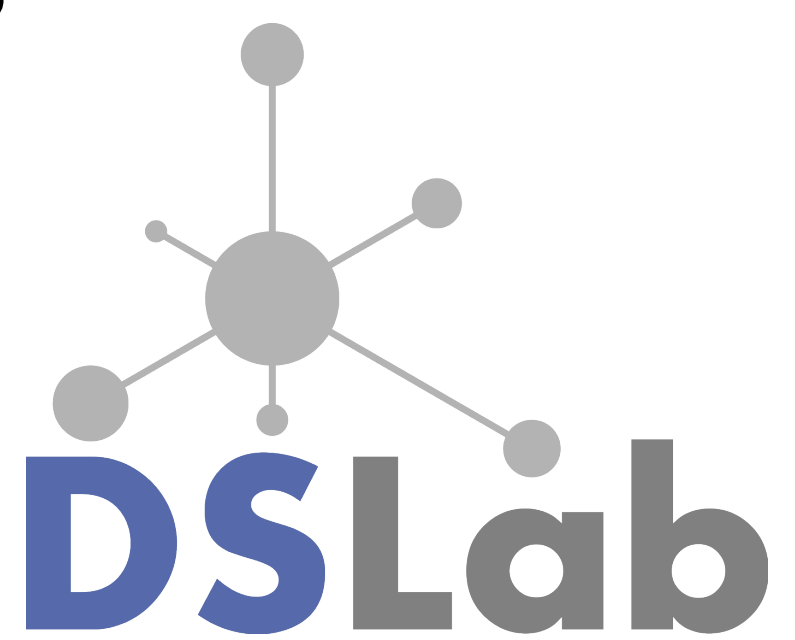




WHAT FACTORS DETERMINE UNEVEN SUBURBANISATION? EXPLAINING URBAN SPRAWL WITH MACHINE LEARNING

HONORATA BOGUSZ*, SZYMON WINNICKI AND PIOTR WÓJCIK
UNIVERSITY OF WARSAW, FACULTY OF ECONOMIC SCIENCES
DATA SCIENCE LAB



BACKGROUND

- suburbanisation is a shift of population from central urban areas into suburbs, resulting in the formation of (sub)urban sprawl
- unrestricted city growth leads to a variety of **formidable consequences**
- in order to execute **adequate social planning**, it is useful to identify factors which push migrants out of the city and pull them to the suburban boroughs

OBJECTIVES

- identify pulling features of boroughs which are key factors in predicting the number of migrants
- find a machine learning algorithm that exhibits the most accurate predictive performance
- uncover non-linear relationships between the number of migrants and the most important regressors

DATA AND METHODS

- Data:** 30 features of boroughs obtained from Polish Statistical Office, Open Street Map, Google Maps, National Electoral Commission, e-podroznik.pl and gratka.pl (most recent observations for year 2018 or 2019)
- Model benchmark:** we base our 7 predictive models on the *gravity model of migration* framework
- Models used:** Ordinary Least Squares, Lasso, Ridge Regression, Elastic Net, Support Vector Regression, Random Forest, Extreme Gradient Boosting
- EDA:** PCA with varimax rotation as a method of choosing variables out of correlated groups for algorithms incapable of variable selection
- Errors:** assessing models' performance by common benchmarks of RMSE, MAE and R^2 on validation and training samples (LOOCV used)
- XAI:** *Permutation-based feature importance* and *Accumulated Local Effects* plots

CONCLUSIONS

- identified important pulling factors (see ALE plots on the right): good amenities, mean relative income, progressiveness of a community, spiritual needs
- identified important pushing factors (see ALE plots on the right): distance and population density (the latter in contrast with the existing literature)
- XGB outdoes OLS by only a slight margin
- assuming linear relationships is infeasible with respect to some regressors - non-linear relationships can be interpreted with ALE plots

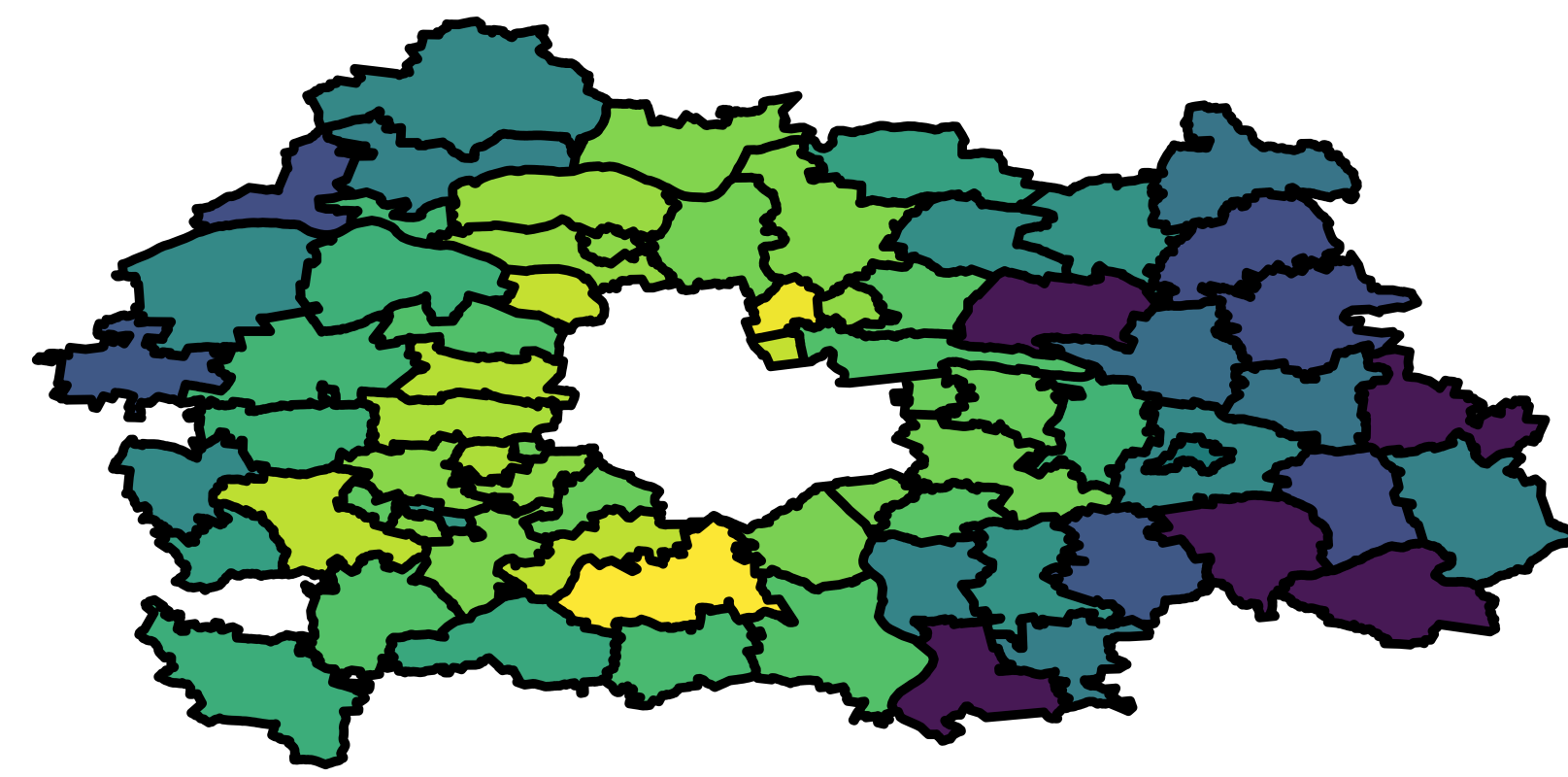
CONTACT INFORMATION

*email h.bogusz@uw.edu.pl
twitter HonorataBogusz

EMPIRICAL VERIFICATION

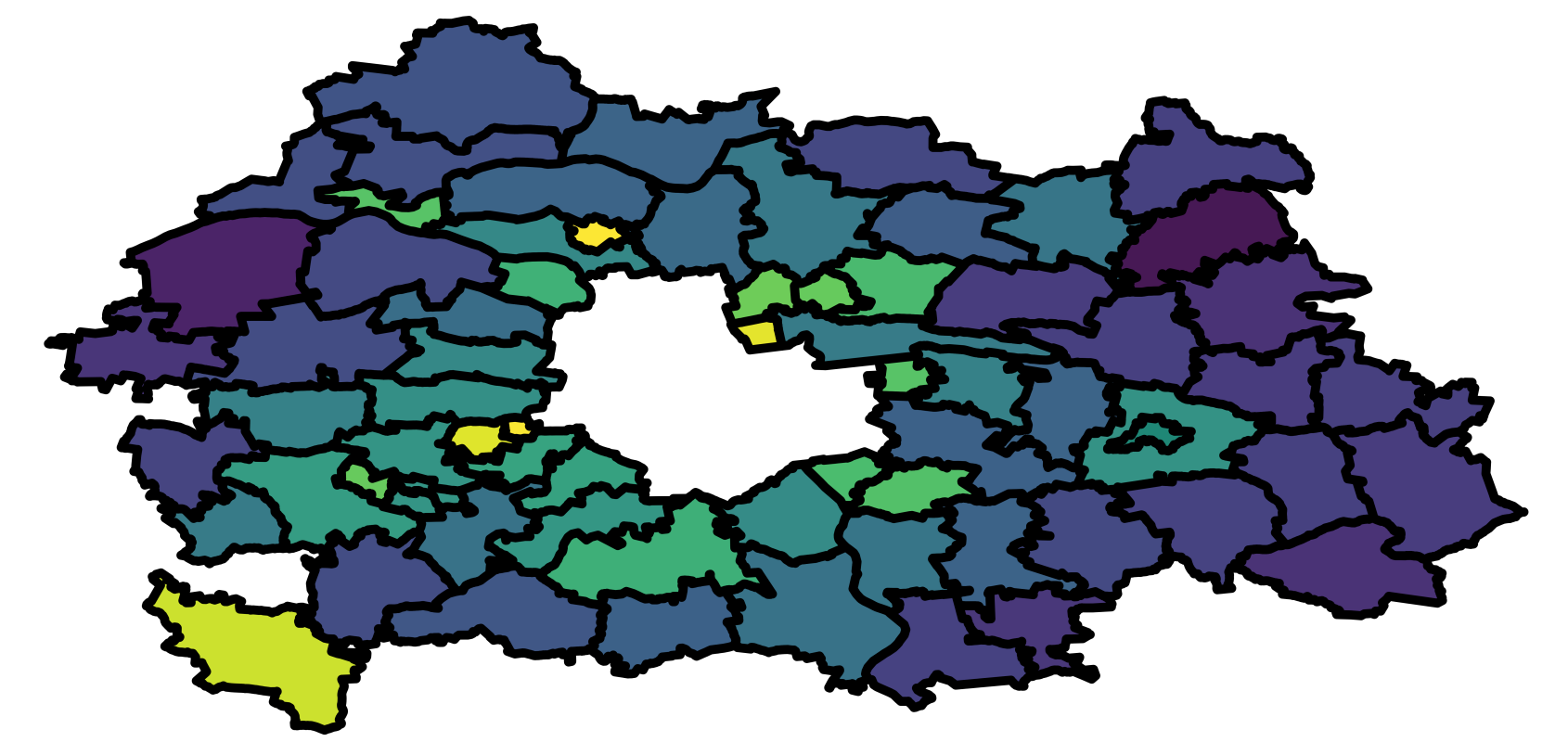
Maps of the key variables of the gravity model of migration, and income per capita (important pulling factor according to the literature) - logarithmed for clearance

a) log (migrants)



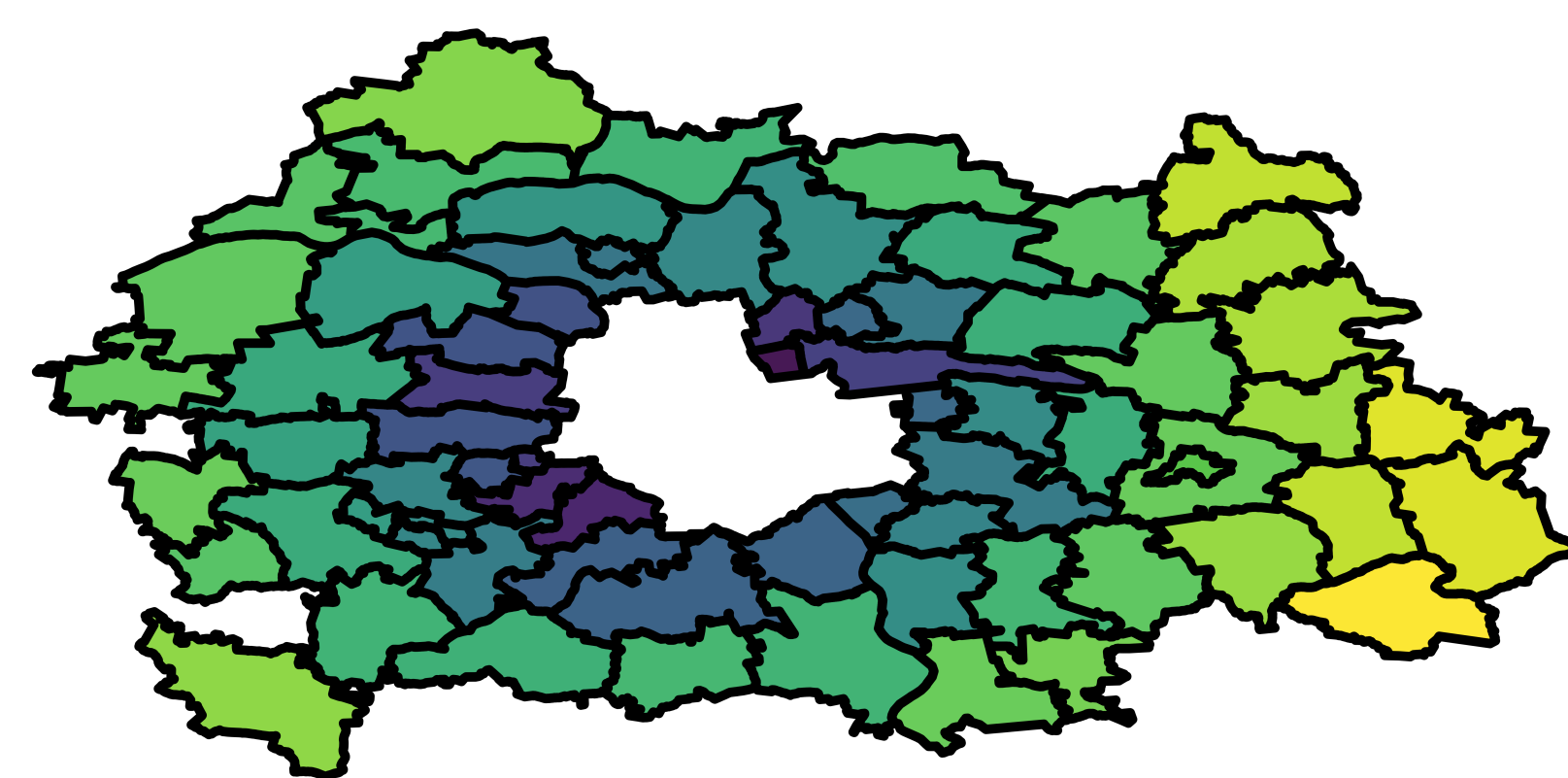
log # of migrants
Data Source: Polish Statistical Office

b) log (population density)



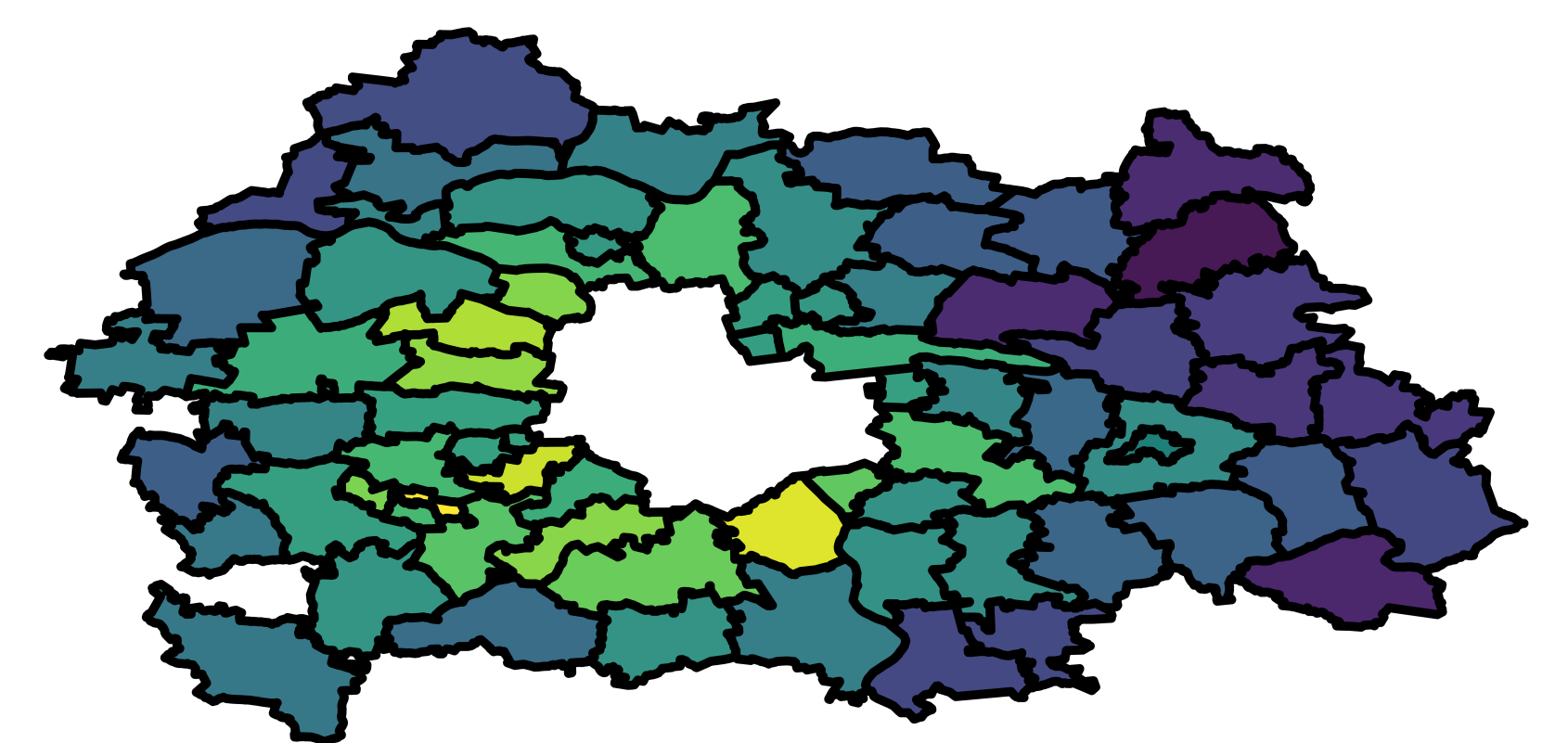
log # of people / km²
Data Source: Polish Statistical Office

c) log (distance)



log km
Data Source: Google Maps

d) log (income per capita)



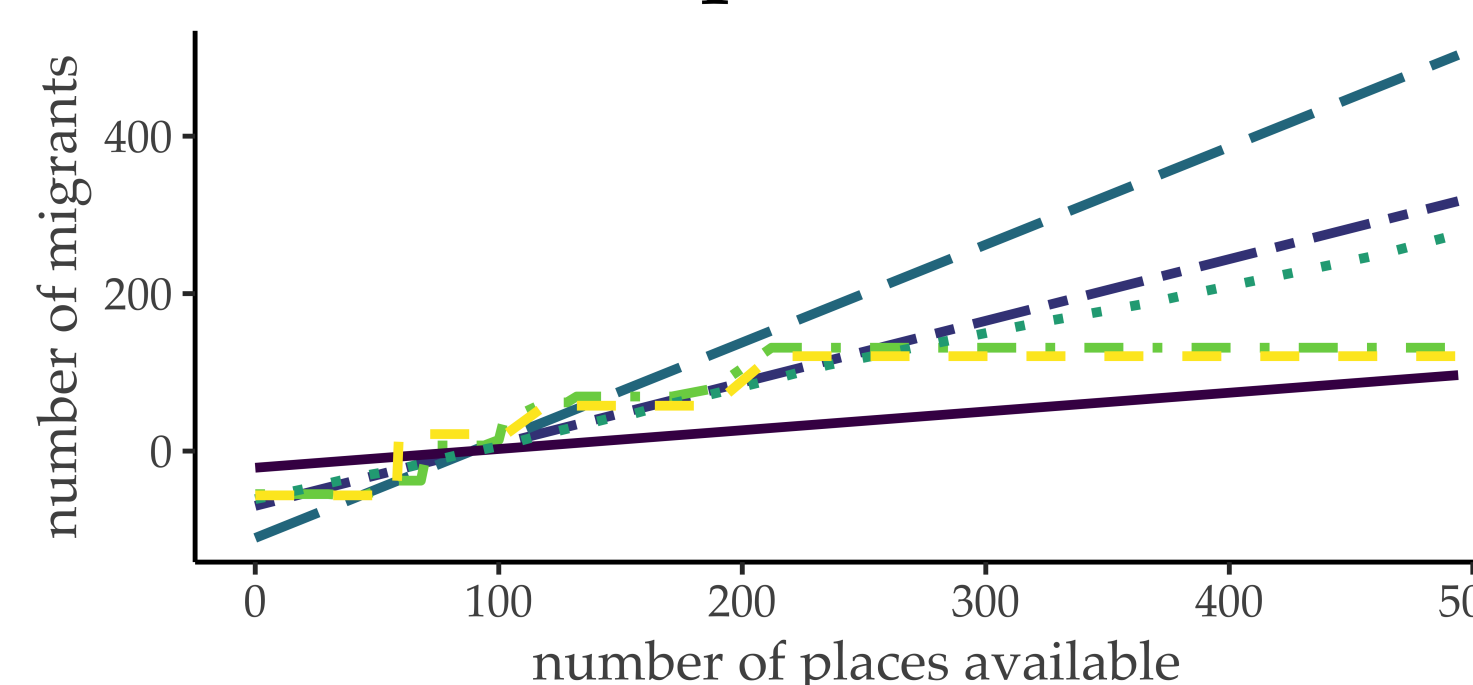
log ratio
Data Source: Polish Statistical Office

Model errors (validation and training sample, calculated by LOOCV)

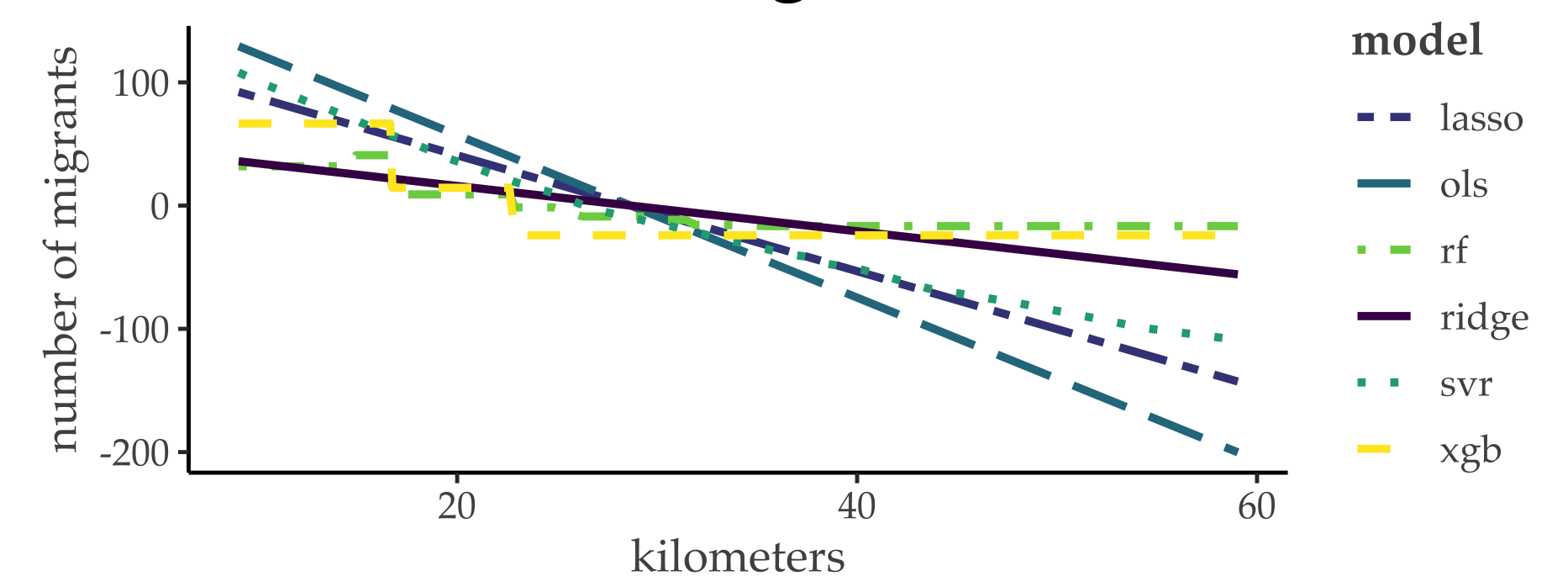
MODEL	VALIDATION			TRAIN		
	RMSE	MAE	R2	RMSE	MAE	R2
OLS	111.12	79.86	0.62	97.55	71.93	0.71
RIDGE	136.15	87.50	0.43	113.59	73.25	0.61
LASSO	124.63	83.92	0.53	97.99	66.81	0.71
SVR	119.15	78.33	0.57	65.63	34.45	0.87
RF	123.05	75.63	0.54	80.29	53.00	0.80
XGB	109.56	70.07	0.63	78.25	45.76	0.81
ELASTIC	124.63	83.92	0.53	97.99	66.81	0.71

ALE plots for 6 most important features as indicated by mean PFI (of all models)

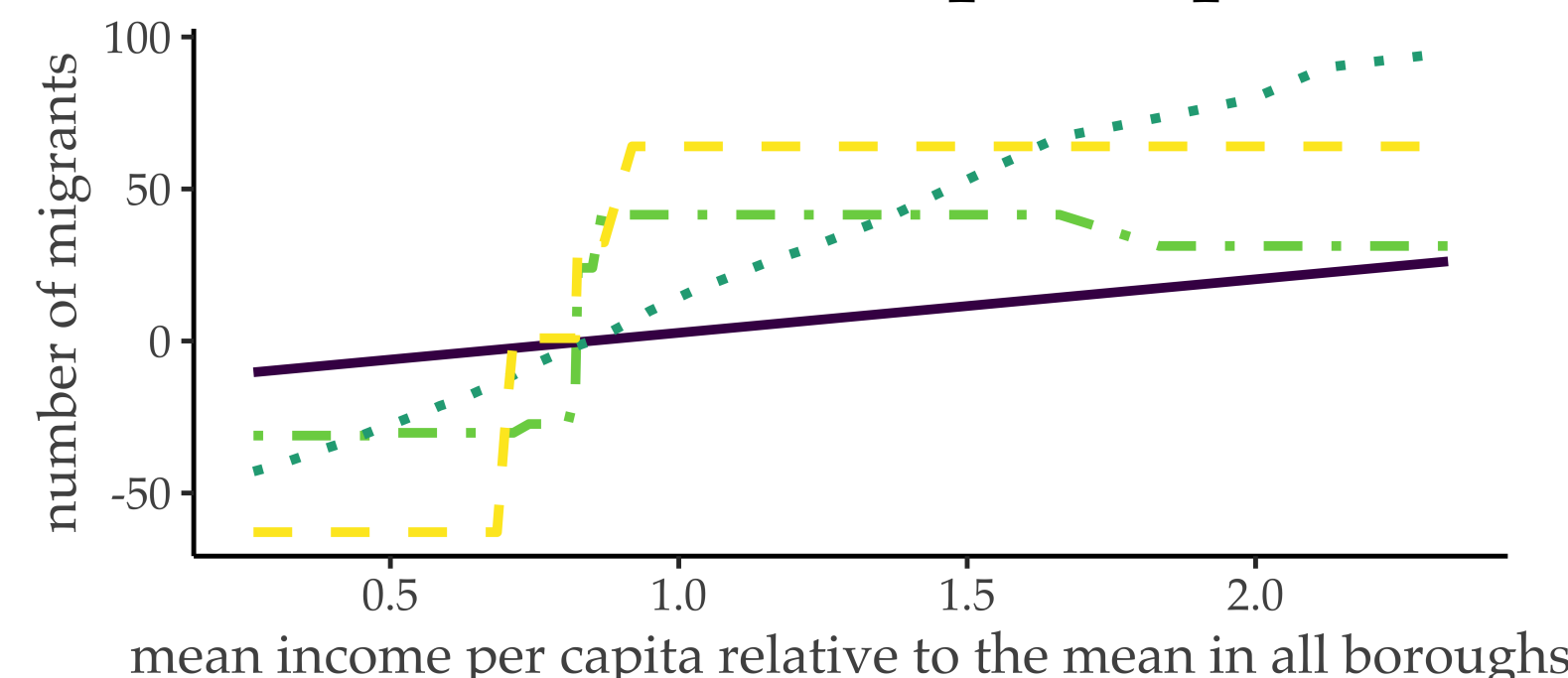
a) number of places in nurseries



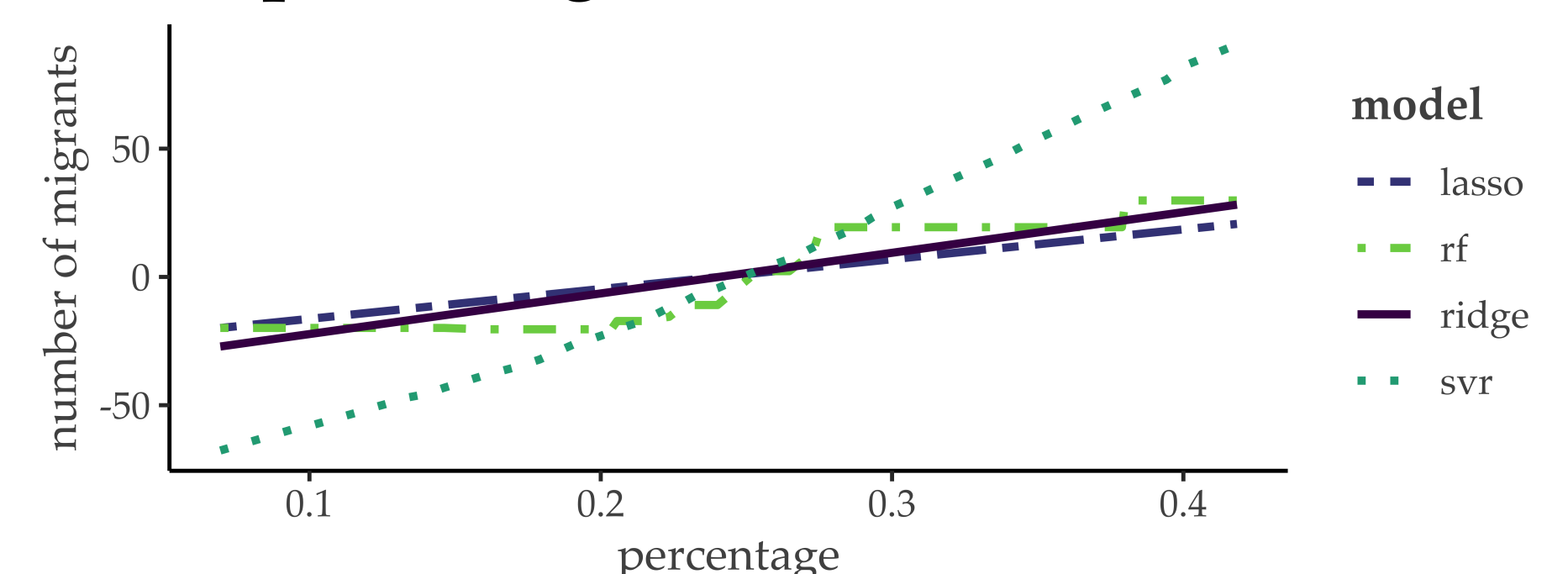
b) distance (straight line)



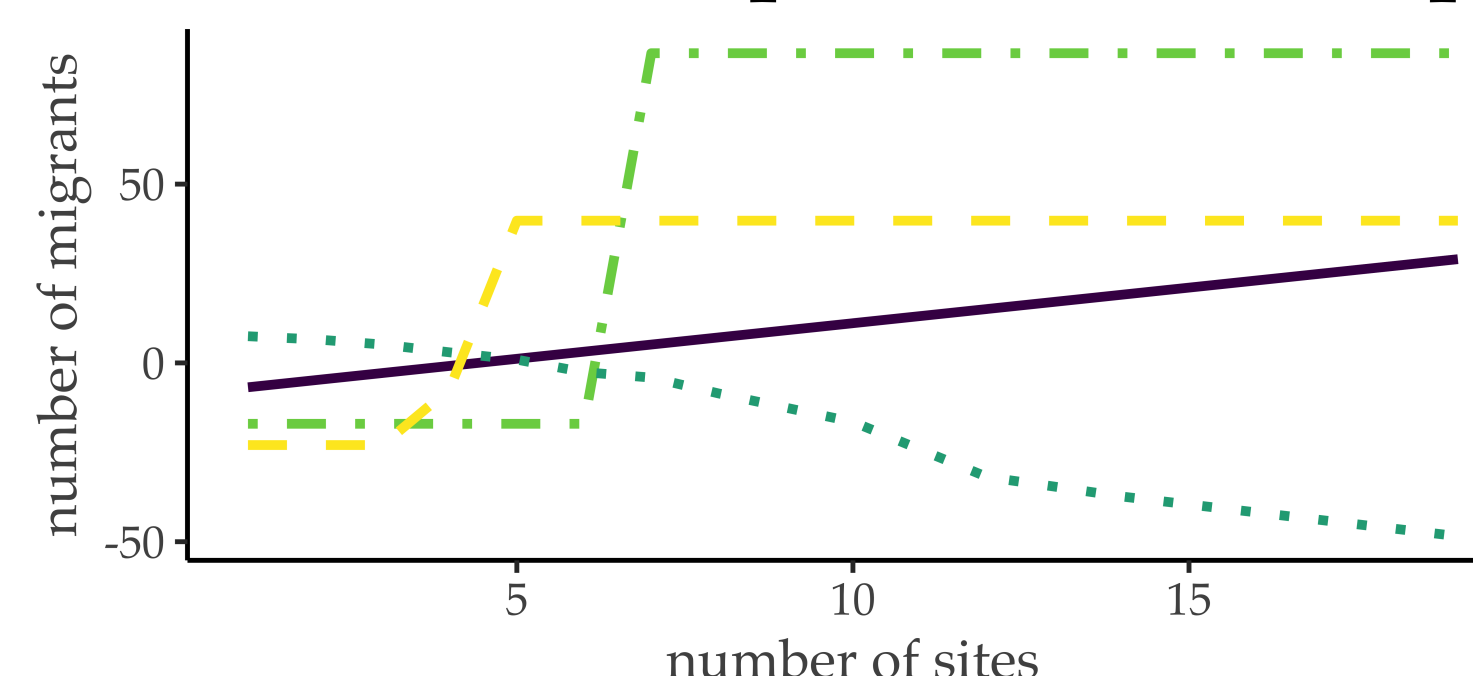
c) relative income per capita



d) percentage of votes for KO in 2019



e) number of places of worship



f) population density

