# Tutorial `reclin`
## A package for probabilistic record linkage and deduplication

Jan van der Laan <dj.vanderlaan@cbs.nl>

# About me

Methodologist/data scientist at Statistics Netherlands (CBS)

Author of a couple of R-packages: `reclin`, `simplermarkdown`, `LaF`, `ldat`, `lvec`

Give multiple courses at CBS: multivariate analysis methods in R, using R in statistical processes.

Current projects: social network analysis, measuring social economic status, measuring segregation.
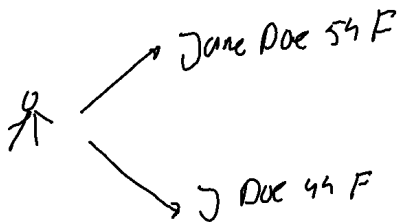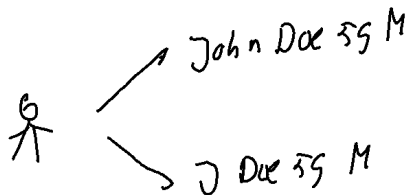
# Outline

— Overview record linkage

— Traditional record linkage / EM

— String similarity / custom similarity functions
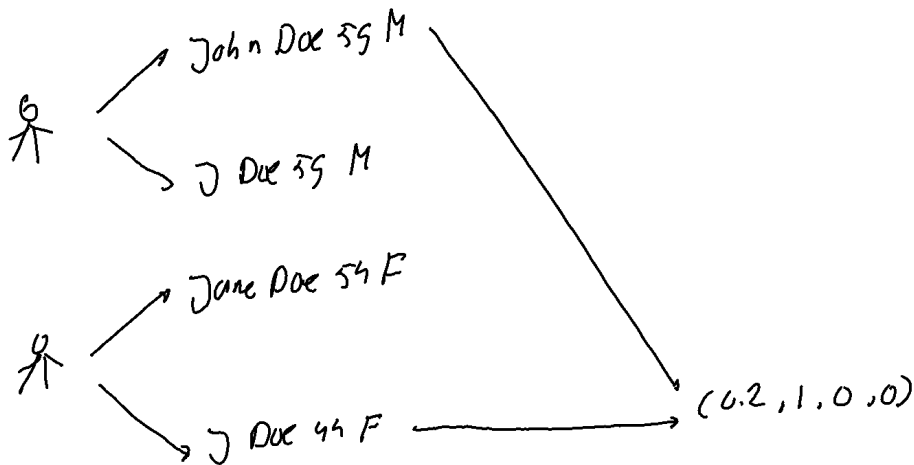
— Machine learning

Each item 10-15 min explanation + 10-15 min working on excercises

# Overview record linkage



John Doe 35 M

J Doe 35 M

Jane Doe 54 F

J Doe 44 F

# Overview record linkage



John Doe 55 M

J Doe 55 M

Jane Doe 54 F

J Doe 44 F

(0.2, 1, 0, 0)

# Overview record linkage



John Doe 55 M    (0.2, 1, 1, 1)

J Doe 55 M    (0.5, 1, 0, 0)

Jane Doe 54 F    (0.2, 1, 0, 1)

J Doe 44 F    (0.2, 1, 0, 0)

# Overview record linkage



John Doe 55 M

J Doe 55 M

Jane Doe 54 F

J Doe 54 F

$(0.2, 1, 1, 1) \longrightarrow 3.2$

$(0.5, 1, 0, 0) \longrightarrow 1.5$

$(0.2, 1, 0, 1) \longrightarrow 1.2$

$(0.2, 1, 0, 0) \longrightarrow 1.2$

# Overview record linkage



John Doe 55 M

J Doe 55 M

Jane Doe 54 F

J Doe 44 F

$(0.2, 1, 1, 1) \rightarrow \underline{3.2}$

$(0.5, 1, 0, 0) \rightarrow 1.5$

$(0.2, 1, 0, 1) \rightarrow \underline{1.2}$

$(0.2, 1, 0, 0) \rightarrow 1.2$

# The record linkage process

1. Generate record pairs
   — Blocking
2. Generate comparison vectors
3. Translate comparison vector in a score measuring likelihood of both records in a pairs belonging to the same object.
   — *Classical* probabilistic linkage using EM
   — Machine learning
   — Simple scoring functions
4. Select pairs with a high enough score
5. Generate linked dataset
   — Force one-to-one linkage

**Enough with the sheets, let's go to R …**

**Contact and more information**
Jan van der Laan <dj.vanderlaan@cbs.nl>
https://cran.r-project.org/package=reclin
https://github.com/djvanderlaan/reclin