

# sdcSpatial: Privacy protected density maps

Edwin de Jonge @edwindjonge

Statistics Netherlands Research & Development  
@edwindjonge

useR! 2019

# sdcSpatial: Privacy protected maps



# sdcSpatial: Privacy protected maps

## Goal:

*Create maps, and ensure that no details are revealed on individuals.*

# sdcSpatial: Privacy protected maps

## sdcSpatial has methods for:

- Creating a raster map: `sdc_raster` for pop density, value density and mean density.
- Which locations are sensitive: `plot_sensitive`, `is_sensitive`
- Adjust raster map to protect data: `protect_smooth`, `protect_quadtree`
- Remove sensitive locations.

# Who am I?

- Statistical consultant, Data Scientist @cbs.nl / Statistics NL
- Expertise:
  - **R programming**
  - Data Cleaning with R
  - **Data visualization**
  - Complex networks analysis
  - @edwindjonge / <https://github.com/edwindj>

# What is SN / CBS?

Statistics Netherlands is producer of all main official statistics in the Netherlands:

- Stats on Demographics, economy (GDP), education, environment, agriculture, Finance etc.
- Part of the European Statistical System, ESS.

## Privacy motivation

- ESS has European Code of Statistical Practice (predates GDPR, European law on Data Protection): no individual information on persons and enterprises may be revealed.

# Sdc in sdcSpatial?

SDC = “Statistical Disclosure Control”

## Collection of statistical methods to:

- Check if density map is safe to be published
- Protect data by slightly altering (aggregated) data
  - adding noise
  - shifting mass
- Most SDC methods operate on records.
- **sdcSpatial works upon location data.**

# Lets create a raster map with sdc\_raster

```
library(sdcSpatial)
unemployed <- sdc_raster( dwellings[c("x", "y")]
                        , dwellings$unemployed
                        , r = 500
                        , min_count = 10 # min support
                        )

print(unemployed)
```

## logical sdc\_raster object:

## resolution: 500 500 , max\_risk: 0.95 , min\_count: 10

## mean sensitivity score [0,1]: 0.4249471

## What is the sensitivity?

Binary score (logical) per raster cell indicating if it's unsafe to publish.

## Calculated:

- Per location  $(x_i, y_i)$  (raster cell)
- Using risk function disclosure\_risk  $r(x, y) \in [0, 1]$ . How



# Type of raster density maps:

Density can be area-based:

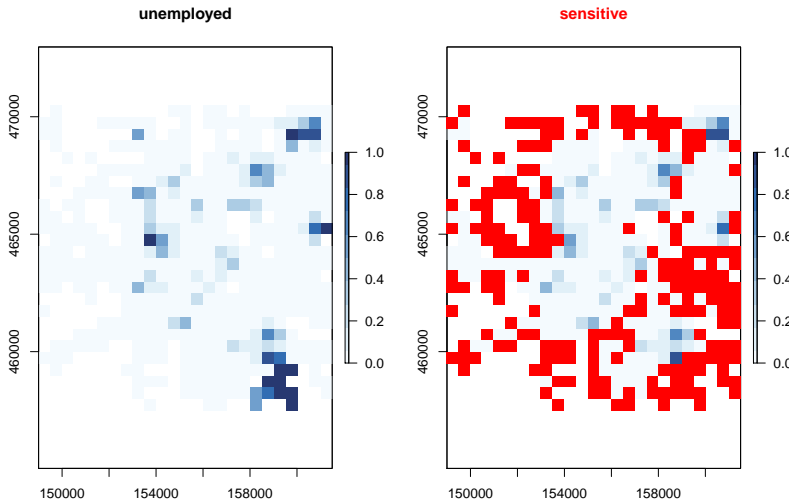
- number of people per square (`unemployed$value$count`): population density
- (total) value per square (`unemployed$value$sum`): number of unemployed per square.

Or density can population based: - Mean value per square (`unemployed$value$mean`): probability of being unemployed per square.

*All types can be valid, but note that (total) value per square strongly interacts with population density. (see <https://xkcd.com/1138>)*

# So let's plot!

```
plot(unemployed)
```

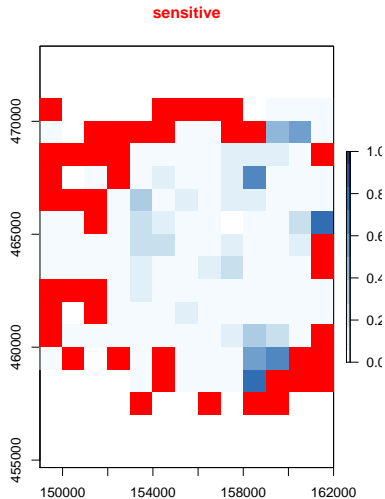
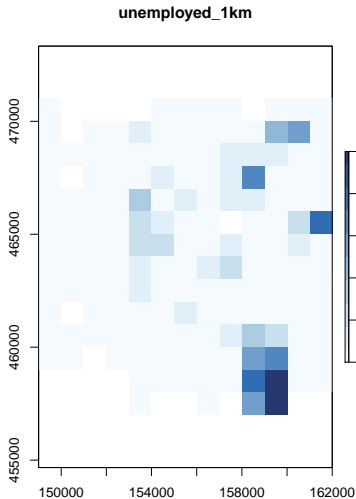


# how to improve?

- Use a coarser raster: `sdcraster`.
- Spatial smoothing: `protect_smooth[^jonge]`.
- Quadtree aggregation: `protect_quadtree`.
- Removing sensitive locations: `remove_sensitive`.

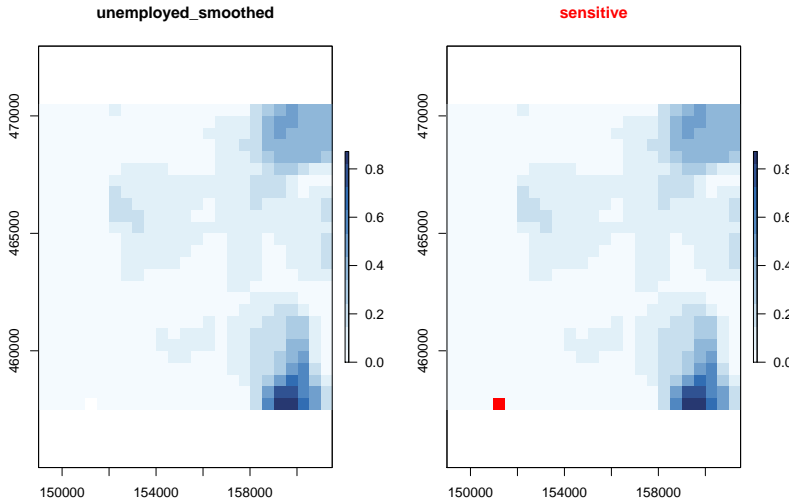
## Option: coarsening

```
unemployed_1km <- sdc_raster(dwellingings[c("x", "y")], dwellingings$unemployed_1km)  
plot(unemployed_1km)
```



# Option Smoothing

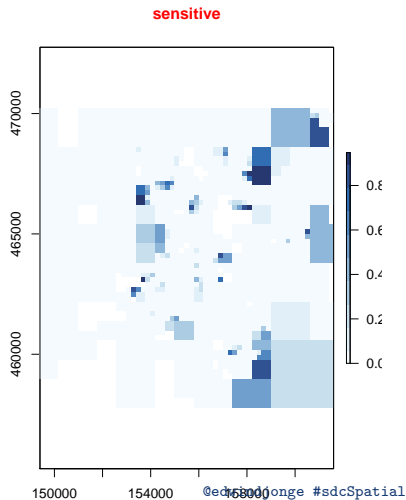
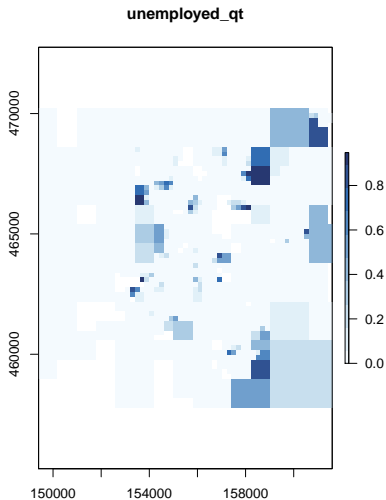
```
unemployed_smoothed <- protect_smooth(unemployed, bw = 150000)  
plot(unemployed_smoothed)
```





# Quadtree

```
unemployed_100m <- sdc_raster(dwellings[c("x","y")], dwellings$unemployed_100m)  
unemployed_qt <- protect_quadtree(unemployed_100m)  
plot(unemployed_qt)
```



# Implementation `sdc_raster`

- built upon excellent `raster`
- Creates internally an `raster::brick` object
- Stores the stats needed for calculating disclosure.





# The end

**Thank you for your attention!**

**Questions?**

**Curious?**

```
install.packages("sdcSpatial")
```

**Feedback and suggestions?**

<https://github.com/edwindj/sdcSpatial/issues>