



# Machine learning techniques to forecast population using Eurostat data: an exploratory study

European Conference on Quality in Official Statistics  
**Q2018**

Krakow, 27th-29th June 2018

**Álvaro Gómez Losada and Néstor Duch-Brown**  
**Joint Research Centre - Seville (Spain)**  
**[alvaro.gomez-losada@ec.europa.eu](mailto:alvaro.gomez-losada@ec.europa.eu)**

## **Contents:**

1. Introduction and aim of this exploratory study.
2. Source of data and input variables.
3. Machine learning methodology.
4. Results.
5. Conclusions: strengths and limitations.

## **1. Introduction**

1. Population projections are accomplished using well-settled, widely-accepted statistical methodologies
2. Machine learning (ML) techniques have emerged as common forecasting approaches in multitude of study fields
3. Typically, ML requires huge amount of data to properly perform
4. Eurostat data are initially not suited to be analyzed by means of ML approaches due to their small size
5. This exploratory study aims to find out how well ML algorithms perform in forecasting population one-year in advance.

## 2. Source of data and input variables

Data were obtained using the *Eurostat* library (Lahti, 2017) from R (RCran, 2017). Target country: *Spain*.

Eurostat's tables and variables:

*demo\_pjan*: population by sex and age (1-99). 2005 to 2016 –**pop**-

*migr\_imm8*: immigration by sex and age (2004 to 2015) –**imm**-.

*migr\_emi2*: emigration by sex and age (2004 to 2015) –**emi**-.

*demo\_r\_frate2*: fertility rate by age (2004 to 2015) –**fer**-.

*demo\_magec*: death by age and sex (2004 to 2015) –**dea**-.

*demo\_nsinrt*: first marriage rate by age and sex (2004 to 2015) –**mar**-.

*demo\_mlexpec*: life expectancy by age and sex (2004 to 2015) –**exp**-.

### **3. Machine Learning methodology**

R software (<https://rstudio.jrc.es/> )

Objective: forecasting the population by age (by year of age)

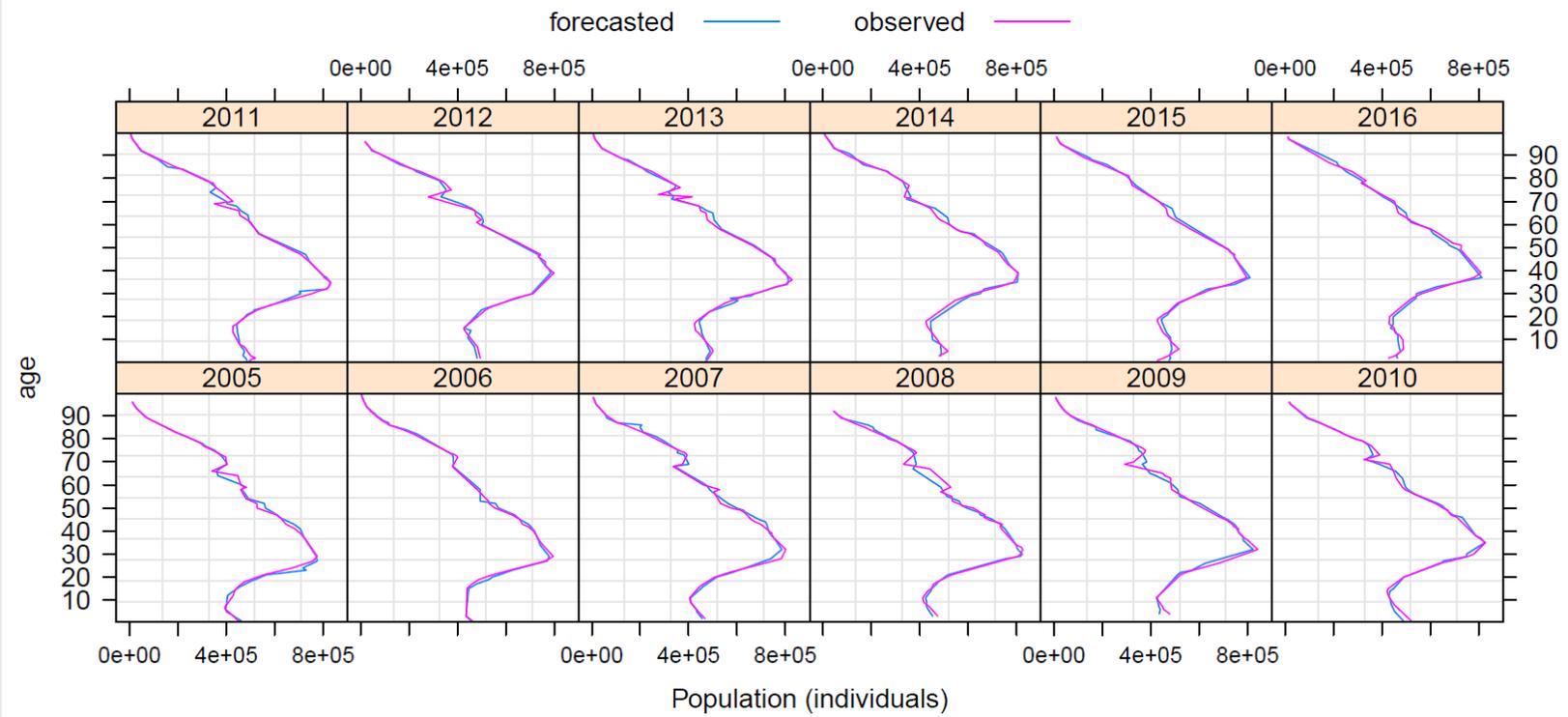
Forecasting approach:

$$\hat{p}_{year} = \hat{f}(pop_{year-1}, imm_{year-1}, emi_{year-1}, fer_{year-1}, dea_{year-1}, mar_{year-1}, exp_{year-1})$$

Algorithm: random forest (Breiman, 2001).

Crossvalidation: 10 hold-out fractions, repeated 3 times

## 4. Results



## 4. Results

Comparison with classical techniques:

Machine learning (ML) vs Arima & exponential smoothing (ES)

Units: population (individuals)

	<b>ML</b>	<b>Arima</b>	<b>ES</b>
RMSE	20318	81010	81247
MAE	14353	60596	62269

RMSE: root mean square deviation

MAE: mean absolute error

## **5. Conclusions**

- 1. Good fit outperforming traditional techniques*
- 2. Promising application of machine learning to traditional statistics*
- 3. Small data sets feed the machine learning approach*
- 4. Very linear relation among variables*
- 5. Further developments will be addressed*



## Joint Research Centre –Seville

Álvaro Gómez Losada  
[alvaro.gomez-losada@ec.europa.eu](mailto:alvaro.gomez-losada@ec.europa.eu)

Néstor Duch-Brown  
[nestor.duch-brown@ec.europa.eu](mailto:nestor.duch-brown@ec.europa.eu)