapter 11.

To calculate Moran's I, we first define a neighbours weights matrix  $W$ .

Global Moran's I test statistic:

$$I = \frac{N \sum_i \sum_j w_{ij} (y_i - \bar{y})(y_j - \bar{y})}{S_0 \sum_i (y_i - \bar{y})^2}, \text{ where } S_0 = \sum_{i=1}^N \sum_{j=1}^N w_{ij}$$

It is often written with deviations  $z$

$$I = \frac{N \sum_i \sum_j w_{ij} (z_i)(z_j)}{S_0 \sum_i (z_i)^2}, \text{ where } S_0 = \sum_{i=1}^N \sum_{j=1}^N w_{ij}$$

Note that in the case of row-standardized weights,  $S_0 = N$ . The  $I$  can be interpreted as: *Relation of the deviation from the mean value between unit*

## 5.1 Global Autocorrelation

*i* and neighbours of unit *i*. Basically, this measures correlation between neighbouring values.

- Negative values: negative autocorrelation
- Around zero: no autocorrelation
- Positive values: positive autocorrelation

To calculate Moran's I, we first need to define the relationship between units. As in the previous example, we define contiguity weights and distance-based weights.

```
# Contiguity (Queens) neighbours weights
queens.nb <- poly2nb(msoa.spdf,
                      queen = TRUE,
                      snap = 1) # we consider points in 1m distance as 'touching'
queens.lw <- nb2listw(queens.nb,
                      style = "W")

# Neighbours within 3km distance
coords <- st_geometry(st_centroid(msoa.spdf))
```

Warning: st\_centroid assumes attributes are constant over geometries

```
dist_3.nb <- dnearneigh(coords,
                           d1 = 0, d2 = 3000)
idw.lw <- nb2listwdist(dist_3.nb,
                        x = coords, # needed for idw
                        type = "idw", # inverse distance weighting
                        alpha = 1, # the decay parameter for distance weighting
                        style = "minmax") # for eigenvalue normalization
```

## 5 Detecting Spatial Dependence

Subsequently, we can calculate the average correlation between neighbouring units.

For contiguity weights, we get:

```
# Global Morans I test of housing values based on contiguity weights
moran.test(msoa.spdf$per_owner, listw = queens.lw, alternative = "two.sided")
```

Moran I test under randomisation

```
data: msoa.spdf$per_owner
weights: queens.lw

Moran I statistic standard deviate = 38.161, p-value < 2.2e-16
alternative hypothesis: two.sided
sample estimates:
Moran I statistic      Expectation      Variance
0.728706855      -0.001018330      0.000365663
```

And for inverse distance weighting, we get:

```
# Global Morans I test of housing values based on idw
moran.test(msoa.spdf$per_owner, listw = idw.lw, alternative = "two.sided")
```

Moran I test under randomisation

```
data: msoa.spdf$per_owner
weights: idw.lw
```

## 5.1 Global Autocorrelation

```
Moran I statistic standard deviate = 65.853, p-value < 2.2e-
16
alternative hypothesis: two.sided
sample estimates:
Moran I statistic      Expectation      Variance
0.6838957350      -0.0010183299      0.0001081719
```

Interpretation: In both cases, we have very strong autocorrelation between neighbouring/closer units (~.7). It barely matters which of the weights matrices we use. This autocorrelation is highly significant. we can thus reject the Null that units are independent of each other (at least at this spatial level and for the share of home owners).

### 5.1.3 Residual-based Moran's I

We can also use the same Moran's I test to inspect spatial autocorrelation in residuals from an estimated linear model.

Let's start with an intercept only model.

```
lm0 <- lm(per_owner ~ 1, msoa.spdf)
lm.morantest(lm0, listw = queens.lw, alternative = "two.sided")
```

```
Global Moran I for regression residuals

data:
model: lm(formula = per_owner ~ 1, data = msoa.spdf)
weights: queens.lw

Moran I statistic standard deviate = 38.177, p-value < 2.2e-
16
```

## 5 Detecting Spatial Dependence

```
alternative hypothesis: two.sided
sample estimates:
Observed Moran I      Expectation      Variance
0.7287068548     -0.0010183299    0.0003653613
```

This is exactly what we have received in the general case of Moran's I.

Now, lets add some predictors. For instance, the distance to the city centre, and the population density may be strongly related to the home ownership rates and explain parts of the spatial dependence.

```
### Distance to city center
# Define centre
centre <- st_as_sf(data.frame(lon = -0.128120855701165,
                                lat = 51.50725909644806),
                      coords = c("lon", "lat"),
                      crs = 4326)
# Reproject
centre <- st_transform(centre, crs = st_crs(msoa.spdf))
# Calculate distance
msoa.spdf$dist_centre <- as.numeric(st_distance(msoa.spdf, centre)) / 1000
# hist(msoa.spdf$dist_centre)

### Run model with predictors
lm1 <- lm(per_owner ~ dist_centre + POPDEN, msoa.spdf)
lm.morantest(lm1, listw = queens.lw, alternative = "two.sided")
```

Global Moran I for regression residuals

```
data:
model: lm(formula = per_owner ~ dist_centre + POPDEN, data = msoa.spdf)
weights: queens.lw
```

## 5.1 Global Autocorrelation

```
Moran I statistic standard deviate = 22.674, p-value < 2.2e-
16
alternative hypothesis: two.sided
sample estimates:
Observed Moran I      Expectation      Variance
0.4298146060          -0.0024065617     0.0003633607
```

There is still considerable auto-correlation in the residuals. However, we have reduced it by a substantial amount with two very simple control variables.

### 5.1.4 Semivariogram

The sample variogram  $\gamma(h)$  for distance intervals  $h_i$  describes the average square difference between the points in this distance interval:

$$\hat{\gamma}(h_i) = \frac{1}{2N(h_i)} \sum_{j=1}^{N(h_i)} (z(s_i) - z(s_i + h'))^2, \quad h_{i,0} \leq h' < h_{i,1}$$

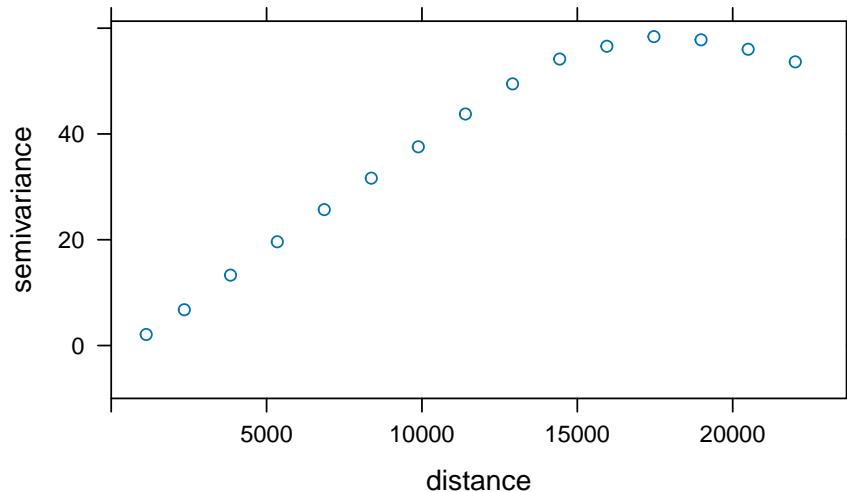
with the number of available pairs  $N(h_i)$  in each distance interval  $h_i$ . Basically, *it is the variance within each distance interval*.

For more information, see for instance the Geospatial Data Science in R by Zia Ahmed or Pebesma and Bivand (2023).

To calculate the empirical semi-vriogram, we can use the package `gstat` with the function `variogram()`.

```
# Variogram No2
v.no2 <- variogram(no2 ~ 1, msoa.spdf)
plot(v.no2, xlim = c(0, 1.075 * max(v.no2$dist)),
      ylim = c(-10, 1.05 * max(v.no2$gamma)))
```

## 5 Detecting Spatial Dependence



Above graphs shows that the variance within each distance interval gradually increases, up to a distance of  $\sim 18\text{km}$ , and then level off at a relative constant level. Lower variances within lower values of distances means that observations are more similar to each other the closer they are.

We can also try to fit a model that resembles the spatial structure. This becomes important when we want to perform spatial interpolation (e.g. to impute missings).

```
# Intial parameter set by eye esitmation
m.no2 <- vgm(60, "Cir", 20000, 0) # Sill, model, range, nugget
# least square fit
m.f.v.no2 <- fit.variogram(v.no2, m.no2)
```

### 5.1 Global Autocorrelation

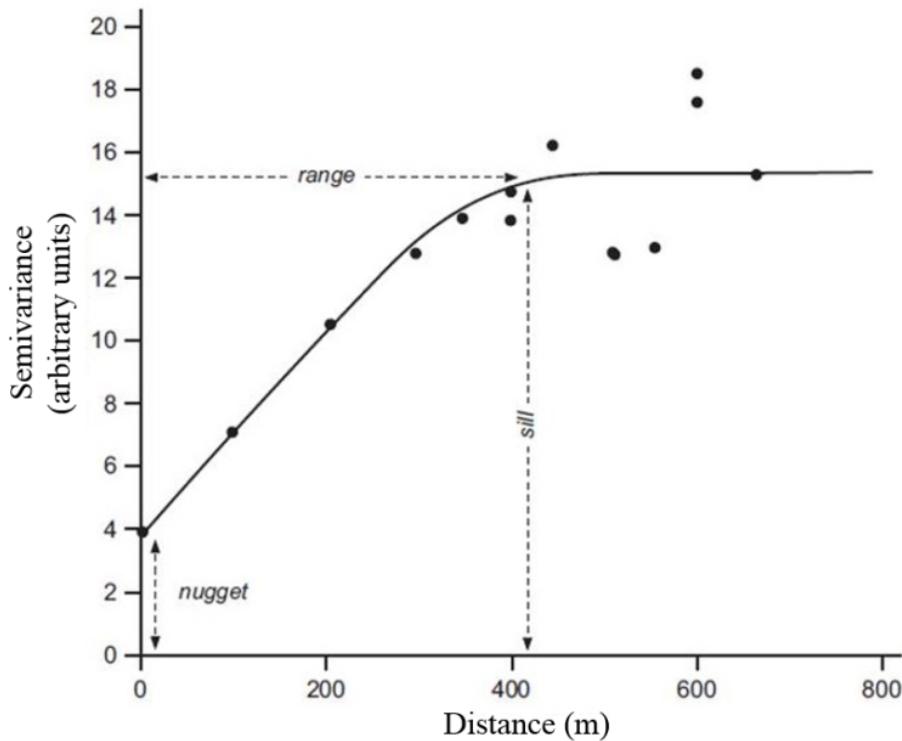
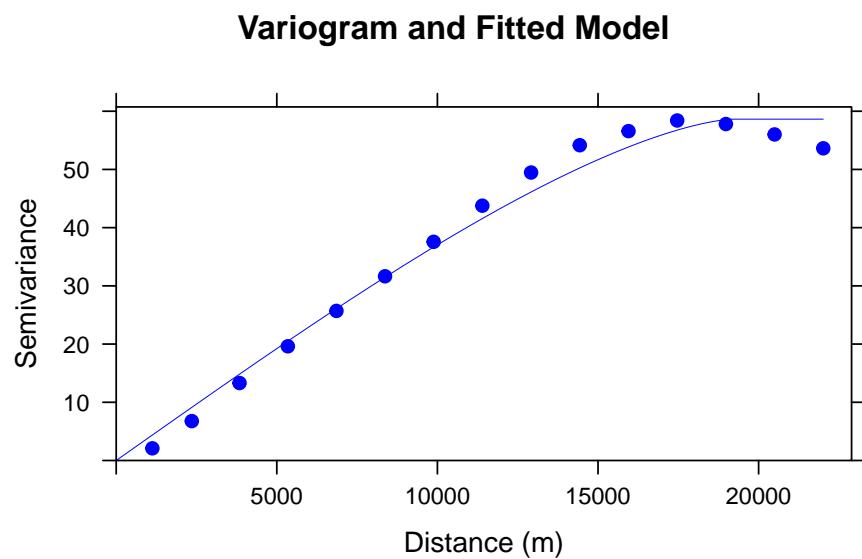


Figure 5.1: Theoretical exponential semi-variogram model. Source: <https://www.aspexit.com/variogram-and-spatial-autocorrelation>

## 5 Detecting Spatial Dependence

```
#### Plot varigram and fitted model:  
plot(v.no2, pl = FALSE,  
      model = m.f.v.no2,  
      col="blue",  
      cex = 0.9,  
      lwd = 0.5,  
      lty = 1,  
      pch = 19,  
      main = "Variogram and Fitted Model",  
      xlab = "Distance (m)",  
      ylab = "Semivariance")
```



### 5.1 Global Autocorrelation

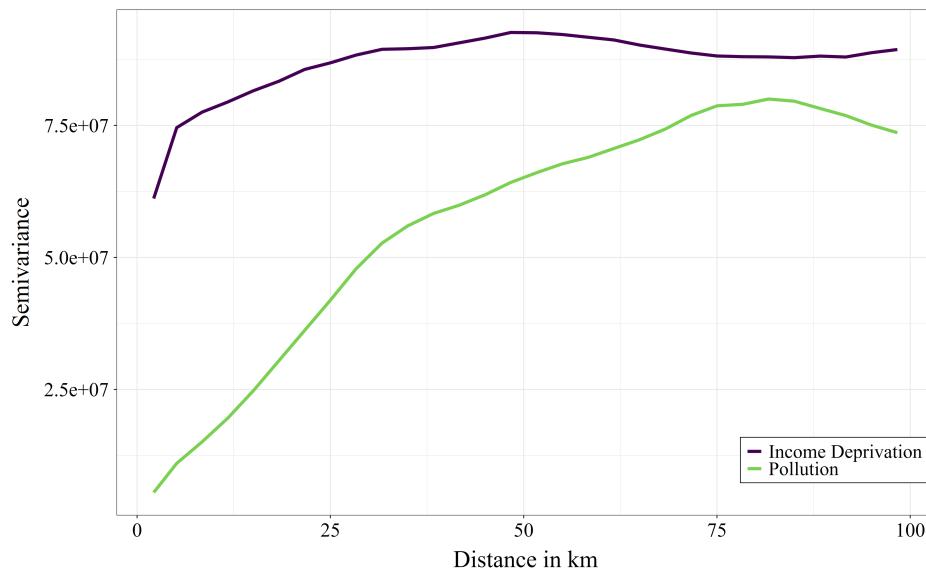


Figure 5.2: Semivariogram of Air pollution and income deprivation in England on the LSOA level for 2019.

## 5 Detecting Spatial Dependence

### 5.1.5 Example

When looking at approx. 10 km distance: the variance in income deprivation is nearly as high when looking at areas within 10km as it would be when looking at areas within 100km distance. This indicates that income deprivation is very local and varies already within smaller areas such as within cities or district. Air pollution, in contrast, has a much lower variance within 10km distances than we would find when looking at the data within 100km distance. This indicates that air pollution has stronger large-scale spatial patterns. When moving locally (e.g. within 10km) to a random location, it would be more difficult to improve in air pollution than it would be to improve in income deprivation.

## 5.2 Local Autocorrelation

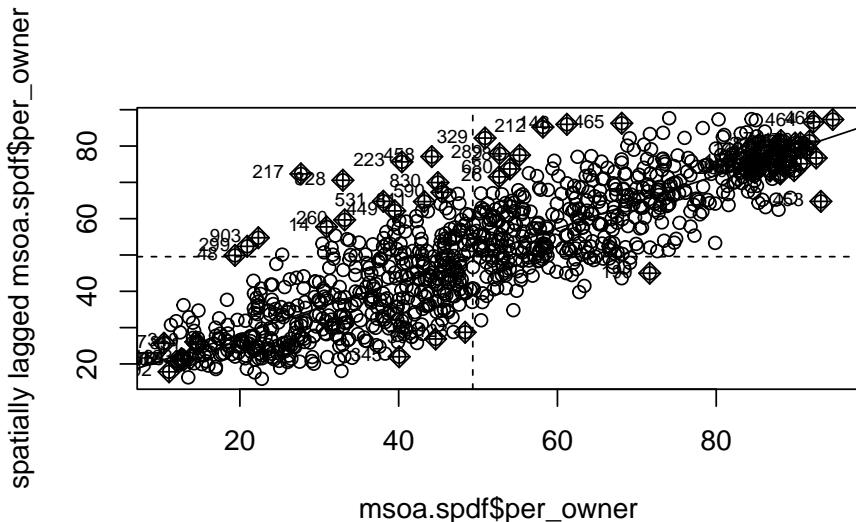
The Global Moran's I statistic above summarizes the spatial pattern by a single value. Although this is helpful to get a feeling of the strength of the general spatial association, it is often more helpful to inspect the spatial pattern in more detail.

The most prominent measure is the Local Indicators of Spatial Association (LISA) (Anselin 1995). LISA measures assess the importance and significance of a statistic at different spatial locations. For more information see for instance the GeoData Materials by Luc Anselin.

For instance, we can use the Moran Plot to identify how single (pairs of) units contribute to the overall dependence.

```
mp <- moran.plot(msoa.spdf$per_owner, queens.lw)
```

## 5.2 Local Autocorrelation



In the lower left corner, we see units with a low-low share of home ownership: focal and neighbouring units have a low share of home owners. In the top right corner, by contrast, we see high-high units.

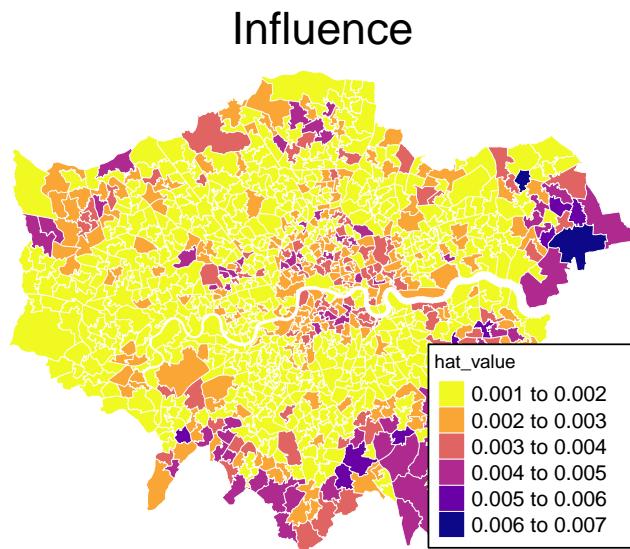
And we can plot influence values on the Overall Moran statistic.

```
msoa.spdf$hat_value <- mp$hat
mp1 <- tm_shape(msoa.spdf) +
  tm_fill(col = "hat_value",
          palette = viridis(n = 10, direction = -1, option = "C"),
          ) +
  tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +
  tm_layout(frame = FALSE,
            legend.frame = TRUE, legend.bg.color = TRUE,
            legend.position = c("right", "bottom"),
```

## 5 Detecting Spatial Dependence

```
legend.outside = FALSE,  
main.title = "Influence",  
main.title.position = "center",  
main.title.size = 1.6,  
legend.title.size = 0.8,  
legend.text.size = 0.8)
```

mp1



### 5.3 Local Moran's I

Local Moran's I is a local version of the overall Moran's I to identify local clusters and local spatial outliers (Anselin 1995). The Local Moran's I is just a local version which is calculated for each location:

### 5.3 Local Moran's $I$

$$\mathbf{I}_i = \frac{z_i \sum_j w_{ij} z_j}{\sum_i (z_i)^2 / (n - 1)}, \text{ where}$$

We use the function `localmoran()` to calculate the local test statistic .

```
loci <- localmoran(msoa.spdf$per_owner, listw = queens.lw)
head(loci)

      Ii          E.Ii       Var.Ii        Z.Ii Pr(z != E(Ii))
1 0.42322928 -1.285364e-04 0.011367934  3.9706976 7.166249e-
05
2 -0.12775982 -2.229957e-05 0.003634711 -2.1187688 3.411001e-
02
3 0.38111534 -6.569549e-04 0.091630752  1.2611995 2.072370e-
01
4 1.02874685 -1.428679e-03 0.279333375  1.9491704 5.127507e-
02
5 0.08553291 -2.108521e-04 0.041275789  0.4220412 6.729949e-
01
6 -0.24014505 -2.228818e-04 0.036321252 -1.2588964 2.080678e-
01
```

It also has an attribute with the Moran plot quadrant of each observation.

```
head(attr(loci, "quadr"))

      mean     median     pysal
1 Low-Low   Low-Low   Low-Low
2 Low-High  Low-High  Low-High
3 High-High High-High High-High
4 High-High High-High High-High
```

## 5 Detecting Spatial Dependence

```
5 High-High High-High High-High  
6 Low-High  Low-High  Low-High
```

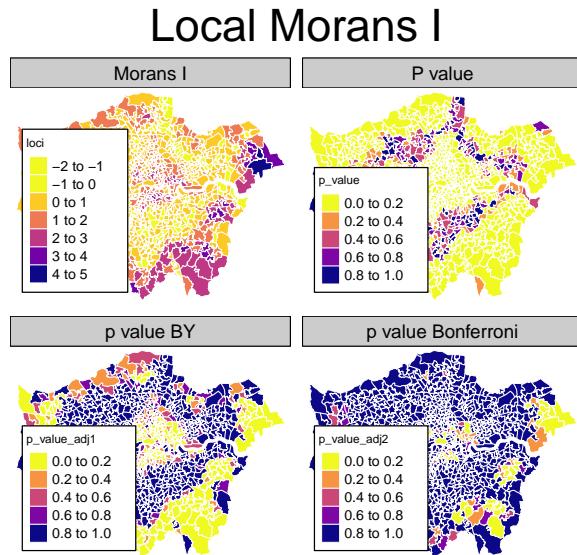
This returns a data.frame with local moran statistic, the expectation of local moran statistic, its variance, and a p value for the statistical significance of each unit. Note that we obviously have a problem of multiple comparisons here and thus may want to correct the significance level, e.g. by Bonferroni adjustment (R. Bivand and Wong 2018).

```
loci.df <- data.frame(loci)  
names(loci.df) <- gsub("\\.", "", names(loci.df))  
msoa.spdf$loci <- loci.df$II  
msoa.spdf$p_value <- loci.df$PrzEIi  
msoa.spdf$p_value_adj1 <- p.adjust(loci.df$PrzEIi, "BY")  
msoa.spdf$p_value_adj2 <- p.adjust(loci.df$PrzEIi, "bonferroni")  
  
mp1 <- tm_shape(msoa.spdf) +  
  tm_fill(col = c("loci", "p_value", "p_value_adj1", "p_value_adj2"),  
          palette = viridis(n = 10, direction = -1, option = "C"),  
          ) +  
  tm_borders(col = "white", lwd = 0.5, alpha = 0.5) +  
  tm_layout(frame = FALSE,  
            legend.frame = TRUE, legend.bg.color = TRUE,  
            legend.position = c("left", "bottom"),  
            legend.outside = FALSE,  
            main.title = "Local Morans I",  
            main.title.position = "center",  
            main.title.size = 1.6,  
            legend.title.size = 0.8,  
            legend.text.size = 0.8,  
            panel.labels = c("Morans I",  
                           "P value",
```

### 5.3 Local Moran's I

```
"p value BY",
 "p value Bonferroni"))
```

```
mp1
```



Something you can often see are so called LISA hotspot maps. They are based on the same idea as the moran plot, and show cluster of high-high and low-low values. We can use the hotspot function to identify the clusters, with a cutoff for singificance and the adjustment for multiple testing.

```
# Calculate clusters
msoa.spdf$lisa_cluster <- hotspot(loci,
                                      "Pr(z != E(Ii))",
                                      cutoff = 0.05,
```

## 5 Detecting Spatial Dependence

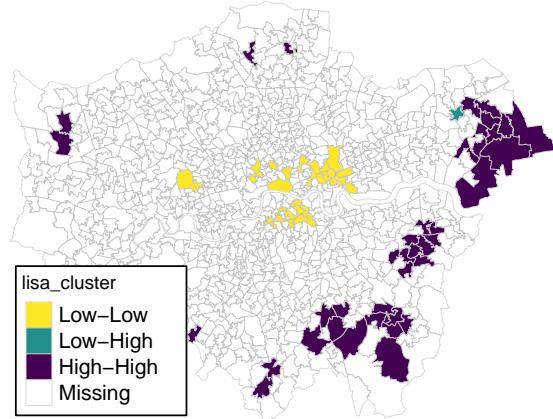
```
quadrant.type = "mean",
p.adjust = "BY")

# Map
mp1 <- tm_shape(msoa.spdf) +
  tm_fill(col = c("lisa_cluster"),
          palette = viridis(n = 3, direction = -1, option = "D"),
          colorNA = "white") +
  tm_borders(col = "grey70", lwd = 0.5, alpha = 0.5) +
  tm_layout(frame = FALSE,
            legend.frame = TRUE, legend.bg.color = TRUE,
            legend.position = c("left", "bottom"),
            legend.outside = FALSE,
            main.title = "Home Ownership \n LISA Clusters p(BY) < 0.05",
            main.title.position = "center",
            main.title.size = 1.6,
            legend.title.size = 0.8,
            legend.text.size = 0.8,)

mp1
```

#### 5.4 Example

### Home Ownership LISA Clusters $p(\text{BY}) < 0.05$



Note that it is not suggested to interpret those cluster as significant in the strict statistical sense. Pebesma and Bivand (2023) suggest to speak of *interesting clusters*. After all, this is an explorative approach. Nevertheless, it can help to identify spatial patterns and clusters.

There are more ways of calculating these hotspot maps and more choices on the cutoffs and calculation of the statistical significance. For more materials see Chapter 15 of Pebesma and Bivand (2023).

## 5.4 Example

### Tate.2021

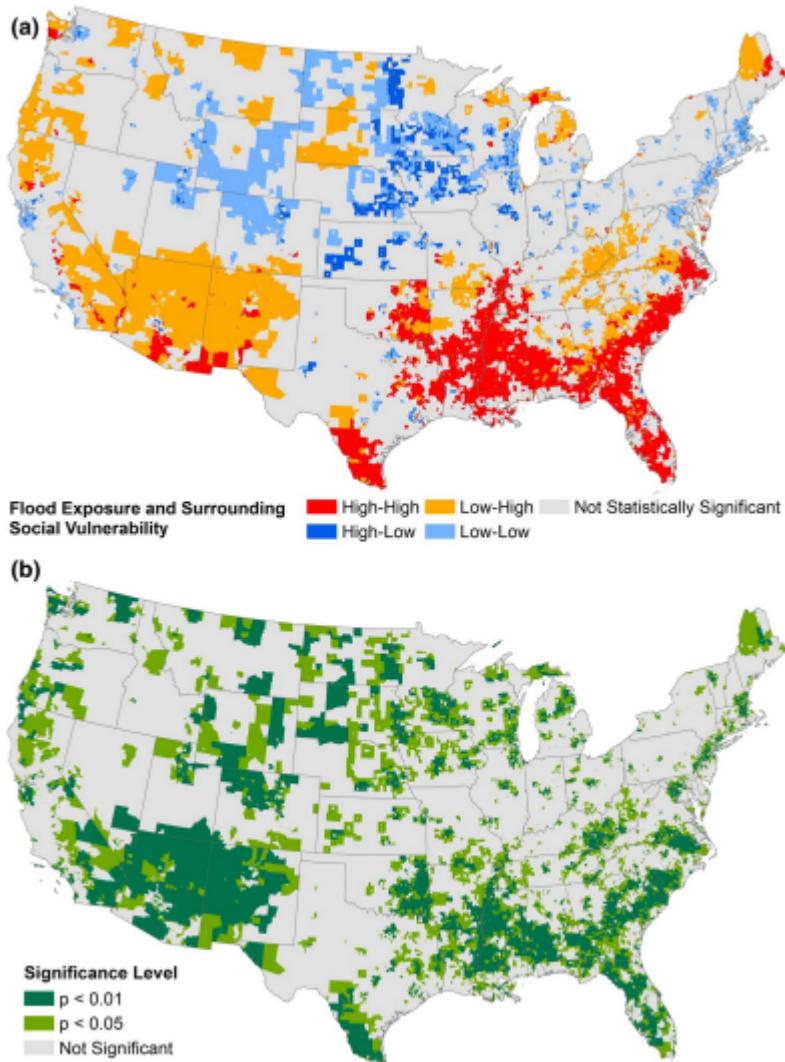
*This study explores the geography of flood exposure and social vulnerability in the conterminous United States based on spatial analysis of fluvial and*

## *5 Detecting Spatial Dependence*

*pluvial flood extent, land cover, and social vulnerability.*

*Mobile homes and racial minorities are most overrepresented in hotspots compared to elsewhere. The results identify priority locations where interventions can mitigate both physical and social aspects of flood vulnerability.*

#### 5.4 Example



**Fig. 4** Bivariate LISA of 100-year flood exposure and surrounding social vulnerability. **a** Cluster map and **b** Cluster significance

## 5.5 Exercises

1. Please calculate a neighbours weights matrix of the nearest 10 neighbours (see `spdep::knearneigh()`), and create a listw object using row normalization.
2. Choose another characteristics from the data (e.g. ethnic groups or house prices) and calculate global Moran's I for it.
3. Produce a LISA cluster map for the characteristic you have chosen.

# 6 Spatial Regression Models

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite") # note: load spd  
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()
```

```
R version 4.4.1 (2024-06-14 ucrt)  
Platform: x86_64-w64-mingw32/x64  
Running under: Windows 11 x64 (build 22631)
```

```
Matrix products: default
```

```
locale:  
[1] LC_COLLATE=English_United Kingdom.utf8  
[2] LC_CTYPE=English_United Kingdom.utf8  
[3] LC_MONETARY=English_United Kingdom.utf8  
[4] LC_NUMERIC=C  
[5] LC_TIME=English_United Kingdom.utf8
```

## 6 Spatial Regression Models

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets   methods    base

other attached packages:
[1] viridisLite_0.4.2 tmap_3.3-4           spatialreg_1.3-
4 Matrix_1.7-0
[5] spdep_1.3-5       spData_2.3.1        mapview_2.11.2     sf_1.0-
16

loaded via a namespace (and not attached):
[1] xfun_0.45          raster_3.6-26        htmlwidgets_1.6.4  lattice_0.22-
6
[5] tools_4.4.1         crosstalk_1.2.1      LearnBayes_2.15.1 parallel_4.4.1
[9] stats4_4.4.1        sandwich_3.1-0       proxy_0.4-27      KernSmooth_2.2
24
[13] satellite_1.0.5    RColorBrewer_1.1-3  leaflet_2.2.2    lifecycle_1.0.4
[17] compiler_4.4.1     deldir_2.0-4       munsell_0.5.1    terra_1.7-
78
[21] codetools_0.2-20   leafsync_0.1.0      stars_0.6-5      htmltools_0.5.8
[25] class_7.3-22       MASS_7.3-60.2      classInt_0.4-
10 lwgeom_0.2-14
[29] wk_0.9.1          abind_1.4-5       boot_1.3-30      multcomp_1.4-
25
[33] nlme_3.1-164       digest_0.6.35      mvtnorm_1.2-
5 splines_4.4.1
[37] fastmap_1.2.0      grid_4.4.1        colorspace_2.1-
0 cli_3.6.2
[41] magrittr_2.0.3      base64enc_0.1-3    dichromat_2.0-
0.1 XML_3.99-0.16.1
[45] survival_3.6-4     leafem_0.2.3       TH.data_1.1-
2 e1071_1.7-14
```

## 6.1 Why do we need spatial regression models

```
[49] scales_1.3.0      sp_2.1-4        rmarkdown_2.27    zoo_1.8-  
12  
[53] png_0.1-8       coda_0.19-4.1   evaluate_0.24.0  knitr_1.47  
[57] tmaptools_3.1-1 s2_1.1.6       rlang_1.1.4     Rcpp_1.0.12  
[61] glue_1.7.0      DBI_1.2.3      rstudioapi_0.16.0 jsonlite_1.8.8  
[65] R6_2.5.1       units_0.8-5
```

### Reload data from previous session

```
load("_data/msoa2_spatial.RData")
```

There are various techniques to model spatial dependence and spatial processes (LeSage and Pace 2009). Here, we will just cover a few of the most common techniques / econometric models. One advantage of the most basic spatial model (SLX) is that this method can easily be incorporated in a variety of other methodologies, such as machine learning approaches.

For more in-depth materials see LeSage and Pace (2009) and Kelejian and Piras (2017). Franzese and Hays (2007), Halleck Vega and Elhorst (2015), LeSage (2014a), Rüttenauer (2022), and Wimpy, Whitten, and Williams (2021) provide article-length introductions. Rüttenauer (2024) is a handbook chapter based on the materials of this workshop.

## 6.1 Why do we need spatial regression models

### 6.1.1 Non-spatial OLS

Let us start with a linear model, where  $\mathbf{y}$  is the outcome or dependent variable ( $N \times 1$ ),  $\mathbf{X}$  are various exogenous covariates ( $N \times k$ ), and  $\epsilon$  ( $N \times 1$ ) is the error term. We are usually interested in the coefficient vector ( $k \times 1$ ) and its insecurity estimates.

## 6 Spatial Regression Models

$$\mathbf{y} = \mathbf{X} +$$

The work-horse for estimating in the social science is the OLS estimator (Wooldridge 2010).

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

### ! OLS assumptions I

1.  $E(\epsilon_i | \mathbf{X}_i) = 0$ : for every value of  $X$ , the average / expectation of the error term equals zero – put differently: the error term is independent of  $X$ ,
2. the observations of the sample are independent and identically distributed (i.i.d),
3. the fourth moments of the variables  $\mathbf{X}_i$  and  $Y_i$  are positive and definite – put differently: extreme values / outliers are very very rare,
4.  $\text{rank}(\mathbf{X}) = K$ : the matrix  $\mathbf{X}$  has full rank – put differently: no perfect multicollinearity between the covariates,

### ! OLS assumptions II

5.  $\text{Var}(\varepsilon|x) = \sigma^2$ : the error terms  $\varepsilon$  are homoskedastic / have the same variance given any value of the explanatory variable,
6.  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ : the error terms  $\varepsilon$  are normally distributed (conditional on the explanatory variables  $X_i$ ).

## 6.1 Why do we need spatial regression models

### 💡 Question

Which of the six assumptions above may be violated by spatial dependence?



### 6.1.2 Problem of ignoring spatial dependence

Does spatial dependence influence the results / coefficient estimates of non-spatial regression models, or in other words: is ignoring spatial dependence harmful?

I've heard different answers, ranging from "It only affects the standard errors" to "it always introduces bias". As so often, the true (or best?) answer is somewhere in the middle: *it depends* (Betz, Cook, and Hollenbach 2020; Cook, Hays, and Franzese 2020; Pace and LeSage 2010; Rüttenauer 2022).

The easiest way to think of it is analogous to the omit variable bias (Betz, Cook, and Hollenbach 2020; Cook, Hays, and Franzese 2020):

## 6 Spatial Regression Models

$$\text{plim } \hat{\beta}_{OLS} = \beta + \gamma \frac{\text{Cov}(\mathbf{x}, \mathbf{z})}{\text{Var}(\mathbf{x})},$$

where  $z$  is some omit variable, and  $\gamma$  is the conditional effect of  $\mathbf{z}$  on  $\mathbf{y}$ . Now imagine that the neighbouring values of the dependent variable  $\mathbf{W}\mathbf{y}$  are autocorrelated to focal unit which we denote with  $\rho > 0$ , and that the covariance between the focal unit's exogenous covariates and  $\mathbf{W}\mathbf{y}$  is not zero. Then we will have an omitted variable bias due to spatial dependence:

$$\text{plim } \hat{\beta}_{OLS} = \beta + \rho \frac{\text{Cov}(\mathbf{x}, \mathbf{W}\mathbf{y})}{\text{Var}(\mathbf{x})} \neq \beta,$$

For completeness, the entire bias is a bit more complicated (Pace and LeSage 2010; Rüttenauer 2022) and looks like:

$$\text{plim } \hat{\beta} = \frac{\sum_{ij} (\mathbf{M}(\delta)\mathbf{M}(\delta)^\top \circ \mathbf{M}(\rho))_{ij}}{\text{tr}(\mathbf{M}(\delta)\mathbf{M}(\delta)^\top)} \beta + \frac{\sum_{ij} (\mathbf{M}(\delta)\mathbf{M}(\delta)^\top \circ \mathbf{M}(\rho)\mathbf{W})_{ij}}{\text{tr}(\mathbf{M}(\delta)\mathbf{M}(\delta)^\top)} \theta,$$

where  $\circ$  denotes the Hadamard product,  $\mathbf{M}(\delta) = (\mathbf{I}_N - \delta\mathbf{W})^{-1}$ , and  $\mathbf{M}(\rho) = (\mathbf{I}_N - \rho\mathbf{W})^{-1}$ .

*(Don't worry, no need to learn by hard!!)*

Essentially, the non-spatial OLS estimator  $\beta_{OLS}$  is biased in the presence of either (Pace and LeSage 2010; Rüttenauer 2022):

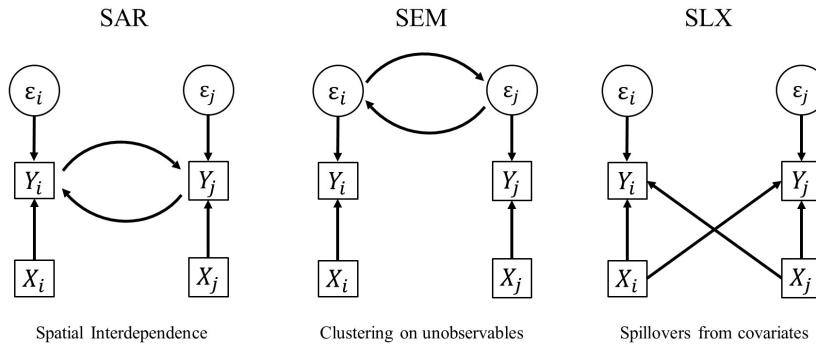
- Spatial autocorrelation in the dependent variable ( $\rho \neq 0$ ) and spatial autocorrelation in the covariate ( $\delta \neq 0$ ). This bias increases with  $\rho$ ,  $\delta$ , and  $\beta$ .

## 6.2 Spatial Regression Models

- Local spatial spillover effects ( $\theta \neq 0$ ) and spatial autocorrelation in the covariate ( $\delta \neq 0$ ). This is analogous to the omitted variable bias resulting from the omission of  $\mathbf{Wx}$ . It increases with  $\theta$  and  $\delta$ , but additionally with  $\rho$  if  $\theta \neq 0$  and  $\delta \neq 0$ .
- An omitted variable and  $E(|\mathbf{x}|) \neq 0$ . This non-spatial omitted variable bias  $\gamma$  is amplified by spatial dependence in the disturbances ( $\lambda$ ) and spatial autocorrelation in the dependent variable ( $\rho$ ), but also increases with positive values of  $\delta$  if either  $\rho \neq 0$  or  $\lambda \neq 0$ . Obviously, it also increases with  $\gamma$ .

## 6.2 Spatial Regression Models

Broadly, spatial dependence or clustering in some characteristic can be the result of three different processes:



Strictly speaking, there are some other possibilities too, such as measurement error or the wrong choice on the spatial level. For instance, imagine we have a city-specific characteristic (e.g. public spending) allocated to neighbourhood units. Obviously, this will introduce heavy autocorrelation on the neighbourhood level by construction.

## 6 Spatial Regression Models

There are three basic ways of incorporating spatial dependence, which then can be further combined. As before, the  $N \times N$  spatial weights matrix  $\mathbf{W}$  defines the spatial relationship between units.

### 6.2.1 Spatial Error Model (SEM)

- Clustering on Unobservables

$$\begin{aligned}\mathbf{y} &= \alpha + \mathbf{X} + \mathbf{u}, \\ \mathbf{u} &= \lambda \mathbf{W} \mathbf{u} +\end{aligned}$$

$\lambda$  denotes the strength of the spatial correlation in the errors of the model:  
*your errors influence my errors.*

- $> 0$ : positive error dependence,
- $< 0$ : negative error dependence,
- $= 0$ : traditional OLS model.

$\lambda$  is defined in the range  $[-1, +1]$ .

### 6.2.2 Spatial Autoregressive Model (SAR)

- Interdependence

$$\mathbf{y} = \alpha + \rho \mathbf{W} \mathbf{y} + \mathbf{X} +$$

$\rho$  denotes the strength of the spatial correlation in the dependent variable (spatial autocorrelation): *your outcome influences my outcome.*

- $> 0$ : positive spatial dependence,
- $< 0$ : negative spatial dependence,

## 6.2 Spatial Regression Models

- $\rho = 0$ : traditional OLS model.

$\rho$  is defined in the range  $[-1, +1]$ .

### 6.2.3 Spatially lagged X Model (SLX)

- Spillovers in Covariates

$$\mathbf{y} = \alpha + \mathbf{X} + \mathbf{W}\mathbf{X} +$$

$\theta$  denotes the strength of the spatial spillover effects from covariate(s) on the dependent variable: *your covariates influence my outcome.*

$\theta$  is basically like any other coefficient from a covariate. It is thus not bound to any range.

Moreover, there are models combining two sets of the above specifications.

### 6.2.4 Spatial Durbin Model (SDM)

- Interdependence
- Spillovers in Covariates

$$\mathbf{y} = \alpha + \rho\mathbf{W}\mathbf{y} + \mathbf{X} + \mathbf{W}\mathbf{X} +$$

## 6 Spatial Regression Models

### 6.2.5 Spatial Durbin Error Model (SDEM)

- Clustering on Unobservables
- Spillovers in Covariates

$$\begin{aligned}\mathbf{y} &= \alpha + \mathbf{X} + \mathbf{WX} + \mathbf{u}, \\ \mathbf{u} &= \lambda \mathbf{Wu} +\end{aligned}$$

### 6.2.6 Combined Spatial Autocorrelation Model (SAC)

- Clustering on Unobservables
- Interdependence

$$\begin{aligned}\mathbf{y} &= \alpha + \rho \mathbf{Wy} + \mathbf{X} + \mathbf{u}, \\ \mathbf{u} &= \lambda \mathbf{Wu} +\end{aligned}$$

### 6.2.7 General Nesting Spatial Model (GNS)

- Clustering on Unobservables
- Interdependence
- Spillovers in Covariates

$$\begin{aligned}\mathbf{y} &= \alpha + \rho \mathbf{Wy} + \mathbf{X} + \mathbf{WX} + \mathbf{u}, \\ \mathbf{u} &= \lambda \mathbf{Wu} +\end{aligned}$$

## 6.2 Spatial Regression Models

### 💡 Manski's reflection problem

The General Nesting Spatial Model (GNS) is only weakly (or not?) identifiable (Gibbons and Overman 2012).

It's analogous to Manski's reflection problem on neighbourhood effects Manski (1993): If people in the same group behave similar, this can be because a) imitating behaviour of the group, b) exogenous characteristics of the group influence the behaviour, and c) members of the same group are exposed to the same external circumstances.  
*We just cannot separate those in observational data.*

Note that all of these models assume different data generating processes (DGP) leading to the spatial pattern. Although there are specifications tests, it is generally not possible to let the data decide which one is the true underlying DGP (Cook, Hays, and Franzese 2020; Rüttenauer 2022). However, there might be theoretical reasons to guide the model specification (Cook, Hays, and Franzese 2020).

Just because SAR is probably the most commonly used model does not make it the best choice. In contrast, various studies (Halleck Vega and Elhorst 2015; Rüttenauer 2022; Wimpy, Whitten, and Williams 2021) highlight the advantages of the relative simple SLX model. Moreover, this specification can basically be incorporated in any other statistical method.

### 6.2.8 A note on missings

Missing values create a problem in spatial data analysis. For instance, in a local spillover model with an average of 10 neighbours, two initial missing values will lead to 20 missing values in the spatially lagged variable. For global spillover models, one initial missing will ‘flow’ through the neighbourhood system until the cutoff point (and create an excess amount of missings).

## *6 Spatial Regression Models*

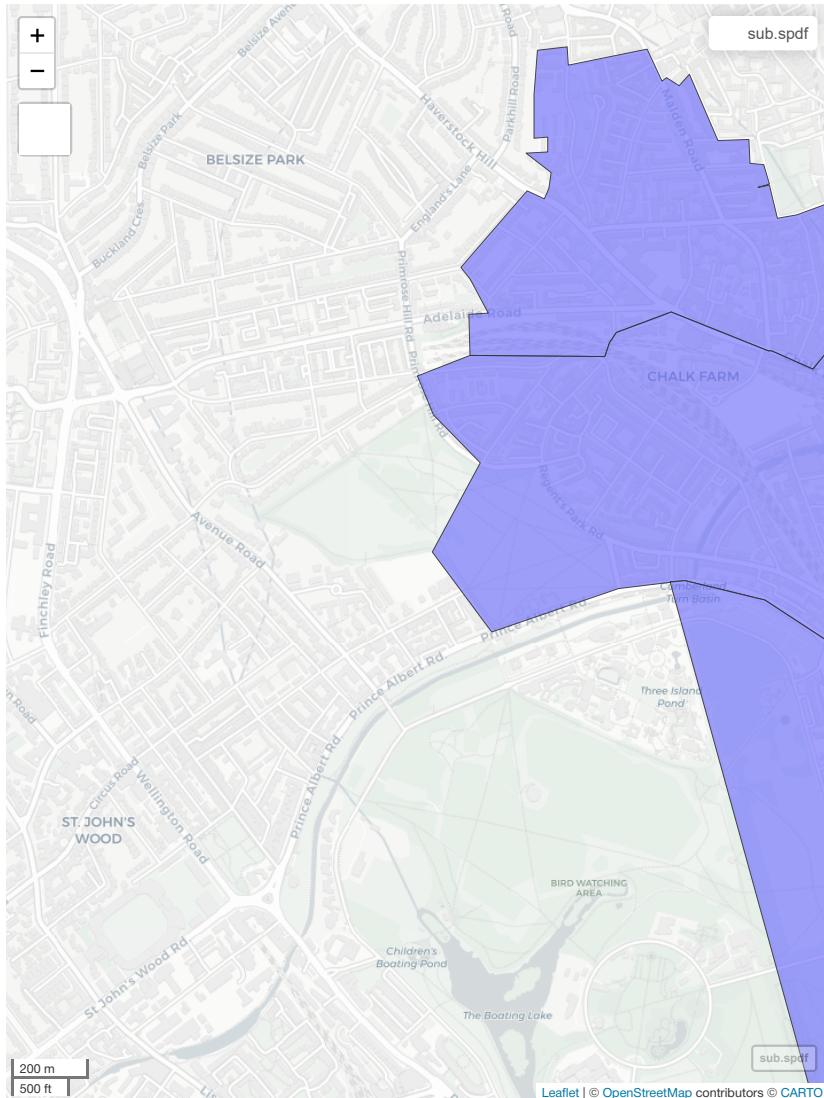
Depending on the data, units with missings can either be dropped and omitted from the initial weights creation, or we need to impute the data first, e.g. using interpolation or Kriging.

### **6.3 Mini Example**

Let's try to make sense of this. We rely on a mini example using a few units in Camden

```
sub.spdf <- msoa.spdf[c(172, 175, 178, 179, 181, 182), ]  
mapview(sub.spdf)
```

### 6.3 Mini Example



## 6 Spatial Regression Models

We then construct queens neighbours, and have a look at the resulting non-normalized matrix  $\mathbf{W}$ .

```
queens.nb <- poly2nb(sub.spdf, queen = TRUE, snap = 1)
W <- nb2mat(queens.nb, style = "B")
W

[,1] [,2] [,3] [,4] [,5] [,6]
172    0    0    1    0    0    0
175    0    0    0    1    0    1
178    1    0    0    1    1    0
179    0    1    1    0    1    1
181    0    0    1    1    0    1
182    0    1    0    1    1    0
attr("call")
nb2mat(neighbours = queens.nb, style = "B")
```

We have selected 6 units. So,  $\mathbf{W}$  is a  $6 \times 6$  matrix. we see that observation 1 has one neighbour: observation 3. Observation 2 has two nieghbours: observation 4 and observation 6. The diagonal is zero: no unit is a neighbour of themselves.

No we row-normalize this matrix.

```
queens.lw <- nb2listw(queens.nb,
                      style = "W")
W_rn <- listw2mat(queens.lw)
W_rn

[,1]      [,2]      [,3]      [,4]      [,5]      [,6]
172 0.0000000 0.0000000 1.0000000 0.0000000 0.0000000 0.0000000
175 0.0000000 0.0000000 0.0000000 0.5000000 0.0000000 0.5000000
```

### 6.3 Mini Example

```
178 0.3333333 0.0000000 0.0000000 0.3333333 0.3333333 0.0000000
179 0.0000000 0.2500000 0.2500000 0.0000000 0.2500000 0.2500000
181 0.0000000 0.0000000 0.3333333 0.3333333 0.0000000 0.3333333
182 0.0000000 0.3333333 0.0000000 0.3333333 0.3333333 0.0000000
```

No every single weight  $w_{ij}$  is divided by the total number of neighbours  $n_i$  of the focal unit. For observation 1, observation 3 is the only neighbour, thus a weight = 1. For observation two, both neighbours have a weight of 1/2. For observation 3 (with three neighbours) each neighbour got a weight of 1/3.

 Question

What happens if we multiply this matrix  $\mathbf{W}$  with a  $N \times 1$  vector  $\mathbf{y}$  or  $\mathbf{x}$ ?

A short reminder on matrix multiplication.

$$\mathbf{W} * \mathbf{y} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} * \begin{bmatrix} y_{11} \\ y_{21} \\ y_{31} \end{bmatrix} = \begin{bmatrix} w_{11}y_{11} + w_{12}y_{21} + w_{13}y_{31} \\ w_{21}y_{11} + w_{22}y_{21} + w_{23}y_{31} \\ w_{31}y_{11} + w_{32}y_{21} + w_{33}y_{31} \end{bmatrix}$$

Each line of  $\mathbf{W} * \mathbf{y}$  just gives a weighted average of the other  $y$ -values  $y_j$  in the sample. In case of the row-normalization, each neighbour gets the same weight  $\frac{1}{n_i}$ . This is simply the mean of  $y_j$  of the neighbours in case of a row-normalized contiguity weights matrix.

Note that the *mean* interpretation is only valid with row-normalization. What would we get with inverse-distance based weights?

Let's look at this in our example

## 6 Spatial Regression Models

```
y <- sub.spdf$med_house_price
x <- sub.spdf$pubs_count

W_rn

[,1]      [,2]      [,3]      [,4]      [,5]      [,6]
172 0.0000000 0.0000000 1.0000000 0.0000000 0.0000000 0.0000000
175 0.0000000 0.0000000 0.0000000 0.5000000 0.0000000 0.5000000
178 0.3333333 0.0000000 0.0000000 0.3333333 0.3333333 0.0000000
179 0.0000000 0.2500000 0.2500000 0.0000000 0.2500000 0.2500000
181 0.0000000 0.0000000 0.3333333 0.3333333 0.0000000 0.3333333
182 0.0000000 0.3333333 0.0000000 0.3333333 0.3333333 0.0000000

y

[1] 376812.5 414625.0 713125.0 322750.0 495000.0 364000.0

x

[1] 1 3 3 1 9 7

W_rn_y <- W_rn %*% y
W_rn_x <- W_rn %*% x
W_rn_y

[,1]
172 713125.0
175 343375.0
178 398187.5
```

### 6.3 Mini Example

```
179 496687.5  
181 466625.0  
182 410791.7
```

```
W_rn_x
```

```
[,1]  
172 3.000000  
175 4.000000  
178 3.666667  
179 5.500000  
181 3.666667  
182 4.333333
```

Let's check if our interpretation is true

```
W_rn_y[1] == y[3]
```

```
[1] TRUE
```

```
W_rn_y[2] == mean(y[c(4, 6)])
```

```
[1] TRUE
```

```
W_rn_y[4] == mean(y[c(2, 3, 5, 6)])
```

```
[1] TRUE
```

## 6.4 Real Example

First, we need the a spatial weights matrix.

```
# Contiguity (Queens) neighbours weights
queens.nb <- poly2nb(msoa.spdf,
                      queen = TRUE,
                      snap = 1) # we consider points in 1m distance as 'tou
queens.lw <- nb2listw(queens.nb,
                      style = "W")
```

We can estimate spatial models using `spatialreg`.

### 6.4.1 SAR

Let's estimate a spatial SAR model using the `lagsarlm()` with contiguity weights. We use median house value as depended variable, and include population density (POPDEN), the air pollution (no2), and the share of ethnic minorities (per\_mixed, per\_asian, per\_black, per\_other).

```
mod_1.sar <- lagsarlm(log(med_house_price) ~ log(no2) + log(POPDEN) +
                        per_mixed + per_asian + per_black + per_other,
                        data = msoa.spdf,
                        listw = queens.lw,
                        Durbin = FALSE) # we could here extend to SDM
summary(mod_1.sar)
```

```
Call:lagsarlm(formula = log(med_house_price) ~ log(no2) + log(POPDEN) +
               per_mixed + per_asian + per_black + per_other, data = msoa.spdf,
               listw = queens.lw, Durbin = FALSE)
```

## 6.4 Real Example

Residuals:

Min	1Q	Median	3Q	Max
-0.5281789	-0.1220524	-0.0099245	0.0992203	1.0936745

Type: lag

Coefficients: (asymptotic standard errors)

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	3.17383180	0.29041604	10.9286	< 2.2e-16
log(no2)	0.39705423	0.04452880	8.9168	< 2.2e-16
log(POPDEN)	-0.05583014	0.01242876	-4.4920	7.055e-06
per_mixed	0.01851577	0.00579832	3.1933	0.001407
per_asian	-0.00228346	0.00045876	-4.9775	6.442e-07
per_black	-0.01263650	0.00100282	-12.6009	< 2.2e-16
per_other	-0.00161419	0.00289082	-0.5584	0.576582

Rho: 0.66976, LR test value: 473.23, p-value: < 2.22e-16

Asymptotic standard error: 0.025311

z-value: 26.461, p-value: < 2.22e-16

Wald statistic: 700.19, p-value: < 2.22e-16

Log likelihood: 196.7203 for lag model

ML residual variance (sigma squared): 0.035402, (sigma: 0.18815)

Number of observations: 983

Number of parameters estimated: 9

AIC: -375.44, (AIC for lm: 95.786)

LM test for residual autocorrelation

test value: 8.609, p-value: 0.0033451

This looks pretty much like a conventional model output, with some additional information: a highly significant `mod_1.sar$rho` of 0.67 indicates strong positive spatial autocorrelation.

Remember that is the coefficient for the term  $\mathbf{y} = \rho \mathbf{W} \mathbf{y}$ .... It is bound to be below 1 for positive autocorrelation.

## 6 Spatial Regression Models

In substantive terms, house prices in the focal unit positively influence house prices in neighbouring units, which again influences house prices among the neighbours of these neighbours, and so on (we'll get back to this).

### ⚠ Warning

The coefficients of covariates in a SAR model are not marginal or partial effects, because of the spillovers and feedback loops in  $\mathbf{y}$  (see below)!

From the coefficient, we can only interpret the direction: there's a positive effect of air pollution and a negative effect of population density, and so on...

### 6.4.2 SEM

SEM models can be estimated using `errorsarlm()`.

```
mod_1.sem <- errorsarlm(log(med_house_price) ~ log(no2) + log(POPDEN) +
                           per_mixed + per_asian + per_black + per_other,
                           data = msoa.spdf,
                           listw = queens.lw,
                           Durbin = FALSE) # we could here extend to SDEM
summary(mod_1.sem)
```

```
Call:errorsarlm(formula = log(med_house_price) ~ log(no2) + log(POPDEN) +
per_mixed + per_asian + per_black + per_other, data = msoa.spdf,
listw = queens.lw, Durbin = FALSE)
```

Residuals:

Min	1Q	Median	3Q	Max
-----	----	--------	----	-----

#### 6.4 Real Example

```
-0.581785 -0.105218 -0.012758  0.094430  0.913425

Type: error
Coefficients: (asymptotic standard errors)
              Estimate Std. Error z value Pr(>|z|)
(Intercept) 12.92801104  0.35239139 36.6865 < 2.2e-16
log(no2)     0.15735296  0.10880727  1.4462 0.1481317
log(POPDEN)  -0.08316270  0.01254315 -6.6301 3.354e-11
per_mixed    -0.03377962  0.00811054 -4.1649 3.115e-05
per_asian    -0.00413115  0.00096849 -4.2656 1.994e-05
per_black    -0.01653816  0.00126741 -13.0488 < 2.2e-16
per_other    -0.01693012  0.00462999 -3.6566 0.0002556

Lambda: 0.88605, LR test value: 623.55, p-value: < 2.22e-16
Asymptotic standard error: 0.015803
z-value: 56.068, p-value: < 2.22e-16
Wald statistic: 3143.6, p-value: < 2.22e-16

Log likelihood: 271.8839 for error model
ML residual variance (sigma squared): 0.026911, (sigma: 0.16405)
Number of observations: 983
Number of parameters estimated: 9
AIC: -525.77, (AIC for lm: 95.786)
```

In this case `mod_1.sem$lambda` gives us the spatial parameter. A highly significant lambda of 0.89 indicates that the errors are highly spatially correlated (e.g. due to correlated unobservables). Again,  $\$ = 1 \$$  would be the maximum.

In spatial error models, we can interpret the coefficients directly, as in a conventional linear model.

## 6 Spatial Regression Models

### 6.4.3 SLX

SLX models can either be estimated with `lmSLX()` directly, or by creating  $\mathbf{WX}$  manually and plugging it into any available model-fitting function.

```
mod_1.slx <- lmSLX(log(med_house_price) ~ log(no2) + log(POPDEN) +
                      per_mixed + per_asian + per_black + per_other,
                      data = msoa.spdf,
                      listw = queens.lw,
                      Durbin = TRUE) # use a formula to lag only specific cov
summary(mod_1.slx)
```

Call:

```
lm(formula = formula(paste("y ~ ", paste(colnames(x)[-1], collapse = "+"))),
   data = as.data.frame(x), weights = weights)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.058e+01	1.539e-01	6.878e+01	0.000e+00
log.no2.	-4.407e-01	1.811e-01	-2.434e+00	1.511e-02
log.POPDEN.	-7.684e-02	1.734e-02	-4.430e+00	1.049e-05
per_mixed	-3.304e-02	1.130e-02	-2.925e+00	3.530e-03
per_asian	-2.381e-03	1.474e-03	-1.615e+00	1.065e-01
per_black	-1.623e-02	1.801e-03	-9.009e+00	1.080e-18
per_other	-2.039e-02	6.564e-03	-3.107e+00	1.947e-03

## 6.4 Real Example

lag.log.no2.	9.936e-01	1.994e-01	4.984e+00	7.384e-07
lag.log.POPDEN.	1.133e-01	2.875e-02	3.939e+00	8.759e-05
lag.per_mixed	1.261e-01	1.429e-02	8.820e+00	5.249e-18
lag.per_asian	-3.828e-03	1.661e-03	-2.305e+00	2.140e-02
lag.per_black	-1.805e-02	2.241e-03	-8.056e+00	2.296e-15
lag.per_other	4.814e-02	7.971e-03	6.039e+00	2.204e-09

In SLX models, we can simply interpret the coefficients of direct and indirect (spatially lagged) covariates.

For instance, lets look at population density:

### 💡 Interpretation SLX

1. A high population density in the focal unit is related to lower house prices (a 1% increase in population density decreases house prices by -0.08%), but
2. A high population density in the neighbouring areas is related to higher house prices (while keeping population density in the focal unit constant). A 1% increase in the *average* population density *across the adjacent neighbourhoods* increases house prices in *the focal unit* by 0.11%)

Potential interpretation: areas with a low population density in central regions of the city (high pop density in surrounding neighbourhoods) have higher house prices. We could try testing this interpretation by including the distance to the city centre as a control.

## 6 Spatial Regression Models

Also note how the air pollution coefficient has changed here, with a negative effect in the focal unit and positive one among the neighbouring units.

An alternative way of estimating the same model is lagging the covariates first.

```
# Loop through vars and create lagged variables
msoa.spdf$log_POPDEN <- log(msoa.spdf$POPDEN)
msoa.spdf$log_no2 <- log(msoa.spdf$no2)
msoa.spdf$log_med_house_price <- log(msoa.spdf$med_house_price)

vars <- c("log_med_house_price", "log_no2", "log_POPDEN",
         "per_mixed", "per_asian", "per_black", "per_other",
         "per_owner", "per_social", "pubs_count")
for(v in vars){
  msoa.spdf[, paste0("w.", v)] <- lag.listw(queens.lw,
                                              var = st_drop_geometry(msoa.spdf))
}

# Alternatively:
w_vars <- create_WX(st_drop_geometry(msoa.spdf[, vars]),
                      listw = queens.lw,
                      prefix = "w")

head(w_vars)

w.log_med_house_price w.log_no2 w.log_POPDEN w.per_mixed w.per_asian
1          12.98382  3.843750    4.662014   4.748368 23.899916
2          12.28730  3.098960    3.300901   3.978275 19.951593
3          12.21207  3.206338    4.009795   3.997487 20.793559
4          12.18176  3.169934    3.630360   2.759082  7.633439
5          12.11159  3.221203    3.993660   3.930061 12.791140
6          12.08393  3.217865    3.876070   3.419488  8.997514
w.per_black w.per_other w.per_owner w.per_social w.pubs_count
```

## 6.4 Real Example

1	7.879758	3.2080074	25.75738	33.85580	8.5454545
2	10.451828	1.6368986	66.42278	15.75042	0.6666667
3	12.965863	1.7526693	58.72637	21.38169	0.2857143
4	12.135478	0.6992118	66.52519	19.70500	0.2000000
5	16.108948	1.3817357	53.05539	29.44022	0.4000000
6	15.312652	0.9611710	59.49460	23.81126	0.1666667

And subsequently we use those new variables in a linear model.

```
mod_1.lm <- lm (log(med_house_price) ~ log(no2) + log(POPDEN) +
  per_mixed + per_asian + per_black + per_other +
  w.log_no2 + w.log_POPDEN +
  w.per_mixed + w.per_asian + w.per_black + w.per_other,
  data = msoa.spdf)
summary(mod_1.lm)
```

Call:

```
lm(formula = log(med_house_price) ~ log(no2) + log(POPDEN) +
  per_mixed + per_asian + per_black + per_other + w.log_no2 +
  w.log_POPDEN + w.per_mixed + w.per_asian + w.per_black +
  w.per_other, data = msoa.spdf)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.50809	-0.16605	-0.01817	0.13055	1.09039

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	10.582440	0.153862	68.779	< 2e-16 ***
log(no2)	-0.440727	0.181063	-2.434	0.01511 *
log(POPDEN)	-0.076840	0.017345	-4.430	1.05e-05 ***
per_mixed	-0.033042	0.011298	-2.925	0.00353 **

## 6 Spatial Regression Models

```
per_asian    -0.002381   0.001474  -1.615   0.10655
per_black    -0.016229   0.001801  -9.009   < 2e-16 ***
per_other    -0.020391   0.006564  -3.107   0.00195 **
w.log_no2     0.993602   0.199370   4.984   7.38e-07 ***
w.log_POPDEN 0.113262   0.028752   3.939   8.76e-05 ***
w.per_mixed   0.126069   0.014294   8.820   < 2e-16 ***
w.per_asian   -0.003828   0.001661  -2.305   0.02140 *
w.per_black   -0.018054   0.002241  -8.056   2.30e-15 ***
w.per_other   0.048139   0.007971   6.039   2.20e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2262 on 970 degrees of freedom
Multiple R-squared:  0.653, Adjusted R-squared:  0.6487
F-statistic: 152.1 on 12 and 970 DF,  p-value: < 2.2e-16
```

Looks pretty similar to `lmSLX()` results, and it should! A big advantage of the SLX specification is that we can use the lagged variables in basically all methods which take variables as inputs, such as non-linear models, matching algorithms, and machine learning tools.

Moreover, using the lagged variables gives a high degree of freedom. For instance, we could (not saying that it necessarily makes sense):

- Use different weights matrices for different variables
- Include higher order neighbours using `nblag()` (with an increasing number of orders we go towards a more global model, but we estimate a coefficient for each spillover, instead of estimating just one)
- Use machine learning techniques to determine the best fitting weights specification.

## 6.4 Real Example

### 6.4.4 SDEM

SDEM models can be estimated using `errorsarlm()` with the additional option `Durbin = TRUE`.

```
mod_1.sdem <- errorsarlm(log(med_house_price) ~ log(no2) + log(POPDEN) +
                           per_mixed + per_asian + per_black + per_other,
                           data = msoa.spdf,
                           listw = queens.lw,
                           Durbin = TRUE) # we could here extend to SDEM
summary(mod_1.sdem)
```

```
Call:errorsarlm(formula = log(med_house_price) ~ log(no2) + log(POPDEN) +
per_mixed + per_asian + per_black + per_other, data = msoa.spdf,
listw = queens.lw, Durbin = TRUE)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.617795	-0.106380	-0.014832	0.095826	0.927446

Type: error

Coefficients: (asymptotic standard errors)

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	10.4422703	0.3652148	28.5921	< 2.2e-16
log(no2)	-0.2057493	0.1264914	-1.6266	0.1038248
log(POPDEN)	-0.0769743	0.0132094	-5.8272	5.635e-09
per_mixed	-0.0222406	0.0079705	-2.7904	0.0052649
per_asian	-0.0037484	0.0010054	-3.7284	0.0001927
per_black	-0.0179751	0.0012383	-14.5161	< 2.2e-16
per_other	-0.0150218	0.0044895	-3.3460	0.0008199
lag.log(no2)	1.0004491	0.1739833	5.7503	8.911e-09
lag.log(POPDEN)	-0.0054241	0.0327802	-0.1655	0.8685763

## 6 Spatial Regression Models

```
lag.per_mixed    0.0669699  0.0169349   3.9545 7.668e-05
lag.per_asian    -0.0018566  0.0015957  -1.1635 0.2446368
lag.per_black    -0.0079949  0.0024833  -3.2195 0.0012842
lag.per_other    0.0273378  0.0087430   3.1268 0.0017671

Lambda: 0.76173, LR test value: 455.7, p-value: < 2.22e-16
Asymptotic standard error: 0.024949
z-value: 30.531, p-value: < 2.22e-16
Wald statistic: 932.15, p-value: < 2.22e-16

Log likelihood: 300.847 for error model
ML residual variance (sigma squared): 0.027504, (sigma: 0.16584)
Number of observations: 983
Number of parameters estimated: 15
AIC: -571.69, (AIC for lm: -117.99)
```

And this SDEM can be interpreted like a combination of SEM and SLX.

First, we still see highly significant auto-correlation in the error term. However, it's lower in magnitude now that we also include the **WX** terms.

Second, the coefficients tell a similar story as in the SLX (use the same interpretation), but some coefficient magnitudes have become smaller.

### 6.4.5 SDM

SDM models can be estimated using `lagsarlm()` with the additional option `Durbin = TRUE`.

```
mod_1.sdm <- lagsarlm(log(med_house_price) ~ log(no2) + log(POPDEN) +
                        per_mixed + per_asian + per_black + per_other,
                        data = msoa.spdf,
                        listw = queens.lw,
```

## 6.4 Real Example

```
Durbin = TRUE) # we could here extend to SDM
summary(mod_1.sdm)

Call:lagsarlm(formula = log(med_house_price) ~ log(no2) + log(POPDEN) +
  per_mixed + per_asian + per_black + per_other, data = msoa.spdf,
  listw = queens.lw, Durbin = TRUE)

Residuals:
    Min          1Q      Median          3Q          Max  
-0.614314 -0.107947 -0.013509  0.092234  0.917398 

Type: mixed
Coefficients: (asymptotic standard errors)
            Estimate Std. Error z value Pr(>|z|)    
(Intercept) 2.7843426 0.2944721 9.4554 < 2.2e-16
log(no2)     -0.3112762 0.1308101 -2.3796 0.0173312
log(POPDEN)   -0.0802866 0.0125213 -6.4120 1.436e-10
per_mixed    -0.0368998 0.0081596 -4.5223 6.118e-06
per_asian     -0.0033726 0.0010636 -3.1711 0.0015189
per_black     -0.0159770 0.0013006 -12.2848 < 2.2e-16
per_other     -0.0209743 0.0047369 -4.4279 9.516e-06
lag.log(no2)  0.4880923 0.1456778 3.3505 0.0008067
lag.log(POPDEN) 0.0781188 0.0207600 3.7629 0.0001679
lag.per_mixed 0.0640880 0.0104646 6.1243 9.110e-10
lag.per_asian  0.0017665 0.0012101 1.4598 0.1443498
lag.per_black  0.0070487 0.0017938 3.9295 8.511e-05
lag.per_other  0.0284822 0.0057774 4.9299 8.226e-07

Rho: 0.73126, LR test value: 501.83, p-value: < 2.22e-16
Asymptotic standard error: 0.025889
z-value: 28.246, p-value: < 2.22e-16
Wald statistic: 797.86, p-value: < 2.22e-16
```

## 6 Spatial Regression Models

```
Log likelihood: 323.9111 for mixed model
ML residual variance (sigma squared): 0.026633, (sigma: 0.1632)
Number of observations: 983
Number of parameters estimated: 15
AIC: -617.82, (AIC for lm: -117.99)
LM test for residual autocorrelation
test value: 36.704, p-value: 1.3747e-09
```

And this SDM can be interpreted like a combination of SAR and SLX.

First, there's still substantial auto-correlation in  $y$ , and this has become even stronger as compared to SAR.

Second, we can interpret the direction of the effect, but we *cannot interpret the coefficient as marginal effects*.

# 7 Spatial Regression Models: Estimation

This section is strongly based on Sarrias (2023), despite being much less detailed than the original.

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite") # note: load spd  
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()  
  
R version 4.4.1 (2024-06-14 ucrt)  
Platform: x86_64-w64-mingw32/x64  
Running under: Windows 11 x64 (build 22631)  
  
Matrix products: default  
  
locale:
```

## 7 Spatial Regression Models: Estimation

```
[1] LC_COLLATE=English_United Kingdom.utf8
[2] LC_CTYPE=English_United Kingdom.utf8
[3] LC_MONETARY=English_United Kingdom.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United Kingdom.utf8

time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets  methods   base

other attached packages:
[1] viridisLite_0.4.2 tmap_3.3-4        spatialreg_1.3-
4 Matrix_1.7-0
[5] spdep_1.3-5     spData_2.3.1       mapview_2.11.2    sf_1.0-
16

loaded via a namespace (and not attached):
[1] xfun_0.45          raster_3.6-26      htmlwidgets_1.6.4 lattice_0.22-
6
[5] tools_4.4.1         crosstalk_1.2.1   LearnBayes_2.15.1 parallel_4.4.1
[9] stats4_4.4.1        sandwich_3.1-0    proxy_0.4-27    KernSmooth_2.2
24
[13] satellite_1.0.5    RColorBrewer_1.1-3 leaflet_2.2.2    lifecycle_1.0.4
[17] compiler_4.4.1     deldir_2.0-4     munsell_0.5.1   terra_1.7-
78
[21] codetools_0.2-20   leafsync_0.1.0   stars_0.6-5    htmltools_0.5.8
[25] class_7.3-22       MASS_7.3-60.2   classInt_0.4-
10 lwgeom_0.2-14
[29] wk_0.9.1          abind_1.4-5    boot_1.3-30   multcomp_1.4-
25
[33] nlme_3.1-164       digest_0.6.35   mvtnorm_1.2-
5      splines_4.4.1
```

```
[37] fastmap_1.2.0      grid_4.4.1        colorspace_2.1-
0   cli_3.6.2
[41] magrittr_2.0.3    base64enc_0.1-3  dichromat_2.0-
0.1 XML_3.99-0.16.1
[45] survival_3.6-4   leafem_0.2.3    TH.data_1.1-
2     e1071_1.7-14
[49] scales_1.3.0     sp_2.1-4       rmarkdown_2.27
12
[53] png_0.1-8        coda_0.19-4.1  evaluate_0.24.0
[57] tmaptools_3.1-1  s2_1.1.6       rlang_1.1.4
[61] glue_1.7.0       DBI_1.2.3     rstudioapi_0.16.0
[65] R6_2.5.1         units_0.8-5  zoo_1.8-
12
[53] png_0.1-8        coda_0.19-4.1  evaluate_0.24.0
[57] tmaptools_3.1-1  s2_1.1.6       rlang_1.1.4
[61] glue_1.7.0       DBI_1.2.3     rstudioapi_0.16.0
[65] R6_2.5.1         units_0.8-5  zoo_1.8-
12
[53] png_0.1-8        coda_0.19-4.1  evaluate_0.24.0
[57] tmaptools_3.1-1  s2_1.1.6       rlang_1.1.4
[61] glue_1.7.0       DBI_1.2.3     rstudioapi_0.16.0
[65] R6_2.5.1         units_0.8-5  jsonlite_1.8.8
```

## Reload data from previous session

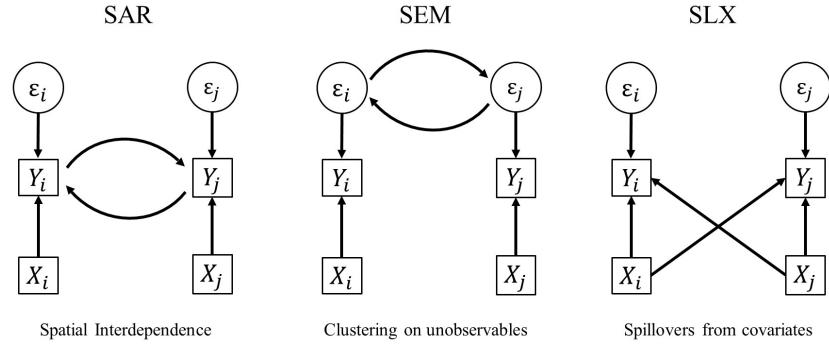
```
load("_data/msoa2_spatial.RData")
```

Note that most of the spatial model specifications can not be estimated by Least Squares (LS), as using (constrained) LS estimators for models containing a spatially lagged dependent variable or disturbance leads to inconsistent results (Anselin and Bera 1998; Franzese and Hays 2007). However, an extensive amount of econometric literature discusses different estimation methods based on (quasi-) maximum likelihood (Anselin 1988; L. Lee 2004; Ord 1975) or instrumental variable approaches using generalized methods of moments (Drukker, Egger, and Prucha 2013; Kelejian and Prucha 1998, 2010), in which the endogenous lagged variables can be instrumented by  $q$  higher order lags of the exogenous regressors ( $\mathbf{X}, \mathbf{WX}, \mathbf{W}^2\mathbf{X}, \dots, \mathbf{W}^q\mathbf{X}$ ) (Kelejian and Prucha 1998).

## 7 Spatial Regression Models: Estimation

### 7.1 Simultaneity bias

Remember what is happening when we estimate a spatial auto-regressive model.



Note the circular process here: My  $X$  influences my  $Y$ , which then influences your  $Y$ , which then influences my  $Y$  again. We write this as

$$\mathbf{y} = \alpha + \rho \mathbf{W} \mathbf{y} + \mathbf{X} + .$$

If we ignore  $\mathbf{X}$  and write the pure auto-regressive term in its reduce form, we get:

$$\mathbf{y} = (\mathbf{I}_n - \rho \mathbf{W})^{-1} \varepsilon,$$

and the spatial lag term is

$$\mathbf{W} \mathbf{y} = \mathbf{W} (\mathbf{I}_n - \rho \mathbf{W})^{-1} \varepsilon.$$

The OLS estimator for the spatial lag term then is

## 7.2 Instrumental variable

$$\hat{\rho}_{OLS} = \left[ \underbrace{(\mathbf{W}\mathbf{y})^\top}_{(1 \times n)} \underbrace{(\mathbf{W}\mathbf{y})}_{(n \times 1)} \right]^{-1} \underbrace{(\mathbf{W}\mathbf{y})^\top}_{(1 \times n)} \underbrace{\mathbf{y}}_{(n \times 1)} .$$

It can then be shown that the OLS estimators equals

$$\hat{\rho}_{OLS} = \rho + [(\mathbf{W}\mathbf{y})^\top (\mathbf{W}\mathbf{y})]^{-1} (\mathbf{W}\mathbf{y})^\top \varepsilon = \rho + \left( \sum_{i=1}^n \mathbf{y}_{Li}^2 \right)^{-1} \left( \sum_{i=1}^n \mathbf{y}_{Li} \epsilon_i \right),$$

with  $\mathbf{y}_{Li}$  defined as the  $i$ th element of the spatial lag operator  $\mathbf{W}\mathbf{y} = \mathbf{y}_L$ . It can further be shown that the second part of the equation  $\neq 0$ , which demonstrates that OLS gives a biased estimate of  $\rho$  (Franzese and Hays 2007; Sarrias 2023).



### Warning

Do not estimate spatial lags of the dependent variable in OLS. It will suffer from simultaneity bias.

## 7.2 Instrumental variable

A potential way of estimating spatial lag /SAR models is 2SLS (Kelejian and Prucha 1998).

We start with our standard model

$$\mathbf{y} = \alpha + \rho \mathbf{W}\mathbf{y} + \mathbf{X} + \varepsilon .$$

## 7 Spatial Regression Models: Estimation

As we have seen above, there is a problem of simultaneity: the “covariate”  $\mathbf{W}\mathbf{y}$  is endogenous. One way of dealing with this endogeneity problem is the Instrumental Variable approach.

So, the question is what are good instruments  $\mathbf{H}$  for  $\mathbf{W}\mathbf{y}$ ? As we have specified the mode, we are sure that  $\mathbf{X}$  determines  $\mathbf{y}$ . Thus, it must be true that  $\mathbf{WX}$  and  $\mathbf{W}^2\mathbf{X}, \dots, \mathbf{W}^l\mathbf{X}$  determines  $\mathbf{W}\mathbf{y}$ .

Note that  $\mathbf{W}^l$  denotes higher orders of  $\mathbf{W}$ . So  $\mathbf{W}^2$  are the second order neighbours (neighbours of neighbours), and  $\mathbf{W}^3$  are the third order neighbours (the neighbours of my neighbour’s neighbours), and so on...

We will discuss this in more detail later, but note for now that the reduced form of the SAR always contains a series of higher order neighbours.

$$(\mathbf{I}_N - \rho \mathbf{W})^{-1} \beta_k = (\mathbf{I}_N + \rho \mathbf{W} + \rho^2 \mathbf{W}^2 + \rho^3 \mathbf{W}^3 + \dots) \beta_k = (\mathbf{I}_N + \sum_{h=1}^{\infty} \rho^h \mathbf{W}^h) \beta_k.$$

Thus, Kelejian and Prucha (1998) suggested to use a set of lagged covariates as instruments for  $\mathbf{WY}$ :

$$\mathbf{H} = \mathbf{X}, \mathbf{WX}, \mathbf{W}^2\mathbf{X}, \dots, \mathbf{W}^l\mathbf{X},$$

where  $l$  is a pre-defined number for the higher order neighbours included. In practice,  $l$  is usually restricted to  $l = 2$ .

This has further been developed by, for instance, using a (truncated) power series as instruments (Kelejian, Prucha, and Yuzefovich 2004):

$$\mathbf{H} = \left[ \mathbf{X}, \mathbf{W} \left( \sum_{l=1}^{\infty} \rho^l \mathbf{W}^l \right) \mathbf{X} \right].$$

We can estimate this using the package **spatialreg** with the function **stsols()**,

## 7.2 Instrumental variable

```
mod_1.sls <- stsls(log(med_house_price) ~ log(no2) + log(POPDEN) +
                     per_mixed + per_asian + per_black + per_other,
                     data = msoa.spdf,
                     listw = queens.lw,
                     robust = TRUE, # heteroskedasticity robust SEs
                     W2X = TRUE) # Second order neighbours are included as instruments (else
summary(mod_1.sls)
```

Call:stsls(formula = log(med\_house\_price) ~ log(no2) + log(POPDEN) +
per\_mixed + per\_asian + per\_black + per\_other, data = msoa.spdf,
listw = queens.lw, robust = TRUE, W2X = TRUE)

Residuals:

Min	1Q	Median	3Q	Max
-0.5464924	-0.1238002	-0.0052299	0.0989150	1.0793093

Coefficients:

	Estimate	HCO	std. Error	z value	Pr(> z )
Rho	0.71004211	0.04678235	15.1776	< 2.2e-16	
(Intercept)	2.73582523	0.50997823	5.3646	8.113e-08	
log(no2)	0.37752751	0.04920257	7.6729	1.688e-14	
log(POPDEN)	-0.05710992	0.01684036	-3.3913	0.0006957	
per_mixed	0.01634307	0.00588488	2.7771	0.0054842	
per_asian	-0.00205426	0.00045905	-4.4750	7.640e-06	
per_black	-0.01166456	0.00128557	-9.0734	< 2.2e-16	
per_other	-0.00280423	0.00332302	-0.8439	0.3987377	

Residual variance (sigma squared): 0.035213, (sigma: 0.18765)

### 7.3 Generalized Method of Moments

Generalized Method of Moments (GMM) provides a way of estimating spatial error / SEM models. A motivation for GMM was that Maximum Likelihood was unfeasible for large samples and its consistency could not be shown. Kelejian and Prucha (1999) thus proposed a Moments estimator for SEM.

We start with the model

$$\begin{aligned}\mathbf{y} &= \alpha + \mathbf{X} + \mathbf{u}, \\ \mathbf{u} &= \lambda \mathbf{W} \mathbf{u} +\end{aligned}$$

The key issue here is to find a consistent estimator for  $\lambda$ . However, we usually do not want to draw inference about  $\lambda$  itself, but only need it to consistently estimate  $\lambda$ . Kelejian and Prucha (1999) thus treat  $\lambda$  as pure nuisance parameter.

In essence, the GMM works as follows (Sarrià 2023):

- 1) First of all obtain a consistent estimate of  $\lambda$ , say  $\tilde{\lambda}$  using either OLS or non-linear least squares (NLS).
- 2) Use this estimate to obtain an estimate of  $\mathbf{u}$ , say  $\hat{\mathbf{u}}$ ,
- 3) Use  $\hat{\mathbf{u}}$ , to estimate  $\lambda$ , say  $\hat{\lambda}$ , using

$$(\hat{\lambda}_{NLS,n}, \hat{\sigma}_{NLS,N}^2) = \operatorname{argmin} \left\{ \frac{1}{n} (\lambda, \sigma^2)^T \frac{1}{n} (\lambda, \sigma^2) : \rho \in [-a, a], \sigma^2 \in [0, b] \right\},$$

- 4) Estimate  $\lambda$  using Equation

### 7.3 Generalized Method of Moments

$$FGLS(\lambda) = \left[ \mathbf{X}^\top (\hat{\lambda})^{-1} \mathbf{X} \right]^{-1} \mathbf{X}^\top (\hat{\lambda})^{-1} \mathbf{y},$$

$$(\lambda) = (\mathbf{I} - \lambda \mathbf{W})^{-1} (\mathbf{I} - \lambda \mathbf{W}^\top)^{-1}$$

For more, see for instance Kelejian and Piras (2017), chapter 2.2.4 or Sarrias (2023).

We can calculate the estimator using `GMerrorsar()` from `spatialreg`.

```
mod_1.gmm <- GMerrorsar(log(med_house_price) ~ log(no2) + log(POPDEN) +
                           per_mixed + per_asian + per_black + per_other,
                           data = msoa.spdf,
                           listw = queens.lw,
                           se.lambda = TRUE) # Provide standard error for lambda
summary(mod_1.gmm)
```

```
Call:GMerrorsar(formula = log(med_house_price) ~ log(no2) + log(POPDEN) +
per_mixed + per_asian + per_black + per_other, data = msoa.spdf,
listw = queens.lw, se.lambda = TRUE)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.8275979	-0.1840855	-0.0096616	0.1610019	1.2270026

Type: GM SAR estimator

Coefficients: (GM standard errors)

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	11.07612114	0.26596129	41.6456	< 2.2e-16
log(no2)	0.67758095	0.08620995	7.8597	3.775e-15
log(POPDEN)	-0.08006377	0.01464953	-5.4653	4.622e-08
per_mixed	-0.01307831	0.00894766	-1.4616	0.1438

## 7 Spatial Regression Models: Estimation

```
per_asian   -0.00521983  0.00090937  -5.7400  9.465e-09
per_black    -0.01957288  0.00134527  -14.5494 < 2.2e-16
per_other     -0.00521695  0.00489760  -1.0652      0.2868

Lambda: 0.69344 (standard error): 0.071248 (z-value): 9.7328
Residual variance (sigma squared): 0.037126, (sigma: 0.19268)
GM argmin sigma squared: 0.037239
Number of observations: 983
Number of parameters estimated: 9
```

## 7.4 Maximum likelihood estimation

### 7.4.1 ML SAR

Maximum Likelihood estimation of spatial models is the most common way of estimation. The procedure to estimate Sar models via ML is based on Ord (1975) and Anselin (1988).

Starting with

$$\mathbf{y} = \rho \mathbf{W}\mathbf{y} + \mathbf{X} + \varepsilon,$$
$$\varepsilon \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 \mathbf{I}_n),$$

and its reduced form

$$\mathbf{y} = (\mathbf{I}_N - \rho \mathbf{W})^{-1}(\mathbf{X} + ),$$
$$\mathbf{y} = \mathbf{A}^{-1}(\mathbf{X} + ),$$

where  $\mathbf{A} = (\mathbf{I}_N - \rho \mathbf{W})$ .

#### 7.4 Maximum likelihood estimation

The ML estimator then chooses the parameters  $\hat{\rho}$ ,  $\hat{\gamma}$ , and  $\hat{\sigma}$  to maximize the probability of fitting the observed sample based on the Likelihood function

$$\begin{aligned}\mathcal{L}(\cdot) &= \log |\mathbf{A}| - \frac{n \log(2\pi)}{2} - \frac{n \log(\sigma^2)}{2} - \frac{1}{2\sigma^2} (\mathbf{A}\mathbf{y} - \mathbf{X})^\top (\mathbf{A}\mathbf{y} - \mathbf{X}) \\ &= \log |\mathbf{A}| - \frac{n \log(2\pi)}{2} - \frac{n \log(\sigma^2)}{2} - \frac{1}{2\sigma^2} [\mathbf{y}^\top \mathbf{A}^\top \mathbf{A}\mathbf{y} - 2(\mathbf{A}\mathbf{y})^\top \mathbf{X} + \mathbf{X}^\top \mathbf{X}],\end{aligned}$$

ML estimation of the SAR works as follows Sarrias (2023):

- 1) Perform the two auxiliary regression of  $\mathbf{y}$  and  $\mathbf{W}\mathbf{y}$  on  $\mathbf{X}$  to obtain the estimators  $\hat{\rho}_O$  and  $\hat{\rho}_L$  as in Equation

$$\begin{aligned}\hat{\rho}_{ML}(\rho) &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - \rho (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}\mathbf{y}, \\ &= \hat{\rho}_O - \hat{\rho}_L.\end{aligned}$$

- 2) Use  $\hat{\rho}_O$  and  $\hat{\rho}_L$  to compute the residuals in Equation

$$\varepsilon_O = \mathbf{y} - \mathbf{X}_0 \quad \text{and} \quad \varepsilon_L = \mathbf{W}\mathbf{y} - \mathbf{X}_L.$$

- 3) By numerical optimization to obtain an estimate of  $\rho$ , maximize the concentrated likelihood given in

$$\ell(\rho) = -\frac{n}{2} - \frac{n}{2} \log(2\pi) - \frac{n}{2} \log \left[ \frac{(\varepsilon_O - \rho\varepsilon_L)^\top (\varepsilon_O - \rho\varepsilon_L)}{n} \right] + \log |\mathbf{I}_n - \rho\mathbf{W}|,$$

- 4) Use the estimate of  $\hat{\rho}$  to plug it back in to the expression for  $\hat{\sigma}^2$

$$\begin{aligned}\hat{\sigma}_{ML}^2(\rho) &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{A}\mathbf{y} \hat{\sigma}_{ML}^2(\rho) = \\ &\frac{(\mathbf{A}\mathbf{y} - \mathbf{X}_{ML})^\top (\mathbf{A}\mathbf{y} - \mathbf{X}_{ML})}{n}\end{aligned}$$

## 7 Spatial Regression Models: Estimation

### 7.4.2 ML SEM

We can also use ML to estimate the spatial error / SEM model of the form

$$\begin{aligned}\mathbf{y} &= \alpha + \mathbf{X} + \mathbf{u}, \\ \mathbf{u} &= \lambda \mathbf{W} \mathbf{u} + \\ \varepsilon &\sim \mathcal{N}(\mathbf{0}_n, \sigma^2 \mathbf{I}_n)\end{aligned}$$

Its reduce for is given by

$$\begin{aligned}\mathbf{y} &= \alpha + \mathbf{X} + (\mathbf{I}_N - \lambda \mathbf{W})^{-1} \cdot \\ \mathbf{y} &= \alpha + \mathbf{X} + \mathbf{B}^{-1}.\end{aligned}$$

where  $\mathbf{B} = (\mathbf{I}_N - \lambda \mathbf{W})$ .

Note that the OLS estimate of the SEM model are unbiased – if there is no omitted variable bias! However, even in that case, they are inefficient if  $\lambda \neq 0$ .

The log-likelihood function is given by

$$\begin{aligned}\ell() &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{(\mathbf{y} - \mathbf{X})^\top (\lambda)(\mathbf{y} - \mathbf{X})}{2\sigma^2} + \log |\mathbf{I}_n - \lambda \mathbf{W}|, \\ (\lambda) &= \mathbf{B}^\top \mathbf{B} = (\mathbf{I}_n - \lambda \mathbf{W})^\top (\mathbf{I}_n - \lambda \mathbf{W})\end{aligned}$$

Based on Anselin and Bera (1998), the ML estimation of SEM follow the procedure (Sarrias 2023):

- 1) Carry out an OLS of  $\mathbf{BX}$  on  $\mathbf{By}$ ; get  $\hat{\alpha}_{OLS}$

## 7.4 Maximum likelihood estimation

- 2) Compute initial set of residuals  $\hat{\epsilon}_{OLS} = \mathbf{By} - \mathbf{BX}_{OLS}$
- 3) Given  $\hat{\epsilon}_{OLS}$ , find  $\hat{\lambda}$  that maximizes the concentrated likelihood

$$\ell(\lambda) = \text{const} + \frac{n}{2} \log \left[ \frac{1}{n} \mathbf{B}^\top \mathbf{B} \right] + \log |\mathbf{B}|.$$

- 4) If the convergence criterion is met, proceed, otherwise repeat steps 1, 2 and 3.
- 5) Given  $\hat{\lambda}$ , estimate  $\gamma(\lambda)$  by GLS and obtain a new vector of residuals,  $\hat{\gamma}(\lambda)$
- 6) Given  $\hat{\gamma}(\lambda)$  and  $\hat{\lambda}$ , estimate  $\hat{\sigma}(\lambda)$ .

The package **spatialreg** Pebesma and Bivand (2023) provides a series of functions to calculate the ML estimator for all spatial models we have considered.

Table from Pebesma and Bivand (2023):

model	model name	maximum likelihood estimation function
SEM	spatial error	<code>errorsarlm(...,</code> <code>Durbin=FALSE)</code>
SEM	spatial error	<code>spautolm(..., family="SAR")</code>
SDEM	spatial Durbin error	<code>errorsarlm(..., Durbin=TRUE)</code>
SLM	spatial lag	<code>lagsarlm(..., Durbin=FALSE)</code>
SDM	spatial Durbin	<code>lagsarlm(..., Durbin=TRUE)</code>
SAC	spatial autoregressive combined	<code>sacsarlm(..., Durbin=FALSE)</code>
GNM	general nested	<code>sacsarlm(..., Durbin=TRUE)</code>

### ML SAR

## 7 Spatial Regression Models: Estimation

```
mod_1.sar <- lagsarlm(log(med_house_price) ~ log(no2) + log(POPDEN) +
                        per_mixed + per_asian + per_black + per_other,
                        data = msoa.spdf,
                        listw = queens.lw,
                        Durbin = FALSE) # we could here extend to SDM
summary(mod_1.sar)

Call:lagsarlm(formula = log(med_house_price) ~ log(no2) + log(POPDEN) +
per_mixed + per_asian + per_black + per_other, data = msoa.spdf,
listw = queens.lw, Durbin = FALSE)

Residuals:
    Min          1Q      Median          3Q          Max
-0.5281789 -0.1220524 -0.0099245  0.0992203  1.0936745

Type: lag
Coefficients: (asymptotic standard errors)
            Estimate Std. Error z value Pr(>|z|)
(Intercept) 3.17383180 0.29041604 10.9286 < 2.2e-16
log(no2)    0.39705423 0.04452880  8.9168 < 2.2e-16
log(POPDEN) -0.05583014 0.01242876 -4.4920 7.055e-06
per_mixed   0.01851577 0.00579832  3.1933  0.001407
per_asian   -0.00228346 0.00045876 -4.9775 6.442e-07
per_black   -0.01263650 0.00100282 -12.6009 < 2.2e-16
per_other   -0.00161419 0.00289082 -0.5584  0.576582

Rho: 0.66976, LR test value: 473.23, p-value: < 2.22e-16
Asymptotic standard error: 0.025311
z-value: 26.461, p-value: < 2.22e-16
Wald statistic: 700.19, p-value: < 2.22e-16

Log likelihood: 196.7203 for lag model
```

## 7.4 Maximum likelihood estimation

```
ML residual variance (sigma squared): 0.035402, (sigma: 0.18815)
Number of observations: 983
Number of parameters estimated: 9
AIC: -375.44, (AIC for lm: 95.786)
LM test for residual autocorrelation
test value: 8.609, p-value: 0.0033451
```

### ML SEM

```
mod_1.sem <- errorsarlm(log(med_house_price) ~ log(no2) + log(POPDEN) +
                           per_mixed + per_asian + per_black + per_other,
                           data = msoa.spdf,
                           listw = queens.lw,
                           Durbin = FALSE) # we could here extend to SDEM
summary(mod_1.sem)
```

```
Call:errorsarlm(formula = log(med_house_price) ~ log(no2) + log(POPDEN) +
                 per_mixed + per_asian + per_black + per_other, data = msoa.spdf,
                 listw = queens.lw, Durbin = FALSE)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.581785	-0.105218	-0.012758	0.094430	0.913425

Type: error

Coefficients: (asymptotic standard errors)

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	12.92801104	0.35239139	36.6865	< 2.2e-16
log(no2)	0.15735296	0.10880727	1.4462	0.1481317
log(POPDEN)	-0.08316270	0.01254315	-6.6301	3.354e-11
per_mixed	-0.03377962	0.00811054	-4.1649	3.115e-05
per_asian	-0.00413115	0.00096849	-4.2656	1.994e-05

7 Spatial Regression Models: Estimation

```
per_black -0.01653816 0.00126741 -13.0488 < 2.2e-16
per_other -0.01693012 0.00462999 -3.6566 0.0002556

Lambda: 0.88605, LR test value: 623.55, p-value: < 2.22e-16
Asymptotic standard error: 0.015803
z-value: 56.068, p-value: < 2.22e-16
Wald statistic: 3143.6, p-value: < 2.22e-16

Log likelihood: 271.8839 for error model
ML residual variance (sigma squared): 0.026911, (sigma: 0.16405)
Number of observations: 983
Number of parameters estimated: 9
AIC: -525.77, (AIC for lm: 95.786)
```

# 8 Exercises II

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite") # note: load spd  
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()
```

```
R version 4.4.1 (2024-06-14 ucrt)  
Platform: x86_64-w64-mingw32/x64  
Running under: Windows 11 x64 (build 22631)
```

```
Matrix products: default
```

```
locale:  
[1] LC_COLLATE=English_United Kingdom.utf8  
[2] LC_CTYPE=English_United Kingdom.utf8  
[3] LC_MONETARY=English_United Kingdom.utf8  
[4] LC_NUMERIC=C  
[5] LC_TIME=English_United Kingdom.utf8
```

## 8 Exercises II

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets  methods   base

other attached packages:
[1] viridisLite_0.4.2 tmap_3.3-4           spatialreg_1.3-
4 Matrix_1.7-0
[5] spdep_1.3-5       spData_2.3.1        mapview_2.11.2    sf_1.0-
16

loaded via a namespace (and not attached):
[1] xfun_0.45          raster_3.6-26        htmlwidgets_1.6.4 lattice_0.22-
6
[5] tools_4.4.1         crosstalk_1.2.1     LearnBayes_2.15.1 parallel_4.4.1
[9] stats4_4.4.1        sandwich_3.1-0      proxy_0.4-27    KernSmooth_2.2
24
[13] satellite_1.0.5    RColorBrewer_1.1-3  leaflet_2.2.2  lifecycle_1.0.4
[17] compiler_4.4.1     deldir_2.0-4       munsell_0.5.1 terra_1.7-
78
[21] codetools_0.2-20   leafsync_0.1.0     stars_0.6-5    htmltools_0.5.8
[25] class_7.3-22      MASS_7.3-60.2     classInt_0.4-
10 lwgeom_0.2-14
[29] wk_0.9.1          abind_1.4-5       boot_1.3-30   multcomp_1.4-
25
[33] nlme_3.1-164       digest_0.6.35    mvtnorm_1.2-
5 splines_4.4.1
[37] fastmap_1.2.0     grid_4.4.1       colorspace_2.1-
0 cli_3.6.2
[41] magrittr_2.0.3     base64enc_0.1-3  dichromat_2.0-
0.1 XML_3.99-0.16.1
[45] survival_3.6-4    leafem_0.2.3     TH.data_1.1-
2 e1071_1.7-14
```

## 8.1 Environmental inequality

```
[49] scales_1.3.0      sp_2.1-4        rmarkdown_2.27    zoo_1.8-  
12  
[53] png_0.1-8       coda_0.19-4.1   evaluate_0.24.0  knitr_1.47  
[57] tmaptools_3.1-1 s2_1.1.6       rlang_1.1.4     Rcpp_1.0.12  
[61] glue_1.7.0      DBI_1.2.3      rstudioapi_0.16.0 jsonlite_1.8.8  
[65] R6_2.5.1       units_0.8-5
```

### Reload data from previous session

```
load("_data/msoa2_spatial.RData")
```

## 8.1 Environmental inequality

How would you investigate the following descriptive research question: Are immigrant minorities in London exposed to higher levels of pollution? Also consider the spatial structure. What's your dependent and what is your independent variable?

### 1) Define a neighbours weights object of your choice

Assume a typical neighbourhood would be 2.5km in diameter

*8 Exercises II*

- 2) Estimate the extent of spatial auto-correlation in air pollution**
  - 3) Estimate a Spatial SAR regression model**
  - 4) Estimate a Spatial SEM regression model**
  - 5) Estimate a Spatial SLX regression model**
  - 6) Estimate a Spatial Durbin regression model**
  - 7) Estimate a Spatial Durbin Error regression model**
- 8.1.1 8) Sneak preview on tomorrow: Which of the spatial model specifications about would you choose / prefer in a real world example?**
- 8.1.2 9) Please calculate the spatially lagged value of the median house price.**
- 8.1.3 10) Can you use the results of the previous task to run a non-linear SLX model, where you predict if an MSOA is within the ulez zone based on the house prices? Can you make sense of the result?**

# 9 Spatial Impacts

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite") # note: load spd  
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()
```

```
R version 4.4.1 (2024-06-14 ucrt)  
Platform: x86_64-w64-mingw32/x64  
Running under: Windows 11 x64 (build 22631)
```

```
Matrix products: default
```

```
locale:  
[1] LC_COLLATE=English_United Kingdom.utf8  
[2] LC_CTYPE=English_United Kingdom.utf8  
[3] LC_MONETARY=English_United Kingdom.utf8  
[4] LC_NUMERIC=C  
[5] LC_TIME=English_United Kingdom.utf8
```

## 9 Spatial Impacts

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets   methods    base

other attached packages:
[1] viridisLite_0.4.2 tmap_3.3-4           spatialreg_1.3-
4 Matrix_1.7-0
[5] spdep_1.3-5       spData_2.3.1        mapview_2.11.2     sf_1.0-
16

loaded via a namespace (and not attached):
[1] xfun_0.45          raster_3.6-26        htmlwidgets_1.6.4  lattice_0.22-
6
[5] tools_4.4.1         crosstalk_1.2.1      LearnBayes_2.15.1 parallel_4.4.1
[9] stats4_4.4.1        sandwich_3.1-0       proxy_0.4-27      KernSmooth_2.2
24
[13] satellite_1.0.5    RColorBrewer_1.1-3  leaflet_2.2.2    lifecycle_1.0.4
[17] compiler_4.4.1     deldir_2.0-4       munsell_0.5.1    terra_1.7-
78
[21] codetools_0.2-20   leafsync_0.1.0      stars_0.6-5      htmltools_0.5.8
[25] class_7.3-22       MASS_7.3-60.2      classInt_0.4-
10 lwgeom_0.2-14
[29] wk_0.9.1          abind_1.4-5       boot_1.3-30      multcomp_1.4-
25
[33] nlme_3.1-164       digest_0.6.35      mvtnorm_1.2-
5 splines_4.4.1
[37] fastmap_1.2.0      grid_4.4.1        colorspace_2.1-
0 cli_3.6.2
[41] magrittr_2.0.3      base64enc_0.1-3    dichromat_2.0-
0.1 XML_3.99-0.16.1
[45] survival_3.6-4     leafem_0.2.3       TH.data_1.1-
2 e1071_1.7-14
```

## 9.1 Coefficient estimates ≠ ‘marginal’ effects

```
[49] scales_1.3.0      sp_2.1-4        rmarkdown_2.27    zoo_1.8-  
12  
[53] png_0.1-8       coda_0.19-4.1   evaluate_0.24.0  knitr_1.47  
[57] tmaptools_3.1-1 s2_1.1.6       rlang_1.1.4     Rcpp_1.0.12  
[61] glue_1.7.0       DBI_1.2.3      rstudioapi_0.16.0 jsonlite_1.8.8  
[65] R6_2.5.1        units_0.8-5
```

### Reload data from previous session

```
load("_data/msoa2_spatial.RData")
```

## 9.1 Coefficient estimates ≠ ‘marginal’ effects



### Warning

Do not interpret coefficients as marginal effects in SAR, SAC, and SDM!!

At first glance, the specifications presented above seem relatively similar in the way of modelling spatial effects. **Yet, they differ in very important aspects.**

First, models with an endogenous spatial term (SAR, SAC, and SDM) assume a very different spatial dependence structure than models with only exogenous spatial terms as SLX and SDEM specifications. While the first three assume **global** spatial dependence, the second two assume **local** spatial dependence (Anselin 2003; Halleck Vega and Elhorst 2015; LeSage and Pace 2009).

Second, the interpretation of the coefficients differs greatly between models with and without endogenous effects. This becomes apparent when consid-

## 9 Spatial Impacts

ering the reduced form of the equations above. Exemplary using the SAR model, the reduced form is given by:

$$\begin{aligned}\mathbf{y} - \rho \mathbf{W} \mathbf{y} &= \mathbf{X} + , \\ (\mathbf{I}_N - \rho \mathbf{W}) \mathbf{y} &= \mathbf{X} + , \\ \mathbf{y} &= (\mathbf{I}_N - \rho \mathbf{W})^{-1}(\mathbf{X} + ),\end{aligned}$$

where  $\mathbf{I}_N$  is an  $N \times N$  diagonal matrix (diagonal elements equal 1, 0 otherwise). This contains no spatially lagged dependent variable on the right-hand side.

If we want to interpret coefficient, we are usually in marginal or partial effects (the association between a unit change in  $X$  and  $Y$ ). We obtain these effects by looking at the first derivative.

When taking the first derivative of the explanatory variable  $\mathbf{x}_k$  from the reduced form in (??) to interpret the partial effect of a unit change in variable  $\mathbf{x}_k$  on  $\mathbf{y}$ , we receive

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}_k} = \underbrace{(\mathbf{I}_N - \rho \mathbf{W})^{-1}}_{N \times N} \beta_k,$$

for each covariate  $k = \{1, 2, \dots, K\}$ . As can be seen, the partial derivative with respect to  $\mathbf{x}_k$  produces an  $N \times N$  matrix, thereby representing the partial effect of each unit  $i$  onto the focal unit  $i$  itself and all other units  $j = \{1, 2, \dots, i-1, i+1, \dots, N\}$ .

Note that the diagonal elements of  $(\mathbf{I}_N - \rho \mathbf{W})^{-1}$  are not zero anymore (as they are in  $\mathbf{W}$ ). Look at the following minimal example:

### 9.1 Coefficient estimates ≠ ‘marginal’ effects

$$\tilde{\mathbf{W}} = \begin{pmatrix} 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \end{pmatrix}, \text{ and normalized } \mathbf{W} = \begin{pmatrix} 0 & 0.5 & 0 & 0.5 & 0 \\ 0.33 & 0 & 0.33 & 0 & 0.33 \\ 0 & 0.5 & 0 & 0.5 & 0 \\ 0.33 & 0 & 0.33 & 0 & 0.33 \\ 0 & 0.5 & 0 & 0.5 & 0 \end{pmatrix}$$

and

$$\rho = 0.6,$$

then

$$\rho \mathbf{W} = \begin{pmatrix} 0 & 0.3 & 0 & 0.3 & 0 \\ 0.2 & 0 & 0.2 & 0 & 0.2 \\ 0 & 0.3 & 0 & 0.3 & 0 \\ 0.2 & 0 & 0.2 & 0 & 0.2 \\ 0 & 0.3 & 0 & 0.3 & 0 \end{pmatrix}.$$

If we want to get the total effect of  $X$  on  $Y$  we need to add the direct association within  $i$  and  $j$  and so on...

$$\begin{aligned} \mathbf{I}_N - \rho \mathbf{W} &= \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 0 & 0.3 & 0 & 0.3 & 0 \\ 0.2 & 0 & 0.2 & 0 & 0.2 \\ 0 & 0.3 & 0 & 0.3 & 0 \\ 0.2 & 0 & 0.2 & 0 & 0.2 \\ 0 & 0.3 & 0 & 0.3 & 0 \end{pmatrix} \\ &= \begin{pmatrix} 1 & -0.3 & 0 & -0.3 & 0 \\ -0.2 & 1 & -0.2 & 0 & -0.2 \\ 0 & 0.3 & 1 & 0.3 & 0 \\ -0.2 & 0 & -0.2 & 1 & -0.2 \\ 0 & -0.3 & 0 & -0.3 & 1 \end{pmatrix}. \end{aligned}$$

## 9 Spatial Impacts

And finally we take the inverse of that

$$\begin{aligned}
 (\mathbf{I}_N - \rho \mathbf{W})^{-1} &= \begin{pmatrix} 1 & -0.3 & 0 & -0.3 & 0 \\ -0.2 & 1 & -0.2 & 0 & -0.2 \\ 0 & 0.3 & 1 & 0.3 & 0 \\ -0.2 & 0 & -0.2 & 1 & -0.2 \\ 0 & -0.3 & 0 & -0.3 & 1 \end{pmatrix}^{-1} \\
 &= \begin{pmatrix} \textcolor{red}{1.1875} & 0.46875 & 0.1875 & 0.46875 & 0.1875 \\ 0.3125 & \textcolor{red}{1.28125} & 0.3125 & 0.28125 & 0.3125 \\ 0.1875 & 0.46875 & \textcolor{red}{1.1875} & 0.46875 & 0.1875 \\ 0.3125 & 0.28125 & 0.3125 & \textcolor{red}{1.28125} & 0.3125 \\ 0.1875 & 0.46875 & 0.1875 & 0.46875 & \textcolor{red}{1.1875} \end{pmatrix}.
 \end{aligned}$$

As you can see,  $(\mathbf{I}_N - \rho \mathbf{W})^{-1}$  has diagonal elements  $> 1$ : these are feedback loops. My  $X$  influences my  $Y$  directly, but my  $Y$  then influences my neighbour's  $Y$ , which then influences my  $Y$  again (also also other neighbour's  $Y$ s). Thus the influence of my  $X$  on my  $Y$  includes a spatial multiplier.

Check yourself:

```

I = diag(5)
rho = 0.6
W = matrix(c(0 , 0.5 , 0 , 0.5 , 0,
            1/3 , 0 , 1/3 , 0 , 1/3,
            0 , 0.5 , 0 , 0.5 , 0,
            1/3 , 0 , 1/3 , 0 , 1/3,
            0 , 0.5 , 0 , 0.5 , 0), ncol = 5, byrow = TRUE)

(IrW = I - rho*W)

[,1] [,2] [,3] [,4] [,5]
[1,] 1.0 -0.3 0.0 -0.3 0.0

```

### 9.1 Coefficient estimates ≠ ‘marginal’ effects

```
[2,] -0.2  1.0 -0.2  0.0 -0.2
[3,]  0.0 -0.3  1.0 -0.3  0.0
[4,] -0.2  0.0 -0.2  1.0 -0.2
[5,]  0.0 -0.3  0.0 -0.3  1.0

# (I - rho*W)^-1
(M = solve(IrW))

[,1]    [,2]    [,3]    [,4]    [,5]
[1,] 1.1875 0.46875 0.1875 0.46875 0.1875
[2,] 0.3125 1.28125 0.3125 0.28125 0.3125
[3,] 0.1875 0.46875 1.1875 0.46875 0.1875
[4,] 0.3125 0.28125 0.3125 1.28125 0.3125
[5,] 0.1875 0.46875 0.1875 0.46875 1.1875
```

The diagonal elements of  $M$  indicate how each unit  $i$  influences itself (change of  $x_i$  on change of  $y_i$ ), and each off-diagonal elements in column  $j$  represents the effect of  $j$  on each other unit  $i$  (change of  $x_j$  on change of  $y_i$ ).

$$\begin{pmatrix} 1.1875 & 0.46875 & 0.1875 & 0.46875 & 0.1875 \\ 0.3125 & 1.28125 & 0.3125 & 0.28125 & 0.3125 \\ 0.1875 & 0.46875 & 1.1875 & 0.46875 & 0.1875 \\ 0.3125 & 0.28125 & 0.3125 & 1.28125 & 0.3125 \\ 0.1875 & 0.46875 & 0.1875 & 0.46875 & 1.1875 \end{pmatrix}.$$

For instance,  $W_{12}$  indicates that unit 2 has an influence of 0.46875 on unit 1. On the other hand,  $W_{53}$  indicates that unit 3 has an influence of magnitude 0.1875 on unit 5.

## 9 Spatial Impacts

### 💡 Question

Why does unit 3 have any effect on unit 5? According to  $\mathbf{W}$  those two units are no neighbours  $w_{53} = 0$ !

## 9.2 Global and local spillovers

The kind of indirect spillover effects in SAR, SAC, and SDM models differs from the kind of indirect spillover effects in SLX and SDEM models: while the first three specifications represent **global spillover effects**, the latter three represent **local spillover effects** (Anselin 2003; LeSage and Pace 2009; LeSage 2014a).

### 9.2.1 Local spillovers

In case of SLX and SDEM the spatial spillover effects can be interpreted as the effect of a one unit change of  $\mathbf{x}_k$  in the spatially weighted neighbouring observations on the dependent variable of the focal unit: the weighted average among neighbours; when using a row-normalised contiguity weights matrix,  $\mathbf{W}\mathbf{x}_k$  is the mean value of  $\mathbf{x}_k$  in the neighbouring units.

Assume we have  $k = 2$  covariates, then

## 9.2 Global and local spillovers

$$\begin{aligned} \underset{N \times N}{\mathbf{W}} \underset{N \times 2}{\mathbf{X}} \underset{2 \times 1}{\mathbf{y}} &= \begin{pmatrix} 0 & 0.5 & 0 & 0.5 & 0 \\ 0.33 & 0 & 0.33 & 0 & 0.33 \\ 0 & 0.5 & 0 & 0.5 & 0 \\ 0.33 & 0 & 0.33 & 0 & 0.33 \\ 0 & 0.5 & 0 & 0.5 & 0 \end{pmatrix} \begin{pmatrix} 3 & 100 \\ 4 & 140 \\ 1 & 200 \\ 7 & 70 \\ 5 & 250 \end{pmatrix} \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \\ &= \begin{pmatrix} 6 & 105 \\ 3 & 190 \\ 6 & 105 \\ 3 & 190 \\ 6 & 105 \end{pmatrix} \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \end{aligned}$$

```
X <- cbind(x1 = c(3,4,1,8,5),
             x2 = c(100,140,200,70,270))
(WX <- W %*% X)
```

```
x1  x2
[1,] 6 105
[2,] 3 190
[3,] 6 105
[4,] 3 190
[5,] 6 105
```

Thus, only direct neighbours – as defined in  $\mathbf{W}$  – contribute to those local spillover effects. The  $\hat{\theta}$  coefficients only estimate how my direct neighbour's  $\mathbf{X}$  values influence my own outcome  $\mathbf{y}$ .

There are no higher order neighbours involved (as long as we do not model them), nor are there any feedback loops due to interdependence.

### 9.2.2 Global spillovers

In contrast, spillover effects in SAR, SAC, and SDM models do not only include direct neighbours but also neighbours of neighbours (second order neighbours) and further higher-order neighbours. This can be seen by rewriting the inverse  $(\mathbf{I}_N - \rho\mathbf{W})^{-1}$  as power series: A power series of  $\sum_{k=0}^{\infty} \mathbf{W}^k$  converges to  $(\mathbf{I} - \mathbf{W})^{-1}$  if the maximum absolute eigenvalue of  $\mathbf{W} < 1$ , which is ensured by standardizing  $\mathbf{W}$ .

$$(\mathbf{I}_N - \rho\mathbf{W})^{-1}\beta_k = (\mathbf{I}_N + \rho\mathbf{W} + \rho^2\mathbf{W}^2 + \rho^3\mathbf{W}^3 + \dots)\beta_k = (\mathbf{I}_N + \sum_{h=1}^{\infty} \rho^h\mathbf{W}^h)\beta_k,$$

where the identity matrix represents the direct effects and the sum represents the first and higher order indirect effects and the above mentioned feedback loops. This implies that a change in one unit  $i$  does not only affect the direct neighbours but passes through the whole system towards higher-order neighbours, where the impact declines with distance within the neighbouring system. Global indirect impacts thus are ‘multiplied’ by influencing direct neighbours as specified in  $\mathbf{W}$  and indirect neighbours not connected according to  $\mathbf{W}$ , with additional feedback loops between those neighbours.

$$\underbrace{(\underbrace{\mathbf{I}_N}_{N \times N} - \underbrace{\rho \mathbf{W}}_{\hat{\equiv} 0.6 N \times N})^{-1}}_{N \times N} \beta_k = \begin{pmatrix} 1.1875 & 0.46875 & 0.1875 & 0.46875 & 0.1875 \\ 0.3125 & 1.28125 & 0.3125 & 0.28125 & 0.3125 \\ 0.1875 & 0.46875 & 1.1875 & 0.46875 & 0.1875 \\ 0.3125 & 0.28125 & 0.3125 & 1.28125 & 0.3125 \\ 0.1875 & 0.46875 & 0.1875 & 0.46875 & 1.1875 \end{pmatrix} (\beta_1 + \beta_2)$$

.

All diagonal elements of  $\text{diag}(\mathbf{W}) = w_{ii} = 0$ . However, diagonal elements of higher order neighbours are not zero  $\text{diag}(\mathbf{W}^2) = \text{diag}(\mathbf{WW}) \neq 0$ .

## 9.2 Global and local spillovers

Intuitively,  $\rho\mathbf{W}$  only represents the effects between direct neighbours (and the focal unit is not a neighbour of the focal unit itself), whereas  $\rho^2\mathbf{W}^2$  contains the effects of second order neighbours, where the focal unit is a second order neighbour of the focal unit itself. Thus,  $(\mathbf{I}_N - \rho\mathbf{W})^{-1}\beta_k$  includes feedback effects from  $\rho^2\mathbf{W}^2$  on (they are part of the direct impacts according to the summary measures below). This is why the diagonal above  $\geq 1$ .

In consequence, local and global spillover effects represent two distinct kinds of spatial spillover effects (LeSage 2014a). The interpretation of local spillover effects is straightforward: it represents the effect of all neighbours as defined by  $\mathbf{W}$  (the average over all neighbours in case of a row-normalised weights matrix).

For instance, the environmental quality in the focal unit itself but also in neighbouring units could influence the attractiveness of a district and its house prices. In this example it seems reasonable to assume that we have local spillover effects: only the environmental quality in directly contiguous units (e.g. in walking distance) is relevant for estimating the house prices.

In contrast, interpreting global spillover effects can be a bit more difficult. Intuitively, the global spillover effects can be seen as a kind of diffusion process. For example, an exogenous event might increase the house prices in one district of a city, thus leading to an adaptation of house prices in neighbouring districts, which then leads to further adaptations in other units (the neighbours of the neighbours), thereby globally diffusing the effect of the exogenous event due to the endogenous term.

Yet, those processes happen over time. In a cross-sectional framework, the global spillover effects are hard to interpret. Anselin (2003) proposes an interpretation as an equilibrium outcome, where the partial impact represents an estimate of how this long-run equilibrium would change due to a change in  $\mathbf{x}_k$  (LeSage 2014a).

### 9.3 Summary impact measures

Note that the derivative in SAR, SAC, and SDM is a  $N \times N$  matrix, returning individual effects of each unit on each other unit, differentiated in *direct, indirect, and total impacts*.

$$(\mathbf{I}_N - \rho \mathbf{W})^{-1} \boldsymbol{\beta} = \begin{pmatrix} 1.1875 & 0.46875 & 0.1875 & 0.46875 & 0.1875 \\ 0.3125 & 1.28125 & 0.3125 & 0.28125 & 0.3125 \\ 0.1875 & 0.46875 & 1.1875 & 0.46875 & 0.1875 \\ 0.3125 & 0.28125 & 0.3125 & 1.28125 & 0.3125 \\ 0.1875 & 0.46875 & 0.1875 & 0.46875 & 1.1875 \end{pmatrix} \boldsymbol{\beta}$$

However, the individual effects (how  $i$  influences  $j$ ) mainly vary because of variation in  $\mathbf{W}$ .

 Do not interpret these as “estimated” individual impacts

We estimate two scalar parameters in a SAR model:  $\beta$  for the direct coefficient and  $\rho$  for the auto-regressive parameter.

All variation in the effects matrix  $(\mathbf{I}_N - \rho \mathbf{W})^{-1}$  comes from the relationship in  $\mathbf{W}$  which we have given a-priori!

Since reporting the individual partial effects is usually not of interest, LeSage and Pace (2009) proposed to average over these effect matrices. While the average diagonal elements of the effects matrix  $(\mathbf{I}_N - \rho \mathbf{W})^{-1}$  represent the so called direct impacts of variable  $x_k$ , the average column-sums of the off-diagonal elements represent the so called indirect impacts (or spatial spillover effects).

direct impacts refer to an average effect of a unit change in  $x_i$  on  $y_i$ , and the indirect (spillover) impacts indicate how a change in  $x_i$ , on average, influences all neighbouring units  $y_j$ .

### 9.3 Summary impact measures

Though previous literature (Halleck Vega and Elhorst 2015; LeSage and Pace 2009) has established the notation of direct and indirect impacts, it is important to note that also the direct impacts comprise a spatial ‘multiplier’ component if we specify an endogenous lagged depended variable, as a change in  $\mathbf{x}_i$  influences  $\mathbf{y}_i$ , which influences  $\mathbf{y}_j$ , which in turn influences  $\mathbf{y}_i$ .

Usually, one should use summary measures to report effects in spatial models (LeSage and Pace 2009). Halleck Vega and Elhorst (2015) provide a nice summary of the impacts for each model:

Model	Direct Impacts	Indirect Impacts	type
OLS/SEM	$\beta_k$	–	–
SAR/SAC	Diagonal elements of $(\mathbf{I} - \rho \mathbf{W})^{-1}$	Off-diagonal elements of $(\mathbf{I} - \rho \mathbf{W})^{-1}$	global
SLX/SDEM	$\beta_k$	$\theta_k$	local
SDM	Diagonal elements of $(\mathbf{I} - \rho \mathbf{W})^{-1} [\beta_k + \mathbf{W}\theta_k]$	Off-diagonal elements of $(\mathbf{I} - \rho \mathbf{W})^{-1} [\beta_k + \mathbf{W}\theta_k]$	global

$$(\mathbf{I}_N - \rho \mathbf{W})^{-1} \beta = \begin{pmatrix} 1.1875 & 0.46875 & 0.1875 & 0.46875 & 0.1875 \\ 0.3125 & 1.28125 & 0.3125 & 0.28125 & 0.3125 \\ 0.1875 & 0.46875 & 1.1875 & 0.46875 & 0.1875 \\ 0.3125 & 0.28125 & 0.3125 & 1.28125 & 0.3125 \\ 0.1875 & 0.46875 & 0.1875 & 0.46875 & 1.1875 \end{pmatrix} \beta$$

The different indirect effects / spatial effects mean conceptually different things:

- Global spillover effects: SAR, SAC, SDM

## 9 Spatial Impacts

- Local spillover effects: SLX, SDEM

**⚠️** Common ratio between direct and indirect impacts in SAR and SAC

Note that impacts in SAR only estimate one single spatial multiplier coefficient. Thus direct and indirect impacts are bound to a common ratio, say  $\phi$ , across all covariates.

if  $\beta_1^{direct} = \phi\beta_1^{indirect}$ , then  $\beta_2^{direct} = \phi\beta_2^{indirect}$ ,  $\beta_k^{direct} = \phi\beta_k^{indirect}$ .

We can calculate these impacts using `impacts()` with simulated distributions, e.g. for the SAR model:

```
mod_1.sar.imp <- impacts(mod_1.sar, listw = queens.lw, R = 300)
summary(mod_1.sar.imp, zstats = TRUE, short = TRUE)
```

```
Impact measures (lag, exact):
      Direct    Indirect      Total
log(no2)   0.447853184  0.754466618  1.202319802
log(POPDEN) -0.062973027 -0.106086209 -0.169059236
per_mixed   0.020884672  0.035182931  0.056067603
per_asian   -0.002575602 -0.004338934 -0.006914536
per_black   -0.014253206 -0.024011369 -0.038264575
per_other   -0.001820705 -0.003067212 -0.004887917
=====
Simulation results ( variance matrix):
=====
Simulated standard errors
      Direct    Indirect      Total
log(no2)   0.0488746577 0.0992947679 0.135731528
log(POPDEN) 0.0138348654 0.0272604116 0.040219045
per_mixed   0.0062340535 0.0114284286 0.017423600
per_asian   0.0005028114 0.0008278141 0.001282866
```

### 9.3 Summary impact measures

```
per_black  0.0009584477 0.0022958909 0.002660797  
per_other   0.0032087463 0.0055369935 0.008732728
```

Simulated z-values:

	Direct	Indirect	Total
log(no2)	9.1699338	7.6050769	8.8654548
log(POPDEN)	-4.6050047	-3.9594280	-4.2677606
per_mixed	3.3619039	3.1048256	3.2393746
per_asian	-5.1308747	-5.2313286	-5.3867117
per_black	-14.9323624	-10.4988735	-14.4378407
per_other	-0.5277511	-0.5269264	-0.5280146

Simulated p-values:

	Direct	Indirect	Total
log(no2)	< 2.22e-16	2.8422e-14	< 2.22e-16
log(POPDEN)	4.1246e-06	7.5129e-05	1.9745e-05
per_mixed	0.00077407	0.0019039	0.0011979
per_asian	2.8840e-07	1.6830e-07	7.1758e-08
per_black	< 2.22e-16	< 2.22e-16	< 2.22e-16
per_other	0.59767208	0.5982447	0.5974892

```
# Alternative with traces (better for large W)  
W <- as(queens.lw, "CsparseMatrix")  
trMatc <- trW(W, type = "mult",  
                 m = 30) # number of powers  
mod_1.sar.imp2 <- impacts(mod_1.sar,  
                           tr = trMatc, # trace instead of listw  
                           R = 300,  
                           Q = 30) # number of power series used for approximation  
summary(mod_1.sar.imp2, zstats = TRUE, short = TRUE)
```

Impact measures (lag, trace):

## 9 Spatial Impacts

	Direct	Indirect	Total
log(no2)	0.447853101	0.754459497	1.202312598
log(POPDEN)	-0.062973015	-0.106085208	-0.169058223
per_mixed	0.020884668	0.035182599	0.056067267
per_asian	-0.002575601	-0.004338893	-0.006914494
per_black	-0.014253203	-0.024011142	-0.038264346
per_other	-0.001820704	-0.003067183	-0.004887888
<hr/>			
Simulation results ( variance matrix):			
<hr/>			
Simulated standard errors			
	Direct	Indirect	Total
log(no2)	0.047110308	0.0924219350	0.127555773
log(POPDEN)	0.014335172	0.0268662539	0.040458337
per_mixed	0.006676480	0.0117103841	0.018188023
per_asian	0.000509860	0.0008152945	0.001284136
per_black	0.001050298	0.0022327528	0.002762871
per_other	0.003306380	0.0056900174	0.008981326
 Simulated z-values:			
	Direct	Indirect	Total
log(no2)	9.5292084	8.1928725	9.4556683
log(POPDEN)	-4.4094068	-3.9827415	-4.2070673
per_mixed	3.1382280	3.0241420	3.0990824
per_asian	-5.0011138	-5.2515590	-5.3198684
per_black	-13.5952749	-10.7839159	-13.8829882
per_other	-0.5551403	-0.5562902	-0.5568004
 Simulated p-values:			
	Direct	Indirect	Total
log(no2)	$< 2.22e-16$	2.2204e-16	$< 2.22e-16$
log(POPDEN)	1.0365e-05	6.8125e-05	2.5871e-05
per_mixed	0.0016997	0.0024934	0.0019412
per_asian	5.7000e-07	1.5082e-07	1.0384e-07

### 9.3 Summary impact measures

```
per_black < 2.22e-16 < 2.22e-16 < 2.22e-16  
per_other  0.5787987  0.5780125  0.5776638
```

The indirect effects in SAR, SAC, and SDM refer to global spillover effects. This means a change of  $x$  in the focal units flows through the entire system of neighbours (direct neighbours, neighbours of neighbours, ...) influencing ‘their  $y$ ’. One can think of this as diffusion or a change in a long-term equilibrium.

*If Log NO<sub>2</sub> increases by one unit, this increases the house price in the focal unit by 0.448 units. Overall, a one unit change in log NO<sub>2</sub> increases the house prices in the entire neighbourhood system (direct and higher order neighbours) by 0.754.*

For SLX models, nothing is gained from computing the impacts, as they equal the coefficients. Again, it’s the effects of direct neighbours only.

```
print(impacts(mod_1.slx, listw = queens.lw))
```

```
Impact measures (SLX, glht):  
          Direct    Indirect      Total  
log(no2)   -0.440727458  0.993602103  0.552874645  
log(POPDEN) -0.076839828  0.113262218  0.036422390  
per_mixed   -0.033042221  0.126068686  0.093026466  
per_asian   -0.002380698 -0.003828126 -0.006208824  
per_black   -0.016229407 -0.018053503 -0.034282910  
per_other   -0.020391354  0.048139008  0.027747654
```

## 9 Spatial Impacts

### 9.4 Examples

#### Boillat, Ceddia, and Bottazzi (2022)

*The paper investigates the effects of protected areas and various land tenure regimes on deforestation and possible spillover effects in Bolivia, a global tropical deforestation hotspot.*

Table 3  
SDM Models (1–5) estimates of ADE, AIE and ATE expressed in terms of elasticity.

	SDM1			SDM2			SDM3		
	ADE	AIE	ATE	ADE	AIE	ATE	ADE	AIE	ATE
TRAVELTIME	-0.15** (0.06)	0.11 (0.19)	-0.04 (0.22)	-0.12* (0.06)	0.18 (0.20)	0.06 (0.22)	-0.01 (0.35)	0.40 (0.40)	0.40** (0.19)
SLOPE	-0.29*** (0.06)	-0.31*** (0.08)	-0.60*** (0.09)	-0.20*** (0.06)	-0.32*** (0.08)	-0.60*** (0.10)	-0.29 (0.26)	-0.32 (0.30)	-0.61*** (0.10)
POPDENS12	-0.004 (0.005)	0.09** (0.04)	0.09** (0.04)	0.00 (0.01)	0.09*** (0.04)	0.09** (0.03)	0.00 (0.08)	0.09 (0.09)	0.10*** (0.03)
WATERACCESS12	0.44*** (0.15)	1.59*** (0.35)	2.03*** (0.35)	0.41** (0.15)	1.44*** (0.35)	1.05*** (0.35)	0.44 (1.04)	1.19 (1.16)	1.63*** (0.32)
POVERTY12	-0.15* (0.08)	-1.08*** (0.20)	-1.22*** (0.22)	-0.13 (0.08)	-1.06*** (0.20)	-1.19*** (0.22)	-0.16 (0.67)	-0.78 (0.73)	-0.93*** (0.21)
PROTECTEDAREA	-0.13*** (0.03)	-0.21*** (0.05)	-0.34*** (0.07)	-0.12*** (0.03)	-0.20*** (0.05)	-0.32*** (0.07)	-0.12 (0.14)	-0.16 (0.14)	-0.23*** (0.06)
NOTITLELAND13	-0.11* (0.06)	0.85*** (0.18)	0.74*** (0.21)	0.17* (0.10)	0.95*** (0.38)	1.12*** (0.38)	0.21 (0.25)	0.25 (0.36)	0.47 (0.40)
INDIGENLAND13	-0.15*** (0.04)	0.22* (0.12)	0.06 (0.13)	-0.01 (0.05)	0.26 (0.19)	0.25 (0.20)	-0.01 (0.05)	-0.01 (0.10)	-0.01 (0.20)
COMMONLAND13	0.01 (0.03)	0.19** (0.09)	0.21** (0.10)	0.12*** (0.04)	0.23 (0.16)	0.35** (0.17)	0.12** (0.05)	0.04 (0.15)	0.15 (0.17)
STATELAND13	-0.20*** (0.06)	0.28* (0.15)	0.08 (0.17)	0.04 (0.09)	0.38 (0.20)	0.42 (0.30)	0.05 (0.13)	-0.12 (0.35)	-0.07 (0.34)
PRIVATELAND13				0.22*** (0.06)	0.11 (0.20)	0.33 (0.22)			
SMALLPROP13							0.24*** (0.03)	0.01 (0.06)	0.26*** (0.07)
MEDIUMPROP13							-0.06 (0.11)	-0.13 (0.16)	-0.19** (0.09)
BUSINESSLAND13							0.03 (0.03)	-0.02 (0.14)	0.01 (0.16)

Standard errors in parentheses.

\*\*\* p < 0.01, \*\* p < 0.05, \* p < 0.1.

*Protected areas – which in Bolivia are all based on co-management schemes - also protect forests in adjacent areas, showing an indirect protective spillover effect. Indigenous lands however only have direct forest protection effects.*

## 9.4 Examples

### Fischer et al. (2009)

*The focus of this paper is on the role of human capital in explaining labor productivity variation among 198 European regions within a regression framework.*

The impact of human capital on regional labor productivity in Europe 105

**Table 2** Direct, indirect and total impact estimates (*t*-statistics in parentheses)

Variables	Spatial Durbin model		
	Mean	Mean	Mean
	direct impact	indirect impact	total impact
Initial labor productivity	0.6677 (27.5716)	0.0683 (1.8992)	0.7361 (26.3921)
Human capital	0.1317 (6.8644)	-0.1968 (-3.7637)	-0.0650 (-1.1847)

*Note:* *t*-statistics based on 10,000 sampled raw parameter estimates of the SDM

*A ceteris paribus increase in the level of human capital is found to have a significant and positive direct impact. But this positive direct impact is offset by a significant and negative indirect (spillover) impact leading to a total impact that is not significantly different from zero.*

*The intuition here arises from the notion that it is relative regional advantages in human capital that matter most for labor productivity, so changing human capital across all regions should have little or no total impact on (average) labor productivity levels.*

### Rüttenauer (2018)

*This study investigates the presence of environmental inequality in Germany - the connection between the presence of foreign-minority population and objectively measured industrial pollution.*

## 9 Spatial Impacts

**Table 3**  
Community-fixed effects estimates. Dependent variable: Industrial air pollution.

	Overall	Urban	Rural	Diff
	M5	M6	M7	(M6-M7)
% Foreigners	0.062*** (0.009)	0.076*** (0.016)	0.035*** (0.007)	0.041* (0.017)
W [% Foreigners]	0.163*** (0.025)	0.391*** (0.070)	0.077*** (0.015)	0.315*** (0.071)
Population	-0.010 (0.007)	-0.011 (0.009)	-0.008 (0.006)	-0.002 (0.011)
W [Population]	-0.039 (0.023)	-0.124*** (0.030)	-0.023 (0.017)	-0.100** (0.034)
% 65 and older	-0.005 (0.004)	-0.001 (0.017)	-0.000 (0.003)	-0.000 (0.017)
W [% 65 and older]	-0.009 (0.007)	-0.101 (0.072)	0.004 (0.005)	-0.105 (0.072)
% Vacant housing	0.014** (0.005)	0.035 (0.027)	0.012** (0.004)	0.023 (0.027)
W [% Vacant housing]	0.021* (0.008)	0.087 (0.092)	0.020** (0.007)	0.068 (0.093)
Living space	-0.012** (0.004)	0.014 (0.022)	-0.008* (0.003)	0.022 (0.022)
W [Living space]	-0.036*** (0.007)	-0.161* (0.070)	-0.022*** (0.006)	-0.139* (0.070)
R <sup>2</sup>	0.018	0.076	0.005	
Adj. R <sup>2</sup>	-0.031	0.067	-0.050	
Num. obs.	93777	9061	84716	

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$ . Standardized coefficients. Cluster robust standard errors in parentheses. W is specified as contiguity neighbours weights matrix. The differences and standard errors in column 4 are obtained by an overall fixed-effects model with interaction terms of the urban dummy and all covariates.

*Results reveal that the share of minorities within a census cell indeed positively correlates with the exposure to industrial pollution. Furthermore, spatial spillover effects are highly relevant: the characteristics of the neighbouring spatial units matter in predicting the amount of pollution. Especially within urban areas, clusters of high minority neighbourhoods are affected by high levels of environmental pollution.*

# 10 Comparing and Selecting Models

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite") # note: load spd  
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()
```

```
R version 4.4.1 (2024-06-14 ucrt)  
Platform: x86_64-w64-mingw32/x64  
Running under: Windows 11 x64 (build 22631)
```

```
Matrix products: default
```

```
locale:  
[1] LC_COLLATE=English_United Kingdom.utf8  
[2] LC_CTYPE=English_United Kingdom.utf8  
[3] LC_MONETARY=English_United Kingdom.utf8  
[4] LC_NUMERIC=C  
[5] LC_TIME=English_United Kingdom.utf8
```

## 10 Comparing and Selecting Models

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets   methods    base

other attached packages:
[1] viridisLite_0.4.2 tmap_3.3-4          spatialreg_1.3-
4 Matrix_1.7-0
[5] spdep_1.3-5       spData_2.3.1        mapview_2.11.2     sf_1.0-
16

loaded via a namespace (and not attached):
 [1] xfun_0.45           raster_3.6-26        htmlwidgets_1.6.4  lattice_0.22-
6
 [5] tools_4.4.1          crosstalk_1.2.1     LearnBayes_2.15.1 parallel_4.4.1
 [9] stats4_4.4.1         sandwich_3.1-0      proxy_0.4-27      KernSmooth_2.2
24
[13] satellite_1.0.5     RColorBrewer_1.1-3  leaflet_2.2.2    lifecycle_1.0.0
[17] compiler_4.4.1      deldir_2.0-4       munsell_0.5.1    terra_1.7-
78
[21] codetools_0.2-20    leafsync_0.1.0     stars_0.6-5     htmltools_0.5.8
[25] class_7.3-22        MASS_7.3-60.2     classInt_0.4-
10 lwgeom_0.2-14
[29] wk_0.9.1            abind_1.4-5       boot_1.3-30     multcomp_1.4-
25
[33] nlme_3.1-164         digest_0.6.35    mvtnorm_1.2-
5 splines_4.4.1
[37] fastmap_1.2.0        grid_4.4.1       colorspace_2.1-
0 cli_3.6.2
[41] magrittr_2.0.3       base64enc_0.1-3  dichromat_2.0-
0.1 XML_3.99-0.16.1
[45] survival_3.6-4      leafem_0.2.3     TH.data_1.1-
2 e1071_1.7-14
```

## 10.1 Specific-to-general

```
[49] scales_1.3.0      sp_2.1-4        rmarkdown_2.27    zoo_1.8-  
12  
[53] png_0.1-8       coda_0.19-4.1   evaluate_0.24.0  knitr_1.47  
[57] tmaptools_3.1-1 s2_1.1.6       rlang_1.1.4      Rcpp_1.0.12  
[61] glue_1.7.0       DBI_1.2.3      rstudioapi_0.16.0 jsonlite_1.8.8  
[65] R6_2.5.1        units_0.8-5
```

### Reload data from previous session

```
load("_data/msoa2_spatial.RData")
```

As we have seen, a variety of spatial model specifications exist that can be used to account for the spatial structure of the data. Thus, selecting the correct model specification remains a crucial task in applied research. For some helpful literature on the issue of model selection see Anselin, Serenini, and Amaral (2024), Halleck Vega and Elhorst (2015), Rüttenauer (2022)

One way of selecting the model specification is the application of empirical specification tests. In general, there are two different strategies: a specific-to-general or a general-to-specific approach (Florax, Folmer, and Rey 2003; Mur and Angulo 2009).

## 10.1 Specific-to-general

The specific-to-general approach is more common in spatial econometrics. This approach starts with the most basic non-spatial model and tests for possible misspecification due to omitted autocorrelation in the error term or the dependent variable.

Anselin et al. (1996) proposed to use Lagrange multiplier (LM) tests for the hypotheses  $H_0: \lambda = 0$  and  $H_0: \rho = 0$ , which are robust against the alternative source of spatial dependence.

### 10.1.1 Lagrange Multiplier Test

We have earlier talked about methods to detect auto-correlation – visualisation and Moran's I. Both methods can tell us that there is spatial autocorrelation. However, both method do not provide any information on why there is autocorrelation. Possible reasons:

- Interdependence ( $\rho$ )
- Clustering on unobservables ( $\lambda$ )
- Spillovers in covariates ()

Lagrange Multiplier test (Anselin et al. 1996):

- (Robust) test for spatial lag dependence  $LM_{\rho}^*$
- (Robust) test for spatial error dependence  $LM_{\lambda}^*$

Robust test for lag dependence:  $H_0: \rho = 0$

$$LM_{\rho}^* = G^{-1} \hat{\sigma}_{\epsilon}^2 \left( \frac{\mathbf{W}^T \mathbf{W} \mathbf{y}}{\hat{\sigma}_{\epsilon}^2} - \frac{\mathbf{W}^T \mathbf{W}}{\hat{\sigma}_{\epsilon}^2} \right)^2 \sim \chi^2$$

Robust test for error dependence:  $H_0: \lambda = 0$

$$LM_{\lambda}^* = \frac{\left( \mathbf{W}^T \mathbf{W} / \hat{\sigma}_{\epsilon}^2 - [T \hat{\sigma}_{\epsilon}^2 (G + T \hat{\sigma}_{\epsilon}^2)^{-1}]^T \mathbf{W} \mathbf{y} / \hat{\sigma}_{\epsilon}^2 \right)^2}{T \left[ 1 - \frac{\hat{\sigma}_{\epsilon}^2}{G + \hat{\sigma}_{\epsilon}^2} \right]} \sim \chi^2$$

with

$$\begin{aligned} G &= (\mathbf{W} \mathbf{X})^T (\mathbf{I} - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) (\mathbf{W} \mathbf{X}) \\ T &= \text{tr}[(\mathbf{W}^T + \mathbf{W}) \mathbf{W}], \end{aligned}$$

where  $\text{tr}(\mathbf{A})$  is the sum of the main diagonal of any square matrix  $\mathbf{A}$ .

To perform the test, we first need to run a conventional OLS model. We take yesterdays example.

## 10.1 Specific-to-general

```
### Load data

load("_data/msoa2_spatial.RData")

### Generate weights
coords <- st_centroid(msoa.spdf)

# Neighbours within 3km distance
dist_15.nb <- dnearneigh(coords, d1 = 0, d2 = 2500)

summary(dist_15.nb)

# There are some empty neighbour sets. Lets impute those with the nearest neighbour.
k2.nb <- knearneigh(coords, k = 1)

# Replace zero
nolink_ids <- which(card(dist_15.nb) == 0)
dist_15.nb[card(dist_15.nb) == 0] <- k2.nb$nn[nolink_ids, ]

summary(dist_15.nb)

# listw object with row-normalization
dist_15.lw <- nb2listw(dist_15.nb, style = "W")

mod_1.lm <- lm(log(no2) ~ per_mixed + per_asian + per_black + per_other +
                 + per_nonUK_EU + per_nonEU + log(POPDEN),
                 data = msoa.spdf)
summary(mod_1.lm)
```

Call:

```
lm(formula = log(no2) ~ per_mixed + per_asian + per_black + per_other +
```

## 10 Comparing and Selecting Models

```
per_nonUK_EU + per_nonEU + log(POPDEN), data = msoa.spdf)

Residuals:
    Min      1Q  Median      3Q     Max 
-0.46035 -0.09063 -0.01010  0.07484  0.90401 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 2.5188244  0.0313089 80.451 < 2e-16 ***
per_mixed   0.0050636  0.0044498  1.138  0.2554    
per_asian    0.0006251  0.0003459  1.807  0.0711 .  
per_black   -0.0009929  0.0007100 -1.398  0.1623    
per_other    0.0141322  0.0022477  6.288 4.86e-10 ***
per_nonUK_EU 0.0132176  0.0011520 11.473 < 2e-16 ***
per_nonEU    0.0026836  0.0010643  2.522  0.0118 *  
log(POPDEN)  0.1391236  0.0079044 17.601 < 2e-16 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1373 on 975 degrees of freedom
Multiple R-squared:  0.6162,    Adjusted R-squared:  0.6134 
F-statistic: 223.6 on 7 and 975 DF,  p-value: < 2.2e-16
```

Subsequently, we run the robust LM test for Lag dependence:

```
test.lag <- lm.RStests(mod_1.lm,
                        listw = dist_15.lw,
                        test = "adjRSlag")
summary(test.lag)

Rao's score (a.k.a Lagrange multiplier) diagnostics for spatial
dependence
data:
```

## 10.1 Specific-to-general

```
model: lm(formula = log(no2) ~ per_mixed + per_asian + per_black +
per_other + per_nonUK_EU + per_nonEU + log(POPDEN), data = msoa.spdf)
test weights: dist_15.lw

      statistic parameter   p.value
adjRSlag     365.69          1 < 2.2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

For the robust LM test on spatial error dependence, we run:

```
test.err <- lm.RStests(mod_1.lm,
                      listw = dist_15.lw,
                      test = "adjRSerr")
summary(test.err)

Rao's score (a.k.a Lagrange multiplier) diagnostics for spatial
dependence
data:
model: lm(formula = log(no2) ~ per_mixed + per_asian + per_black +
per_other + per_nonUK_EU + per_nonEU + log(POPDEN), data = msoa.spdf)
test weights: dist_15.lw

      statistic parameter   p.value
adjRSerr    324.32          1 < 2.2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

In fact, you can compare them all at once:

```
test.all <- lm.RStests(mod_1.lm,
                      listw = dist_15.lw,
```

## 10 Comparing and Selecting Models

```
test = "all")
summary(test.all)

Rao's score (a.k.a Lagrange multiplier) diagnostics for spatial
dependence
data:
model: lm(formula = log(no2) ~ per_mixed + per_asian + per_black +
per_other + per_nonUK_EU + per_nonEU + log(POPDEN), data = msoa.spdf)
test weights: dist_15.lw

      statistic parameter   p.value
RSerr      1850.70      1 < 2.2e-16 ***
RSlag      1892.06      1 < 2.2e-16 ***
adjRSerr    324.32      1 < 2.2e-16 ***
adjRSlag    365.69      1 < 2.2e-16 ***
SARMA     2216.38      2 < 2.2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

In both cases, the LM test indicates the presence of significance lag dependence and significance error dependence.

 Take the LM test with a grain of salt

**DO NOT** use the above results as a reason to say: “well, there is Lag dependence and there is Error dependence. So, I will use a SAC-like specification”. This would usually be a bad choice, as the LM tests assume the absence od any SLX terms.

### 10.1.2 Problem

The specific-to-general approach based on the robust LM test offers a good performance in distinguishing between SAR, SEM, and non-spatial OLS (Florax, Folmer, and Rey 2003).

Still, in their original paper, Anselin et al. (1996) already note the declining power of the robust  $\text{LM}_\lambda$  test for spatial error dependence with increasing autocorrelation in the dependent variable (indicating some uncertainty under a SAC-like DGP).

Mur and Angulo (2009) demonstrate strong drawbacks of the specific-to-general approach under non-optimal conditions like heteroscedasticity or endogeneity.

Moreover, the test disregard the presence of spatial dependence from local spillover effects ( $\theta$  is assumed to be zero), as resulting from an SLX-like process. Cook, Hays, and Franzese (2020), for instance, show theoretically that an SLX-like dependence structure leads to the rejection of both hypotheses  $H_0: \lambda = 0$  and  $H_0: \rho = 0$ , though no autocorrelation is present (Elhorst and Halleck Vega 2017; Rüttenauer 2022).

#### 💡 Potential solution?

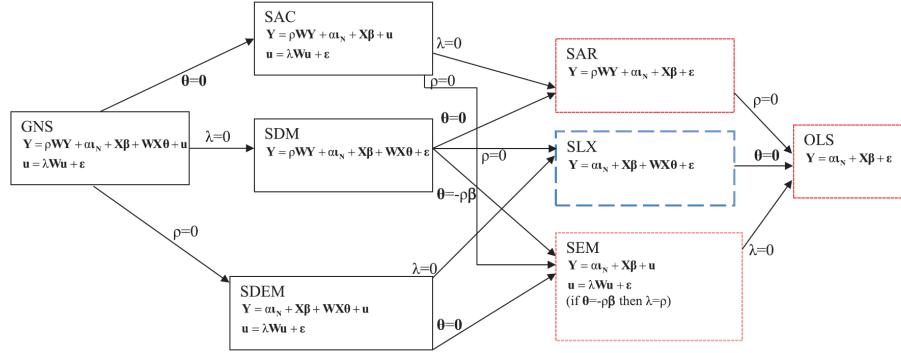
In a recent preprint, Anselin, Serenini, and Amaral (2024) proposed a two-step test approach (“STGE-Pre”). First, we need to test the presence of any SLX (**WX**) term. If we can omit those terms, we proceed with the robust LM lag test and the robust LM error test.

## 10.2 General-to-specific approach

The general-to-specific approach depicts the opposite method of specification search. This approach starts with the most general model and stepwise

## 10 Comparing and Selecting Models

imposes restrictions on the parameters of this general model.



Note: GNS = general nesting spatial model, SAC = spatial autoregressive combined model, SDM = spatial Durbin model, SDEM = spatial Durbin error model, SAR = spatial autoregressive model, SLX = spatial lag of  $\mathbf{X}$  model, SEM = spatial error model, OLS = ordinary least squares model.

Figure 10.1: Halleck Vega and Elhorst (2015): Nesting of different Spatial Econometric Model Specifications

In theory, we would

- 1) start with a GNS specification and
- 2) subsequently restrict the model to simplified specifications based on the significance of parameters in the GNS.

The problem with this strategy is that the GNS is only weakly identified and, thus, is of little help in selecting the correct restrictions (Burridge, Elhorst, and Zigova 2016).

The most intuitive alternative would be to start with one of the two-source models SDM, SDEM, or SAC. This, however, bears the risk of imposing the wrong restriction in the first place (Cook, Hays, and Franzese 2020). Furthermore, Cook, Hays, and Franzese (2020) show that more complicated restrictions are necessary to derive all single-source models from SDEM or SAC specifications. Anselin, Serenini, and Amaral (2024) argue that the Ge

### 10.3 General advice?

LeSage and Pace (2009), LeSage (2014a), Elhorst (2014) argue that there are strong analytical reasons to restrict the model specifications to a subset, as the SDM subsumes the SLX and SAR model, and the SDEM subsumes SLX and SEM.

It is easily observed that SDM reduces to SLX if  $\rho = 0$  and to SAR if  $= 0$ , while the SDEM reduces to SLX if  $\lambda = 0$  and to SEM if  $= 0$ . Less intuitively, (Anselin 1988) has also shown that the SDM subsumes the SEM. Therefore, we can express the reduced form and rearrange terms:

$$\begin{aligned}\mathbf{y} &= \mathbf{X} + (\mathbf{I}_N - \lambda \mathbf{W})^{-1} \\ (\mathbf{I}_N - \lambda \mathbf{W})\mathbf{y} &= (\mathbf{I}_N - \lambda \mathbf{W})\mathbf{X} + \\ (\mathbf{I}_N - \lambda \mathbf{W})\mathbf{y} &= \mathbf{X} - \lambda \mathbf{WX} + \\ \mathbf{y} &= (\mathbf{I}_N - \lambda \mathbf{W})^{-1}(\mathbf{X} + \mathbf{WX} + ).\end{aligned}$$

Thus, the SEM constitutes a special case of an SDM with the relative simple restriction  $= -\lambda$ , meaning direct and indirect effects are constrained to a common factor (Anselin 1988, 2003).

The fact that SDM subsumes SAR, SLX, and SEM leads to the conclusion that applied research should only consider SDM and SDEM as model specifications (LeSage 2014a). Especially in the case of a likely omitted variable bias, (LeSage and Pace 2009, ~68) argue in favour of using the SDM.

Nonetheless, others propose to use the SLX specification as point of departure (Gibbons and Overman 2012; Halleck Vega and Elhorst 2015). First, scholars have argued that SAC and SDM models are only weakly identified in practice (Gibbons and Overman 2012; Pinkse and Slade 2010). Second, the global spillover specification in SAR, SAC, and SDM often seems to be theoretically implausible.

And finally:

## 10.4 Design and Theory

Some argue that the best way of choosing the appropriate model specification is to exclude one or more sources of spatial dependence – autocorrelation in the dependent variable, autocorrelation in the disturbances, or spatial spillover effects of the covariates – by design Gibbons, Overman, and Patacchini (2015).

**Natural experiments** are probably the best way of making one or more sources of spatial dependence unlikely, thereby restricting the model alternatives to a subset of all available models. However, the opportunities to use natural experiments are restricted in social sciences, making it a favourable but often impractical way of model selection.

Cook, Hays, and Franzese (2020) and Rüttenauer (2022) argue that theoretical considerations should guide the model selection.

- 1) Rule out some sources of spatial dependence by theory, and thus restrict the specifications to a subset (*Where does the spatial dependence come from?*),
- 2) Theoretical mechanisms may guide the choice of either global or local spillover effects.

## 10.5 Monte Carlo simulation

This section discusses results from Rüttenauer (2022). The aim: how do different spatial models perform under different scenarios?

The DGP of the Monte Carlo simulation follows a GNS, where  $k$  and  $\epsilon$  are independent and randomly distributed  $\mathcal{N}(0, \sigma_v^2)$  and  $\mathcal{N}(0, \sigma_\epsilon^2)$  with

## 10.5 Monte Carlo simulation



Jeffrey Wooldridge  
@jmwooldridge

...

In 2018 I was invited to give a talk at SOCHER in Chile, to give my opinions about using spatial methods for policy analysis. I like the idea of putting in spatial lags of policy variables to measure spillovers. Use fixed effects with panel data, compute fully robust ses.

[Tweet übersetzen](#)



1:30 vorm. · 10. März 2021

48 Retweets 7 Zitate 316 „Gefällt mir“-Angaben 69 Lesezeichen



Twittiere deine Antwort!

[Antworten](#)



Jeffrey Wooldridge @jmwooldridge · 10. März 2021

...

For the life of me, I couldn't figure out how putting in spatial lags of Y had any value. After preparing a course in July 2020, I was even more negative about this practice. It seems an unnecessary complication developed by theorists.

Figure 10.2: “I will use spatial lags of X, not spatial lags of Y”, J. Wooldridge on twitter

## 10 Comparing and Selecting Models

a mean of zero, and  $\mathbf{x}_k$  is the  $k$ th column-vector of  $\mathbf{X}$  for  $k = 1, \dots, K$  covariates ( $K$  is fixed at 2 in the simulations). The parameter  $\rho$  represents the autocorrelation in the dependent variable,  $\lambda$  the autocorrelation in the disturbances, and  $\delta_k$  the autocorrelation in covariate  $k$ .

$$\begin{aligned}\mathbf{y} &= \rho \mathbf{W} \mathbf{y} + \mathbf{X} + \mathbf{W} \mathbf{X} + \mathbf{u}, \\ \mathbf{u} &= \lambda \mathbf{W} \mathbf{u} + \mathbf{X} + , \\ \mathbf{x}_k &= \delta_k \mathbf{W} \mathbf{x}_k + .\end{aligned}$$

The parameter-vector  $\beta$  specifies the correlation between  $\mathbf{x}$  and the disturbance vector  $\mathbf{u}$ , thereby defining the strength of an omitted variable bias. In reduced form, this DGP can be written as

$$\begin{aligned}\mathbf{y} &= (\mathbf{I}_N - \rho \mathbf{W})^{-1} [(\mathbf{I}_N - \delta_k \mathbf{W})^{-1} \beta_k \\ &\quad + \mathbf{W} (\mathbf{I}_N - \delta_k \mathbf{W})^{-1} \theta_k \\ &\quad + (\mathbf{I}_N - \lambda \mathbf{W})^{-1} ((\mathbf{I}_N - \delta_k \mathbf{W})^{-1} \gamma_k + )].\end{aligned}$$

The parameter vector  $\beta$  was fixed at  $= (0.2 \ 0.5)^\top$ , and the noise parameters were fixed at  $\sigma_v^2, \sigma_\varepsilon^2 = 1$  for all trials. All other parameters vary between the following two options for each parameter (vector):

- $\rho \in \{0, 0.5\}$ ,
- $\lambda \in \{0, 0.5\}$ ,
- $\delta \in \left\{ (0 \ 0)^\top, (0.4 \ 0.7)^\top \right\}$ ,
- $\theta \in \left\{ (0 \ 0)^\top, (0.1 \ 0.8)^\top \right\}$ ,
- $\gamma \in \left\{ (0 \ 0)^\top, (0.3 \ 0)^\top \right\}$ ,

leading to a total of 32 distinct combinations. Note that this selection of parameters intentionally violates the common ratio assumption between direct and indirect effects, as this should be a more common case in practical research. All combinations were simulated in 1000 trials, with the same starting seed for each combination. If you're interested in the simulations, see replication code on Github.

## 10.5 Monte Carlo simulation

### 10.5.1 Without omitted variable bias

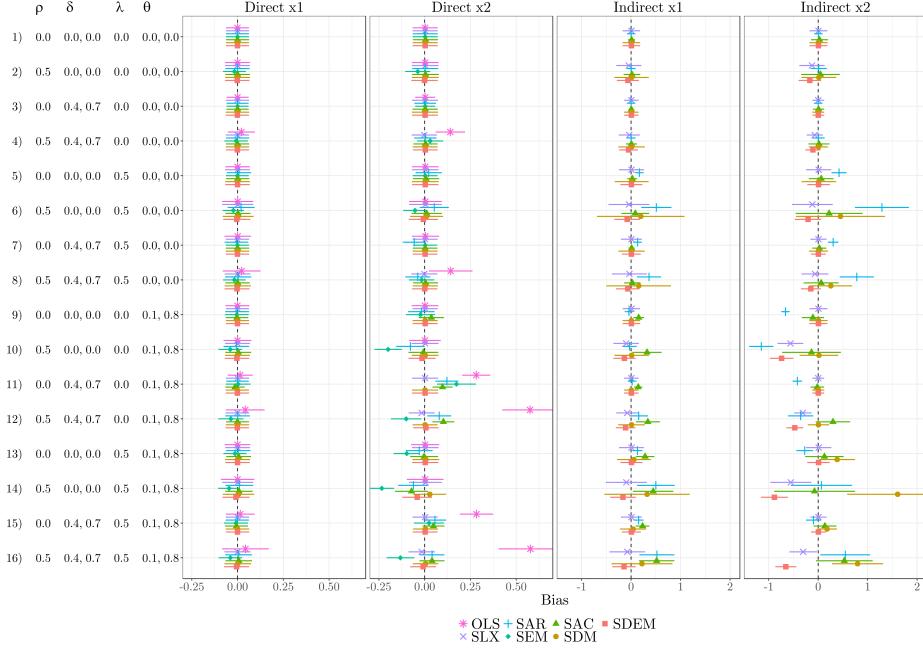


Figure 10.3: Bias of impacts and 95% confidence interval of empirical standard deviation without omv:  $\rho = (0.2, 0.5)^\top$ ,  $\delta = (0, 0)^\top$ .  $\rho$  = autocorrelation in the dependent variable ( $\mathbf{W}\mathbf{y}$ );  $\delta$  = autocorrelation in the covariates ( $\mathbf{x}_k = f(\mathbf{W}\mathbf{x}_k)$ );  $\lambda$  = autocorrelation in the disturbances ( $\mathbf{W}\mathbf{u}$ );  $\theta$  = spatial spillover effects of covariates ( $\mathbf{W}\mathbf{X}$ );  $\phi$  = strength of omv.

SLX, SDM, and SDEM all provide quite accurate estimates of the direct impacts (most visible in column 2). SAR, SEM, and SAC, in contrast, yield some drawbacks: especially in the presence of local spillover effects, these three specifications are biased (see lower part). Furthermore, SAR and SEM suffer from bias if autocorrelation in the disturbance and autocorrelation in the dependent variable are present simultaneously (see line 6

## 10 Comparing and Selecting Models

and 8). Though SLX is downwardly biased in case of autocorrelation in the dependent variable and the covariates (e.g. line 12 and 16), and SDM as well as SDEM yield some bias in case of a GNS-like process (line 14 and 16), those biases are rather moderate. This indicates that SLX, SDM, and SDEM are most robust against misspecification regarding the direct impacts.

Several differences exist regarding the indirect impacts. Most obviously, the often used SAR specification suffers from considerable bias: it overestimates indirect impacts in case of autocorrelation in the disturbances, and offers biased estimates if local spillover effects exist (which are not restricted to a common ratio). The latter also applies to SAC: though SAC offers relatively accurate estimates for  $\mathbf{x}_2$ , it overestimates indirect impacts for  $\mathbf{x}_1$ .

Regarding the remaining three specifications – SLX, SDM, and SDEM – conclusions are less obvious. SDM and SDEM suffer from large bias for high values of  $\theta_k$  (see  $\mathbf{x}_2$ ) if the DGP follows a GNS-like process (line 14 and 16): SDM overestimates the indirect impacts, while SDEM underestimates the indirect impacts. In addition, SDM performs badly if the true DGP is SDEM (line 13), and SDEM performs badly if the true DGP is SDM (line 10), whereas the bias increases with higher values of  $\theta_k$  in both cases. Similar to SDEM, SLX underestimates the indirect impacts in presence of global spillovers / autocorrelation in the dependent variable.

### 10.5.2 With omitted variable bias

### 10.5.3 Indirect impacts if DGP = GNS

Below an illustration about the indirect impacts, if the spatial process is a combination of

- 1) Clustering on Unobservables
- 2) Interdependence (in the outcome)

## 10.5 Monte Carlo simulation

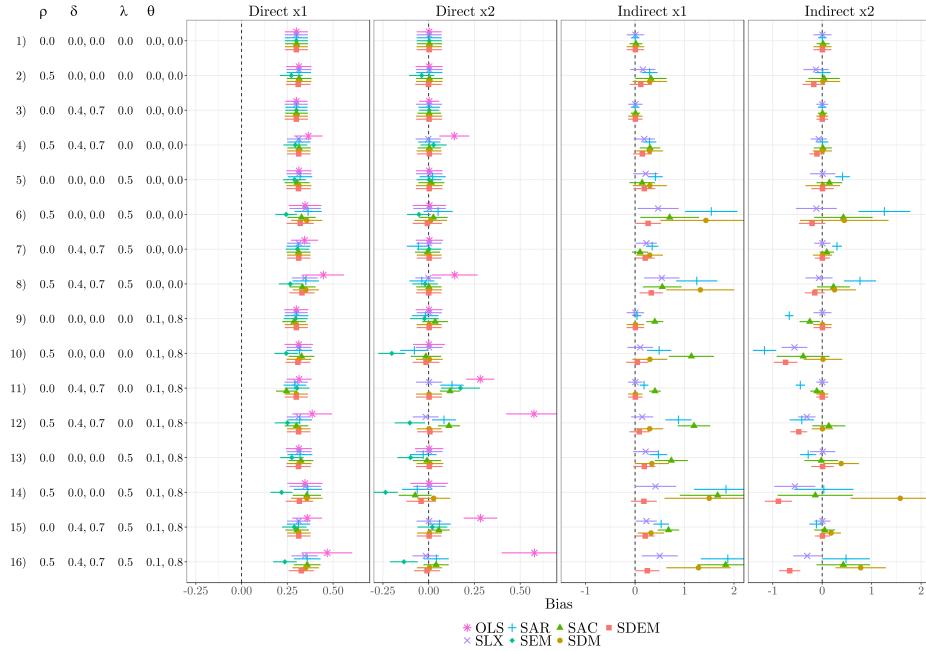


Figure 10.4: Bias of impacts and 95% confidence interval of empirical standard deviation with  $\text{omv} = (0.2, 0.5)^\top$ ,  $= (0.3, 0)^\top$ .  $\rho$  = autocorrelation in the dependent variable ( $\mathbf{W}\mathbf{y}$ ); = autocorrelation in the covariates ( $\mathbf{x}_k = f(\mathbf{W}\mathbf{x}_k)$ );  $\lambda$  = autocorrelation in the disturbances ( $\mathbf{W}\mathbf{u}$ ); = spatial spillover effects of covariates ( $\mathbf{W}\mathbf{X}$ ); = strength of  $\text{omv}$ .

## 10 Comparing and Selecting Models

### 3) Spillovers in Covariates

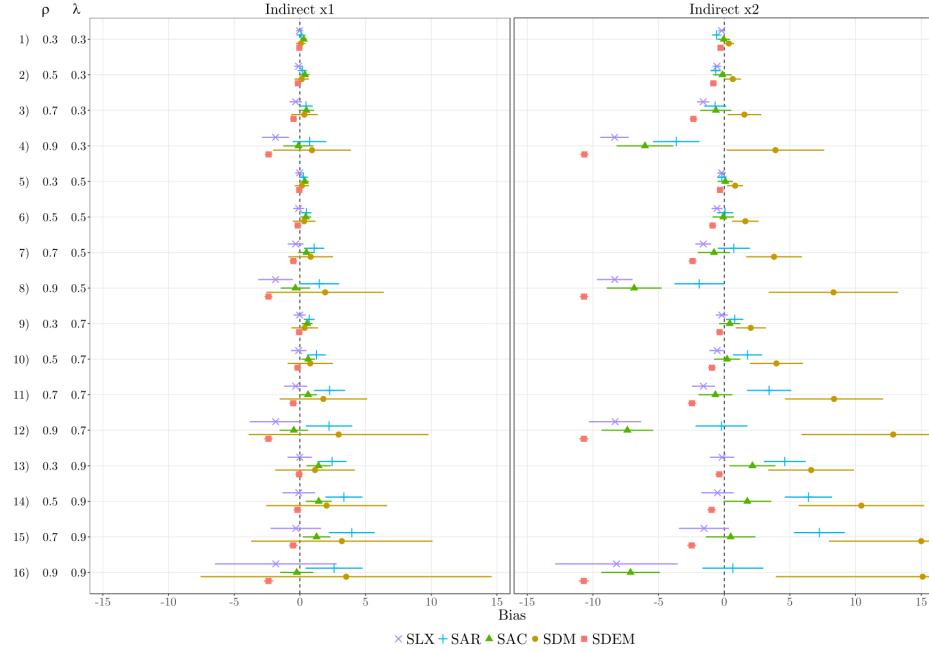


Figure 10.5: Bias of indirect impacts and 95% confidence interval of empirical standard deviation for different strengths of autocorrelation:  $\rho = (0.2, 0.5)^\top$ ,  $\lambda = (0, 0)^\top$ ,  $\lambda = (0, 0)^\top$ ,  $\lambda = (0.1, 0.8)^\top$ .  $\rho$  = autocorrelation in the dependent variable ( $\mathbf{W}\mathbf{y}$ );  $\lambda$  = autocorrelation in the covariates ( $\mathbf{x}_k = f(\mathbf{W}\mathbf{x}_k)$ );  $\lambda$  = autocorrelation in the disturbances ( $\mathbf{W}\mathbf{u}$ );  $\rho$  = spatial spillover effects of covariates ( $\mathbf{W}\mathbf{X}$ );  $\lambda$  = strength of omv.

First, in a GNS-like situation, the bias in SDM grows with increasing autocorrelation in  $\mathbf{y}$  ( $\rho$ ) and increasing autocorrelation in the disturbances ( $\lambda$ ).

Second, the bias in SLX and SDEM increases with higher values of  $\rho$ , but

## 10.6 Example: House prices in London

is unaffected from the strength of  $\lambda$ .

Third, though SLX and SDEM suffer from the same problem, the bias from omitting global autocorrelation is less severe in SLX than in SDEM.

Thus, the SLX outperforms SDEM. Furthermore, SLX outperforms SDM in most situations; only if the autocorrelation in the dependent variable is much stronger than the autocorrelation in the disturbances ( $\rho = 0.9$ ,  $\lambda = 0.3$ ), SDM yields lower bias than SLX. Note that the SAC yields relatively low biases for the indirect impacts in GNS-like processes, but at the same time produces relative large biases in the direct impacts.

## 10.6 Example: House prices in London

The example is taken from Rüttenauer (2024).

As an example to compare the different spatial model specifications, we estimate the effect of local characteristics such as green space and public transport connectivity on the median house price. The relation between environmental characteristics and housing choice and prices has been investigated in several studies (Anselin and Lozano-Gracia 2008; Kley and Dovbischuk 2021; Liebe, van Cranenburgh, and Chorus 2023). The data for the current example was retrieved from the London Datastore<sup>1</sup>, the 2011 Census<sup>2</sup> and OpenStreetMaps and combined at the Middle Layer Super Output Areas (MSOA). There are 983 MSOAs in London with an average population size of around 8,000 residents. The script for compiling and preparing the data can be found in the Supplementary Materials. All data preparation and analysis were performed with the statistical software R. For a comprehensive overview of spatial software see R. Bivand, Millo, and Piras (2021) or Pebesma and Bivand (2023).

---

<sup>1</sup>For house prices, see: <https://data.london.gov.uk/dataset/average-house-prices>. For London accessibility scores see: <https://data.london.gov.uk/dataset/public-transport-accessibility-levels>

<sup>2</sup>For UK demographics, see: [https://www.nomisweb.co.uk/sources/census\\_2011](https://www.nomisweb.co.uk/sources/census_2011)

## 10 Comparing and Selecting Models

```
# Load the packages
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "texreg", "extrafont",
        "ggplot2", "ggthemes", "rmapshaper", "viridis", "gridExtra") # n
lapply(pkgs, require, character.only = TRUE)

# Load the data
load("_data/msoa3_spatial.RData")

loadfonts()

# Create Contiguity (Queens) neighbours weights
queens.nb <- poly2nb(msoa.spdf,
                      queen = TRUE,
                      snap = 1) # we consider points in 1m distance as 'touching'
queens.lw <- nb2listw(queens.nb,
                      style = "W")

#### Plot house prices

# Get some larger scale boundaries
borough.spdf <- st_read(dsn = paste0("_data", "/statistical-gis-boundaries"),
                         layer = "London_Borough_Excluding_MHW" # Note: no file
                         )
# transform to only inner lines
borough_inner <- rmapshaper::ms_innerlines(borough.spdf)
borough_inner <- borough.spdf

# Plot with inner lines
msoa.spdf$med_house_price_ln <- log(msoa.spdf$med_house_price)
msoa.spdf$pt_access_index_ln <- log(msoa.spdf$pt_access_index)
```

## 10.6 Example: House prices in London

```
gp <- ggplot(msoa.spdf)+  
  geom_sf(aes(fill = med_house_price_ln))+  
  scale_fill_viridis_c(option = "A")  
  geom_sf(data = borough_inner, color = "gray92", fill = NA)+  
  coord_sf(datum = NA)+  
  theme_map()  
  labs(fill = "log price")  
  theme(plot.title = element_text(hjust = 0.5))+  
  ggtitle("Median house price")  
  
gp2 <- ggplot(msoa.spdf)+  
  geom_sf(aes(fill = pt_access_index_ln))+  
  scale_fill_viridis_c(option = "D")  
  geom_sf(data = borough_inner, color = "gray92", fill = NA)+  
  coord_sf(datum = NA)+  
  theme_map()  
  labs(fill = "log index")  
  theme(plot.title = element_text(hjust = 0.5))+  
  ggtitle("Public transport access")  
  
cairo_ps(file = paste("fig/", "Maps.eps", sep=""), width = 10, height = 4,  
        bg = "white", family = "Times New Roman")  
par(mar = c(0, 0, 0, 0))  
par(mfrow = c(1, 1), oma = c(0, 0, 0, 0))  
grid.arrange(gp, gp2, ncol = 2)  
dev.off()  
  
jpeg(file = paste("fig/", "Maps.jpeg", sep=""), width = 10, height = 4,  
     units = "in", res = 300, type = "cairo",  
     bg = "white", family = "Times New Roman")  
par(mar = c(0, 0, 0, 0))  
par(mfrow = c(1, 1), oma = c(0, 0, 0, 0))
```

## 10 Comparing and Selecting Models

```
grid.arrange(gp, gp2, ncol = 2)
dev.off()
```

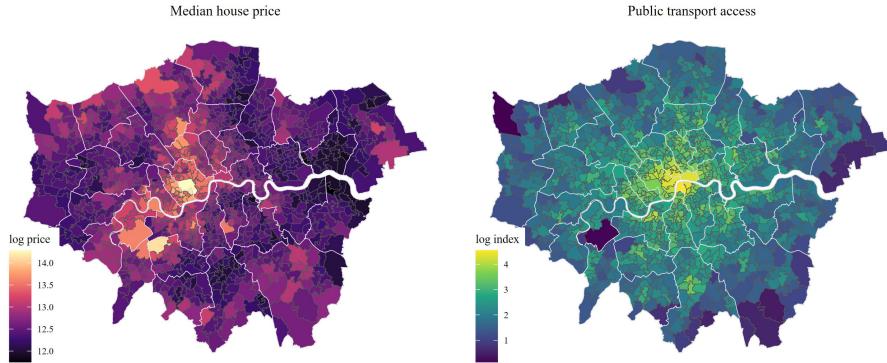


Figure 10.6: Spatial distribution of log-transformed median house prices and transport accessibility across London.

Figure 10.6 shows an unclassified choropleth map of house prices and public transport access across London, both log-scaled for mapping. As we would expect, both indicators follow a relatively strong spatial pattern of positive autocorrelation: house prices first decrease with increasing distance to the centre, and then seem to increase again in suburban areas. Moreover, there seems to be a pattern of higher prices towards the west and particularly high prices around Hyde Park. Public transport accessibility steadily decreases with distance to the city centre. Spatial regression models thus seem to be important here for two reasons: a) observations are not independent of each other but follow clear spatial patterns, and b) surrounding / adjacent urban characteristics likely play a role for housing demand and prices in the focal unit as well.

In Table 1, we regress the median house price in 2011 on the area (in  $\text{km}^2$ ) covered by green space according to OpenStreetMaps, an index of public transport access (ranging from 0-low accessibility to 100-high accessibility), and several population characteristics from the 2011 census

## 10.6 Example: House prices in London

such as population density, the percent of non-UK residents and the percent of social housing. Reported are results from (1) non-spatial OLS, (2) Spatial Autoregressive (SAR), (3) Spatial Error Model (SEM), (4) Spatial Lag of X (SLX), (5) Spatial Durbin Model (SDM), (6) and Spatial Durbin Error Model (SDEM). All variables were standardized before estimation, and we thus interpret coefficients in standard deviations. Note that we do not estimate results for Spatial Autoregressive Combined (SAC) models because of its severe drawbacks for applied research (LeSage 2014a).

```
# Specifcy variables and formula
fm <- med_house_price ~ park_kmsq + pt_access_index + POPDEN + per_nonUK + per_social

# Standardize variables
vars <- all.vars(fm)
msoa_sd.spdf <- msoa.spdf
for(v in vars){
  msoa_sd.spdf[, v] <- as.numeric(scale(msoa_sd.spdf[, v, drop = TRUE]))
}

# Estimate the models
mod_1.ols <- lm(fm, data = msoa_sd.spdf)

# Spatial autoregressive model
mod_1.sar <- lagsarlm(fm,
                        data = msoa_sd.spdf,
                        listw = queens.lw,
                        Durbin = FALSE) # we could here extend to SDM

# Spatial error model
mod_1.sem <- errorsarlm(fm,
                           data = msoa_sd.spdf,
                           listw = queens.lw,
                           Durbin = FALSE) # we could here extend to SDEM
```

## 10 Comparing and Selecting Models

```
# SLX
mod_1.slx <- lmSLX(fm,
                      data = msoa_sd.spdf,
                      listw = queens.lw,
                      Durbin = TRUE) # use a formula to lag only specific cov

# Spatial Durbin
mod_1.sdm <- lagsarlm(fm,
                        data = msoa_sd.spdf,
                        listw = queens.lw,
                        Durbin = TRUE) # we could here extend to SDM

# Spatial Durbun Error
mod_1.sdem <- errorsarlm(fm,
                           data = msoa_sd.spdf,
                           listw = queens.lw,
                           Durbin = TRUE) # we could here extend to SDEM

### Coefficient Output
# Get AIC and N for all models to get common gof stats
aic.l <- sapply(list(mod_1.ols, mod_1.sar, mod_1.sem, mod_1.slx, mod_1.sdm),
                 FUN = function(x) AIC(x))

n.l <- sapply(list(mod_1.ols, mod_1.sar, mod_1.sem, mod_1.slx, mod_1.sdm),
              FUN = function(x) length(residuals(x)))

# Create table
mod_1.slx.lm <- mod_1.slx
class(mod_1.slx.lm) <- "lm" # only for gofs
# screenreg(list(mod_1.ols, mod_1.sar, mod_1.sem, mod_1.slx.lm, mod_1.sdm,
#                 custom.coef.map = list('(Intercept)' = '(Intercept)',
```

## 10.6 Example: House prices in London

```
#             'park_kmsq' = 'Green space',
#             'pt_access_index' = 'Public transport access',
#             'POPDEN' = 'Population density',
#             'per_nonUK' = '% non-UK',
#             'per_social' = '% social housing',
#             'lag.park_kmsq' = 'W Green space',
#             'lag.pt_access_index' = 'W Public transport access',
#             'lag.POPDEN' = 'W Population density',
#             'lag.per_nonUK' = 'W % non-UK',
#             'lag.per_social' = 'W % social housing',
#             'rho' = 'rho',
#             'lambda' = 'lambda'),
# custom.model.names = c("OLS", "SAR", "SEM", "SLX", "SDM", "SDEM"),
# dcolumn = TRUE, caption.above = TRUE, digits = 3,
# caption = "Spatial regression models. Outcome variable: median house price.",
# include.nobs = FALSE,
# include.loglik = FALSE,
# include.aic = FALSE,
# include.lr = TRUE,
# include.wald = FALSE,
# include.fstatistic = FALSE,
# include.rmse = FALSE,
# custom.gof.rows = list('Num. obs.' = n.l,
#                       'AIC' = aic.l),
# reorder.gof = c(1, 3:6, 2))

wordreg(list(mod_1.ols, mod_1.sar, mod_1.sem, mod_1.slx.lm, mod_1.sdm, mod_1.sdem),
        file = "fig/Regression.doc",
        custom.coef.map = list('(Intercept)' = '(Intercept)',
                               'park_kmsq' = 'Green space',
                               'pt_access_index' = 'Public transport access',
                               'POPDEN' = 'Population density',
```

## 10 Comparing and Selecting Models

```
'per_nonUK' = 'Percent non-UK',
'per_social' = 'Percent social housing',
'lag.park_kmsq' = 'W Green space',
'lag.pt_access_index' = 'W Public transp
'lag.POPDEN' = 'W Population density',
'lag.per_nonUK' = 'W Percent non-UK',
'lag.per_social' = 'W Percent social hous
'rho' = 'rho',
'lambda' = 'lambda'),
custom.model.names = c("OLS", "SAR", "SEM", "SLX", "SDM", "SDEM"
dcolumn = TRUE, caption.above = TRUE, digits = 3,
caption = "Spatial regression models. Outcome variable: median h
include.nobs = FALSE,
include.loglik = FALSE,
include.aic = FALSE,
include.lr = TRUE,
include.wald = FALSE,
include.fstatistic = FALSE,
include.rmse = FALSE,
custom.gof.rows = list('Num. obs.' = n.l,
'AIC' = aic.l),
reorder.gof = c(1, 3:6, 2))
```

```
texreg(list(mod_1.ols, mod_1.sar, mod_1.sem, mod_1.slx.lm, mod_1.sdm, mod_
file = "fig/Regression.tex",
custom.coef.map = list('(Intercept)' = '(Intercept)',
'park_kmsq' = 'Green space',
'pt_access_index' = 'Public transport ac
'POPDEN' = 'Population density',
'per_nonUK' = 'Percent non-UK',
'per_social' = 'Percent social housing',
'lag.park_kmsq' = 'W Green space',
```

## 10.6 Example: House prices in London

```
'lag.pt_access_index' = 'W Public transport access',
'lag.POPDEN' = 'W Population density',
'lag.per_nonUK' = 'W Percent non-UK',
'lag.per_social' = 'W Percent social housing',
'rho' = 'rho',
'lambda' = 'lambda'),
custom.model.names = c("OLS", "SAR", "SEM", "SLX", "SDM", "SDEM"),
dcolumn = TRUE, caption.above = TRUE, digits = 3, use.packages = FALSE,
caption = "Spatial regression models. Outcome variable: median house price.",
include.nobs = FALSE, fontsize = "scriptsize",
include.loglik = FALSE,
include.aic = FALSE,
include.lr = TRUE,
include.wald = FALSE,
include.fstatistic = FALSE,
include.rmse = FALSE,
custom.gof.rows = list('Num. obs.' = n.1,
                      'AIC' = aic.l),
reorder.gof = c(1, 3:6, 2))
```

Compared to results from conventional non-spatial models, Table 1 comes with several additions: First, variables starting with a “W” (or “lag”) indicate the spatially lagged variable or in the case of row-normalized weights matrices the average value of the respective variable across the local neighbours. Moreover, there are two auto-regressive parameters: “rho” for the estimated auto-correlation in the dependent variable and “lambda” for the estimated auto-correlation in the error term. In case of the SAR, a highly significant  $\hat{\rho}$  coefficient of 0.786 indicates strong positive spatial auto-correlation in the median house price: the house price in adjacent areas positively impacts the focal house prices. A  $\hat{\lambda}$  of 0.89 in the SEM however indicates that there is very strong spatial auto-correlation among the (remaining) error variance. The likelihood ratio test in the goodness-of-fit statistics are highly significant in both cases, rejecting the

## 10 Comparing and Selecting Models

Table 10.1: Spatial regression models. Outcome variable: median house price.

	OLS	SAR	SEM	SLX	SDM	SD
(Intercept)	0.000 (0.027)	-0.012 (0.017)	0.022 (0.139)	0.007 (0.024)	0.002 (0.015)	0.000 (0.015)
Green space	0.204*** (0.029)	0.133*** (0.018)	0.100*** (0.015)	0.136*** (0.026)	0.106*** (0.016)	0.106*** (0.016)
Public transport access	0.366*** (0.033)	0.097*** (0.021)	-0.054 (0.033)	-0.152** (0.054)	-0.100** (0.034)	-0.100** (0.034)
Population density	0.189*** (0.037)	0.055* (0.023)	-0.094*** (0.027)	-0.112* (0.044)	-0.111*** (0.028)	-0.111*** (0.028)
Percent non-UK	-0.033 (0.033)	-0.050* (0.020)	-0.250*** (0.033)	-0.235*** (0.053)	-0.262*** (0.033)	-0.262*** (0.033)
Percent social housing	-0.402*** (0.032)	-0.202*** (0.020)	-0.260*** (0.022)	-0.306*** (0.035)	-0.266*** (0.022)	-0.266*** (0.022)
W Green space			0.249*** (0.040)	-0.029 (0.026)	0.000 (0.026)	0.000 (0.026)
W Public transport access			0.696*** (0.069)	0.239*** (0.045)	0.000 (0.045)	0.000 (0.045)
W Population density			0.455*** (0.065)	0.136*** (0.041)	0.000 (0.041)	0.000 (0.041)
W Percent non-UK			0.304*** (0.066)	0.300*** (0.042)	0.000 (0.042)	0.000 (0.042)
W Percent social housing			-0.352*** (0.053)	0.119*** (0.035)	-0.029 (0.035)	-0.029 (0.035)
Num. obs.	983	983	983	983	983	983
R <sup>2</sup>	0.263			0.439		
Adj. R <sup>2</sup>	0.259			0.433		
LR test: statistic		789.480	934.291		732.684	695.121
LR test: p-value		0.000	0.000		0.000	0.000
AIC	2502.492	1715.012	1570.201	2244.135	1513.451	1551.121

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

## 10.6 Example: House prices in London

NULL of no spatial auto-autocorrelation.

Given the strong positive auto-correlation in the dependent variable in SAR and SDM, we cannot directly interpret the coefficients as marginal effects. Similar to auto-regressive temporal models, we need to account for the spatial multiplier effect. For SEM, SLX and SDEM, we could directly interpret the coefficients of Table 1. However, we plot the impacts of all five models in Figure Figure 10.7 for reasons of comparison. Note that SEM only has direct and no indirect impacts.

```
# Get direct and indirect impacts
mod.1 <- list(mod_1.sar, mod_1.slx, mod_1.sdm, mod_1.sdem)
imp.1 <- vector(mode = "list", length = length(mod.1))

for(i in 1:length(mod.1)){
  imp.1[[i]] <- spatialreg::impacts(mod.1[[i]], listw = queens.lw, R = 600)
}

# Add SEM

# Extract summary measures
extract.imp <- function(x){
  s <- summary(x, zstats = TRUE, short = TRUE)
  names <- attr(x, "bnames")
  l <- length(names)
  effs <- c("Direct", "Indirect", "Total")
  if(attr(x, "type") == "lag" | attr(x, "type") == "mixed"){
    coefs = c(s$res$direct, s$res$indirect, s$res$total)
  }else{
    coefs = c(s$impacts$direct, s$impacts$indirect, s$impacts$total)
  }
  df <- data.frame(var = rep(names, 3),
```

## 10 Comparing and Selecting Models

```
eff = rep(effs, each = 1),
coef = coefs,
se = c(s$semat[, 1], s$semat[, 2], s$semat[, 3]),
pval = c(s$pzmat[, 1], s$pzmat[, 2], s$pzmat[, 3])
)
}

mods <- c("SAR", "SLX", "SDM", "SDEM")
for(i in 1:length(imp.l)){
  tmp <- extract.imp(imp.l[[i]])
  tmp$mod <- mods[i]
  if(i == 1){
    imp.res <- tmp
  }else{
    imp.res <- rbind(imp.res, tmp)
  }
}

# Add SEM
sem.coefs <- summary(mod_1.sem)$Coef[,-3]
colnames(sem.coefs) <- c("coef", "se", "pval")
sem.df <- data.frame(var = rownames(sem.coefs),
                      eff = "Direct",
                      sem.coefs,
                      mod = "SEM")

imp.res <- rbind(imp.res, sem.df[-1, ])

### Plot the effects

# Coef Labels
```

## 10.6 Example: House prices in London

```
imp.res$lab <- as.character(sprintf("%.3f", round(imp.res$coef, 3)))  
  
# # Add stars  
# imp.res$lab[imp.res$pval <= 0.1 & imp.res$pval > 0.05] <- paste0(imp.res$lab[imp.res$pval  
# imp.res$lab[imp.res$pval <= 0.05 & imp.res$pval > 0.01] <- paste0(imp.res$lab[imp.res$pval  
# imp.res$lab[imp.res$pval <= 0.01 & imp.res$pval > 0.001] <- paste0(imp.res$lab[imp.res$pval  
# imp.res$lab[imp.res$pval <= 0.001] <- paste0(imp.res$lab[imp.res$pval <= 0.001], "***")  
  
# # Get rid of leading zero  
# imp.res$lab <- gsub("0\\.", " \\. ", imp.res$lab)  
  
# Confidence intervals  
interval2 <- -qnorm((1-0.95)/2) # 95% multiplier  
imp.res$lb <- imp.res$coef - imp.res$se * interval2  
imp.res$ub <- imp.res$coef + imp.res$se * interval2  
  
# Rename variables  
names <- list('park_kmsq' = 'Green space',  
              'pt_access_index' = 'Public transport access',  
              'POPDEN' = 'Population density',  
              'per_nonUK' = 'Percent non-UK',  
              'per_social' = 'Percent social housing')  
  
imp.res$var <- factor(imp.res$var, levels = rev(names(names)), labels = rev(names))  
imp.res$mod <- factor(imp.res$mod, levels = rev(c("SAR", "SEM", "SLX", "SDM", "SDEM")))  
  
# Plot  
zp_all <- ggplot(imp.res[imp.res$eff != "Total", ], aes(colour = mod, shape = mod, fill = mod)) +  
  facet_grid(. ~ eff, scales = "free_x") +  
  geom_hline(yintercept = 0, colour = scales::alpha("black", 0.3), lty = 2) +  
  geom_pointrange(aes(x = var, y = coef, ymin = lb, ymax = ub),  
                  lwd = 0.7, position = position_dodge(width = 1/1.4),  
                  size = 1.5)
```

## 10 Comparing and Selecting Models

```
fill = "black") +  
  geom_text(aes(label = lab,  
                x = var,  
                y = coef),  
            size = 3.0, show.legend = FALSE,  
            vjust = -0.35, hjust = -0.035, position = position_dodge(width =  
              coord_flip() + theme_bw() +  
              scale_x_discrete(expand = c(0.1, 0.1)) +  
              theme(legend.title = element_blank()) +  
              labs(y = "Impacts on house price", x = "") +  
              scale_shape_manual(values = rev(c(18, 19, 17, 15, 25))) +  
              scale_color_viridis_d() +  
              scale_fill_viridis_d() +  
              theme(text = element_text(size = 13),  
                    legend.position = "bottom",  
                    legend.background = element_blank(),  
                    legend.box.background = element_rect(colour = "black"),  
                    legend.key = element_blank(),  
                    axis.text.y = element_text(colour = "black", size = 13),  
                    axis.title.x = element_text(colour = "black", size = 13),  
                    axis.text.x = element_text(colour = "black", size = 13),  
                    strip.background = element_blank(),  
                    strip.text = element_text(size = 13, colour = "black",  
                                              angle = 90)),  
            ) +  
  ggtitle(element_blank()) +  
  guides(colour = guide_legend(override.aes = list(linetype = 0), reverse = T))  
  
  cairo_ps(file = paste("fig/", "Coefplot.eps", sep = ""), width = 8, height = 6,  
           bg = "white", family = "Times New Roman")  
  par(mar = c(0, 0, 0, 0))  
  par(mfrow = c(1, 1), oma = c(0, 0, 0, 0))
```

## 10.6 Example: House prices in London

```
zp_all  
dev.off()  
  
jpeg(file = paste("fig/", "Coefplot.jpeg", sep=""), width = 8, height = 6,  
      units = "in", res = 300, type = "cairo",  
      bg = "white", family = "Times New Roman")  
par(mar = c(0, 0, 0, 0))  
par(mfrow = c(1, 1), oma = c(0, 0, 0, 0))  
zp_all  
dev.off()
```

We start with the results of the SAR model in Figure Figure 10.7. A one standard-deviation increase of green space in the focal unit is associated with a 0.161 standard deviation increase in house prices within the same spatial unit. However, there are also highly significant diffusion processes. This increase in green space in the focal unit will also increase house prices in neighbouring units and the neighbours of these neighbours. This indirect impact will add up to a 0.458 standard deviation increase in house prices across neighbouring units connected through the spatial weights system. Similarly, an increase in public transport accessibility is associated with a 0.118 standard-deviation higher median house price in the unit itself and an additional 0.458 deviation increase diffusing though the neighbouring regions. Note that direct and indirect effects are bound to a common ration, as SAR only estimates one single spatial parameter  $\hat{\rho}$ . In our case, every indirect impact equals approximately 2.83 times the direct impact. This is a very restrictive conditions and a severe drawback of the SAR model.

The SLX - similar to SAR - estimates a positive impact of green space in the focal but also in adjacent neighbourhoods on house prices in the focal unit. A one standard deviation in the focal unit is associated with 0.136 standard-deviations higher house price in the focal unit. If green spaces in adjacent neighbourhoods increase on average by one standard deviation, this would increase house prices in the focal unit by 0.249 standard deviations. Note

## 10 Comparing and Selecting Models

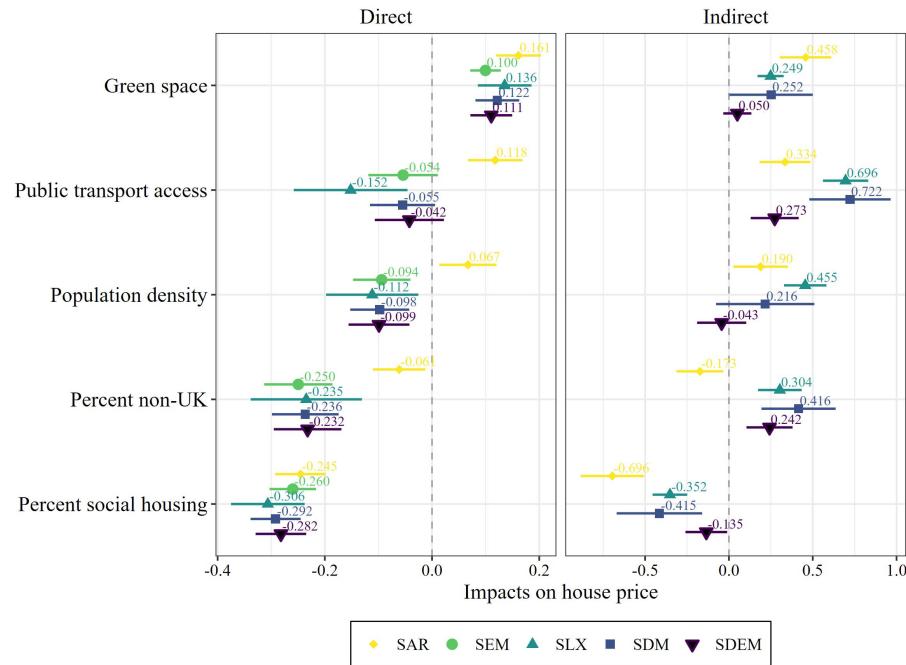


Figure 10.7: Direct and indirect impacts from spatial regression models.  
 Dependent variable: mean house prices. All variables are standardised.

## *10.6 Example: House prices in London*

that the SLX tells a different story about the effect of public transport access than SAR: there is a negative direct and a very strong and positive indirect effect. A one standard deviation increase in public transport access in the focal unit is associated with -0.152 standard deviations lower house prices. In contrast, more public transport in the local surrounding (the average neighbours) is associated with 0.696 standard deviations higher prices. This is in line with the idea that public transport facilities are usually not particularly attractive: it is good to have them close but not too close. The same is true for population density: it is good to live in a broader area with high population density as indicated by the indirect impacts (probability indicating high centrality), but the local neighbourhood should have a low population density as indicated by the negative direct impact.

We could go further with the other models. However, interpretation in SDM follows the same logic as SAR, and interpretation in SDEM aligns to SLX. Interpretation in SEM is analogous to non-spatial OLS, as there are no indirect impacts. Moreover, it is important to keep in mind that the indirect impacts are summary measures which sum over all impacts from or onto neighbouring regions. The indirect public transport effect of 0.696 in SLX would occur if the average public transport access across neighbours would increase by one standard deviation. This only occurs if all neighbours would simultaneously increase public transport access by one standard deviation.



# 11 Exercises III

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite",
         "plm", "lfe", "splm", "SDPDmod")
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()

R version 4.4.1 (2024-06-14 ucrt)
Platform: x86_64-w64-mingw32/x64
Running under: Windows 11 x64 (build 22631)

Matrix products: default

locale:
[1] LC_COLLATE=English_United Kingdom.utf8
[2] LC_CTYPE=English_United Kingdom.utf8
[3] LC_MONETARY=English_United Kingdom.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United Kingdom.utf8
```

### 11 Exercises III

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets  methods   base

other attached packages:
[1] SDPDmod_0.0.5      splm_1.6-5        lfe_3.0-0       plm_2.6-
4
[5] viridisLite_0.4.2   tmap_3.3-4       spatialreg_1.3-
4
Matrix_1.7-0
[9] spdep_1.3-5        spData_2.3.1     mapview_2.11.2   sf_1.0-
16

loaded via a namespace (and not attached):
[1] fastmap_1.2.0      leaflet_2.2.2     TH.data_1.1-2
[4] dotCall64_1.1-1    fixest_0.12.1     XML_3.99-0.16.1
[7] digest_0.6.35      lifecycle_1.0.4   dreamerr_1.4.0
[10] LearnBayes_2.15.1  survival_3.6-4   terra_1.7-78
[13] magrittr_2.0.3     compiler_4.4.1   rlang_1.1.4
[16] tools_4.4.1       collapse_2.0.14  knitr_1.47
[19] htmlwidgets_1.6.4   sp_2.1-4        classInt_0.4-10
[22] RColorBrewer_1.1-3 multcomp_1.4-25  abind_1.4-5
[25] KernSmooth_2.23-24 numDeriv_2016.8-1.1 leafsync_0.1.0
[28] grid_4.4.1        stats4_4.4.1    xtable_1.8-4
[31] e1071_1.7-14      leafem_0.2.3    colorspace_2.1-0
[34] scales_1.3.0      MASS_7.3-60.2   dichromat_2.0-
0.1
[37] cli_3.6.2         mvtnorm_1.2-5   rmarkdown_2.27
[40] miscTools_0.6-28   generics_0.1.3   RSpectra_0.16-1
[43] rstudioapi_0.16.0  tmaptools_3.1-1  bdsmatrix_1.3-7
[46] DBI_1.2.3         proxy_0.4-27   stringr_1.5.1
[49] splines_4.4.1     stars_0.6-5    parallel_4.4.1
```

## 11.1 Environmental inequality (continued)

```
[52] s2_1.1.6           stringmagic_1.1.2    base64enc_0.1-3
[55] boot_1.3-30        sandwich_3.1-0      jsonlite_1.8.8
[58] Formula_1.2-5     crosstalk_1.2.1   units_0.8-5
[61] spam_2.10-0       glue_1.7.0        lwgeom_0.2-14
[64] codetools_0.2-20  stringi_1.8.4    deldir_2.0-4
[67] raster_3.6-26     lmtest_0.9-40    munsell_0.5.1
[70] htmltools_0.5.8.1 satellite_1.0.5  R6_2.5.1
[73] wk_0.9.1          maxLik_1.5-2.1  Rdpack_2.6
[76] evaluate_0.24.0   lattice_0.22-6  rbibutils_2.2.16
[79] png_0.1-8         class_7.3-22    Rcpp_1.0.12
[82] coda_0.19-4.1    nlme_3.1-164   xfun_0.45
[85] zoo_1.8-12
```

### Reload data from previous session

```
load("_data/msoa2_spatial.RData")
```

## 11.1 Environmental inequality (continued)

Let's use the same neighbours weights definition as before:

```
coords <- st_centroid(msoa.spdf)
```

```
Warning: st_centroid assumes attributes are constant over geometries
```

```
# Neighbours within 3km distance
dist_15.nb <- dnearneigh(coords, d1 = 0, d2 = 2500)

summary(dist_15.nb)
```

### 11 Exercises III

```
Neighbour list object:  
Number of regions: 983  
Number of nonzero links: 15266  
Percentage nonzero weights: 1.579859  
Average number of links: 15.53001  
4 regions with no links:  
158 463 478 505  
6 disjoint connected subgraphs  
Link number distribution:  
  
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25  
4 5 9 23 19 26 36 31 53 39 61 63 59 48 42 35 24 31 28 30 27 26 25 19 38 29  
26 27 28 29 30 31 32 33 34  
32 38 26 16 20 10 8 1 2  
5 least connected regions:  
160 469 474 597 959 with 1 link  
2 most connected regions:  
565 567 with 34 links
```

```
# There are some empty neighbour sets. Lets impute those with the nearest  
k2.nb <- knearneigh(coords, k = 1)
```

```
# Replace zero  
nolink_ids <- which(card(dist_15.nb) == 0)  
dist_15.nb[card(dist_15.nb) == 0] <- k2.nb$nn[nolink_ids, ]  
  
summary(dist_15.nb)
```

```
Neighbour list object:  
Number of regions: 983  
Number of nonzero links: 15270  
Percentage nonzero weights: 1.580273  
Average number of links: 15.53408
```

### 11.1 Environmental inequality (continued)

```
2 disjoint connected subgraphs
Link number distribution:

 1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26
 9  9 23 19 26 36 31 53 39 61 63 59 48 42 35 24 31 28 30 27 26 25 19 38 29 32
27 28 29 30 31 32 33 34
38 26 16 20 10  8  1  2
9 least connected regions:
158 160 463 469 474 478 505 597 959 with 1 link
2 most connected regions:
565 567 with 34 links
```

```
# listw object with row-normalization
dist_15.lw <- nb2listw(dist_15.nb, style = "W")
```

and estimate the spatial SAR model:

```
mod_1.sar <- lagsarlm(log(no2) ~ per_mixed + per_asian + per_black + per_other
+ per_nonUK_EU + per_nonEU + log(POPDEN),
data = msoa.spdf,
listw = dist_15.lw,
Durbin = FALSE) # we could here extend to SDM
summary(mod_1.sar)
```

```
Call:lagsarlm(formula = log(no2) ~ per_mixed + per_asian + per_black +
per_other + per_nonUK_EU + per_nonEU + log(POPDEN), data = msoa.spdf,
listw = dist_15.lw, Durbin = FALSE)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.2140485	-0.0267085	-0.0021421	0.0238337	0.3505513

### 11 Exercises III

```
Type: lag
Coefficients: (asymptotic standard errors)
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -1.7004e-02 1.8122e-02 -0.9383 0.348110
per_mixed    3.4376e-04 1.4758e-03  0.2329 0.815810
per_asian   -8.5205e-05 1.1494e-04 -0.7413 0.458507
per_black   -4.2754e-04 2.3468e-04 -1.8218 0.068484
per_other    1.9693e-03 7.4939e-04  2.6279 0.008591
per_nonUK_EU 8.9027e-04 3.9638e-04  2.2460 0.024703
per_nonEU    1.8460e-03 3.5159e-04  5.2506 1.516e-07
log(POPDEN)  1.8650e-02 2.7852e-03  6.6963 2.138e-11

Rho: 0.9684, LR test value: 2002.5, p-value: < 2.22e-16
Asymptotic standard error: 0.0063124
      z-value: 153.41, p-value: < 2.22e-16
Wald statistic: 23535, p-value: < 2.22e-16

Log likelihood: 1562.401 for lag model
ML residual variance (sigma squared): 0.0020568, (sigma: 0.045352)
Number of observations: 983
Number of parameters estimated: 10
AIC: -3104.8, (AIC for lm: -1104.3)
LM test for residual autocorrelation
test value: 108.97, p-value: < 2.22e-16
```

#### 1) Please calculate the true multiplier matrix of this SAR model.

The multiplier matrix is given by  $(\mathbf{I}_N - \rho\mathbf{W})^{-1}$ .

## 11.2 Inkar data: the effect of regional characteristics on life expectancy

- 2) Create an N x N effects matrix for the effect of the non-EU citizens. What is the effect of unit 6 on unit 10? Why is this larger than the effect of unit 5 on unit 8?**
- 3) Calculate and interpret the summary impact measures of the SAR model.**
- 4) Is SAR the right model choice or would you rather estimate a different model? Please run a Durbin model and calculate its impact summary measures**
- 5) Please repeat with a Durbin Error model. Why are the impacts here identical to the coefficients?**

## 11.2 Inkar data: the effect of regional characteristics on life expectancy

Below, we read and transform some characteristics of the INKAR data on German counties.

```
load("_data/inkar2.Rdata")
```

Variables are

Variable	Description
“Kennziffer”	ID
“Raumeinheit”	Name
“Aggregat”	Level
“year”	Year
“poluation_density”	Population Density
“median_income”	Median Household income (only for 2020)

### 11 Exercises III

Variable	Description
“gdp_in1000EUR”	Gross Domestic Product in 1000 euros
“unemployment_rate”	Unemployment rate
“share_longterm_unemployed”	Share of longterm unemployed (among unemployed)
“share_working_industry”	Share of employees in undustrial sector
“share_foreigners”	Share of foreign nationals
“share_college”	Share of school-finishers with college degree
“recreational_space”	Recreational space per inhabitant
“car_density”	Density of cars
“life_expectancy”	Life expectancy

### 11.3 County shapes

```
kreise.spdf <- st_read(dsn = "_data/vg5000_ebenen_1231",
                         layer = "VG5000_KRS")

Reading layer `VG5000_KRS' from data source
`C:\work\Lehre\Geodata_Spatial_Regression\_data\vg5000_ebenen_1231'
  using driver `ESRI Shapefile'
Simple feature collection with 400 features and 24 fields
Geometry type: MULTIPOLYGON
Dimension:      XY
Bounding box:  xmin: 280353.1 ymin: 5235878 xmax: 921261.6 ymax: 6101302
Projected CRS: ETRS89 / UTM zone 32N
```

## 11.4 Estimate an FE model with SLX specification

### 1) Please map the life expectancy across Germany

- a) Merge data with the shape file (as with conventional data)
- b) Create a map of life-expectancy

### 2) Choose some variables that could predict life expectancy. See for instance the following paper.

### 3) Generate a neighbours object (e.g. the 10 nearest neighbours).

### 4) Estimate a cross-sectional spatial model for the year 2020 and calculate the impacts.

### 5) Calculate the spatial lagged variables for your covariates (e.g. use `create_WX()`, which needs a non-spatial df as input) .

### 6) Can you run a spatial machine learning model? (for instance, using `randomForest`)?

You could even go further and use higher order neighbours (e.g. `nblag(queens.nb, maxlag = 3)`) to check the importance of direct neighbours and the neighbours neighbours and so on ...

## 11.4 Estimate an FE model with SLX specification

- a) Loops over years to generate WX
- b) Estimate a twoways FE SLX panel model
- c) Estimate a twoways FE SAR panel model (use `spml()`)

*11 Exercises III*

- d) Estimate the summary impacts.

# 12 Spatio-temporal models

## Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite",
         "plm", "splm", "SDPDmod")
lapply(pkgs, require, character.only = TRUE)
```

## Session info

```
sessionInfo()

R version 4.4.1 (2024-06-14 ucrt)
Platform: x86_64-w64-mingw32/x64
Running under: Windows 11 x64 (build 22631)

Matrix products: default

locale:
[1] LC_COLLATE=English_United Kingdom.utf8
[2] LC_CTYPE=English_United Kingdom.utf8
[3] LC_MONETARY=English_United Kingdom.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United Kingdom.utf8
```

## 12 Spatio-temporal models

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets   methods    base

other attached packages:
[1] SDPDmod_0.0.5      splm_1.6-5       plm_2.6-4      viridisLite_0.4.2
[5] tmap_3.3-4         spatialreg_1.3-4 Matrix_1.7-0    spdep_1.3-
5
[9] spData_2.3.1       mapview_2.11.2    sf_1.0-16

loaded via a namespace (and not attached):
[1] fastmap_1.2.0        leaflet_2.2.2      TH.data_1.1-2
[4] dotCall64_1.1-1      fixest_0.12.1     XML_3.99-0.16.1
[7] digest_0.6.35        lifecycle_1.0.4   dreamerr_1.4.0
[10] LearnBayes_2.15.1   survival_3.6-4   terra_1.7-78
[13] magrittr_2.0.3       compiler_4.4.1   rlang_1.1.4
[16] tools_4.4.1         collapse_2.0.14  knitr_1.47
[19] htmlwidgets_1.6.4    sp_2.1-4        classInt_0.4-10
[22] RColorBrewer_1.1-3  multcomp_1.4-25  abind_1.4-5
[25] KernSmooth_2.23-24 numDeriv_2016.8-1.1 leafsync_0.1.0
[28] grid_4.4.1          stats4_4.4.1     xtable_1.8-4
[31] lfe_3.0-0           e1071_1.7-14   leafem_0.2.3
[34] colorspace_2.1-0    scales_1.3.0    MASS_7.3-60.2
[37] dichromat_2.0-0.1   cli_3.6.2       mvtnorm_1.2-5
[40] rmarkdown_2.27       miscTools_0.6-28 generics_0.1.3
[43] RSpectra_0.16-1     rstudioapi_0.16.0 tmaptools_3.1-1
[46] bdsmatrix_1.3-7     DBI_1.2.3       proxy_0.4-27
[49] stringr_1.5.1       splines_4.4.1   stars_0.6-5
[52] parallel_4.4.1      s2_1.1.6       stringmagic_1.1.2
[55] base64enc_0.1-3     boot_1.3-30    sandwich_3.1-0
[58] jsonlite_1.8.8      Formula_1.2-5  crosstalk_1.2.1
```

## 12.1 Static panel data models

```
[61] units_0.8-5          spam_2.10-0        glue_1.7.0
[64] lwgeom_0.2-14        codetools_0.2-20    stringi_1.8.4
[67] deldir_2.0-4         raster_3.6-26      lmtest_0.9-40
[70] munsell_0.5.1        htmltools_0.5.8.1   satellite_1.0.5
[73] R6_2.5.1              wk_0.9.1           maxLik_1.5-2.1
[76] Rdpack_2.6            evaluate_0.24.0    lattice_0.22-6
[79] rbibutils_2.2.16     png_0.1-8          class_7.3-22
[82] Rcpp_1.0.12           coda_0.19-4.1     nlme_3.1-164
[85] xfun_0.45             zoo_1.8-12
```

Elhorst (2014) provides a comprehensive introduction to spatial panel data methods. Article length introduction to spatial panel data models (FE and RE) can be found in Elhorst (2012), Millo and Piras (2012) and Croissant and Millo (2018). LeSage (2014b) discusses Bayesian panel data methods.

Note that we will only discuss some basics here, as the complete econometrics of these models and their estimation strategy become insanely complicated (L. Lee and Yu 2010).

### 12.1 Static panel data models

The idea behind a static panel data with auto-regressive term is similar to the cross sectional situation (Millo and Piras 2012).

$$\mathbf{y} = \rho(\mathbf{I}_T \otimes \mathbf{W}_N)\mathbf{y} + \mathbf{X} + \mathbf{u}.$$

where  $\otimes$  is the Kronecker product (block-wise multiplication).

## 12 Spatio-temporal models

$$\underbrace{\mathbf{I}_T \otimes \mathbf{W}_N}_{\substack{T \times T \\ NT \times NT}} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} \begin{bmatrix} v_1 w_1 & v_1 w_2 & \cdots & v_1 w_m \\ v_2 w_1 & v_2 w_2 & \cdots & v_2 w_m \\ \vdots & \vdots & \ddots & \vdots \\ v_n w_1 & v_n w_2 & \cdots & v_n w_m \end{bmatrix}$$

$$\begin{pmatrix} \begin{bmatrix} v_1 w_1 & v_1 w_2 & \cdots & v_1 w_m \\ v_2 w_1 & v_2 w_2 & \cdots & v_2 w_m \\ \vdots & \vdots & \ddots & \vdots \\ v_n w_1 & v_n w_2 & \cdots & v_n w_m \end{bmatrix} & 0 & \cdots & 0 \\ 0 & \begin{bmatrix} v_1 w_1 & v_1 w_2 & \cdots & v_1 w_m \\ v_2 w_1 & v_2 w_2 & \cdots & v_2 w_m \\ \vdots & \vdots & \ddots & \vdots \\ v_n w_1 & v_n w_2 & \cdots & v_n w_m \end{bmatrix} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \begin{bmatrix} v_1 w_1 & v_1 w_2 & \cdots & v_1 w_m \\ v_2 w_1 & v_2 w_2 & \cdots & v_2 w_m \\ \vdots & \vdots & \ddots & \vdots \\ v_n w_1 & v_n w_2 & \cdots & v_n w_m \end{bmatrix} \end{pmatrix}$$

Here we model only spatial dependence within each cross-section and multiply the same spatial weights matrix  $T$  times. Off block-diagonal elements are all zero. So there is no spatial dependence that goes across time.

The error term can be decomposed into two parts:

$$\mathbf{u} = (\mathbf{1}_T \otimes \mathbf{I}_N) + ,$$

where  $\mathbf{1}_T$  is a  $T \times 1$  vector of ones,  $\mathbf{I}_N$  an  $N \times N$  identity matrix,  $\mathbf{v}$  is a vector of time-invariant individual specific effects (not spatially autocorrelated).

We could obviously extent the specification to allow for error correlation by specifying

## 12.1 Static panel data models

$$= \lambda(\mathbf{I}_T \otimes \mathbf{W}_N) + .$$

The individual effects can be treated as fixed or random.

### Fixed Effects

In the FE model, the individual specific effects are treated as fixed. If we re-write the equation above, we derive at the well-known fixed effects formula with an additional spatial autoregressive term:

$$y_{it} = \rho \sum_{j=1}^N w_{ij} y_{jt} + \mathbf{x}_{it} + \mu_i + \nu_{it},$$

where  $\mu_i$  denote the individual-specific fixed effects.

As with the standard spatial lag model, we cannot rely on the OLS estimator because of the simultaneity problem. The coefficients are thus estimated by maximum likelihood (Elhorst 2014).

### Random Effects

In the RE model, the individual specific effects are treated as components of the error  $\mu \sim \text{IID}(o, \sigma_\mu^2)$ . The model can then be written as

$$\begin{aligned} \mathbf{y} &= \rho(\mathbf{I}_T \otimes \mathbf{W}_N)\mathbf{y} + \mathbf{X} + \mathbf{u}, \\ \mathbf{u} &= (\mathbf{I}_T \otimes \mathbf{I}_N) + [\mathbf{I}_T \otimes (\mathbf{I}_N - \lambda \mathbf{W}_N)]^{-1}. \end{aligned}$$

As with the conventional random effects model, we make the strong assumption that the unobserved individual effects are uncorrelated with the covariates  $\mathbf{X}$  in the model.

## 12.2 Dynamic panel data models

Relying on panel data and repeated measures over time, comes with an additional layer of dependence / autocorrelation between units. We have spatial dependence (with its three potential sources), and we have temporal/serial dependence within each unit over time.

A general dynamic model would account for all sources of potential dependence, including combinations (Elhorst 2012). The most general model can be written as:

$$\begin{aligned}\mathbf{y}_t &= \tau\mathbf{y}_{t-1} + \rho(\mathbf{I}_T \otimes \mathbf{W}_N)\mathbf{y}_t + \gamma(\mathbf{I}_T \otimes \mathbf{W}_N)\mathbf{y}_{t-1} \\ &\quad + \mathbf{X} + (\mathbf{I}_T \otimes \mathbf{W}_N)\mathbf{X} + \mathbf{u}_t, \\ \mathbf{u}_t &= +(\mathbf{I}_T \otimes \mathbf{I}_N) + \varepsilon_t, \\ \varepsilon_t &= \psi_{t-1} + \lambda(\mathbf{I}_T \otimes \mathbf{W}_N) + \nu_t,\end{aligned}$$

Where  $\mathbf{X}$  could further contain time-lagged covariates. Compared to the static spatial panel model, we have introduced temporal dependency in the outcome  $\tau\mathbf{y}_{t-1}$  and the spatially lagged outcome  $\gamma(\mathbf{I}_T \otimes \mathbf{W}_N)\mathbf{y}_{t-1}$ , and in the error term  $\psi_{t-1}$ .

### 12.2.1 Impacts in spatial panel models

Note that similar to the distinction between local and global spillovers, we now have to distinguish between short-term and long-term effects. A change in  $\mathbf{X}_t$  influences focal  $Y$  and neighbour's  $Y$  but also contemporaneous  $Y$  and future  $Y$ .

While the short-term effects are the known impacts

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}_k} = (\mathbf{I} - \rho \mathbf{W}_{NT})^{-1} [\beta_k + \mathbf{W}_{NT} \theta_k].$$

### 12.3 Example: Local employment impacts of immigration

The long-term impacts, by contrast, additionally account for the effect multiplying through time

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}_k} = [(1 - \tau)\mathbf{I} - (\rho + \gamma)\mathbf{W}_{NT}]^{-1} [\beta_k + \mathbf{W}_{NT}\theta_k].$$

For more information see Elhorst (2012).

**Table 1** Short-term, long-term, direct and indirect (spatial spillover) effects of different models

Type of model	Short-term direct effect	Short-term indirect effect	Long-term direct effect	Long-term indirect effect	Shortcoming
0. Static spatial Durbin model	–	–	$[(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^d$	$[(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^{\text{sum}}$	No short-term effects
1. Error terms lagged in space and/or in time	–	–	$\beta_{1k}$	–	No short-term effects No indirect effects
2. Dynamic model + spatial error corr.	$\beta_{1k}$	–	$\beta_{1k}/(1 - \tau)$	–	No indirect effects
3. Dynamic spatial Durbin model	$[(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^d$	$[(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^{\text{sum}}$	$[(1 - \tau)I - (\delta + \eta)W]^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^d$	$[(1 - \tau)I - (\delta + \eta)W]^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^{\text{sum}}$	Parameters not identified
4. $\beta_2 = 0$	$[(I - \delta W)^{-1}(\beta_{1k}I_N)]^d$	$[(I - \delta W)^{-1}(\beta_{1k}I_N)]^{\text{sum}}$	$[(1 - \tau)I - (\delta + \eta)W]^{-1}(\beta_{1k}I_N)]^d$	$[(1 - \tau)I - (\delta + \eta)W]^{-1}(\beta_{1k}I_N)]^{\text{sum}}$	Ratio ind./dir effects the same for every X
5. $\delta = 0$	$[(\beta_{1k}I_N + \beta_{2k}W)]^d$	$[(\beta_{1k}I_N + \beta_{2k}W)]^{\text{sum}}$	$[(1 - \tau)I - \eta W]^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^d$	$[(1 - \tau)I - \eta W]^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^{\text{sum}}$	No short-term global indirect effects
6. $\eta = -\tau\delta$	$[(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^d$	$[(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^{\text{sum}}$	$\left[\frac{1}{1-\tau}(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)\right]^d$	$\left[\frac{1}{1-\tau}(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)\right]^{\text{sum}}$	Ratio ind./dir effects constant over time
7. $\eta = 0$	$[(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^d$	$[(I - \delta W)^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^{\text{sum}}$	$[(1 - \tau)I - \delta W]^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^d$	$[(1 - \tau)I - \delta W]^{-1}(\beta_{1k}I_N + \beta_{2k}W)]^{\text{sum}}$	–

Figure 12.1: Summary impact measures in dynamic spatial panel models (Elhorst 2012)

### 12.3 Example: Local employment impacts of immigration

Fingleton, Olner, and Pryce (2020): Estimating the local employment impacts of immigration: A dynamic spatial panel model. Urban Studies, 57(13), 2646–2662. <https://doi.org/10.1177/0042098019887916>

*This paper highlights a number of important gaps in the UK evidence base on the employment impacts of immigration, namely: (1) the lack of research*

## 12 Spatio-temporal models

*on the local impacts of immigration – existing studies only estimate the impact for the country as a whole; (2) the absence of long-term estimates – research has focused on relatively short time spans – there are no estimates of the impact over several decades, for example; (3) the tendency to ignore spatial dependence of employment which can bias the results and distort inference – there are no robust spatial econometric estimates we are aware of.*

*We illustrate our approach with an application to London and find that no migrant group has a statistically significant long-term negative effect on employment. EU migrants, however, are found to have a significant positive impact, which may have important implications for the Brexit debate. Our approach opens up a new avenue of inquiry into subnational variations in the impacts of immigration on employment.*

### 12.4 Estimation in R

To estimate spatial panel models in R, we can use the `splm` package of Millo and Piras (2012).

We use a standard example with longitudinal data from the `plm` package here.

```
data(Produc, package = "plm")
data(usaww)

head(Produc)

  state year region    pcap    hwy   water   util     pc    gsp    emp
1 ALABAMA 1970       6 15032.67 7325.80 1655.68 6051.20 35793.80 28418 1010.5
2 ALABAMA 1971       6 15501.94 7525.94 1721.02 6254.98 37299.91 29375 1021.9
3 ALABAMA 1972       6 15972.41 7765.42 1764.75 6442.23 38670.30 31303 1072.3
4 ALABAMA 1973       6 16406.26 7907.66 1742.41 6756.19 40084.01 33430 1135.5
```

## 12.4 Estimation in R

**Table 2.** Parameter estimates and elasticities controlling for selection effects.

Variable	Parameter	Estimate	St. error	t ratio	Long-run total elasticity
$\ln E_{t-1}$	$\gamma$	0.2937	0.1041	2.821	
$\mathbf{W}_N \ln E_t$	$\rho_1$	0.6007	0.05877	10.22	
$\ln \text{Irish}_t$	$\beta_1$	-0.02064	0.02428	-0.8503	-0.0597
$\ln \text{Indian}_t$	$\beta_2$	-0.02757	0.02187	-1.261	-0.0797
$\ln \text{Pakistani}_t$	$\beta_3$	-0.03552	0.01915	-1.855	-0.1027
$\ln \text{European}_t$	$\beta_4$	0.0970	0.01497	6.48	0.2805
$\ln \text{RoW}_t$	$\beta_5$	0.03834	0.01609	2.383	0.1109
$\ln \text{UK}_t$	$\beta_6$	0.2337	0.05815	4.019	0.6759
$\ln \text{UnemploymentLQ}_t$	$\beta_7$	-0.1201	0.03638	-3.302	
$\mathbf{W}_N \ln E_{t-1}$	$\theta$	-0.2402	0.1106	-2.172	
	$\rho_2$	-0.2558	0.031753	-8.0988 <sup>2</sup>	
	$\sigma^2_\varepsilon$	0.0583			
	$\sigma^2_\mu$	0.4921			
Stationarity conditions					
$\rho_1 + \theta$		0.36058			
$\gamma + (\rho_1 + \theta)e_{\max}$		0.65425			
$\gamma + (\rho_1 + \theta)e_{\min}$		0.069051			
$\rho_1 - \theta$		0.84089			
$\gamma - (\rho_1 - \theta)e_{\max}$		-0.54722			
$\gamma - (\rho_1 - \theta)e_{\min}$		0.8175			
Max eig of $\mathbf{B}_N^{-1} \mathbf{C}_N$		0.32256			

Notes: Given a bootstrap sampling distribution, the GM estimation method for  $\rho_2$  is used to obtain 100 estimates under the null of zero error dependence and the mean and variance of the null distribution used to calculate the t ratio.

Figure 12.2: Impacts on employment, Fingleton, Olner, and Pryce (2020)

## 12 Spatio-temporal models

```
5 ALABAMA 1974      6 16762.67 8025.52 1734.85 7002.29 42057.31 33749 1169.8
6 ALABAMA 1975      6 17316.26 8158.23 1752.27 7405.76 43971.71 33604 1155.4
    unemp
1   4.7
2   5.2
3   4.7
4   3.9
5   5.5
6   7.7
```

```
usaww[1:10, 1:10]
```

	ALABAMA	ARIZONA	ARKANSAS	CALIFORNIA	COLORADO	CONNECTICUT	DELAWARE
ALABAMA	0.0 0.0000000	0	0.0	0.0	0	0	0
ARIZONA	0.0 0.0000000	0	0.2	0.2	0	0	0
ARKANSAS	0.0 0.0000000	0	0.0	0.0	0	0	0
CALIFORNIA	0.0 0.3333333	0	0.0	0.0	0	0	0
COLORADO	0.0 0.1428571	0	0.0	0.0	0	0	0
CONNECTICUT	0.0 0.0000000	0	0.0	0.0	0	0	0
DELAWARE	0.0 0.0000000	0	0.0	0.0	0	0	0
FLORIDA	0.5 0.0000000	0	0.0	0.0	0	0	0
GEORGIA	0.2 0.0000000	0	0.0	0.0	0	0	0
IDAHO	0.0 0.0000000	0	0.0	0.0	0	0	0
	FLORIDA	GEORGIA	IDAHO				
ALABAMA	0.25	0.25	0				
ARIZONA	0.00	0.00	0				
ARKANSAS	0.00	0.00	0				
CALIFORNIA	0.00	0.00	0				
COLORADO	0.00	0.00	0				
CONNECTICUT	0.00	0.00	0				
DELAWARE	0.00	0.00	0				
FLORIDA	0.00	0.50	0				
GEORGIA	0.20	0.00	0				

## 12.4 Estimation in R

```
IDAHO      0.00    0.00      0
```

**Produc** contains data on US States Production - a panel of 48 observations from 1970 to 1986. **usaww** is a spatial weights matrix of the 48 continental US States based on the queen contiguity relation.

Let start with an FE SEM model, using function **spml()** for maximum likelihood estimation of static spatial panel models.

```
# Gen listw object
usalw <- mat2listw(usaww, style = "W")

# Spec formula
fm <- log(gsp) ~ log(pcap) + log(pc) + log(emp) + unemp

### Estimate FE SEM model
semfe.mod <- spml(formula = fm, data = Produc,
                    index = c("state", "year"), # ID column
                    listw = usalw,             # listw
                    model = "within",         # one of c("within", "random", "pooling").
                    effect = "individual",   # type of fixed effects
                    lag = FALSE,              # spatial lag of Y
                    spatial.error = "b",       # "b" (Baltagi), "kkp" (Kapoor, Kelejian and Pr
                    method = "eigen",          # estimation method, for large data e.g. ("spam"
                    na.action = na.fail,      # handling of missings
                    zero.policy = NULL)        # handling of missings

summary(semfe.mod)
```

Spatial panel fixed effects error model

Call:

## 12 Spatio-temporal models

```
spml(formula = fm, data = Produc, index = c("state", "year"),
      listw = usalw, na.action = na.fail, model = "within", effect = "individual",
      lag = FALSE, spatial.error = "b", method = "eigen", zero.policy = NULL)

Residuals:
    Min. 1st Qu. Median 3rd Qu. Max.
-0.1246945 -0.0237699 -0.0034993 0.0170886 0.1882224

Spatial error parameter:
    Estimate Std. Error t-value Pr(>|t|)
rho 0.557401 0.033075 16.853 < 2.2e-16 ***
               Estimate Std. Error t-value Pr(>|t|)
log(pcap) 0.0051438 0.0250109 0.2057 0.83705
log(pc)    0.2053026 0.0231427 8.8712 < 2e-16 ***
log(emp)   0.7822540 0.0278057 28.1328 < 2e-16 ***
unemp     -0.0022317 0.0010709 -2.0839 0.03717 *
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

A RE SAR model, by contrast, can be estimated using the following options:

```
### Estimate an RE SAR model
sarre.mod <- spml(formula = fm, data = Produc,
                    index = c("state", "year"), # ID column
                    listw = usalw,             # listw
                    model = "random",          # one of c("within", "random", "random")
                    effect = "individual",    # type of fixed effects
                    lag = TRUE,                # spatial lag of Y
                    spatial.error = "none",    # "b" (Baltagi), "kkr" (Kapoor,
                    method = "eigen",          # estimation method, for large d
```

## 12.4 Estimation in R

```
    na.action = na.fail,      # handling of missings
    zero.policy = NULL)      # handling of missings

summary(sarre.mod)
```

ML panel with spatial lag, random effects

Call:

```
spreml(formula = formula, data = data, index = index, w = listw2mat(listw),
       w2 = listw2mat(listw2), lag = lag, errors = errors, cl = cl,
       method = "eigen", zero.policy = .2)
```

Residuals:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.38	1.57	1.70	1.70	1.80	2.13

Error variance parameters:

	Estimate	Std. Error	t-value	Pr(> t )
phi	21.3175	8.2929	2.5706	0.01015 *

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	Pr(> t )
lambda	0.161615	0.029042	5.5648	2.625e-08 ***

Coefficients:

	Estimate	Std. Error	t-value	Pr(> t )							
(Intercept)	1.65814987	0.15071855	11.0016	< 2.2e-16 ***							
log(pcap)	0.01294505	0.02493997	0.5190	0.6037							
log(pc)	0.22555375	0.02163422	10.4258	< 2.2e-16 ***							
log(emp)	0.67081074	0.02642113	25.3892	< 2.2e-16 ***							
unemp	-0.00579716	0.00089175	-6.5009	7.984e-11 ***							
---											
Signif. codes:	0	'***'	0.001	'**'	0.01	'*'	0.05	'.'	0.1	' '	1

## 12 Spatio-temporal models

Note that Millo and Piras (2012) use a different notation, namely  $\lambda$  for lag dependence, and  $\rho$  for error dependence....

Again, we have to use an additional step to get impacts for SAR-like models.

```
# Number of years
T <- length(unique(Produc$year))

# impacts
sarre.mod.imp <- impacts(sarre.mod,
                           listw = usalw,
                           time = T)
summary(sarre.mod.imp)

Impact measures (lag, trace):
      Direct      Indirect       Total
log(pcap) 0.013028574 0.002411880 0.015440454
log(pc)    0.227009032 0.042024438 0.269033470
log(emp)   0.675138835 0.124983264 0.800122098
unemp     -0.005834562 -0.001080108 -0.006914669
=====
Simulation results ( variance matrix):
Direct:

Iterations = 1:200
Thinning interval = 1
Number of chains = 1
Sample size per chain = 200

1. Empirical mean and standard deviation for each variable,
plus standard error of the mean:
```

## 12.4 Estimation in R

	Mean	SD	Naive SE	Time-series SE
log(pcap)	0.015148	0.0233842	1.654e-03	0.0014375
log(pc)	0.227589	0.0209440	1.481e-03	0.0013352
log(emp)	0.675362	0.0257468	1.821e-03	0.0018206
unemp	-0.005867	0.0008737	6.178e-05	0.0000552

2. Quantiles for each variable:

	2.5%	25%	50%	75%	97.5%
log(pcap)	-0.027689	-0.001279	0.015377	0.030439	0.064451
log(pc)	0.190018	0.213978	0.228329	0.239585	0.270045
log(emp)	0.615845	0.658976	0.678070	0.693001	0.718042
unemp	-0.007506	-0.006454	-0.005874	-0.005316	-0.004153

=====

Indirect:

```
Iterations = 1:200
Thinning interval = 1
Number of chains = 1
Sample size per chain = 200
```

1. Empirical mean and standard deviation for each variable,  
plus standard error of the mean:

	Mean	SD	Naive SE	Time-series SE
log(pcap)	0.002893	0.0046021	3.254e-04	2.814e-04
log(pc)	0.042586	0.0107716	7.617e-04	7.617e-04
log(emp)	0.126014	0.0279297	1.975e-03	1.975e-03
unemp	-0.001100	0.0003142	2.222e-05	2.222e-05

2. Quantiles for each variable:

	2.5%	25%	50%	75%	97.5%

## 12 Spatio-temporal models

```
log(pcap) -0.005854 -0.0003247 0.002498 0.0055664 0.0116637
log(pc)    0.024515  0.0349812 0.041803 0.0488222 0.0628312
log(emp)   0.076839  0.1060017 0.126582 0.1457542 0.1732680
unemp     -0.001771 -0.0012790 -0.001067 -0.0008751 -0.0005793
```

```
=====
```

Total:

```
Iterations = 1:200
Thinning interval = 1
Number of chains = 1
Sample size per chain = 200
```

1. Empirical mean and standard deviation for each variable,  
plus standard error of the mean:

	Mean	SD	Naive SE	Time-series SE
log(pcap)	0.018041	0.027857	1.970e-03	1.708e-03
log(pc)	0.270174	0.027756	1.963e-03	1.963e-03
log(emp)	0.801376	0.040368	2.854e-03	2.579e-03
unemp	-0.006966	0.001097	7.758e-05	6.966e-05

2. Quantiles for each variable:

	2.5%	25%	50%	75%	97.5%
log(pcap)	-0.033584	-0.001613	0.018430	0.03578	0.072426
log(pc)	0.220933	0.250828	0.268104	0.28661	0.331620
log(emp)	0.717083	0.772866	0.801234	0.82941	0.874786
unemp	-0.009268	-0.007654	-0.006954	-0.00625	-0.004821

There is an alternative by using the package **SDPDmod** by Rozeta Simonovska (see vignette).

## 12.4 Estimation in R

```
### FE SAR model
sarfe.mod2 <- SDPDm(formula = fm,
                      data = Produc,
                      W = usaww,
                      index = c("state","year"), # ID
                      model = "sar",           # on of c("sar","sdm"),
                      effect = "individual",   # FE structure
                      dynamic = FALSE,         # time lags of the dependet variable
                      LYtrans = TRUE)          # Lee-Yu transformation (bias correction)

summary(sarfe.mod2)
```

sar panel model with individual fixed effects

Call:

```
SDPDm(formula = fm, data = Produc, W = usaww, index = c("state",
  "year"), model = "sar", effect = "individual", dynamic = FALSE,
  LYtrans = TRUE)
```

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	Pr(> t )
rho	0.27856	0.02400	11.607	< 2.2e-16 ***

Coefficients:

	Estimate	Std. Error	t-value	Pr(> t )
log(pcap)	-0.0468700	0.0262162	-1.7878	0.0738 .
log(pc)	0.1859579	0.0237252	7.8380	4.578e-15 ***
log(emp)	0.6230728	0.0305554	20.3916	< 2.2e-16 ***
unemp	-0.0044701	0.0008917	-5.0130	5.359e-07 ***
---				

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

And subsequently, we can calculate the impacts of the model.

## 12 Spatio-temporal models

```
# Impacts
sarfe.mod2.imp <- impactsSDPDm(sarfe.mod2,
                                    NSIM = 200, # N simulations
                                    sd = 12345) # seed

summary(sarfe.mod2.imp)
```

Impact estimates for spatial (static) model

Direct:

	Estimate	Std. Error	t-value	Pr(> t )
log(pcap)	-0.04548734	0.02599122	-1.7501	0.0801 .
log(pc)	0.18811130	0.02383583	7.8920	2.975e-15 ***
log(emp)	0.63774644	0.03027064	21.0682	< 2.2e-16 ***
unemp	-0.00459922	0.00089695	-5.1276	2.935e-07 ***

Indirect:

	Estimate	Std. Error	t-value	Pr(> t )
log(pcap)	-0.01635496	0.00940025	-1.7398	0.08189 .
log(pc)	0.06724002	0.00990916	6.7856	1.156e-11 ***
log(emp)	0.22817291	0.02236760	10.2010	< 2.2e-16 ***
unemp	-0.00164883	0.00037103	-4.4440	8.831e-06 ***

Total:

	Estimate	Std. Error	t-value	Pr(> t )
log(pcap)	-0.0618423	0.0352523	-1.7543	0.07938 .
log(pc)	0.2553513	0.0313969	8.1330	4.188e-16 ***
log(emp)	0.8659194	0.0371907	23.2832	< 2.2e-16 ***
unemp	-0.0062481	0.0012311	-5.0752	3.871e-07 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

## 12.5 Example: Industrial facilities and municipal income

Note: I did not manage to estimate a dynamic panel model with SDPDm.

### 12.5 Example: Industrial facilities and municipal income

Rüttenauer and Best (2021): Environmental Inequality and Residential Sorting in Germany: A Spatial Time-Series Analysis of the Demographic Consequences of Industrial Sites. *Demography*, 58(6), 2243–2263.  
<https://doi.org/10.1215/00703370-9563077>

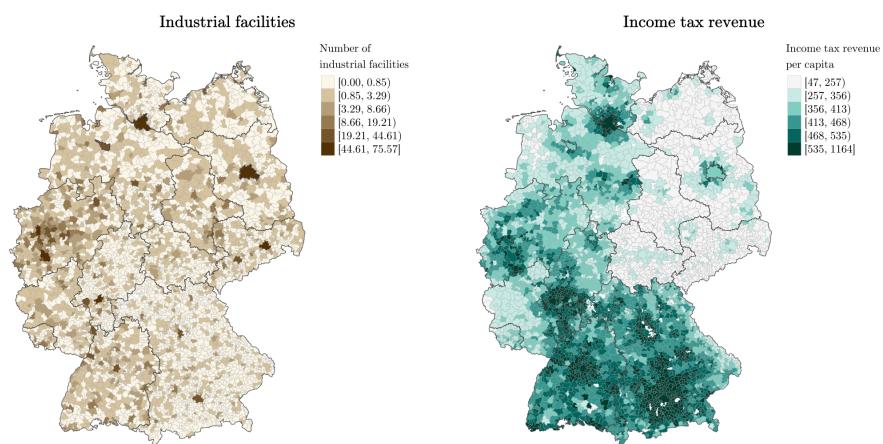


Figure 12.3: Spatial distribution of industrial facilities and income tax revenue per municipality for 2015.

## 12 Spatio-temporal models

Table 3: Fixed effects individual slopes (FEIS) estimator Dep. var.: Income tax revenue

	Overall (1)	West (2)	East (3)	Overall (4)	West (5)	East (6)
Number facilities <sub>t-1</sub>	-0.052** (0.020)	-0.080** (0.026)	0.006 (0.015)			
W Number facilities <sub>t-1</sub>	0.047 (0.044)	-0.033 (0.063)	0.213*** (0.046)			
Relative rank <sub>t-1</sub>				-0.027* (0.011)	-0.031* (0.015)	-0.010 (0.009)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.736	0.716	0.942	0.729	0.710	0.933
Adj. R <sup>2</sup>	0.736	0.716	0.942	0.729	0.710	0.933
Num. obs.	40095	31374	8721	40095	31374	8721
Num. groups: id	4455	3486	969	4455	3486	969

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$ , † $p < 0.1$ . Standardized coefficients. Cluster robust standard errors in parentheses. W is the spatially lagged coefficient. Controls: % aged 18 or younger, % aged 65 or above, population density, population density<sup>2</sup>, % of foreigners, year dummies (except year, all additionally included as spatial lag in models 1 - 3). Slopes for FEIS: trade tax revenue per capita, trade tax revenue per capita<sup>2</sup>.

# 13 Other Models

## 13.0.1 Required packages

```
pkgs <- c("sf", "mapview", "spdep", "spatialreg", "tmap", "viridisLite", "Gwmodel") # not  
lapply(pkgs, require, character.only = TRUE)
```

## 13.0.2 Session info

```
sessionInfo()
```

```
R version 4.4.1 (2024-06-14 ucrt)  
Platform: x86_64-w64-mingw32/x64  
Running under: Windows 11 x64 (build 22631)
```

```
Matrix products: default
```

```
locale:  
[1] LC_COLLATE=English_United Kingdom.utf8  
[2] LC_CTYPE=English_United Kingdom.utf8  
[3] LC_MONETARY=English_United Kingdom.utf8  
[4] LC_NUMERIC=C  
[5] LC_TIME=English_United Kingdom.utf8
```

### 13 Other Models

```
time zone: Europe/Berlin
tzcode source: internal

attached base packages:
[1] stats      graphics   grDevices utils      datasets   methods    base

other attached packages:
[1] GWmodel_2.3-2     Rcpp_1.0.12       sp_2.1-4        robustbase_0.99-
2
[5] viridisLite_0.4.2 tmap_3.3-4       spatialreg_1.3-
4 Matrix_1.7-0
[9] spdep_1.3-5      spData_2.3.1      mapview_2.11.2   sf_1.0-
16

loaded via a namespace (and not attached):
[1] xfun_0.45          raster_3.6-26      htmlwidgets_1.6.4 lattice_0.22-
6
[5] tools_4.4.1         crosstalk_1.2.1   LearnBayes_2.15.1 parallel_4.4.1
[9] stats4_4.4.1        sandwich_3.1-0    spacetime_1.3-
1 proxy_0.4-27
[13] DEoptimR_1.1-3     xts_0.14.0       KernSmooth_2.23-
24 satellite_1.0.5
[17] RColorBrewer_1.1-3 leaflet_2.2.2    lifecycle_1.0.4   FNN_1.1.4
[21] compiler_4.4.1     deldir_2.0-4     munsell_0.5.1   terra_1.7-
78
[25] codetools_0.2-20   leafsync_0.1.0   stars_0.6-5     htmltools_0.5.8
[29] class_7.3-22       MASS_7.3-60.2   classInt_0.4-
10 lwgeom_0.2-14
[33] wk_0.9.1          abind_1.4-5     boot_1.3-30    multcomp_1.4-
25
[37] nlme_3.1-164       digest_0.6.35   mvtnorm_1.2-
5 splines_4.4.1
[41] fastmap_1.2.0      grid_4.4.1      colorspace_2.1-
0 cli_3.6.2
```

### 13.1 Geographically weighted regression

```
[45] magrittr_2.0.3      base64enc_0.1-3    dichromat_2.0-
0.1 XML_3.99-0.16.1  leafem_0.2.3      TH.data_1.1-
[49] survival_3.6-4    e1071_1.7-14     rmarkdown_2.27   zoo_1.8-12
2           scales_1.3.0      evaluate_0.24.0  knitr_1.47      png_0.1-
[53] 8                  rlang_1.1.4       glue_1.7.0      tmaptools_3.1-
[57] coda_0.19-4.1      jsonlite_1.8.8   R6_2.5.1       DBI_1.2.3
1           [61] s2_1.1.6          intervals_0.15.4
[65] rstudioapi_0.16.0
[69] units_0.8-5
```

#### 13.0.3 Reload data from previous session

```
load("_data/msoa2_spatial.RData")
```

## 13.1 Geographically weighted regression

Does the relation between  $y$  and  $x$  vary depending on the region we are looking at? With geographically weighted regressions (GWR), we can exploit the spatial heterogeneity in relations / coefficients.

GWR (Brunsdon, Fotheringham, and Charlton 1996; Gollini et al. 2015) is mainly an explorative tool for spatial data analysis in which we estimate an equation at different geographical points. For  $L$  given locations across London, we receive  $L$  different coefficients.

$$\hat{\beta}_l = (\mathbf{X}^\top \mathbf{M}_l \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M}_l \mathbf{Y},$$

## 13 Other Models

The  $N \times N$  matrix  $\mathbf{M}_l$  defines the weights at each local point  $l$ , assigning higher weights to closer units. The local weights are determined by a kernel density function with a pre-determined bandwidth  $b$  around each point (either a fixed distance or an adaptive k nearest neighbours bandwidth). Models are estimated via `gwr.basic()` or `gwr.robust()` of the `GWmodel` package.

```
# Search for the optimal bandwidth
set.seed(123)
hv_1.bw <- bw.gwr(log(med_house_price) ~ log(no2) + log(POPDEN) + pubs_cou-
                     data = as_Spatial(msoa.spdf),
                     kernel = "boxcar",
                     adaptive = TRUE)
```

```
Adaptive bandwidth: 615 CV score: 117.989
Adaptive bandwidth: 388 CV score: 107.5287
Adaptive bandwidth: 247 CV score: 89.99347
Adaptive bandwidth: 160 CV score: 76.23795
Adaptive bandwidth: 106 CV score: 66.39574
Adaptive bandwidth: 73 CV score: 62.89816
Adaptive bandwidth: 52 CV score: 59.46008
Adaptive bandwidth: 39 CV score: 56.70472
Adaptive bandwidth: 31 CV score: 54.97107
Adaptive bandwidth: 26 CV score: 53.27627
Adaptive bandwidth: 23 CV score: 54.23635
Adaptive bandwidth: 28 CV score: 54.47944
Adaptive bandwidth: 25 CV score: 52.5378
Adaptive bandwidth: 24 CV score: 53.74594
Adaptive bandwidth: 25 CV score: 52.5378
```

```
hv_1.bw
```

### 13.1 Geographically weighted regression

[1] 25

```
### GWR
hv_1.gwr <- gwr.robust(log(med_house_price) ~ log(no2) + log(POPDEN) + pubs_count,
                           data = as_Spatial(msoa.spdf),
                           kernel = "boxcar",
                           adaptive = TRUE,
                           bw = hv_1.bw,
                           longlat = FALSE)
print(hv_1.gwr)
```

```
*****
*          Package    GWmodel           *
*****
Program starts at: 2024-06-29 17:44:40.25414
Call:
gwr.basic(formula = formula, data = data, bw = bw, kernel = kernel,
adaptive = adaptive, p = p, theta = theta, longlat = longlat,
dMat = dMat, F123.test = F123.test, cv = T, W.vect = W.vect,
parallel.method = parallel.method, parallel.arg = parallel.arg)

Dependent (y) variable: med_house_price
Independent variables: no2 POPDEN pubs_count
Number of data points: 983
*****
*          Results of Global Regression           *
*****
Call:
lm(formula = formula, data = data)

Residuals:
```

### 13 Other Models

```
      Min       1Q   Median     3Q    Max
-0.90930 -0.24801 -0.05018  0.20925  1.49660

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 10.630835  0.194980 54.523 < 2e-16 ***
log(no2)     0.716116  0.074916  9.559 < 2e-16 ***
log(POPDEN) -0.111104  0.023101 -4.809 1.75e-06 ***
pubs_count   0.008513  0.004209  2.023  0.0434 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 0.359 on 979 degrees of freedom
Multiple R-squared:  0.1177
Adjusted R-squared:  0.115
F-statistic: 43.55 on 3 and 979 DF,  p-value: < 2.2e-16
***Extra Diagnostic information
Residual sum of squares: 126.1498
Sigma(hat): 0.3585988
AIC: 781.3976
AICc: 781.459
BIC: -142.6963
*****
*          Results of Geographically Weighted Regression      *
*****
*****Model calibration information*****
Kernel function: boxcar
Adaptive bandwidth: 25 (number of nearest neighbours)
Regression points: the same locations as observations are used.
Distance metric: Euclidean distance metric is used.

*****Summary of GWR coefficient estimates:*****
      Min.    1st Qu.   Median   3rd Qu.    Max.

```

### 13.1 Geographically weighted regression

```
Intercept      -5.4436033 10.0869061 13.1177690 15.5252567 26.8582
log(no2)       -4.0174929 -0.7502331  0.0833824  1.0580475  6.2823
log(POPDEN)    -0.8610728 -0.3139087 -0.1368702 -0.0200140  0.4031
pubs_count     -0.3421595 -0.0252434 -0.0034784  0.0192272  0.2222
*****Diagnostic information*****
Number of data points: 983
Effective number of parameters (2trace(S) - trace(S'S)): 136.1695
Effective degrees of freedom (n-2trace(S) + trace(S'S)): 846.8305
AICc (GWR book, Fotheringham, et al. 2002, p. 61, eq 2.33): -
133.0433
AIC (GWR book, Fotheringham, et al. 2002, GWR p. 96, eq. 4.22): -
316.0801
BIC (GWR book, Fotheringham, et al. 2002, GWR p. 61, eq. 2.34): -
496.9587
Residual sum of squares: 36.33063
R-square value: 0.7459107
Adjusted R-square value: 0.7050051
*****Program stops at: 2024-06-29 17:44:47.805478
```

The results give a range of coefficients for different locations. Let's map those individual coefficients.

```
# Spatial object
gwr.spdf <- st_as_sf(hv_1.gwr$SDF)
gwr.spdf <- st_make_valid(gwr.spdf)

# Map
tmap_mode("view")

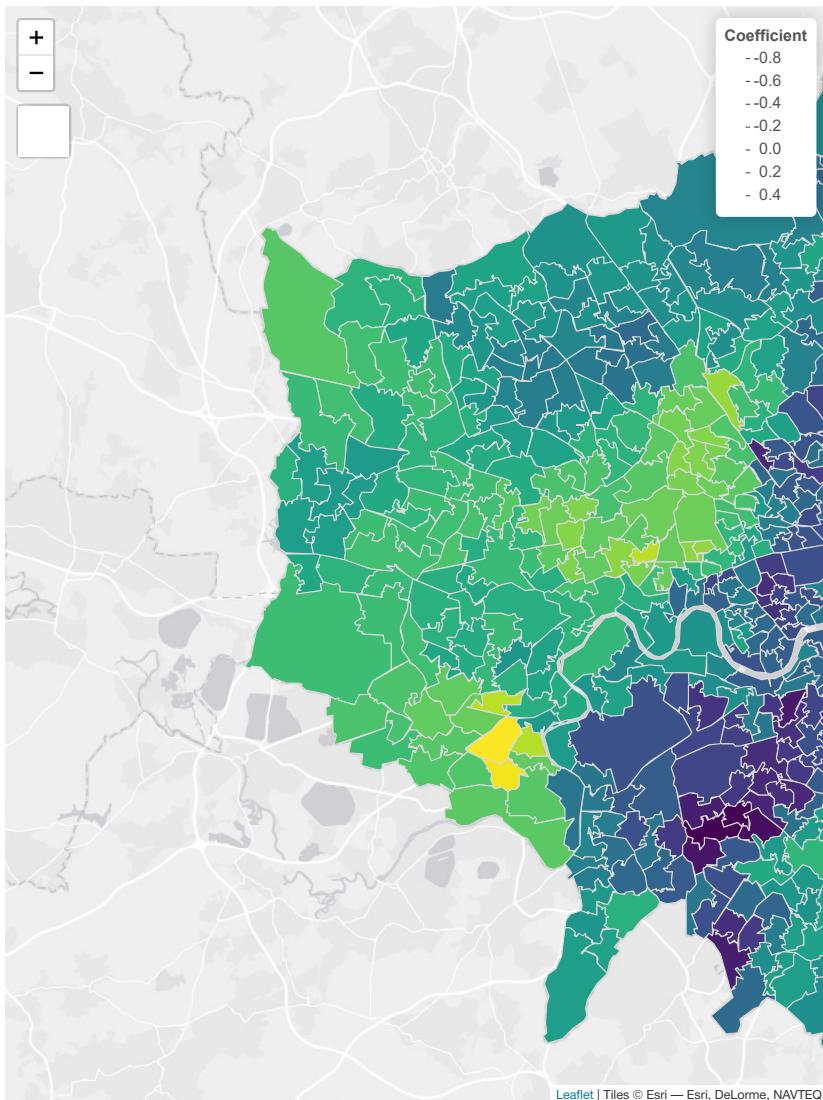
tmap mode set to interactive viewing
```

### 13 Other Models

```
mp2 <- tm_shape(gwr.spdf) +  
  tm_fill(col = "log(POPDEN)",  
          style = "cont",  
          # n = 8,  
          title = "Coefficient",  
          palette = viridis(100),  
          midpoint = TRUE, stretch.palette = TRUE) +  
  tm_borders(col = "grey85") +  
  tm_layout(legend.frame = TRUE, legend.bg.color = TRUE,  
            #legend.position = c("right", "bottom"),  
            legend.outside = TRUE,  
            main.title = "Coefficient of log population density",  
            main.title.position = "center",  
            title.snap.to.legend = TRUE)
```

mp2

### 13.1 Geographically weighted regression



## 13 Other Models

Just from looking at the map, there may be a connection with the underground network - the effect of population density on house values seems to be stronger / more positive where underground connection is weaker?!

### 13.2 Non-Linear Models

Models with endogenous regressors (SAR)

- In the literature: mostly spatial probit considered
- Spatial logit rather uncommon (non normally distributed errors)

Issues with non-linear spatial models

- Estimation: with dependent observations, we need to maximize one  $n$ -dimensional (log-)likelihood instead of a product of  $n$  independent distributions
- Estimation challenging and computationally intense
- Hard to interpret due to non-linear effects in non-linear models

Elhorst et al. (2017), Franzese, Hays, and Cook (2016)

#### 13.2.1 Problem with non-linear models

Spatial-SAR-Probit

$$\mathbf{y}^* = \rho \mathbf{W} \mathbf{y}^* + \mathbf{X} + \mathbf{y}_i = \{1 \text{ if } y_i^* > 0; 0 \text{ if } y_i^* \leq 0\}$$

or in reduced form:

$$\mathbf{y}^* = (\mathbf{I} - \rho \mathbf{W})^{-1} \mathbf{X} + \mathbf{u}, \quad \mathbf{u} = (\mathbf{I} - \rho \mathbf{W})^{-1}, \text{ with } \mathbf{u} \sim MVN(0, (\mathbf{I} - \rho \mathbf{W})^\top (\mathbf{I} - \rho \mathbf{W})^{-1})$$

## 13.2 Non-Linear Models

- Probability  $y^*$  is a latent variable, not observed
- We only observe binary outcome  $y_i$
- $\text{Cov}(y_i, y_j)$  is not the same as  $\text{Cov}(y_i^*, y_j^*)$
- Error term is heteroskedastic and spatially correlated

Probability

$$\text{Prob}[y^* > 0] = \text{Prob}[(\mathbf{I} - \rho \mathbf{W})^{-1} \mathbf{X} + (\mathbf{I} - \rho \mathbf{W})^{-1}] = \text{Prob}[(\mathbf{I} - \rho \mathbf{W})^{-1} < (\mathbf{I} - \rho \mathbf{W})^{-1} \mathbf{X}]$$

or in using the observed outcome:

$$\text{Prob}[y_i = 1 | \mathbf{X}] = \text{Prob}[u_i < [(\mathbf{I} - \rho \mathbf{W})^{-1} \mathbf{X}]_i] = \{[(\mathbf{I} - \rho \mathbf{W})^{-1} \mathbf{X}]_i / \sigma_{ui}\}$$

- $\{\}$  is an n-dimensional cumulative-normal distribution
- $\sigma_{ui}$  equals  $(\mathbf{I} - \rho \mathbf{W})^\top (\mathbf{I} - \rho \mathbf{W})^{-1}_{ii}$ , not constant
- no analytical solution

### 13.2.2 Estimation

Estimation methods for Spatial-SAR Probit / Logit

- Expectation Maximization (McMillen 1992).
- (Linearized) Generalized Methods of Moments (Klier and McMillen 2008).
- Recursive Importance Sampling (Beron and Vijverberg 2004).
- Maximum Simulated Likelihood RIS (Franzese, Hays, and Cook 2016)

## 13 Other Models

- Bayesian approach with Markov Chain Monte Carlo simulations (LeSage and Pace 2009): R package **spatialprobit**

Note that it can be hard to interpret the results. As in the linear case, it is necessary to compute the impacts. However, the ‘marginal’ effects may vary with values of the independent variables and the location (Lacombe and LeSage 2018).

### 13.2.3 Suggestion

In case you are not familiar with the econometric estimation methods and spatial linear models with AR term.

If necessary, I would recommend using `spatialprobit` relying on Bayesian MCMC, e.g. 7500 and 2500).

- So far, no ‘best practice’ guide
- No systematic comparison of estimation methods
- In R: Only **spatialprobit** provides impact measures?
- Hard to interpret results

Work-around: If the specification is theoretical plausible, using SLX probit / logit might be a practical solution!

## References

- Anselin, Luc. 1988. *Spatial Econometrics: Methods and Models*. Studies in Operational Regional Science. Dordrecht: Kluwer.
- . 1995. “Local Indicators of Spatial Association-LISA.” *Geographical Analysis* 27 (2): 93–115. <https://doi.org/10.1111/j.1538-4632.1995.tb00338.x>.
- . 2003. “Spatial Externalities, Spatial Multipliers, and Spatial Econometrics.” *International Regional Science Review* 26 (2): 153–66. <https://doi.org/10.1177/0160017602250972>.
- Anselin, Luc, and Anil K. Bera. 1998. “Spatial Dependence in Linear Regression Models with an Introduction to Spatial Econometrics.” In *Handbook of Applied Economic Statistics*, edited by Aman Ullah and David E. A. Giles, 237–89. New York: Dekker.
- Anselin, Luc, Anil K. Bera, Raymond Florax, and Mann J. Yoon. 1996. “Simple Diagnostic Tests for Spatial Dependence.” *Regional Science and Urban Economics* 26 (1): 77–104. [https://doi.org/10.1016/0166-0462\(95\)02111-6](https://doi.org/10.1016/0166-0462(95)02111-6).
- Anselin, Luc, and Nancy Lozano-Gracia. 2008. “Errors in Variables and Spatial Effects in Hedonic House Price Models of Ambient Air Quality.” *Empirical Economics* 34 (1): 5–34. <https://doi.org/10.1007/s00181-007-0152-3>.
- Anselin, Luc, Renan Serenini, and Pedro Amaral. 2024. “Spatial Econometric Model Specification Search: Another Look.” <https://doi.org/10.13140/RG.2.2.10650.86721>.
- Appelhans, Tim, Florian Detsch, Chritoph Reudenbach, and Stefan Woel-lauer. 2021. “Mapview: Interactive Viewing of Spatial Data in R.”
- Beron, Kurt J., and Wim P. M. Vijverberg. 2004. “Probit in a Spatial Con-

## References

- text: A Monte Carlo Analysis.” In *Advances in Spatial Econometrics: Methodology, Tools and Applications*, edited by Luc Anselin, Florax, Raymond J. G. M., and Sergio J. Rey, 169–95. Berlin, Heidelberg: Springer Berlin Heidelberg. [https://doi.org/10.1007/978-3-662-05617-2\\_8](https://doi.org/10.1007/978-3-662-05617-2_8).
- Betz, Timm, Scott J. Cook, and Florian M. Hollenbach. 2020. “Spatial Interdependence and Instrumental Variable Models.” *Political Science Research and Methods* 8 (4): 646–61. <https://doi.org/10.1017/psrm.2018.61>.
- Bivand, Roger S., and Colin Ruedel. 2018. “Rgeos: Interface to Geometry Engine - Open Source ('GEOS').”
- Bivand, Roger, Giovanni Millo, and Gianfranco Piras. 2021. “A Review of Software for Spatial Econometrics in R.” *Mathematics* 9 (11): 1276. <https://doi.org/10.3390/math911276>.
- Bivand, Roger, and Gianfranco Piras. 2015. “Comparing Implementations of Estimation Methods for Spatial Econometrics.” *Journal of Statistical Software* 63 (18): 1–36. <https://doi.org/10.18637/jss.v063.i18>.
- Bivand, Roger, and David W. S. Wong. 2018. “Comparing Implementations of Global and Local Indicators of Spatial Association.” *TEST* 27 (3): 716–48. <https://doi.org/10.1007/s11749-018-0599-x>.
- Boillat, Sébastien, M. Graziano Ceddia, and Patrick Bottazzi. 2022. “The Role of Protected Areas and Land Tenure Regimes on Forest Loss in Bolivia: Accounting for Spatial Spillovers.” *Global Environmental Change* 76 (September): 102571. <https://doi.org/10.1016/j.gloenvcha.2022.102571>.
- Brunsdon, Chris, A. Stewart Fotheringham, and Martin E. Charlton. 1996. “Geographically Weighted Regression: A Method for Exploring Spatial Nonstationarity.” *Geographical Analysis* 28 (4): 281–98. <https://doi.org/10.1111/j.1538-4632.1996.tb00936.x>.
- Burridge, Peter, J. Paul Elhorst, and Katarina Zigova. 2016. “Group Interaction in Research and the Use of General Nesting Spatial Models.” In *Spatial Econometrics: Qualitative and Limited Dependent Variables*, edited by Badi H. Baltagi, James P. LeSage, and R. Kelley Pace, 37:223–58. Advances in Econometrics. Emerald Group Publishing

## References

- Limited. <https://doi.org/10.1108/S0731-905320160000037016>.
- Cliff, Andrew, and Keith Ord. 1972. “Testing for Spatial Autocorrelation Among Regression Residuals.” *Geographical Analysis* 4 (3): 267–84. <https://doi.org/10.1111/j.1538-4632.1972.tb00475.x>.
- Cook, Scott J., Jude C. Hays, and Robert J. Franzese. 2020. “Model Specification and Spatial Interdependence.” In *The Sage Handbook of Research Methods in Political Science and International Relations*, edited by Luigi Curini and Robert Franzese, 1st ed, 730–47. Thousand Oaks: SAGE Inc.
- Croissant, Yves, and Giovanni Millo, eds. 2018. “Spatial Panels.” In *Panel Data Econometrics with R*, 245–84. Chichester, UK: John Wiley & Sons, Ltd. <https://doi.org/10.1002/9781119504641.ch10>.
- Drukker, David M., Peter Egger, and Ingmar R. Prucha. 2013. “On Two-Step Estimation of a Spatial Autoregressive Model with Autoregressive Disturbances and Endogenous Regressors.” *Econometric Reviews* 32 (5-6): 686–733. <https://doi.org/10.1080/07474938.2013.741020>.
- Elhorst, J. Paul. 2012. “Dynamic Spatial Panels: Models, Methods, and Inferences.” *Journal of Geographical Systems* 14 (1): 5–28. <https://doi.org/10.1007/s10109-011-0158-4>.
- . 2014. *Spatial Econometrics: From Cross-Sectional Data to Spatial Panels*. SpringerBriefs in Regional Science. Berlin and Heidelberg: Springer. <https://doi.org/10.1007/978-3-642-40340-8>.
- Elhorst, J. Paul, and S. Halleck Vega. 2017. “The SLX Model: Extensions and the Sensitivity of Spatial Spillovers to W.” *Papeles de Economía Española* 152: 34–50.
- Elhorst, J. Paul, Pim Heijnen, Anna Samarina, and Jan P. A. M. Jacobs. 2017. “Transitions at Different Moments in Time: A Spatial Probit Approach.” *Journal of Applied Econometrics* 32 (2): 422–39. <https://doi.org/10.1002/jae.2505>.
- Fingleton, Bernard, Daniel Olner, and Gwilym Pryce. 2020. “Estimating the Local Employment Impacts of Immigration: A Dynamic Spatial Panel Model.” *Urban Studies* 57 (13): 2646–62. <https://doi.org/10.1177/0042098019887916>.
- Fischer, Manfred M., Monika Bartkowska, Aleksandra Riedl, Sascha Sar-

## References

- dadvar, and Andrea Kunnert. 2009. “The Impact of Human Capital on Regional Labor Productivity in Europe.” *Letters in Spatial and Resource Sciences* 2 (2-3): 97–108. <https://doi.org/10.1007/s12076-009-0027-7>.
- Florax, Raymond, Hendrik Folmer, and Sergio J. Rey. 2003. “Specification Searches in Spatial Econometrics: The Relevance of Hendry’s Methodology.” *Regional Science and Urban Economics* 33 (5): 557–79. [https://doi.org/10.1016/S0166-0462\(03\)00002-4](https://doi.org/10.1016/S0166-0462(03)00002-4).
- Franzese, Robert J., and Jude C. Hays. 2007. “Spatial Econometric Models of Cross-Sectional Interdependence in Political Science Panel and Time-Series-Cross-Section Data.” *Political Analysis* 15 (2): 140–64. <https://doi.org/10.1093/pan/mpm005>.
- Franzese, Robert J., Jude C. Hays, and Scott J. Cook. 2016. “Spatial- and Spatiotemporal-Autoregressive Probit Models of Interdependent Binary Outcomes.” *Political Science Research and Methods* 4 (01): 151–73. <https://doi.org/10.1017/psrm.2015.14>.
- Gibbons, Steve, and Henry G. Overman. 2012. “Mostly Pointless Spatial Econometrics?” *Journal of Regional Science* 52 (2): 172–91. <https://doi.org/10.1111/j.1467-9787.2012.00760.x>.
- Gibbons, Steve, Henry G. Overman, and Eleonora Patacchini. 2015. “Spatial Methods.” In *Handbook of Regional and Urban Economics*, edited by Gilles Duranton, J. Vernon Henderson, and William C. Strange, 5:115–68. Amsterdam: Elsevier. <https://doi.org/10.1016/B978-0-444-59517-1.00003-9>.
- Gollini, Isabella, Binbin Lu, Martin Charlton, Christopher Brunsdon, and Paul Harris. 2015. “GWmodel : An R Package for Exploring Spatial Heterogeneity Using Geographically Weighted Models.” *Journal of Statistical Software* 63 (17). <https://doi.org/10.18637/jss.v063.i17>.
- Gräler, Benedikt, Edzer Pebesma, and Gerard Heuvelink. 2016. “Spatio-Temporal Interpolation Using Gstat.” *The R Journal* 8 (1): 204–18.
- Halleck Vega, Solmaria, and J. Paul Elhorst. 2015. “The SLX Model.” *Journal of Regional Science* 55 (3): 339–63. <https://doi.org/10.1111/jors.12188>.
- Kelejian, Harry H., and Gianfranco Piras. 2017. *Spatial Econometrics*. Elsevier. <https://doi.org/10.1016/C2016-0-04332-2>.

## References

- Kelejian, Harry H., and Ingmar R. Prucha. 1998. "A Generalized Spatial Two-Stage Least Squares Procedure for Estimating a Spatial Autoregressive Model with Autoregressive Disturbances." *The Journal of Real Estate Finance and Economics* 17 (1): 99–121. <https://doi.org/10.1023/A:1007707430416>.
- . 1999. "A Generalized Moments Estimator for the Autoregressive Parameter in a Spatial Model." *International Economic Review* 40 (2): 509–33. <https://doi.org/10.1111/1468-2354.00027>.
- . 2010. "Specification and Estimation of Spatial Autoregressive Models with Autoregressive and Heteroskedastic Disturbances." *Journal of Econometrics* 157 (1): 53–67. <https://doi.org/10.1016/j.jeconom.2009.10.025>.
- Kelejian, Harry H., Ingmar R. Prucha, and Yevgeny Yuzefovich. 2004. "Instrumental Variable Estimation of a Spatial Autoregressive Model with Autoregressive Disturbances: Large and Small Sample Results." In *Spatial and Spatiotemporal Econometrics*, edited by James P. LeSage and R. Kelley Pace, 163–98. Advances in Econometrics. Amsterdam and Boston: Elsevier.
- Kley, Stefanie, and Tetiana Dovbushchuk. 2021. "How a Lack of Green in the Residential Environment Lowers the Life Satisfaction of City Dwellers and Increases Their Willingness to Relocate." *Sustainability* 13 (7): 3984. <https://doi.org/10.3390/su13073984>.
- Klier, Thomas, and Daniel P. McMillen. 2008. "Clustering of Auto Supplier Plants in the United States: Generalized Method of Moments Spatial Logit for Large Samples." *Journal of Business & Economic Statistics* 26 (4): 460–71.
- Lacombe, Donald J., and James P. LeSage. 2018. "Use and Interpretation of Spatial Autoregressive Probit Models." *The Annals of Regional Science* 60 (1): 1–24. <https://doi.org/10.1007/s00168-015-0705-x>.
- Lee, Barrett A., Sean F. Reardon, Glenn Firebaugh, Chad R. Farrell, Stephen A. Matthews, and David O'Sullivan. 2008. "Beyond the Census Tract: Patterns and Determinants of Racial Segregation at Multiple Geographic Scales." *American Sociological Review* 73 (5): 766–91. <https://doi.org/10.1177/000312240807300504>.

## References

- Lee, Lung-fei. 2004. “Asymptotic Distributions of Quasi-Maximum Likelihood Estimators for Spatial Autoregressive Models.” *Econometrica* 72 (6): 1899–1925.
- Lee, Lung-fei, and Jihai Yu. 2010. “Estimation of Spatial Autoregressive Panel Data Models with Fixed Effects.” *Journal of Econometrics* 154 (2): 165–85. <https://doi.org/10.1016/j.jeconom.2009.08.001>.
- LeSage, James P. 2014a. “What Regional Scientists Need to Know about Spatial Econometrics.” *The Review of Regional Studies* 44 (1): 13–32. <https://doi.org/https://dx.doi.org/10.2139/ssrn.2420725>.
- . 2014b. “Spatial Econometric Panel Data Model Specification: A Bayesian Approach.” *Spatial Statistics* 9 (August): 122–45. <https://doi.org/10.1016/j.spasta.2014.02.002>.
- LeSage, James P., and R. Kelley Pace. 2009. *Introduction to Spatial Econometrics*. Statistics, Textbooks and Monographs. Boca Raton: CRC Press.
- . 2014. “The Biggest Myth in Spatial Econometrics.” *Econometrics* 2 (4): 217–49. <https://doi.org/10.3390/econometrics2040217>.
- Liebe, Ulf, Sander van Cranenburgh, and Caspar Chorus. 2023. “Maximizing Utility or Avoiding Losses? Uncovering Decision Rule-Heterogeneity in Sociological Research with an Application to Neighbourhood Choice.” *Sociological Methods & Research*, July, 00491241231186657. <https://doi.org/10.1177/00491241231186657>.
- Lovelace, Robin, Jakub Nowosad, and Jannes Muenchow. 2019. *Geocomputation with R*. 1st ed. Chapman & Hall/CRC the R Series. Boca Raton: Chapman & Hall/CRC.
- Manski, Charles F. 1993. “Identification of Endogenous Social Effects: The Reflection Problem.” *The Review of Economic Studies* 60 (3): 531–42. <https://doi.org/10.2307/2298123>.
- McMillen, Daniel P. 1992. “Probit with Spatial Autocorrelation.” *Journal of Regional Science* 32 (3): 335–48. <https://doi.org/10.1111/j.1467-9787.1992.tb00190.x>.
- Millo, Giovanni, and Gianfranco Piras. 2012. “Splm: Spatial Panel Data Models in R.” *Journal of Statistical Software* 47 (1). <https://doi.org/10.18637/jss.v047.i01>.

## References

- Mohai, Paul, and Robin Saha. 2007. “Racial Inequality in the Distribution of Hazardous Waste: A National-Level Reassessment.” *Social Problems* 54 (3): 343–70. <https://doi.org/10.1525/sp.2007.54.3.343>.
- Moran, P. A. P. 1950. “Notes on Continuous Stochastic Phenomena.” *Biometrika* 37 (1/2): 17. <https://doi.org/10.2307/2332142>.
- Mur, Jesús, and Ana Angulo. 2009. “Model Selection Strategies in a Spatial Setting: Some Additional Results.” *Regional Science and Urban Economics* 39 (2): 200–213. <https://doi.org/10.1016/j.regsciurbeco.2008.05.018>.
- Neumayer, Eric, and Thomas Plümper. 2016. “W.” *Political Science Research and Methods* 4 (01): 175–93. <https://doi.org/10.1017/psrm.2014.40>.
- Ord, John Keith. 1975. “Estimation Methods for Models of Spatial Interaction.” *Journal of the American Statistical Association* 70 (349): 120–26. <https://doi.org/10.2307/2285387>.
- Pace, R. Kelley, and James P. LeSage. 2010. “Omitted Variable Biases of OLS and Spatial Lag Models.” In *Progress in Spatial Analysis*, edited by Antonio Páez, Julie Gallo, Ron N. Buliung, and Sandy Dall’erba, 17–28. Berlin and Heidelberg: Springer.
- Pebesma, Edzer. 2018. “Simple Features for R: Standardized Support for Spatial Vector Data.” *The R Journal* 10 (1): 439. <https://doi.org/10.32614/RJ-2018-009>.
- Pebesma, Edzer, and Roger Bivand. 2023. *Spatial Data Science: With Applications in R*. First. Boca Raton: Chapman and Hall/CRC. <https://doi.org/10.1201/9780429459016>.
- Pinkse, Joris, and Margaret E. Slade. 2010. “The Future of Spatial Econometrics.” *Journal of Regional Science* 50 (1): 103–17. <https://doi.org/10.1111/j.1467-9787.2009.00645.x>.
- Rüttenauer, Tobias. 2018. “Neighbours Matter: A Nation-wide Small-area Assessment of Environmental Inequality in Germany.” *Social Science Research* 70: 198–211. <https://doi.org/10.1016/j.ssresearch.2017.11.009>.
- . 2022. “Spatial Regression Models: A Systematic Comparison of Different Model Specifications Using Monte Carlo Experiments.”

## References

- Sociological Methods & Research* 51 (2): 728–59. <https://doi.org/10.1177/0049124119882467>.
- . 2024. “Spatial Data Analysis.” arXiv. <https://arxiv.org/abs/2402.09895>.
- Rüttenauer, Tobias, and Henning Best. 2021. “Environmental Inequality and Residential Sorting in Germany: A Spatial Time-Series Analysis of the Demographic Consequences of Industrial Sites.” *Demography* 58 (6): 2243–63. <https://doi.org/10.1215/00703370-9563077>.
- Sarriás, Mauricio. 2023. *Intermediate Spatial Econometrics with Applications in R*.
- Tennekes, Martijn. 2018. “Tmap : Thematic Maps in R.” *Journal of Statistical Software* 84 (6). <https://doi.org/10.18637/jss.v084.i06>.
- Tobler, Waldo R. 1970. “A Computer Movie Simulating Urban Growth in the Detroit Region.” *Economic Geography* 46: 234–40. <https://doi.org/10.2307/143141>.
- Ward, Michael Don, and Kristian Skrede Gleditsch. 2008. *Spatial Regression Models*. Vol. 155. Quantitative Applications in the Social Sciences. Thousand Oaks: Sage.
- Wimpy, Cameron, Guy D. Whitten, and Laron K. Williams. 2021. “X Marks the Spot: Unlocking the Treasure of Spatial-X Models.” *The Journal of Politics* 83 (2): 722–39. <https://doi.org/10.1086/710089>.
- Wong, David. 2009. “The Modifiable Areal Unit Problem (MAUP).” In *The Sage Handbook of Spatial Analysis*, edited by A. Stewart Fotheringham and Peter Rogerson, 105–24. Los Angeles and London: Sage.
- Wooldridge, Jeffrey M. 2010. *Econometric Analysis of Cross Section and Panel Data*. Cambridge, Mass.: MIT Press.