# 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, using the excellent `raster` package by Hijmans (2019).
- 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.

**Motivation for `sdcSpatial`**

- 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`

```r
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:**

a) Per location $(x_i, y_i)$ (raster cell)

b) Using risk function `disclosure_risk` $r(x, y) \in [0, 1]$. How

# Type of raster density maps:
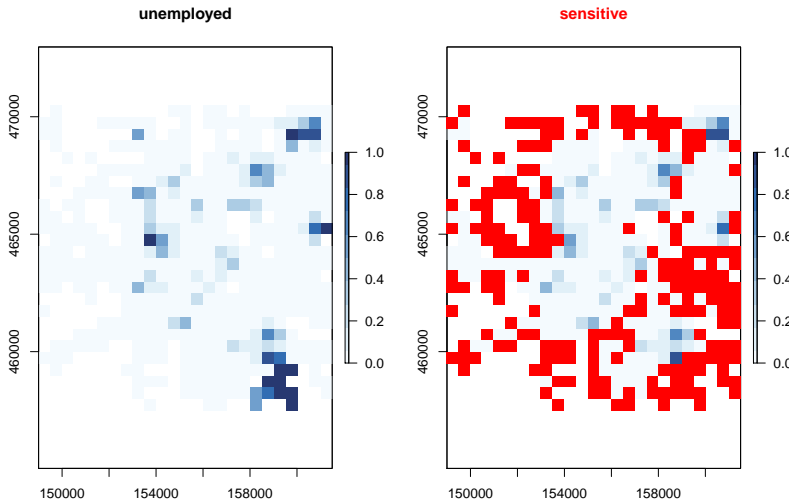
Density can be area-based:

- 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!
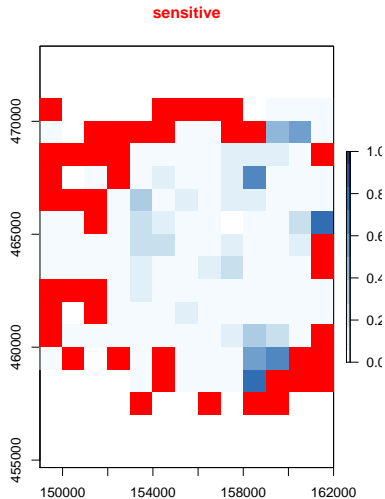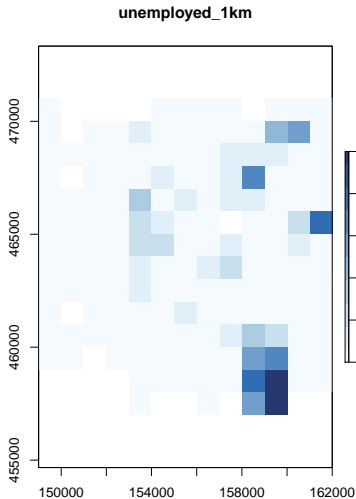
`plot(unemployed)`

# how to improve?

**Options:**

a) Use a coarser raster: `sdc_raster`.
b) Spatial smoothing: `protect_smooth` method by Wolf and Jonge (2018), Jonge and Wolf (2016).
c) Quadtree aggregation: `protect_quadtree` method by Suñé et al. (2017).
d) Removing sensitive locations: `remove_sensitive`.

# Option: coarsening

```
unemployed_1km <- sdc_raster(dwellings[c("x", "y")], dwell:
plot(unemployed_1km)
```

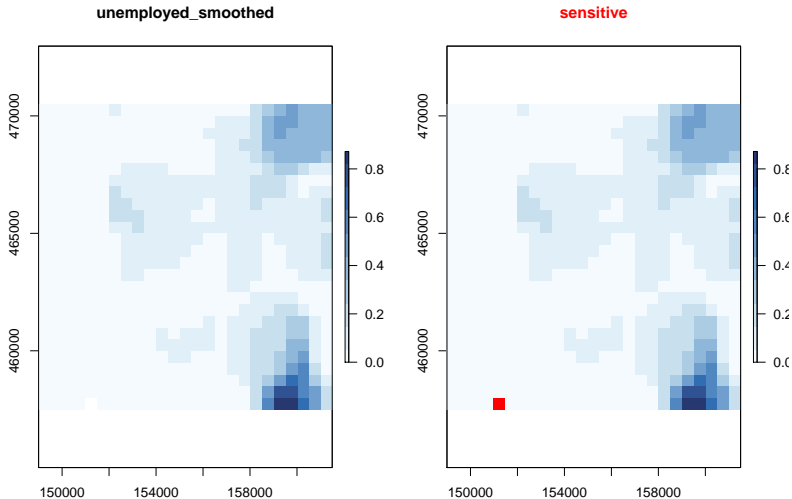# Option: Coarsening

**Pros**
- Simple and easy explainable

**Cons**
- Detailed spatial patterns are removed
- visually unattractive: "Blocky"

# Option: KDE-smoothing

```
unemployed_smoothed <- protect_smooth(unemployed, bw = 1500
plot(unemployed_smoothed)
```
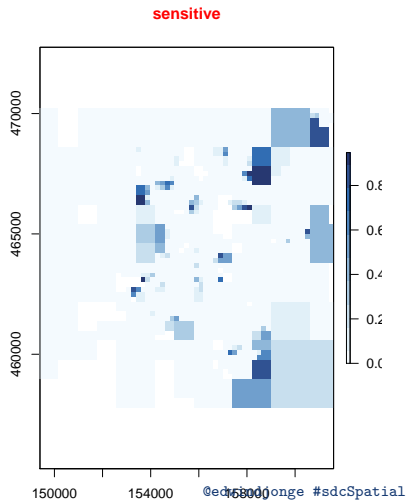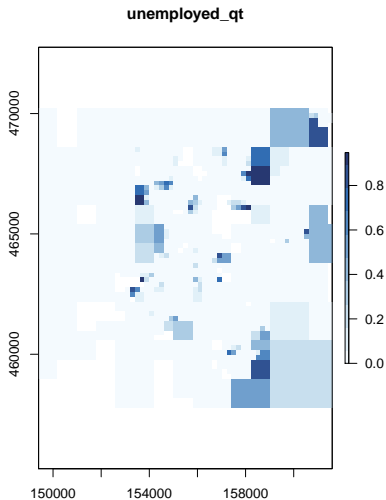
# Options: KDE-smoothing

## Pro's
- Often enhances spatial pattern visualisation.
- Makes it a density map, e.g. contour map

## Con's
- Does not remove all sensitive values (depends on bandwidth `bw`)
- A fixed band width is used for all locations: may removed detailed patterns
  spatial processes often have location dependent band widths.
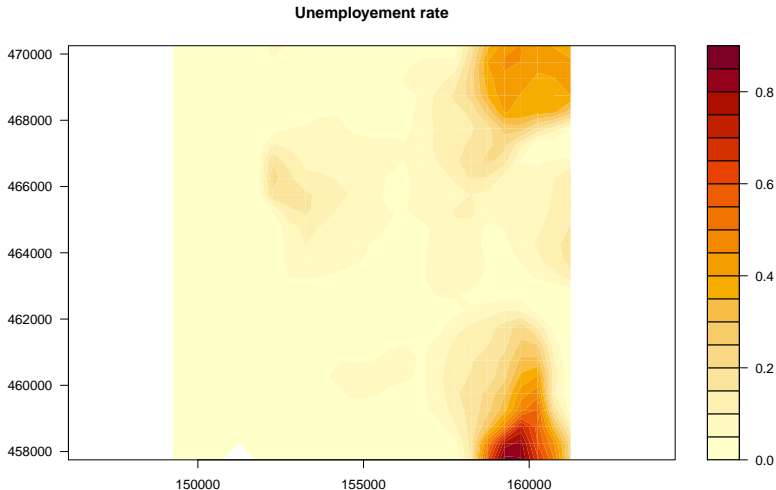  ($=$ future work)

## Option: Quadtree

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

# Publish: visual interpolation

- To make the raster visually more attractive, but still safe:

```
raster::filledContour(mean(unemployed_smoothed), main="Unemployement rate")
```



Unemployement rate

# Option: Quadtree

**Pro**
- Adapts to data density
- Adjusts until no sensitive data is left.

**Cons**
- Visually: "Blocky" / "Mondrian-like" result.

**Cons**

# The end

**Thank you for your atention!**

**Questions?**

**Curious?**
```r
install.packages("sdcSpatial")
```

**Feedback and suggestions?**
https://github.com/edwindj/sdcSpatial/issues

# References

Hijmans, Robert J. 2019. *Raster: Geographic Data Analysis and Modeling*. https://CRAN.R-project.org/package=raster.

Jonge, Edwin de, and Peter-Paul de Wolf. 2016. "Spatial Smoothing and Statistical Disclosure Control." In *Privacy in Statistical Databases*, edited by Josep Domingo-Ferrer and Mirjana Pejić-Bach, 107–17. Springer.

Suñé, E., C. Rovira, D. Ibáñez, and M. Farré. 2017. "Statistical Disclosure Control on Visualising Geocoded Population Data Using Quadtrees." http://nt17.pg2.at/data/x_abstracts/x_abstract_286.docx.

Wolf, Peter-Paul de, and Edwin de Jonge. 2018. "Spatial Smoothing and Statistical Disclosure Control." In *Privacy in Statistical Databases - Psd 2018*, edited by Josep Domingo-Ferrer and Francisco Montes Suay. Springer.