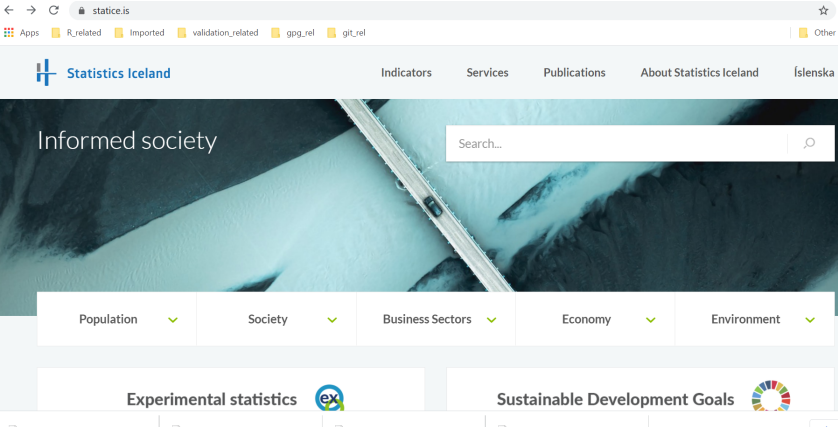# Correcting for population overestimates by using statistical classification methods

Violeta Calian and Margherita Zuppardo
Statistics Iceland

NTTS-2021

# About us

## Context and focus

- ▶ Solving an official statistics problem, **using** *R*
- ▶ Local context: basic principles at Statistics Iceland, while increasing **efficiency**
  - ▶ reproducibility, flexibility
  - ▶ transparency, peer review, collaboration
  - ▶ based on scientific/statistical knowledge
- ▶ General context: evolution of goals in official statistics:
  - ▶ description
  - ▶ modeling, estimation, **prediction**
  - ▶ reporting **uncertainty**
  - ▶ up-to-date!

# Formulation of the problem

### The social statistics problem:

▶ Overestimating population size of a given country/area at a given time

▶ Main cause: de-registration issues
▶ Impact: bias in estimates of demographic/social measures

▶ Input data: from multiple registers / databases (i.e. large set of attributes)
▶ To predict: status as present/absent
▶ Status of solutions across NSIs: *fuzzy* and *SOL*

## Examples of NSIs' solutions

- ► indices on SOL (Estonia)
- ► scores on SOL (Sweden)
- ► few classification models (Iceland, 2011 Census: cumulative link for ordinal response)

## While in the whole wide world:

- ► deep learning
- ► ensemble learning
- ► ML applications for: face/speech recognition, financial and medical fields, . . .

## Specific questions:

- ► *training* data and issues with the *status* variable: delay and noise
- ► a completely different formulation of the problem: macro vs micro
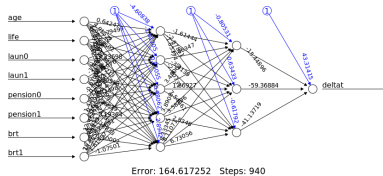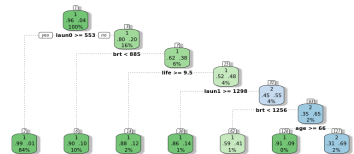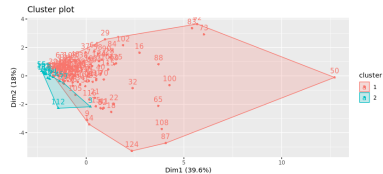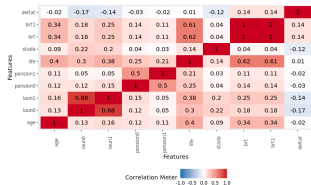- ► interpretability and solutions

# R / open source code:

Key role in modernizing official statistics

illustrated by the present work

developed for the purpose of Census 2021

but also for routine population estimates

# The main R - packages

# Results





Cluster plot





Error: 164.617252   Steps: 940

Next?

https://github.com/violetacln/SLOPA

violeta.calian@hagstofa.is and margherita.zuppardo@hagstofa.is

Thank you!