

Introduction to Spatial Data in R

V. Gómez-Rubio

Departamento de Matemáticas
Escuela de Ingenieros Industriales de Albacete
U. de Castilla-La Mancha

Imperial, 16 May 2014

based on work by Roger S. Bivand, Edzer Pebesma and H. Rue

Why spatial data in R?

- What is R, and why should we pay the price of using it?
- How does the community around R work, what are its shared principles?
- How does applied spatial data analysis fit into R?
- But I have a non-standard research question ...

Where do we find spatial problems?

Geography How are settlements located according to the presence of natural resources, mountains, rivers, etc?

Econometrics Where are flats more expensive in a city?

Ecology How are different species of trees distributed in a forest?

Epidemiology How does the risk of suffering from a particular disease change with location? Is high risk linked to the presence of some pollution sources?

Environmetrics How can we produce an estimation of the pollution in the air from samples obtained at a set of stations?

Public Policy Where is unemployment higher? What regions should benefit from certain types of policies?

Applied spatial data analysis with R

- R can be used to tackle most of these problems, at least initially...
- Packages for importing commonly encountered spatial data formats
- Range of contributed packages in spatial statistics and increasing awareness of the importance of spatial data analysis in the broader community
- Current contributed packages with spatial statistics applications (see R spatial projects):

point patterns: **spatstat**, **VR:spatial**, **spIancs**;

geostatistics: **gstat**, **geoR**, **geoRglm**, **fields**, **spBayes**, **RandomFields**,
VR:spatial, **sgeostat**, **vardiag**;

lattice/area data: **spdep**, **DCluster**, **spgwr**, **ade4**.

modelling tools: **mgcv**, **INLA**, **R2WinBUGS**, **R2BayesX**.

A John Snow example

Even though we know that John Snow already had a working hypothesis about Cholera epidemics, his data remain interesting, especially if we use a GIS (GRASS) to find the street distances from mortality dwellings to the Broad Street pump:

```
v.digit -n map=vsnow4 bgcmd="d.rast map=snow"  
v.to.rast input=vsnow4 output=rfsnow use=val value=1  
r.buffer input=rfsnow output=buff2 distances=4  
r.cost -v input=buff2 output=snowcost_not_broad \  
    start_points=vpump_not_broad  
r.cost -v input=buff2 output=snowcost_broad start_points=vpump_broad
```

We have two raster layers of cost distances along the streets, one distances from the Broad Street pump, the other distances from other pumps.

Cholera mortalities, Soho

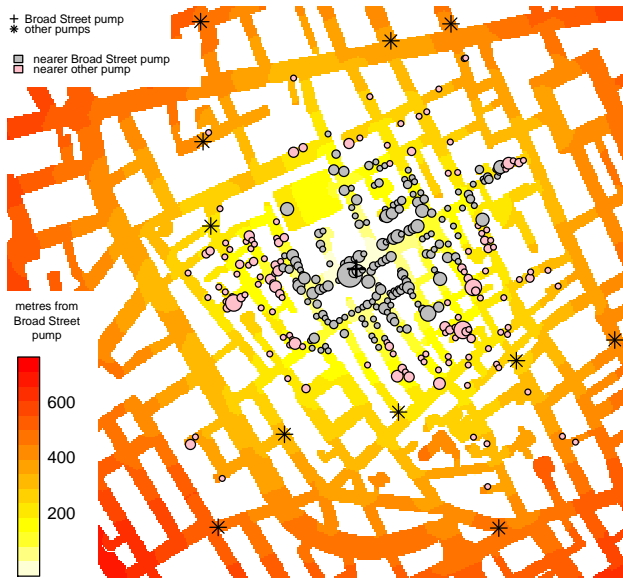
We have read from GRASS into R a point layer of mortalities (counts per address) called `death`, the two distance cost raster layers, and the point locations of the pumps. Overlaying the addresses on the raster, we can pick off the street distances from each address to the nearest pump, and create a new variable `b_nearer`. Using this variable, we can tally the mortalities by nearest pump:

```
> o <- overlay(sohoSG, deaths)
> deaths <- spCbind(deaths, as(o, "data.frame"))
> deaths$b_nearer <- deaths$snowcost_broad < deaths$snowcost_not_broad
> by(deaths$Num_Cases, deaths$b_nearer, sum)
```

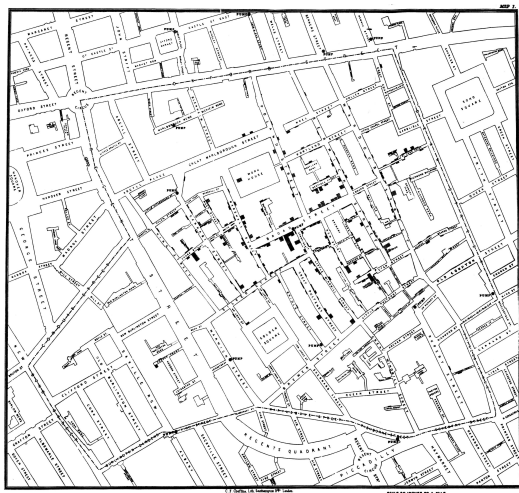
```
deaths$b_nearer: FALSE
[1] 221
```

```
-----
deaths$b_nearer: TRUE
[1] 357
```

Cholera mortalities, Soho



Cholera mortalities, Soho



Source: Wikipedia

Analysing Spatial Data in R

- The background for the tutorial is provided in the R News note by Roger Bivand and Edzer Pebesma, November 2005, pp. 9–13 and the book by Bivand et al. (2008)
- First we'll look at the representation of spatial data in R, with stress on the classes provided in the **sp** package
- After that, we'll see how to read and write spatial data in commonly-used exchange formats, and how to handle coordinate reference systems
- An introduction to spatial analysis using some spatial econometrics will be given and disease mapping models will come next
- Then we will see how to work with point data and analyse point patterns
- We'll show how analysis packages for geostatistics are being adapted to use these representations directly
- Finally, we will move on to the spatio-temporal case

Workshop infrastructure

- Task views are one of the nice innovations on CRAN that help navigate in the jungle of contributed packages — the Spatial task view is a useful resource
- The task view is also a point of entry to the Rgeo website hosted off CRAN, and updated quite often; it tries to mention in more detail contributed packages for spatial data analysis
- It also provides a link to the **sp** development area on Sourceforge, with CVS access to **sp**
- Finally, it links to the R-sig-geo mailing list, which is the preferred place to ask questions about analysing spatial/geographical data
- Additional resources can be found at web site related to the book by Bivand et al. (2008):
<http://www.asdar-book.org>

Analysing Spatial Data in R: Representing Spatial Data

Let's start with 3 examples...

- Point patterns** Location of the starting point of tornados in the US in 2009. Where do tornados tend to appear more often?
- Geostatistics** Study of the distribution of heavy metals around the Meuse river (in the border between Belgium and the Netherlands)
- Lattice Data** Analysis of the cases of sudden infant death syndrome in North Carolina and its possible link to the ethnic distribution of the population

What type of data do we need?

- Point patterns** Coordinates of the points and, possibly, a boundary to bound the study region.
Sometimes a data.frame with more information related to each tornado (state, date, time, EF scale, Economic Loss, etc.)
- Geostatistics** Coordinates of the sampling points plus levels of heavy metals at those points.
Possibly, several layers describing the type of terrain
- Lattice Data** Boundaries for each area in the study region.
Attached data to each area may be available as well (for example, population, etc.)

What type of analysis do we need?

- Point patterns** Estimates of the spatial distribution of tornados. A surface with the probability of occurrence is often used.
- Geostatistics** Methods for predicting the concentration of heavy metals over the study region (usually, a grid is used). Common methods include interpolation, kriging, and others.
- Lattice Data** Estimates of some parameter of interest for each area. These are often based on linear models (LMs, GLMs, GLMMs, GAMs, etc.)

Object framework

- To begin with, all contributed packages for handling spatial data in R had different representations of the data.
- This made it difficult to exchange data both within R between packages, and between R and external file formats and applications.
- The result has been an attempt to develop shared classes to represent spatial data in R, allowing some shared methods and many-to-one, one-to-many conversions.
- Bivand and Pebesma chose to use new-style classes (S4) to represent spatial data, and are confident that this choice was justified.

Spatial objects

- The foundation object is the `Spatial` class, with just two slots (new-style class objects have pre-defined components called slots)
- The first is a bounding box, and is mostly used for setting up plots
- The second is a CRS class object defining the coordinate reference system, and may be set to `CRS(as.character(NA))`, its default value.
- Operations on `Spatial*` objects should update or copy these values to the new `Spatial*` objects being created

Spatial points

- The most basic spatial data object is a point, which may have 2 or 3 dimensions
- A single coordinate, or a set of such coordinates, may be used to define a `SpatialPoints` object; coordinates should be of mode `double` and will be promoted if not
- The points in a `SpatialPoints` object may be associated with a row of attributes to create a `SpatialPointsDataFrame` object
- The coordinates and attributes may, but do not have to be keyed to each other using ID values

Tornado Data 2009

- We will use some Tornado data to show the analysis of point patterns
- These data have been obtained from the *Storm Prediction Center*^a
- Tornado data from 1955 until 2009 are available
- In addition to the coordinates, we have a wealth of related information for each tornado

^a<http://www.spc.noaa.gov/wcm/index.html#data>



Spatial points

The Tornado data are provided in a cvs file that we can read to make a `SpatialPoints` object.

```
> library(sp)
> d<-read.csv(file="datasets/2009_torn.csv", header=FALSE)
> #Names obtained from accompanying description file
> names(d)<-c("Number", "Year", "Month", "Day", "Date",
+   "Time", "TimeZone", "State", "FIPS", "StateNumber",
+   "EFscale", "Injuries", "Fatalities", "Loss", "Closs",
+   "SLat", "SLon", "ELat", "ELon", "Length",
+   "Width", "NStates", "SNumber", "SG", "1FIPS",
+   "2FIPS", "3FIPS", "4FIPS")
> coords <- SpatialPoints(d[, c("SLon", "SLat")])
> summary(coords)
```

Object of class `SpatialPoints`

Coordinates:

```
      min max
SLon -158.064  0
SLat   0.000 49
Is projected: NA
proj4string : [NA]
Number of points: 1182
```

Spatial points

Now we'll add the original data frame to make a `SpatialPointsDataFrame` object. Many methods for standard data frames just work with `SpatialPointsDataFrame` objects.

```
> storn <- SpatialPointsDataFrame(coords, d)
> names(storn)
```

[1] "Number"	"Year"	"Month"	"Day"	"Date"
[6] "Time"	"TimeZone"	"State"	"FIPS"	"StateNumber"
[11] "EFscale"	"Injuries"	"Fatalities"	"Loss"	"Closs"
[16] "SLat"	"SLon"	"ELat"	"ELon"	"Length"
[21] "Width"	"NStates"	"SNumber"	"SG"	"1FIPS"
[26] "2FIPS"	"3FIPS"	"4FIPS"		

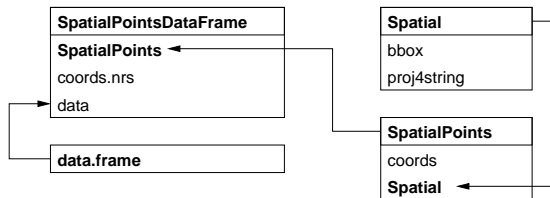
```
> summary(storn$Fatalities)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.00000	0.00000	0.00000	0.02538	0.00000	8.00000

```
> table(storn$Month)
```

1	2	3	4	5	6	7	8	9	10	11	12
6	38	117	234	201	274	125	60	8	66	3	50

Spatial points classes and their slots



Spatial lines and polygons

- A `Line` object is just a spaghetti collection of 2D coordinates; a `Polygon` object is a `Line` object with equal first and last coordinates
- A `Lines` object is a list of `Line` objects, such as all the contours at a single elevation; the same relationship holds between a `Polygons` object and a list of `Polygon` objects, such as islands belonging to the same county
- `SpatialLines` and `SpatialPolygons` objects are made using lists of `Lines` or `Polygons` objects respectively
- `SpatialLinesDataFrame` and `SpatialPolygonsDataFrame` objects are defined using `SpatialLines` and `SpatialPolygons` objects and standard data frames, and the ID fields are here required to match the data frame row names

Spatial lines: Tornado trajectories

- The Tornado data includes starting and ending points of the tornado
- Although we know that tornados do not follow a straight line, a line can be used to represent the path that the tornado followed

```
> sl<-lapply(unique(d$Number), function(X){
+   dd<-d[which(d$Number==X),c("SLon", "SLat", "ELon", "ELat")]
+
+
+   L<-lapply(1:nrow(dd), function(i){
+     Line(matrix(as.numeric(dd[i,]),ncol=2, byrow=TRUE))
+   })
+   Lines(L, ID=as.character(X))
+ })
> Tl<-SpatialLines(sl)
> summary(Tl)
```

Object of class SpatialLines

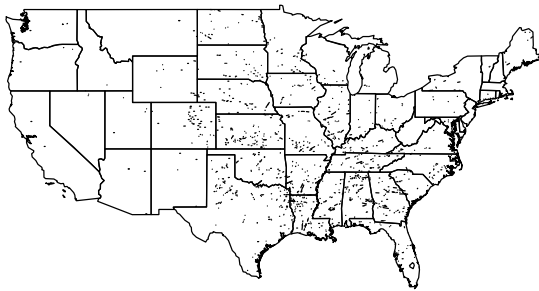
Coordinates:

```
      min max
x -158.064  0
y   0.000 49
Is projected: NA
proj4string : [NA]
```

Spatial polygons: US states boundaries

The *Storm Prediction Center* also provides maps with the boundaries of the US states. These can be used to place data into context by displaying the starting points of the tornados over a map of some of the US states:

```
> load("datasets/statesth.RData")  
> plot(statesth)  
> plot(T1, add=TRUE)
```



Spatial lines

There is a helper function `contourLines2SLDF` to convert the list of contours returned by `contourLines` into a `SpatialLinesDataFrame` object. This example shows how the data slot row names match the ID slot values of the set of `Lines` objects making up the `SpatialLinesDataFrame`, note that some `Lines` objects include multiple `Line` objects:

```
> library(maptools)

> volcano_sl <- ContourLines2SLDF(contourLines(volcano))
> row.names(slot(volcano_sl, "data"))

[1] "C_1" "C_2" "C_3" "C_4" "C_5" "C_6" "C_7" "C_8" "C_9" "C_10"

> sapply(slot(volcano_sl, "lines"), function(x) slot(x, "ID"))

[1] "C_1" "C_2" "C_3" "C_4" "C_5" "C_6" "C_7" "C_8" "C_9" "C_10"

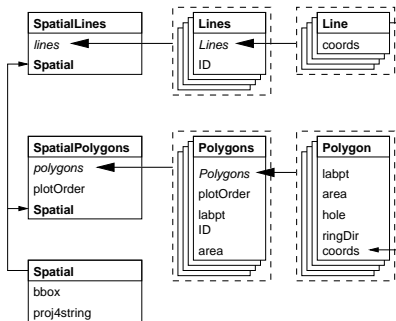
> sapply(slot(volcano_sl, "lines"), function(x) length(slot(x, "Lines")))

[1] 3 4 1 1 1 2 2 3 2 1

> volcano_sl$level

[1] 100 110 120 130 140 150 160 170 180 190
Levels: 100 110 120 130 140 150 160 170 180 190
```

Spatial Polygons classes and slots



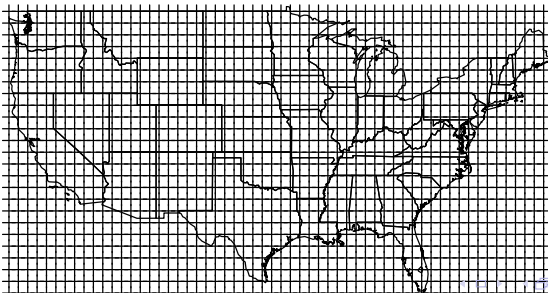
Spatial grids and pixels

- There are two representations for data on regular rectangular grids (oriented N-S, E-W): `SpatialPixels` and `SpatialGrid`
- `SpatialPixels` are like `SpatialPoints` objects, but the coordinates have to be regularly spaced; the coordinates are stored, as are grid indices
- `SpatialPixelsDataFrame` objects only store attribute data where it is present, but need to store the coordinates and grid indices of those grid cells
- `SpatialGridDataFrame` objects do not need to store coordinates, because they fill the entire defined grid, but they need to store NA values where attribute values are missing

Spatial grids

In a point pattern analysis the intensity of the underlying process is often estimated in the study region. This often requires using a grid so that the spatial intensity is computed. A grid can be defined as follows:

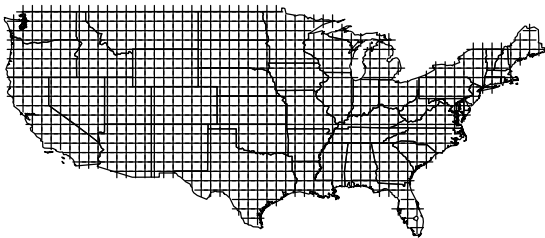
```
> #Step of the grid
> h<-1
> xrange<-diff(bbox(statesth)[1,])
> yrange<-diff(bbox(statesth)[2,])
> nx<-ceiling( xrange/h )
> ny<-ceiling(yrange/h)
> grdtop<-GridTopology(cellcentre.offset=bbox(statesth)[,1],
+   cellsize=c(h,h), cells.dim=c(nx,ny))
> grd<-SpatialGrid(grdtop, proj4string=CRS("+proj=longlat") )
>
> plot(grd)
> plot(statesth, add=TRUE)
```



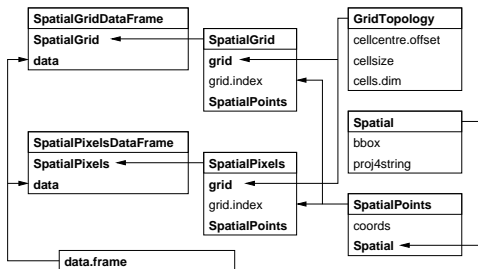
Spatial pixels

Although this is a convenient way of creating and using grids, some of the points fall in the middle of the ocean, far from any US state. We could combine the US boundaries and the grid to keep only the points that are inside the US boundaries in a `SpatialPixels` object:

```
> grdidx<-overlay(grd, statesth)
> grd2<-SpatialPixels(SpatialPoints(coordinates(grd))[!is.na(grdidx),] )
> proj4string(grd)<-CRS("+proj=longlat")
>
> plot(grd2)
> plot(statesth, add=TRUE)
```



Spatial grid and pixels classes and their slots



Spatial classes provided by **sp**

This table summarises the classes provided by **sp**, and shows how they build up to the objects of most practical use, the `Spatial*DataFrame` family objects:

data type	class	attributes	extends
points	<code>SpatialPoints</code>	none	<code>Spatial</code>
points	<code>SpatialPointsDataFrame</code>	<code>data.frame</code>	<code>SpatialPoints</code>
pixels	<code>SpatialPixels</code>	none	<code>SpatialPoints</code>
pixels	<code>SpatialPixelsDataFrame</code>	<code>data.frame</code>	<code>SpatialPixels</code> <code>SpatialPointsDataFrame</code>
full grid	<code>SpatialGrid</code>	none	<code>SpatialPixels</code>
full grid	<code>SpatialGridDataFrame</code>	<code>data.frame</code>	<code>SpatialGrid</code>
line	<code>Line</code>	none	
lines	<code>Lines</code>	none	Line list
lines	<code>SpatialLines</code>	none	<code>Spatial</code> , Lines list
lines	<code>SpatialLinesDataFrame</code>	<code>data.frame</code>	<code>SpatialLines</code>
polygon	<code>Polygon</code>	none	Line
polygons	<code>Polygons</code>	none	Polygon list
polygons	<code>SpatialPolygons</code>	none	<code>Spatial</code> , Polygons list
polygons	<code>SpatialPolygonsDataFrame</code>	<code>data.frame</code>	<code>SpatialPolygons</code>

Methods provided by **sp**

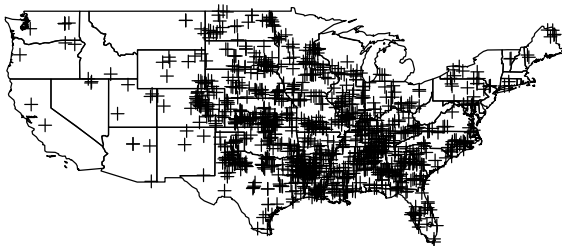
This table summarises the methods provided by **sp**:

method	what it does
[select spatial items (points, lines, polygons, or rows/cols from a grid) and/or attributes variables
\$, \$<-, [[, [[<-	retrieve, set or add attribute table columns
spsample	sample points from a set of polygons, on a set of lines or from a gridded area
bbox	get the bounding box
proj4string	get or set the projection (coordinate reference system)
coordinates	set or retrieve coordinates
coerce	convert from one class to another
overlay	combine two different spatial objects

Overlying tornados and US states

The Tornado data comprises tornados occurred in Puerto Rico, Alaska and other regions or states. In order to select only those tornados found in the main continental region of the US, we can do an overlay:

```
> sidx<-overlay(storn, statesth)  
> storn2<-storn[!is.na(sidx),]  
>  
  
> plot(storn2)  
> plot(statesth, add=TRUE)
```



Using `Spatial` family objects

- Very often, the user never has to manipulate `Spatial` family objects directly, as we have been doing here, because methods to create them from external data are also provided
- Because the `Spatial*DataFrame` family objects behave in most cases like data frames, most of what we are used to doing with standard data frames just works — like `"["` or `$` (but no `merge`, etc., yet)
- These objects are very similar to typical representations of the same kinds of objects in geographical information systems, so they do not suit spatial data that is not geographical (like medical imaging) as such
- They provide a standard base for analysis packages on the one hand, and import and export of data on the other, as well as shared methods, like those for visualisation we turn to now

Analysing Spatial Data in R: Vizualising Spatial Data

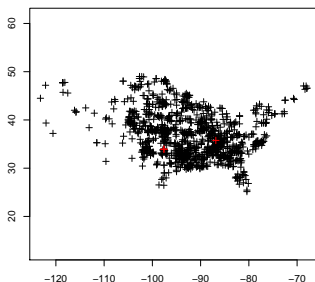
Vizualising Spatial Data

- Displaying spatial data is one of the chief reasons for providing ways of handling it in a statistical environment
- Of course, there will be differences between analytical and presentation graphics here as well — the main point is to get a usable display quickly, and move to presentation quality cartography later
- In general, maintaining aspect is vital, and that can be done in both base and lattice graphics in R (note that both **sp** and **maps** display methods for spatial data with geographical coordinates “stretch” the y-axis)
- We'll look at the basic methods for displaying spatial data in **sp**; other packages have their own methods, but the next unit will show ways of moving data from them to **sp** classes

Just spatial objects

- There are base graphics plot methods for the key `Spatial*` classes, including the `Spatial` class, which just sets up the axes
- In base graphics, additional plots can be added by overplotting as usual, and the `locator()` and `identify()` functions work as expected
- In general, most `par()` options will also work, as will the full range of graphics devices (although some copying operations may disturb aspect)
- First we will display the positional data of the objects discussed in the first unit

Plotting a SpatialPoints object

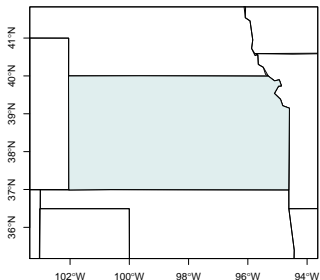


While plotting the `SpatialPoints` object would have called the `plot` method for `Spatial` objects internally to set up the axes, we start by doing it separately:

```
> library(sp)
> plot(as(storn2, "Spatial"), axes=TRUE)
> plot(storn2, add=TRUE)
> plot(storn2[storn2$EFscale==4,], col="red", add=TRUE)
>
```

Then we plot the points with the default plotting character, and subset, overplotting points with EF scale of 4 in red, using the `[]` method

Plotting a SpatialPolygons object

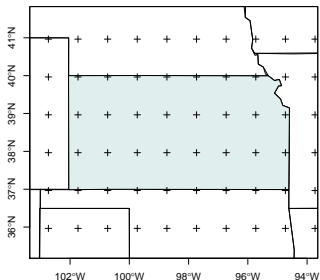


In plotting the `SpatialPolygons` object, we use the `xlim=` and `ylim=` arguments to restrict the display area to match the soil sample points.

```
> #plot(rivers, axes=TRUE, col="azure1", ylim=c(329400, 329400))
> library(sp)
> kansas<-statesth[statesth$NAME=="Kansas",]
> plot(statesth, axes=TRUE, xlim=c(-103,-94), ylim=c(36,41))
> plot(statesth[statesth$NAME=="Kansas",], col="azure2",
> box())
```

If the `axes=` argument is `FALSE` or omitted, no axes are shown — the default is the opposite from standard base graphics plot methods

Plotting a SpatialPixels object



Both SpatialPixels and SpatialGrid objects are plotted like SpatialPoints objects, with plotting characters

```
> #plot(rivers, axes=TRUE, col="azure1", ylim=c(329400, 329500))
> library(sp)
> kansas<-statesth[statesth$NAME=="Kansas",]
> plot(statesth, axes=TRUE, xlim=c(-103,-94), ylim=c(36,41), col="white",
> plot(statesth[statesth$NAME=="Kansas",], col="azure2",
> box()
> #coords <- SpatialPixels(SpatialPoints(meuse.grid[, c("x", "y")],
> #meusegrid <- SpatialPixelsDataFrame(coords, meuse.grid)
> #
> #plot(meusegrid, add=TRUE, col="grey60", cex=0.15)
>
> plot(grd2, add=TRUE)
```

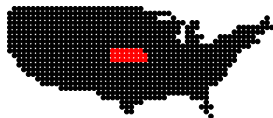
While points, lines, and polygons are often plotted without attributes, this is rarely the case for gridded objects

Including attributes

- To include attribute values means making choices about how to represent their values graphically, known in some GIS as symbology
- It involves choices of symbol shape, colour and size, and of which objects to differentiate
- When the data are categorical, the choices are given, unless there are so many different categories that reclassification is needed for clear display
- Once we've looked at some examples, we'll go on to see how class intervals may be chosen for continuous data

Points of the grid inside kansas

● OFF KANSAS
● IN KANSAS

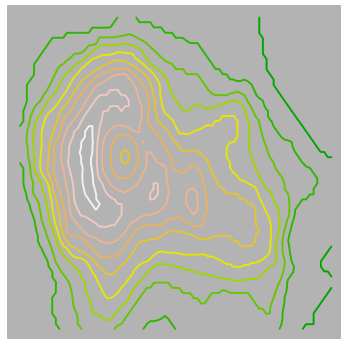


We will usually need to get the category levels and match them to colours (or plotting characters) “by hand”

```
> #meuse1$ffreq1 <- as.numeric(meuse1$ffreq)
> #plot(meuse1, col=meuse1$ffreq1, pch=19)
> kidx<-overlay(grd2, kansas)
> grd2df<-SpatialPointsDataFrame(grd2, data.frame(KANSAS=
> #levels(grd2df$KANSAS)<-c("", "KANSAS")
>
> plot(grd2df, col=grd2df$KANSAS, pch=19)
> labs <- c("OFF KANSAS", "IN KANSAS")
> cols <- 1:2
> legend("topleft", legend=labs, col=cols, pch=19, bty="n")
```

It is also typical that the legend() involves more code than everything else together, but very often the same vectors are used repeatedly and can be assigned just once

Coloured contour lines



Here again, the values are represented as a categorical variable, and so do not require classification

```
> library(maptools)
> volcano_sl <- ContourLines2SLDF(contourLines(volcano))
> volcano_sl$level1 <- as.numeric(volcano_sl$level)
> pal <- terrain.colors(nlevels(volcano_sl$level))
> plot(volcano_sl, bg="grey70", col=pal[volcano_sl$level1])
```

Using class membership for colour palette look-up is a very typical idiom, so that the `col=` argument is in fact a vector of colour values

Class intervals

- Class intervals can be chosen in many ways, and some have been collected for convenience in the **classInt** package
- The first problem is to assign class boundaries to values in a single dimension, for which many classification techniques may be used, including pretty, quantile, natural breaks among others, or even simple fixed values
- From there, the intervals can be used to generate colours from a colour palette, using the very nice `colorRampPalette()` function
- Because there are potentially many alternative class memberships even for a given number of classes (by default from `nclass.Sturges`), choosing a communicative set matters

Class intervals

We will try just two styles, quantiles and Fisher-Jenks natural breaks for five classes, among the many available. They yield quite different impressions, as we will see:

```
> storn2$LLoss<-log(storn2$Loss+.0001)
> library(classInt)
> library(RColorBrewer)
> pal <- brewer.pal(3, "Blues")
> q5 <- classIntervals(storn2$LLoss, n=5, style="quantile")
> q5

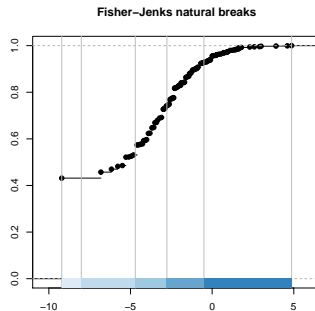
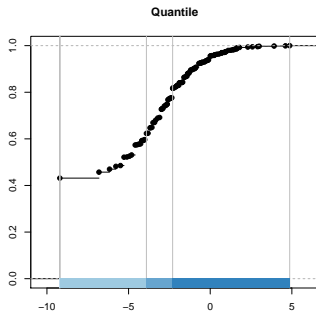
style: quantile
  one of 8,495,410 possible partitions of this variable into 5 classes
  [-9.21034,-9.21034)  [-9.21034,-9.21034)  [-9.21034,-3.907035)
                        0                        0                        697
[-3.907035,-2.301586)  [-2.301586,4.867535]
                        210                        261

> fj5 <- classIntervals(storn2$LLoss, n=5, style="fisher")
> fj5

style: fisher
  one of 8,495,410 possible partitions of this variable into 5 classes
  [-9.21034,-8.011393)  [-8.011393,-4.705556)  [-4.705556,-2.771788)
                        504                        116                        245
[-2.771788,-0.5023957)  [-0.5023957,4.867535]
                        218                        85

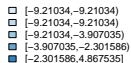
> plot(q5, pal=pal)
> plot(fj5, pal=pal)
```

Class interval plots

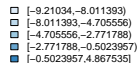


Two versions of the (log)-losses caused by tornadoes

Quantile



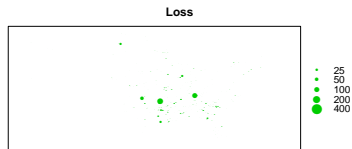
Fisher-Jenks natural breaks



Lattice graphics

- Lattice graphics will only come into their own later on, when we want to plot several variables with the same scale together for comparison
- The workhorse method is `spplot`, which can be used as an interface to the underlying `xypplot` or `levelplot` methods, or others as suitable; overplotting must be done in the single call to `spplot` — see gallery
- It is often worthwhile to load the **lattice** package so as to have direct access to its facilities
- Please remember that lattice graphics are displayed on the current graphics device by default only in interactive sessions — in loops or functions, they must be explicitly `print`'ed

Bubble plots



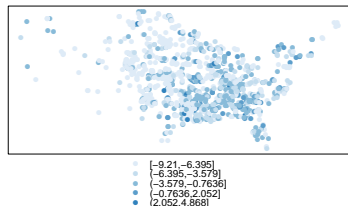
Bubble plots are a convenient way of representing the attribute values by the size of a symbol

```
> library(lattice)
```

```
> print(bubble(storn2, "Loss", maxsize=2, key.entries=25*
```

As with all lattice graphics objects, the function can return an object from which symbol sizes can be recovered

Level plots



The use of lattice plotting methods yields easy legend generation, which is another attraction

```
> bpal <- colorRampPalette(pal)(6)
> print( splot(storn2, "LLoss", col.regions=bpal, cuts=5
```

Here we are showing the distances from the river of grid points in the study area; we can also pass in intervals chosen previously

More realism

- So far we have just used canned data and spatial objects rather than anything richer
- The vizualisation methods are also quite flexible — both the base graphics and lattice graphics methods can be extensively customised
- It is also worth recalling the range of methods available for **sp** objects, in particular the `overlay` and `spsample` methods with a range of argument signatures
- These can permit further flexibility in display, in addition to their primary uses

Analysing Spatial Data in R: Accessing spatial data

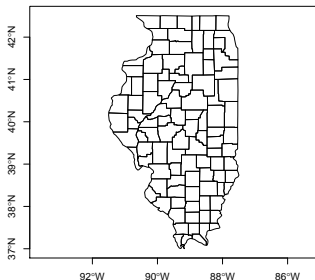
Introduction

- Having described how spatial data may be represented in R, and how to visualise these objects, we need to move on to accessing user data
- There are quite a number of packages handling and analysing spatial data on CRAN, and others off-CRAN, and their data objects can be converted to or from **sp** object form
- We need to cover how coordinate reference systems are handled, because they are the foundation for spatial data integration
- Both here, and in relation to reading and writing various file formats, things have advanced a good deal since the R News note

Creating objects within R

- As mentioned previously, **maptools** includes `ContourLines2SLDF()` to convert contour lines to `SpatialLinesDataFrame` objects
- **maptools** also allows lines or polygons from **maps** to be used as **sp** objects
- **maptools** can export **sp** objects to **PBSmapping**
- **maptools** uses **gpclib** to check polygon topology and to dissolve polygons
- **maptools** converts some **sp** objects for use in **spatstat**
- **maptools** can read GSHHS high-resolution shoreline data into `SpatialPolygon` objects

Using **maps** data: Illinois counties



There are number of valuable geographical databases in map format that can be accessed directly — beware of IDs!

```
> library(maptools)
> library(maps)
> ill <- map("county", regions="illinois", plot=FALSE, fill=TRUE)
> IDs <- sub("^illinois", "", ill$names)
> ill_sp <- map2SpatialPolygons(ill, IDs, CRS("+proj=longlat"))
> plot(ill_sp, axes=TRUE)
```

Coordinate reference systems

- Coordinate reference systems (CRS) are at the heart of geodetics and cartography: how to represent a bumpy ellipsoid on the plane
- We can speak of geographical CRS expressed in degrees and associated with an ellipse, a prime meridian and a datum, and projected CRS expressed in a measure of length, and a chosen position on the earth, as well as the underlying ellipse, prime meridian and datum.
- Most countries have multiple CRS, and where they meet there is usually a big mess — this led to the collection by the European Petroleum Survey Group (EPSG, now Oil & Gas Producers (OGP) Surveying & Positioning Committee) of a geodetic parameter dataset

Coordinate reference systems

- The EPSG list among other sources is used in the workhorse PROJ.4 library, which as implemented by Frank Warmerdam handles transformation of spatial positions between different CRS
- This library is interfaced with R in the **rgdal** package, and the CRS class is defined partly in **sp**, partly in **rgdal**
- A CRS object is defined as a character NA string or a valid PROJ.4 CRS definition
- The validity of the definition can only be checked if **rgdal** is loaded

Here: neither here nor there

In a Dutch navigation example, a chart position in the ED50 datum has to be compared with a GPS measurement in WGS84 datum right in front of the jetties of IJmuiden, both in geographical CRS. Using the `spTransform` method makes the conversion, using EPSG and external information to set up the ED50 CRS. The difference is about 124m; lots of details about CRS in general can be found in [Grids & Datums](#).

```
> library(rgdal)

> ED50 <- CRS(paste("+init=epsg:4230", "+towgs84=-87,-96,-120,0,0,0,0"))
> IJ.east <- as(char2dms("4d31\ '00\"E"), "numeric")
> IJ.north <- as(char2dms("52d28\ '00\"N"), "numeric")
> IJ.ED50 <- SpatialPoints(cbind(x=IJ.east, y=IJ.north), ED50)
> res <- spTransform(IJ.ED50, CRS("+proj=longlat +datum=WGS84"))
> spDistsN1(coordinates(IJ.ED50), coordinates(res), longlat=TRUE)*1000

[1] 124.0994
```

CRS are muddled

- If you think CRS are muddled, you are right, like time zones and daylight saving time in at least two dimensions
- But they are the key to ensuring positional interoperability, and “mashups” — data integration using spatial position as an index must be able to rely on data CRS for integration integrity
- The situation is worse than TZ/DST because there are lots of old maps around, with potentially valuable data; finding correct CRS values takes time
- On the other hand, old maps and odd choices of CRS origins can have their charm ...

Reading vectors

- GIS vector data are points, lines, polygons, and fit the equivalent **sp** classes
- There are a number of commonly used file formats, all or most proprietary, and some newer ones which are partly open
- GIS are also handing off more and more data storage to DBMS, and some of these now support spatial data formats
- Vector formats can also be converted outside R to formats that are easier to read

Reading vectors

- GIS vector data can be either topological or spaghetti — legacy GIS was topological, desktop GIS spaghetti
- **sp** classes are not bad spaghetti, but no checking of lines or polygons is done for errant topology
- A topological representation in principal only stores each point once, and builds arcs (lines between nodes) from points, polygons from arcs — GRASS 6 has a nice topological model
- Only **RArcInfo** tries to keep some traces of topology in importing legacy ESRI ArcInfo binary vector data (or e00 format data) — **maps** uses topology because that was how things were done then

Reading shapefiles

- The ESRI ArcView and now ArcGIS standard(ish) format for vector data is the shapefile, with at least a DBF file of data, an SHP file of shapes, and an SHX file of indices to the shapes; an optional PRJ file is the CRS
- Many shapefiles in the wild do not meet the ESRI standard specification, so hacks are unavoidable unless a full topology is built
- Both **maptools** and **shapefiles** contain functions for reading and writing shapefiles; they cannot read the PRJ file, but do not depend on external libraries
- There are many valid types of shapefile, but they sometimes occur in strange contexts — only some can be happily represented in R so far

Reading shapefiles: maptools

```
> library(maptools)
> #list.files()
> getinfo.shape("datasets/s_01au07.shp")

Shapefile type: Polygon, (5), # of Shapes: 57

> US <- readShapePoly("datasets/s_01au07.shp")
```

There are `readShapePoly`, `readShapeLines`, and `readShapePoints` functions in the `maptools` package, and in practice they now handle a number of infelicities. They do not, however, read the CRS, which can either be set as an argument, or updated later with the `proj4string` method

Reading vectors: rgdal

```
> US1 <- readOGR(dsn="datasets", layer="s_01au07")
```

```
OGR data source with driver: ESRI Shapefile
```

```
Source: "datasets", layer: "s_01au07"
```

```
with 57 features and 5 fields
```

```
Feature type: wkbPolygon with 2 dimensions
```

```
> cat(strwrap(proj4string(US1)), sep="\n")
```

```
+proj=longlat +datum=NAD83 +no_defs +ellps=GRS80
```

```
+towgs84=0,0,0
```

Using the OGR vector part of the Geospatial Data Abstraction Library lets us read shapefiles like other formats for which drivers are available. It also supports the handling of CRS directly, so that if the imported data have a specification, it will be read. OGR formats differ from platform to platform — the next release of rgdal will include a function to list available formats. Use FWTools to convert between formats.

Reading rasters

- There are very many raster and image formats; some allow only one band of data, others think data bands are RGB, while yet others are flexible
- There is a simple `readAsciiGrid` function in **maptools** that reads ESRI Arc ASCII grids into `SpatialGridDataFrame` objects; it does not handle CRS and has a single band
- Much more support is available in **rgdal** in the `readGDAL` function, which — like `readOGR` — finds a usable driver if available and proceeds from there
- Using arguments to `readGDAL`, subregions or bands may be selected, which helps handle large rasters

Reading rasters: rgdal

```
> getGDALDriverNames()$name
```

[1] AAIGrid	ACE2	ADRG	AIG	AirSAR
[6] BAG	BIGGIF	BLX	BMP	BSB
[11] BT	CEOS	COASP	COSAR	CPG
[16] CTG	DIMAP	DIPEX	DODS	DOQ1
[21] DOQ2	DTED	EOGRID	ECRGTOT	EHdr
[26] EIR	ELAS	ENVI	EPSILON	ERS
[31] ESAT	FAST	FIT	FujiBAS	GenBin
[36] GFF	GIF	GMT	GRASSASCIIGrid	GRIB
[41] GS7BG	GSAG	GSBG	GSC	GTiff
[46] GTX	GXF	HDF4	HDF4Image	HDF5
[51] HDF5Image	HF2	HFA	HTTP	IDA
[56] ILWIS	INGR	ISIS2	ISIS3	JAXAPALSAR
[61] JDEM	JPEG	JPEG2000	KMLSUPEROVERLAY	L1B
[66] LAN	LCP	Leveller	LOSLAS	MEM
[71] MFF	MFF2	MSGN	NDF	netCDF
[76] NGSGEID	NITF	NTv2	NWT_GRC	NWT_GRD
[81] OGD1	OZI	PAux	PCIDSK	PCRaster
[86] PDF	PDS	PNG	PNM	PostGISRaster
[91] R	Rasterlite	RIK	RMF	RPFTOT
[96] RS2	RST	SAGA	SAR_CEOS	SDTS
[101] SGI	SNODAS	SRP	SRTMHGT	Terragen
[106] TIL	TSX	USGSDEM	VRT	WCS
[111] WMS	XPM	XYZ	ZMap	

114 Levels: AAIGrid ACE2 ADRG AIG AirSAR BAG BIGGIF BLX BMP BSB BT CEOS COASP COSAR CPG ... ZMap

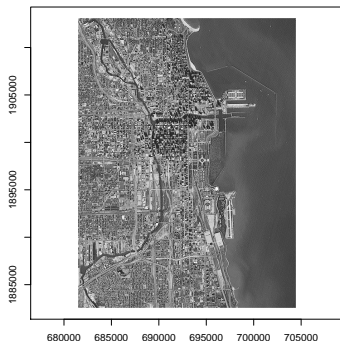
```
> list.files()
```

```
[1] "SP27GTIF.TIF"
```

```
> SP27GTIF <- readGDAL("SP27GTIF.TIF")
```

SP27GTIF.TIF has GDAL driver Gtiff
and has 929 rows and 699 columns

Reading rasters: rgdal



This is a single band GeoTiff, mostly showing downtown Chicago; a lot of data is available in geotiff format from US public agencies, including Shuttle radar topography mission seamless data — we'll get back to this later

```
> image(SP27GTIF, col=grey(1:99/100), axes=TRUE)
```

Reading rasters: rgdal

```
> summary(SP27GTIF)
```

```
Object of class SpatialGridDataFrame
```

```
Coordinates:
```

```
      min      max  
x 681480 704407.2  
y 1882579 1913050.0
```

```
Is projected: TRUE
```

```
proj4string :
```

```
[+proj=tmerc +lat_0=36.66666666666666 +lon_0=-88.33333333333333 +k=0.9999749999999999  
+x_0=152400.3048006096 +y_0=0 +datum=NAD27 +units=us-ft +no_defs +ellps=clrk66  
+nadgrids=@conus,@alaska,@ntv2_0.gsb,@ntv1_can.dat]
```

```
Grid attributes:
```

```
  cellcentre.offset cellsize cells.dim  
x      681496.4      32.8      699  
y      1882595.2      32.8      929
```

```
Data attributes:
```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
	4.0	78.0	104.0	115.1	152.0	255.0

Writing objects

- In **rgdal**, `writeGDAL` can write for example multi-band GeoTiffs, but there are fewer write than read drivers; in general CRS and georeferencing are supported — see `gdalDrivers`
- The **rgdal** function `writeOGR` can be used to write vector files, including those formats supported by drivers, including now KML — see `ogrDrivers`
- External software (including different versions) tolerate output objects in varying degrees, quite often needing tricks - see mailing list archives
- In **maptools**, there are functions for writing **sp** objects to shapefiles — `writePolyShape`, etc., as Arc ASCII grids — `writeAsciiGrid`, and for using the R PNG graphics device for outputting image overlays for Google Earth

GIS interfaces

- GIS interfaces can be as simple as just reading and writing files — loose coupling, once the file formats have been worked out, that is
- Loose coupling is less of a burden than it was with smaller, slower machines, which is why the **GRASS** 5 interface was tight-coupled, with R functions reading from and writing to the GRASS database directly
- The GRASS 6 interface **spgrass6** on CRAN also runs R within GRASS, but uses intermediate temporary files; the package is under development but is quite usable
- Use has been made of COM and Python interfaces to ArcGIS; typical use is by loose coupling except in highly customised work situations
- Carson Farmer has developed a plug-in for QGIS (manageR) to provide a bridge between R and QGIS