

Margherita Zuppardo

Violeta Calian

Ómar Harðarson



MACHINE LEARNING METHODS

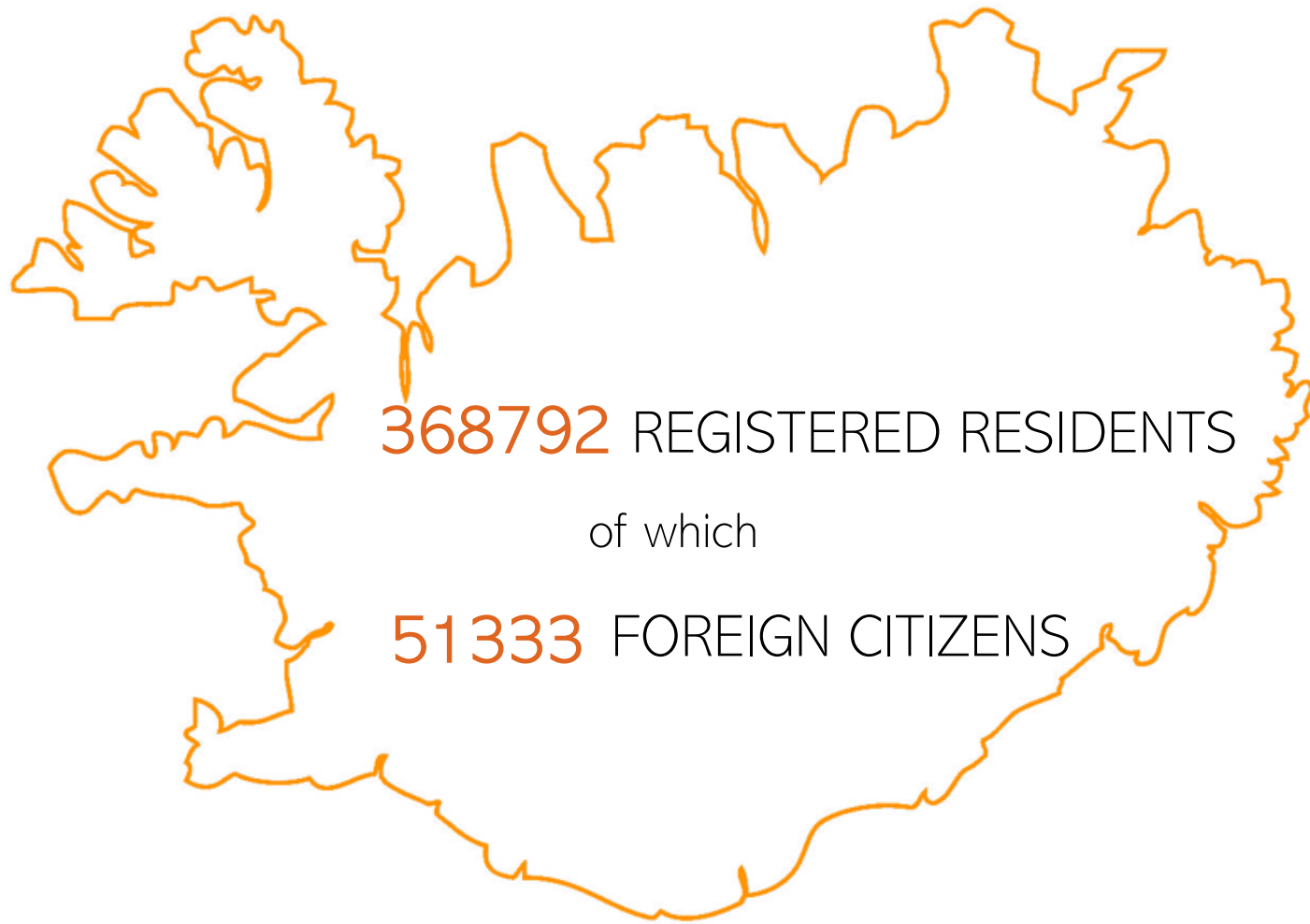
FOR ESTIMATING THE CENSUS POPULATION




Statistics Iceland

NSM 2022


2021 DIGITAL CENSUS



HOW MANY ARE ACTUALLY HERE?

 **People** Change of address → **Moving from Iceland**

Moving from Iceland

Notify a change of address 

Attn.

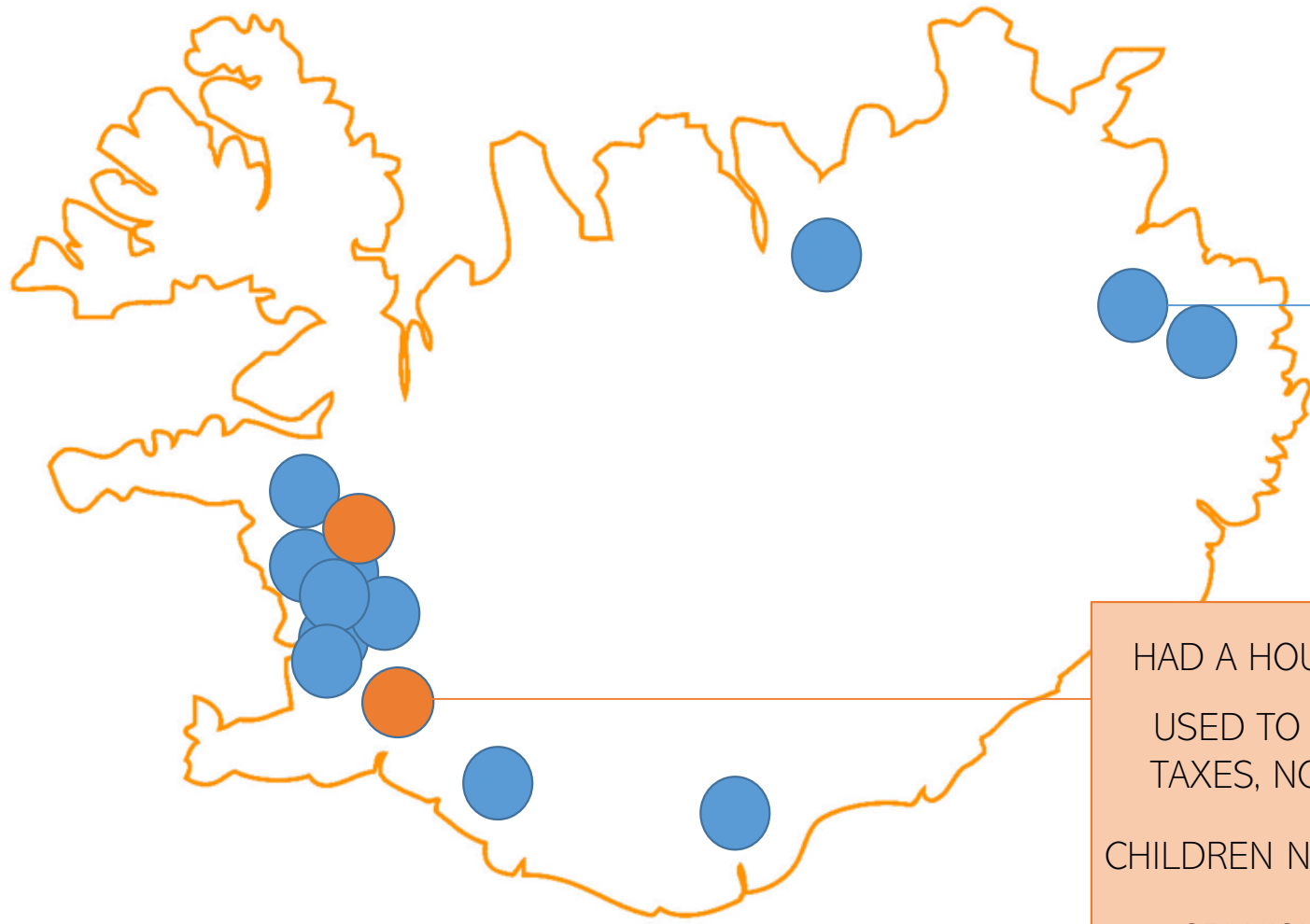
Transfers of legal domicile from Iceland **must be registered within 7 days.**

IN 2020:

7994 DE-REGISTRATIONS

5833 FROM NON-CITIZENS

SIGNS OF LIFE



HAD A HOUSE, SOLD IT
USED TO PAY INCOME
TAXES, NOT ANYMORE
CHILDREN NOT IN SCHOOL
SPANISH SPOUSE
ETC

OUT

OWNS A HOUSE
WORKS IN ICELAND
HAS CHILDREN IN SCHOOL
ETC

IN

BINARY CLASSIFICATION!

SOLUTION

- TRAIN MACHINE LEARNING CLASSIFICATION ALGORITHMS
- CHOOSE THE BEST ALGORITHM BASED ON PERFORMANCE MEASURES
- FIND THE OPTIMUM REGIME OF THIS ALGORITHM

STEP 1: BUILDING A DATAFRAME

SOURCES:

LFS survey 2014-2018

NATIONAL REGISTERS

TRAINING TABLE:

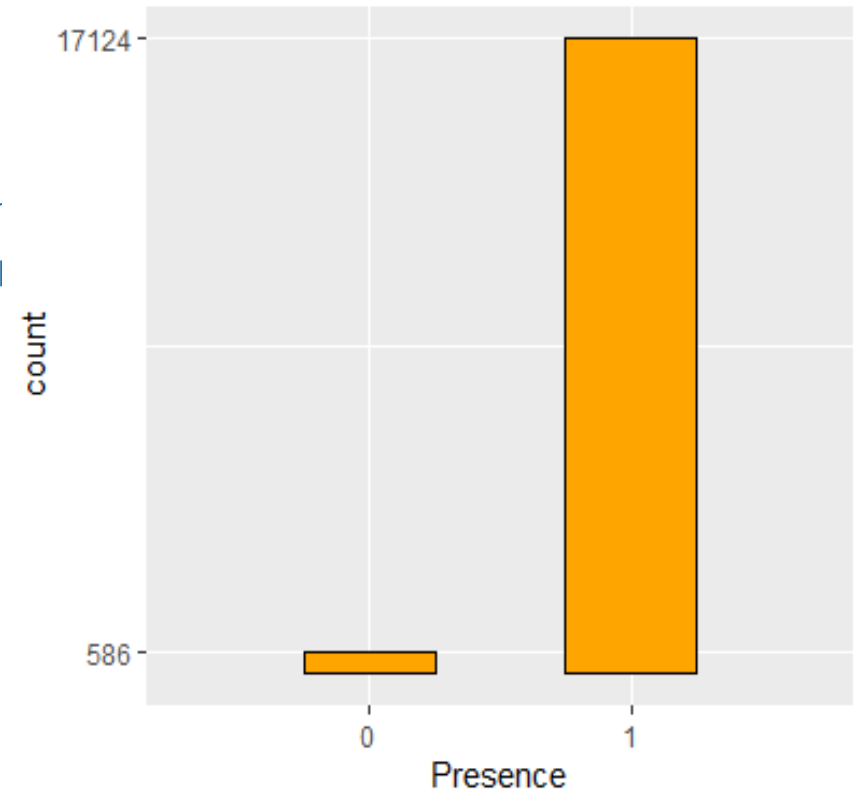
17710 ROWS

21 COLUMNS

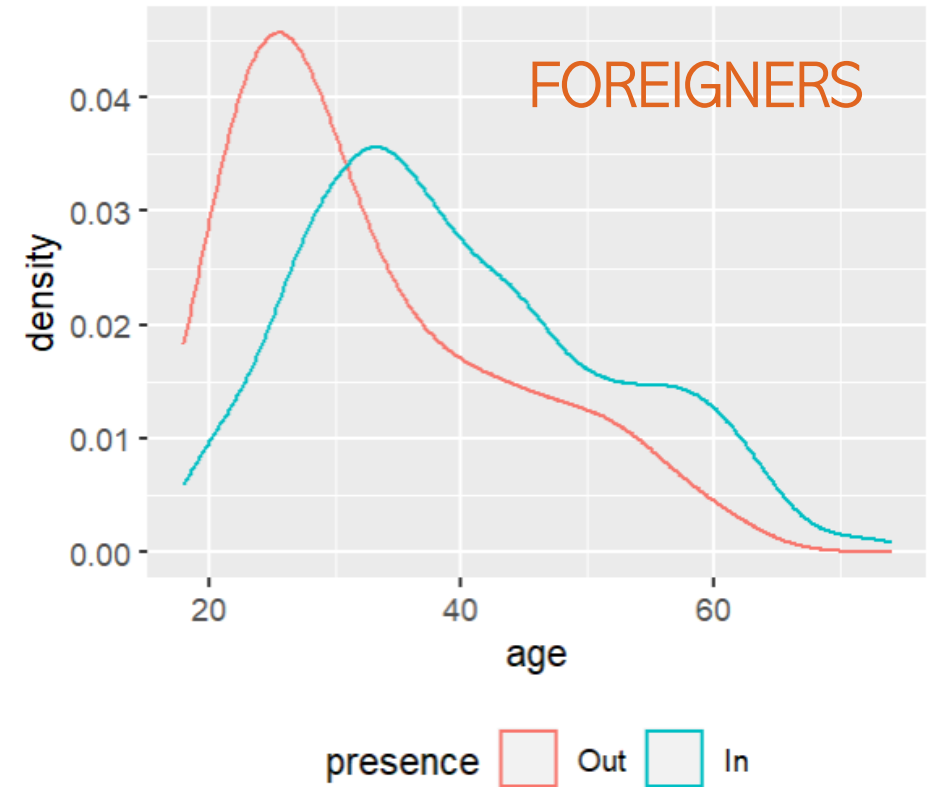
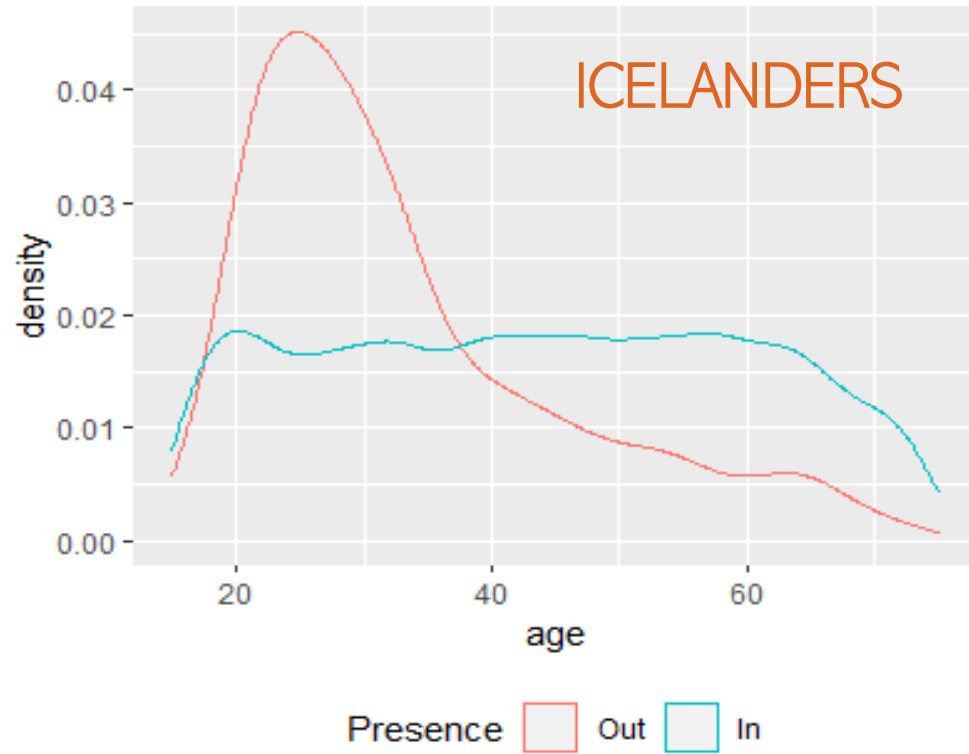
Presence

sex
country_of_birth
citizenship
region
ever_abroad
married
children
home_owner
student_abroad

age
income_1yr
income_2yr
months_worked_1yr
months_worked_2yr
year_max_income
year_max_brt
n_skoli
time_in_Iceland
n_movement3yr
n_changes



PREDICTORS



6 Available Models

The models below are available in `train`. The code behind these protocols can be obtained using the function `getModelInfo` or by going to the [github repository](#).

Show entries

Search:

Model	method	Value	Type	Libraries	Tuning Parameters
AdaBoost Classification Trees	adaboost		Classification	fastAdaboost	nIter, method
AdaBoost.M1	AdaBoost.M1		Classification	adabag, plyr	mfinal, maxdepth, coeflearn
Adaptive Mixture Discriminant Analysis	amdai		Classification	adaptDA	model
Adjacent Categories Probability Model for Ordinal Data	vglmAdjCat		Classification	VGAM	parallel, link
Bagged AdaBoost	AdaBag		Classification	adabag, plyr	mfinal, maxdepth
Bagged FDA using gCV Pruning	bagFDAGCV		Classification	earth	degree
Bagged Flexible Discriminant Analysis	bagFDA		Classification	earth, mda	degree, nprune
Binary Discriminant Analysis	binda		Classification	binda	lambda.freqs
Boosted Classification Trees	ada		Classification	ada, plyr	iter, maxdepth, nu
Boosted Logistic Regression	LogitBoost		Classification	caTools	nIter
C4.5-like Trees	J48		Classification	RWeka	C, M
C5.0	C5.0		Classification	C50, plyr	trials, model, winnow

CHOOSING A MODEL

Cumulative Probability Model for Ordinal Data	vglmCumulative	Classification	VGAM	parallel, link
DeepBoost	deepboost	Classification	deepboost	num_iter, tree_depth, beta, lambda, loss_type
Diagonal Discriminant Analysis	dda	Classification	sparsediscrim	model, shrinkage
Distance Weighted Discrimination with Polynomial Kernel	dwdPoly	Classification	kerndwd	lambda, qval, degree, scale
Distance Weighted Discrimination with Radial Basis Function Kernel	dwdRadial	Classification	kernlab, kerndwd	lambda, qval, sigma
Factor-Based Linear Discriminant Analysis	RFlda	Classification	HiDimDA	q
Flexible Discriminant Analysis	fda	Classification	earth, mda	degree, nprune
Fuzzy Rules Using Chi's Method	FRBCS.CHI	Classification	frbs	num.labels, type.mf
Fuzzy Rules Using Genetic Cooperative-Competitive Learning and Ensemble	FH.GBML	Classification	frbs	max.num.rule, popu.size, max.gen
Fuzzy Rules Using the Structural Learning Algorithm on	SLAVE	Classification	frbs	num.labels, max.iter, max.gen

<https://topepo.github.io/caret/>

CHOOSING A MODEL

	TRUE IN	TRUE OUT
PREDICTED IN	TP	FN
PREDICTED OUT	FP	TN

$$\text{Accuracy} = \frac{TN + TP}{N_{tot}}$$

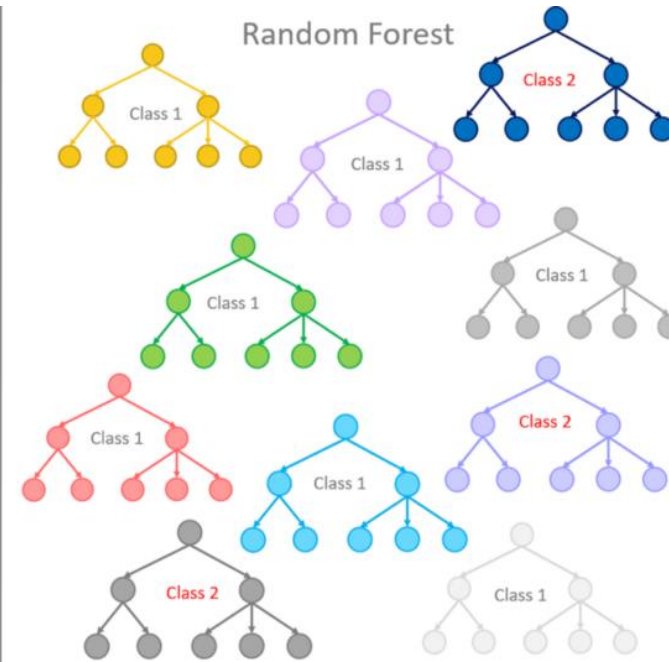
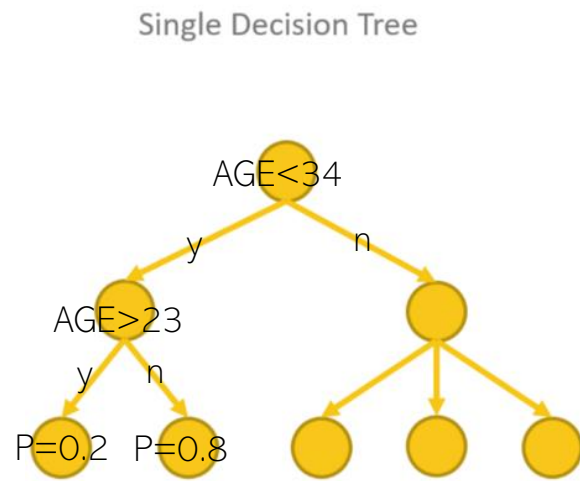
$$\text{Specificity} = \frac{TN}{TN + FP}$$

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

$$\text{Population error} = \frac{|FP - FN|}{N_{tot}}$$

Method	Accuracy (%)	Specificity (%)	Sensitivity (%)	Population error (%)
Register data	<96.7	0.0	100.0	>3.3
Logistic regression	96.8	14.2	99.6	2.6
Decision tree	97.0	16.4	99.8	2.3
Neural network	96.9	17.5	99.6	1.6
AdaBoost	95.1	26.3	97.5	0.0
Random forest (untuned)	97.0	25.1	99.5	2.1
Optimized RF (final model)	96.0	48.0	98.0	0.04
Latest results (revised data)	96.7	56.0	98.2	0.2

RANDOM FOREST



OUTPUT: PROBABILITY TO BE
INSIDE THE COUNTRY

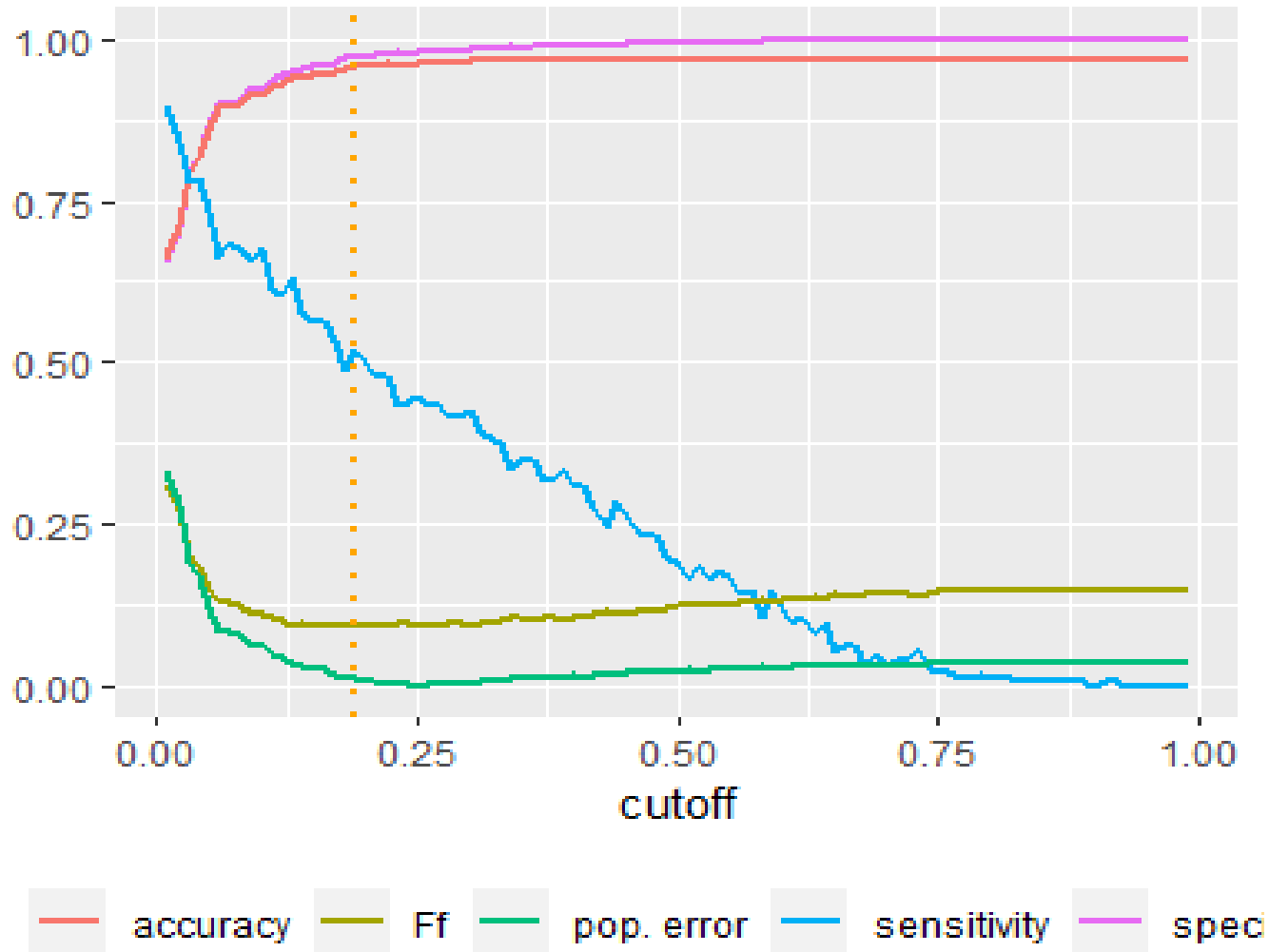
<https://rosaria-silipo.medium.com/>

MODEL TUNING

WE TUNE THE PROBABILITY CUTOFF

Specificity > 98%

WE OPTIMIZE EVERYTHING ELSE

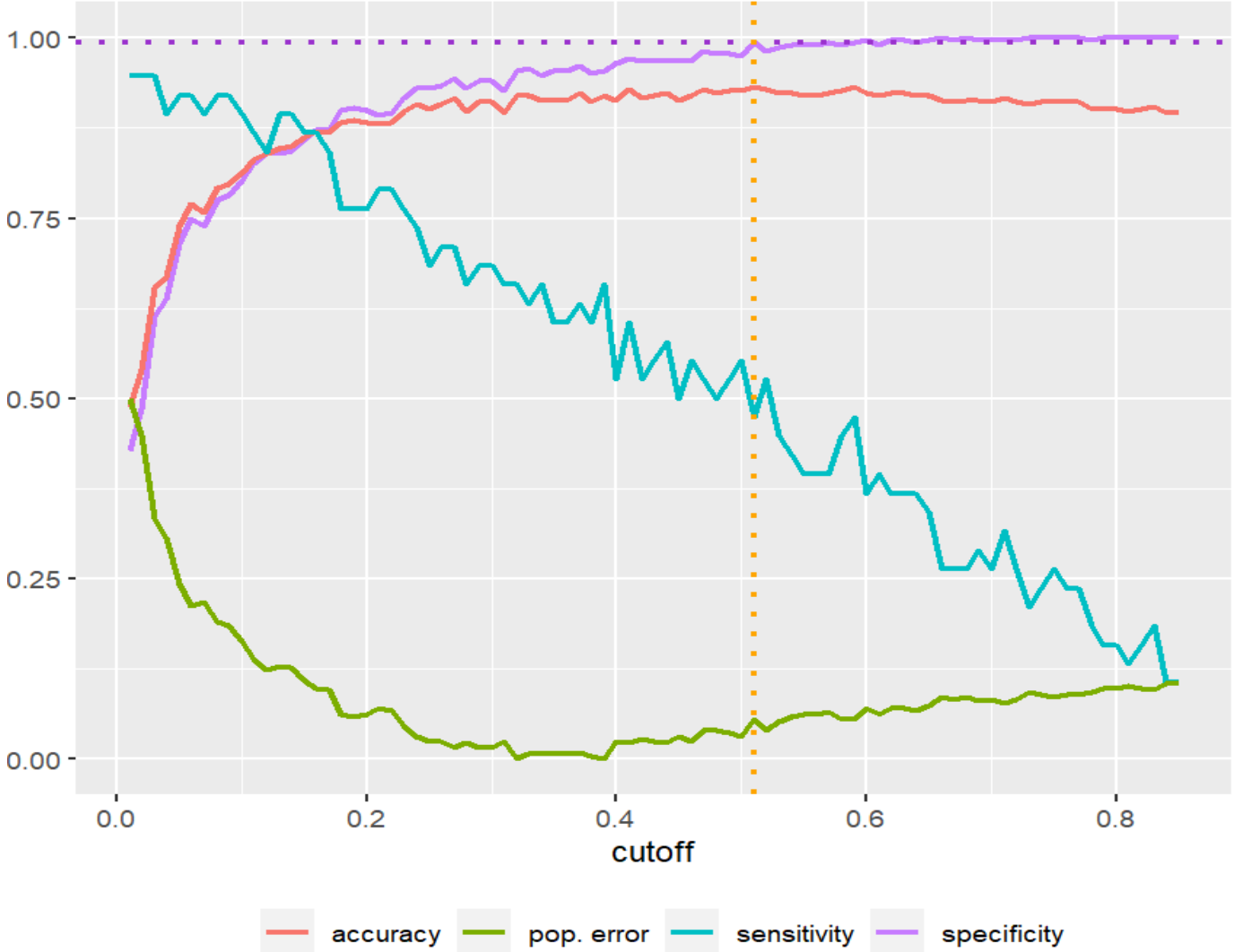


96.0	48.0	98.0	0.04
Accuracy	Sensitivity	Specificity	Pop. error

FOREIGN DATA

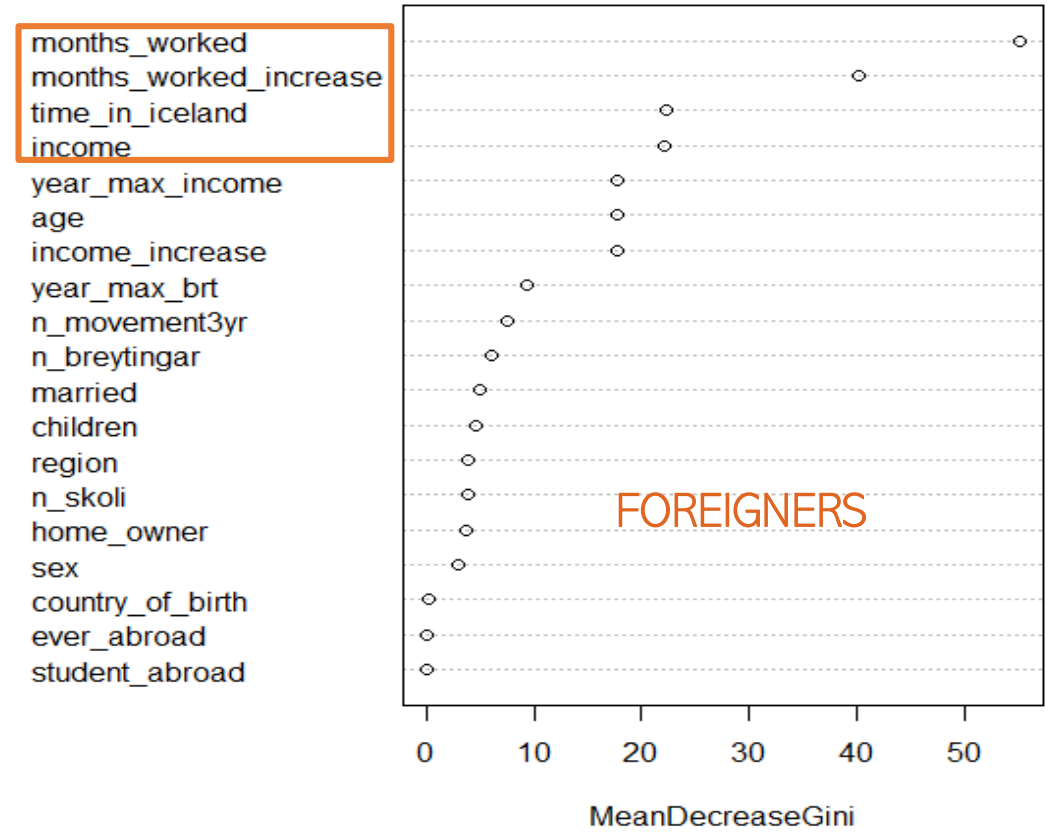
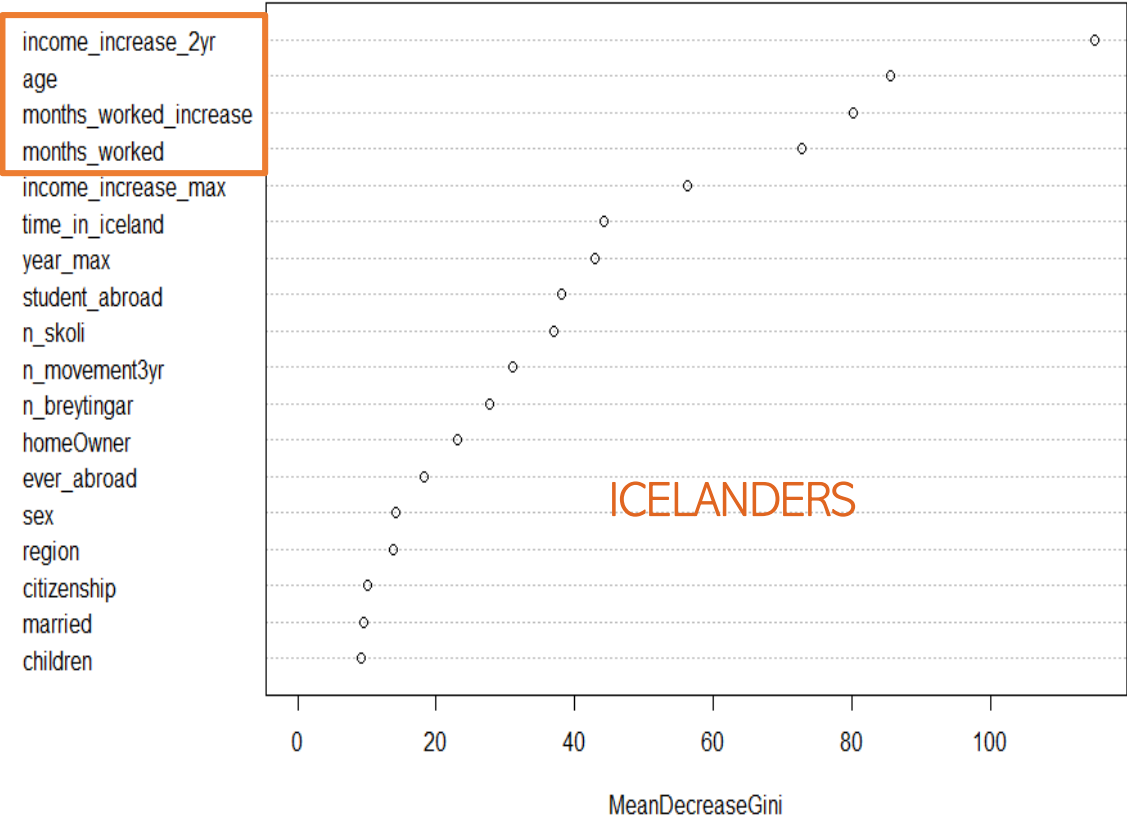
Specificity > 98%

MUCH NOISIER!



92.0	44.6	98.2	0.05
Accuracy	Sensitivity	Specificity	Pop. error

IMPORTANCE OF PREDICTORS



368 792 REGISTERED RESIDENTS

7100 PEOPLE OUT

OUT OF WHICH

3770 FOREGNEIRS

CENSUS RESULTS

	IN	OUT
In school 6-15 of age	45886	808
Not in school 6-15	537	8710
Working in November	181904	2325
Having a car	153608	1915

IN CONCLUSION

WE TRAINED A RANDOM FOREST MODEL TO PREDICT THE ACTUAL 2021 POPULATION IN ICELAND

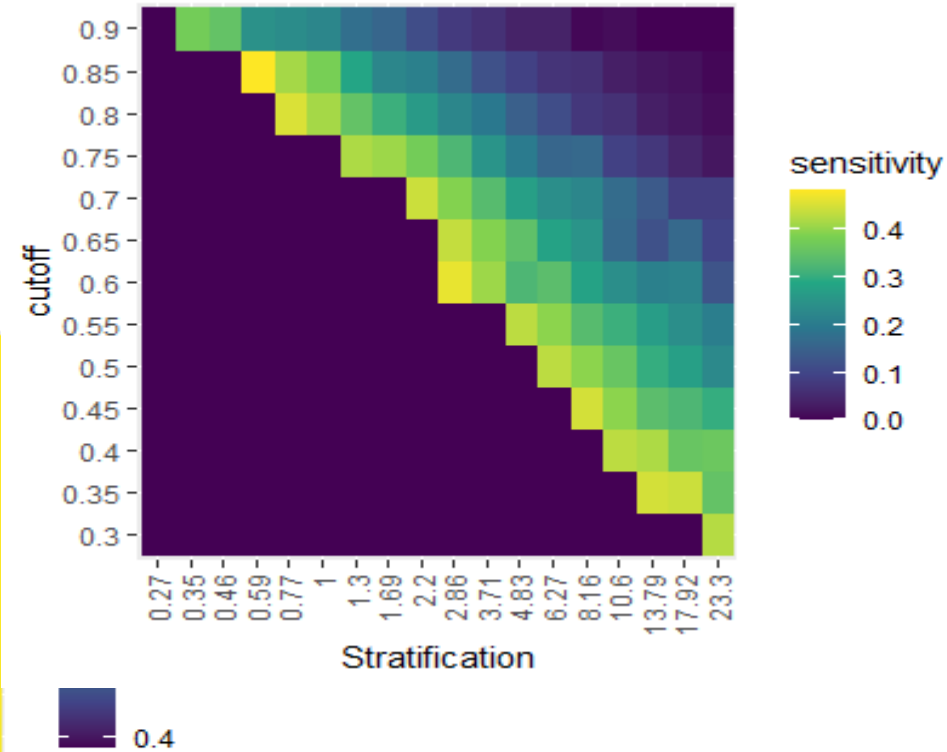
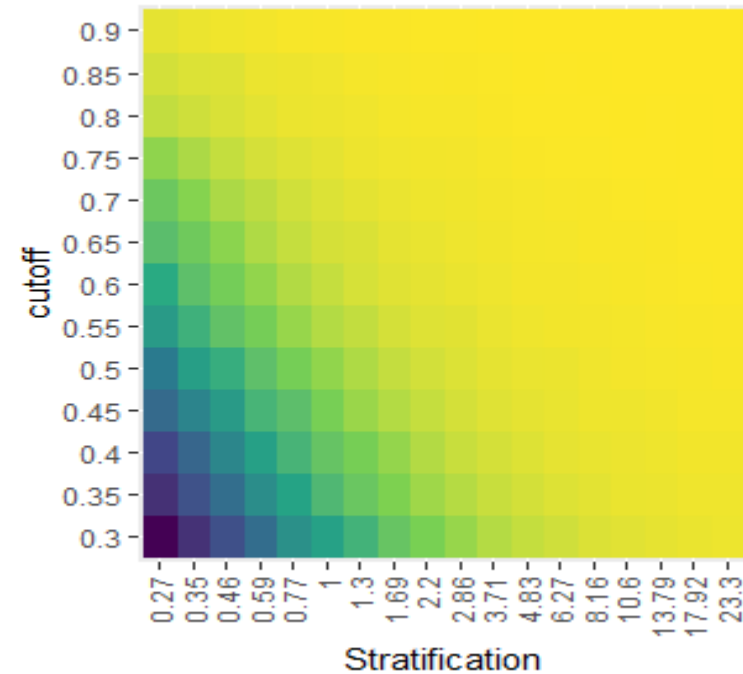
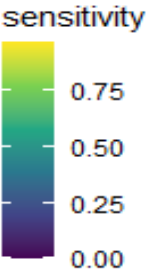
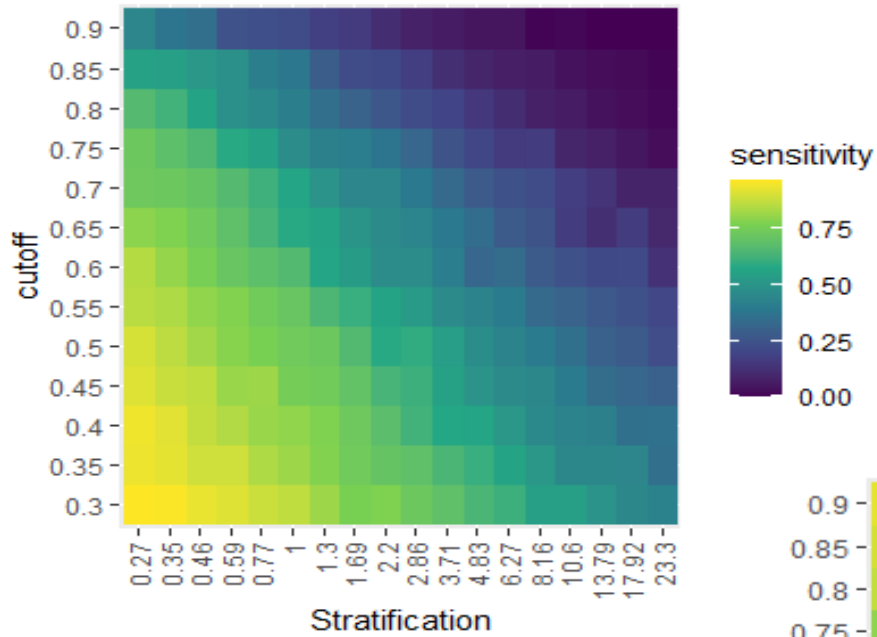
THE MODEL IS ACCURATE BUT 'SAFE' (HIGH SENSITIVITY)

IT CAN BE IMPROVED BY INCREASING THE QUALITY OF THE DATA, AND MAY BE USED REGULARLY IN THE FUTURE TO ESTIMATE THE ICELANDIC POPULATION

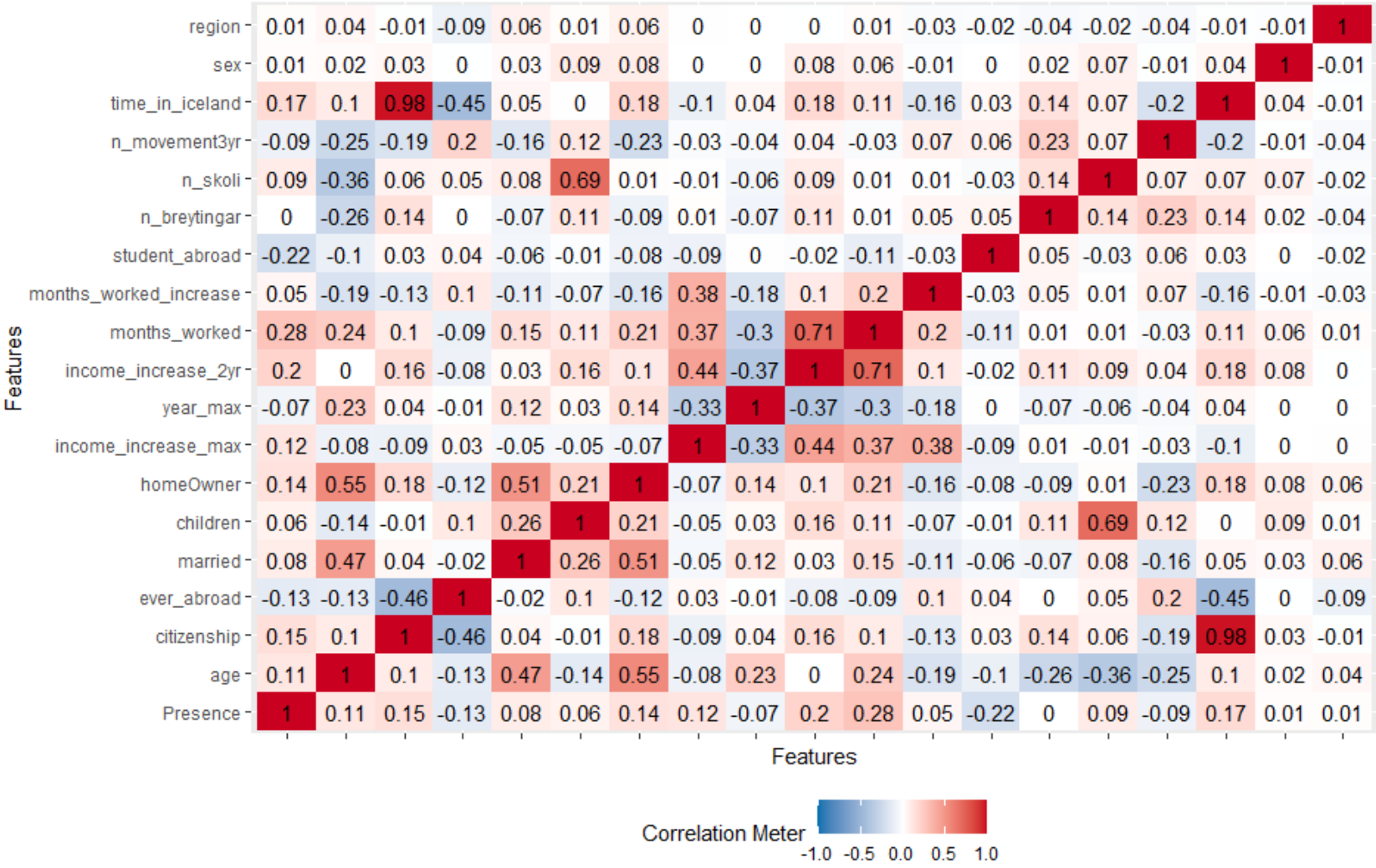


Statistics Iceland

APPENDIX: DATA STRATIFICATION



CORRELATIONS



GINI INDEX

$$Gini\ Index = 1 - \sum (P(x=k))^2$$