

# Spatio-temporal objects to proxy a PostgreSQL table



ifgi  
Institute for Geoinformatics  
University of Münster



Edzer Pebesma

September 25, 2015

## Abstract

This vignette describes and implements a class that proxies data sets in a PostgreSQL database with classes in the `spacetime` package. This might allow access to data sets too large to fit into R memory.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Setting up a database</b>	<b>2</b>
<b>3</b>	<b>A proxy class</b>	<b>3</b>
<b>4</b>	<b>Selection based on time period and/or region</b>	<b>4</b>
<b>5</b>	<b>Closing the database connection</b>	<b>8</b>
<b>6</b>	<b>Limitations and alternatives</b>	<b>8</b>

## 1 Introduction

Massive data are difficult to analyze with R, because R objects reside in memory. Spatio-temporal data easily become massive, either because the spatial domain contains a lot of information (satellite imagery), or many time steps are available (high resolution sensor data), or both. This vignette shows how data residing in a data base can be read into R using spatial or temporal selection.

In case the commands are not evaluated because CRAN packages cannot access an external data base, a document with evaluated commands is found [here](#).

This vignette was run using the following libraries:

```
> library(RPostgreSQL)
```

```
> library(sp)
> library(spacetime)
```

## 2 Setting up a database

We will first set the characteristics of the database<sup>1</sup>

```
> dbname = "postgis"
> user = "edzer"
> password = "pw"
> #password = ""
```

Next, we will create a driver and connect to the database:

```
> drv <- dbDriver("PostgreSQL")
> con <- dbConnect(drv, dbname=dbname, user=user, password=password,
+ host='localhost', port='5432')
```

It should be noted that these first two commands are specific to PostgreSQL; from here on, commands are generic and should work for any database connector that uses the interface of package DBI.

We now remove a set of tables (if present) so they can be created later on:

```
> dbRemoveTable(con, "rural_attr")

[1] TRUE

> dbRemoveTable(con, "rural_space")

[1] TRUE

> dbRemoveTable(con, "rural_time")

[1] TRUE

> dbRemoveTable(con, "space_select")

[1] TRUE
```

Now we will create the table with spatial features (observation locations). For this, we need the `rgdal` function `writeOGR`, which by default creates an index on the geometry:

```
> data(air)
> rural = STFDF(stations, dates, data.frame(PM10 = as.vector(air)))
> rural = as(rural, "STSDf")
> p = rural@sp
> sp = SpatialPointsDataFrame(p, data.frame(geom_id=1:length(p)))
> library(rgdal)
> OGRstring = paste("PG:dbname=", dbname, " user=", user,
+ " password=", password, " host=localhost", sep = "")
> print(OGRstring)
```

---

<sup>1</sup>It is assumed that the database is *spatially enabled*, i.e. it understands how simple features are stored. The standard for this from the open geospatial consortium is described [here](#).

```
[1] "PG:dbname=postgis user=edzer password=pw host=localhost"
```

```
> writeOGR(sp, OGRstring, "rural_space", driver = "PostgreSQL")
```

In case you have problems replicating this, verify that your `rgdal` installation provides the PostgreSQL driver, e.g. by checking that

```
> subset(ogrDrivers(), name == "PostgreSQL")$write
```

```
[1] TRUE
```

prints a `TRUE`, and not a `logical(0)`.

Second, we will write the table with times to the database, and create an index to time:

```
> df = data.frame(time = index(rural@time), time_id = 1:nrow(rural@time))
> dbWriteTable(con, "rural_time", df)
```

```
[1] TRUE
```

```
> idx = "create index time_idx on rural_time (time);"
> dbSendQuery(con, idx)
```

```
<PostgreSQLResult:(7580,0,21)>
```

Finally, we will write the full attribute data table to PostgreSQL, along with its indexes to the spatial and temporal tables:

```
> idx = rural@index
> names(rural@data) = "pm10" # lower case
> df = cbind(data.frame(geom_id = idx[,1], time_id = idx[,2]), rural@data)
> dbWriteTable(con, "rural_attr", df)
```

```
[1] TRUE
```

### 3 A proxy class

The following class has as components a spatial and temporal data structure, but no spatio-temporal attributes (they are assumed to be the most memory-hungry). The other slots refer to the according tables in the PostgreSQL database, the name(s) of the attributes in the attribute table, and the database connection.

```
> setClass("ST_PG", contains = "ST",
+         # slots = c(space_table = "character",
+         representation(space_table = "character",
+         time_table = "character",
+         attr_table = "character",
+         attr = "character",
+         con = "PostgreSQLConnection"))
```

Next, we will create an instance of the new class:

```

> rural_proxy = new("ST_PG",
+   #ST(rural@sp, rural@time, rural@endTime),
+   as(rural, "ST"),
+   space_table = "rural_space",
+   time_table = "rural_time",
+   attr_table = "rural_attr",
+   attr = "pm10",
+   con = con)

```

## 4 Selection based on time period and/or region

The following two helper functions create a character string with an SQL command that for a temporal or spatial selection:

```

> .SqlTime = function(x, j) {
+   stopifnot(is.character(j))
+   require(xts)
+   t = .parseISO8601(j)
+   t1 = paste("'", t$first.time, "'", sep = "")
+   t2 = paste("'", t$last.time, "'", sep = "")
+   what = paste("geom_id, time_id", paste(x@attr, collapse = ","), sep = ", ")
+   paste("SELECT", what, "FROM", x@attr_table, "AS a JOIN", x@time_table,
+       "AS b USING (time_id) WHERE b.time >= ", t1, "AND b.time <=", t2, ";")
+ }
> .SqlSpace = function(x, i) {
+   stopifnot(is(i, "Spatial"))
+   writeOGR(i, OGRstring, "space_select", driver = "PostgreSQL")
+   what = paste("geom_id, time_id", paste(x@attr, collapse = ","), sep = ", ")
+   paste("SELECT", what, "FROM", x@attr_table,
+       "AS a JOIN (SELECT p.wkb_geometry, p.geom_id FROM",
+       x@space_table, " AS p, space_select AS q",
+       "WHERE ST_Intersects(p.wkb_geometry, q.wkb_geometry))",
+       "AS b USING (geom_id);")
+ }

```

The following selection method selects a time period only, as defined by the methods in package `xts`. A time period is defined as a valid ISO8601 string, e.g. 2005-05 is the full month of May for 2005.

```

> setMethod("[", "ST_PG", function(x, i, j, ... , drop = TRUE) {
+   stopifnot(missing(i) != missing(j)) # either of them present
+   if (missing(j))
+       sql = .SqlSpace(x,i)
+   else
+       sql = .SqlTime(x,j)
+   print(sql)
+   df = dbGetQuery(x@con, sql)
+   STSDF(x@sp, x@time, df[x@attr], as.matrix(df[c("geom_id", "time_id"))))
+ })

[1] "["

```

```

> pm10_20050101 = rural_proxy[, "2005-01-01"]

[1] "SELECT geom_id, time_id, pm10 FROM rural_attr AS a JOIN rural_time AS b USING (time_id)"

> summary(pm10_20050101)

Object of class STSDF
  with Dimensions (s, t, attr): (70, 4383, 1)
[[Spatial:]]
Object of class SpatialPoints
Coordinates:
      min      max
coords.x1 6.28107 14.78617
coords.x2 47.80847 54.92497
Is projected: FALSE
proj4string :
[+init=epsg:4326 +proj=longlat +ellps=WGS84 +datum=WGS84 +no_defs
+towgs84=0,0,0]
Number of points: 70
[[Temporal:]]
      Index      timeIndex
Min.   :1998-01-01  Min.   : 1
1st Qu.:2000-12-31  1st Qu.:1096
Median :2004-01-01  Median :2192
Mean   :2004-01-01  Mean   :2192
3rd Qu.:2006-12-31  3rd Qu.:3288
Max.   :2009-12-31  Max.   :4383
[[Data attributes:]]
      pm10
Min.   : 1.792
1st Qu.: 8.167
Median :12.833
Mean   :15.641
3rd Qu.:20.750
Max.   :45.375

> summary(rural[, "2005-01-01"])

Object of class SpatialPointsDataFrame
Coordinates:
      min      max
coords.x1 6.28107 14.78617
coords.x2 47.80847 54.92497
Is projected: FALSE
proj4string :
[+init=epsg:4326 +proj=longlat +ellps=WGS84 +datum=WGS84 +no_defs
+towgs84=0,0,0]
Number of points: 45
Data attributes:
      pm10
Min.   : 1.792

```

```

1st Qu.: 8.167
Median :12.833
Mean   :15.641
3rd Qu.:20.750
Max.    :45.375

> pm10_NRW = rural_proxy[DE_NUTS1[10,],]

[1] "SELECT geom_id, time_id, pm10 FROM rural_attr AS a JOIN (SELECT p.wkb_geometry, p.geoname_id FROM rural_attr AS p WHERE p.time_id = 10) AS b ON a.geom_id = b.wkb_geometry"

> summary(pm10_NRW)

Object of class STSDF
  with Dimensions (s, t, attr): (70, 4383, 1)
[[Spatial:]]
Object of class SpatialPoints
Coordinates:
      min      max
coords.x1 6.28107 14.78617
coords.x2 47.80847 54.92497
Is projected: FALSE
proj4string :
[+init=epsg:4326 +proj=longlat +ellps=WGS84 +datum=WGS84 +no_defs
+towgs84=0,0,0]
Number of points: 70
[[Temporal:]]
      Index      timeIndex
Min.   :1998-01-01 Min.    : 1
1st Qu.:2000-12-31 1st Qu.:1096
Median :2004-01-01 Median :2192
Mean   :2004-01-01 Mean   :2192
3rd Qu.:2006-12-31 3rd Qu.:3288
Max.   :2009-12-31 Max.   :4383
[[Data attributes:]]
      pm10
Min.   : 1.667
1st Qu.: 10.177
Median : 14.875
Mean   : 17.651
3rd Qu.: 22.000
Max.   :161.305

> summary(rural[DE_NUTS1[10,],])

Object of class STSDF
  with Dimensions (s, t, attr): (6, 3591, 1)
[[Spatial:]]
Object of class SpatialPoints
Coordinates:
      min      max
coords.x1 6.28107 8.950597

```

```

coords.x2 50.65324 51.862000
Is projected: FALSE
proj4string :
[+init=epsg:4326 +proj=longlat +ellps=WGS84 +datum=WGS84 +no_defs
+towgs84=0,0,0]
Number of points: 6
[[Temporal:]]
      Index      timeIndex
Min.   :2000-01-01  Min.    : 731
1st Qu.:2002-07-30  1st Qu.:1672
Median :2005-01-31  Median :2588
Mean   :2005-01-22  Mean    :2580
3rd Qu.:2007-07-17  3rd Qu.:3486
Max.   :2009-12-31  Max.    :4383
[[Data attributes:]]
      pm10
Min.   : 1.667
1st Qu.: 10.177
Median : 14.875
Mean   : 17.651
3rd Qu.: 22.000
Max.   :161.305

```

Clearly, the temporal and spatial components are not subsetted, so do not reflect the actual selection made; the attribute data however do; the following selection step “cleans” the unused features/times:

```

> dim(pm10_NRW)

      space      time variables
      70      4383           1

> pm10_NRW = pm10_NRW[T,]
> dim(pm10_NRW)

      space      time variables
      6      3591           1

```

Comparing sizes, we see that the selected object is smaller:

```

> object.size(rural)

3083192 bytes

> object.size(pm10_20050101)

103776 bytes

> object.size(pm10_NRW)

1090928 bytes

```

## 5 Closing the database connection

The following commands close the database connection and release the driver resources:

```
> dbDisconnect(con)

[1] TRUE

> dbUnloadDriver(drv)

[1] TRUE
```

## 6 Limitations and alternatives

The example code in this vignette is meant as an example and is not meant as a full-fledged database access mechanism for spatio-temporal data bases. In particular, the selection here can do only *one* of spatial locations (entered as features) or time periods. If database access is only based on time, a spatially enabled database (such as PostGIS) would not be needed.

For massive databases, data would typically not be loaded into the database from R first, but from somewhere else.

An alternative to access from R large, possibly massive spatio-temporal data bases for the case where the data base is accessible through a sensor observation service (SOS) is provided by the R package [sos4R](#), which is also on CRAN.