ESRA Conference 2019 Predictive Modeling and Machine Learning in Survey Research Zagreb, 18 Jul 2019



# UPDATE ON MACHINE LEARNING IN DATA ANALYSIS FOR SOCIAL RESEARCH

Arne Bethmann (MEA) Jonas Beste (IAB) Giuseppe Casalicchio (LMU) Bernd Bischl (LMU)





This project has received funding from the European Union under grant greement VS/2018/0285 and the European Union's Horizon 2020 research and innovation programme under grant greements No 676536, No 654221 SPONSORED BY THE









### **Machine Learning Tasks**

	Numerical outcome	Categorical outcome
<b>Supervised</b> outcome observed	<ul> <li>Regression task</li> <li>Predict the average fiancial wealth of migrants and non-migrants</li> <li>Predict blood measures from dried blood spots</li> </ul>	<ul> <li>Classification tasks</li> <li>Predict propensity scores to analyse the difference in health between migrants and non- migrants</li> <li>Analyze the effect of education on old age poverty</li> </ul>
<b>Unsupervised</b> outcome unobserved	<ul> <li>Factor analysis tasks</li> <li>Predict individual factor scores on the EURO-D scale</li> <li>Predict movement intensity from accelerometry data</li> </ul>	<ul> <li>Clustering tasks</li> <li>Find similar groups of labour market trajectories among SHARE individuals</li> <li>Find typology of SHARE interviewers in order to identify suspicious clusters</li> </ul>

## Some algorithms/models



	Numerical outcome	outcome Categorical outcome	
<b>Supervised</b> outcome observed	<ul> <li>Regression tasks</li> <li>Ordinary Least Squares</li> <li>Linear Models</li> <li>Splines, Lasso, Ridge</li> <li>Regression Trees, Random Forests</li> <li>Support Vector Machines</li> <li>k-Nearest-Neighbors</li> <li>Neural Nets, "Deep learning"</li> </ul>	<ul> <li>Classification tasks</li> <li>Generalized Linear Models (Logistic or probit,)</li> <li>Splines, Lasso, Ridge</li> <li>Classification Trees, Random Forests</li> <li>Support Vector Machines</li> <li>k-Nearest-Neighbors</li> <li>Neural Nets, "Deep learning"</li> <li>Naive bayes</li> </ul>	
<b>Unsupervised</b> outcome unobserved	<ul> <li>Factor analysis tasks</li> <li>Principal Components Analysis (PCA)</li> <li>Exploratory Factor Analysis (EFA)</li> <li>Confirmatory Factor Analysis (CFA)</li> </ul>	<ul> <li>Clustering tasks</li> <li>Hierarchical clustering (k- Means, Linkage, Ward)</li> <li>Model based (Finite mixtures, )</li> <li>Latent Class Analysis (for ordinal outomes)</li> </ul>	

... among many others

## Estimating effects as a regression tasks

- Estimate effect of variable on outcome
- Produce easily interpretable statistics
- Applicable to any (ML) model that yields individual predictions (i.e. model agnostic)



# **Regression tasks**



#### **Example data**





#### **Example learner/model: regression tree**





- 1. Estimate (complex, multivariate) model
- 2. Set all cases in sample to first value for variable (e. g. gender = 0 "female")
- 3. Predict response for all cases
- 4. Set all cases in sample to next value for variable (e. g. gender = 1 "male")
- 5. Predict response for all cases again
- 6. Calculate AME / APE as mean difference between two predictions for all cases

# Average marginal effects



	Linear Model	<b>Regression Tree</b>
Gender (male vs. female)		
AME / APE	3.61€	3.68€
95 % CI Lower Bound	2.43€	2.42€
95 % CI Upper Bound	4.65€	4.91€
Age (per year)		
AME / APE	0.32€	0.75€
95 % CI Lower Bound	-0.05€	0.09€
95 % CI Upper Bound	0.66€	1.32€

Data: ISSP Germany 2012; Cls: bootstrap estimates



 "ame" R package (Average Marginal Effects)

Giuseppe Casalicchio (Comp. Statistics, LMU) https://github.com/compstat-Imu/ame

- Integrates with mlr "Machine Learning in R" – a very comprehensive ML framework Bernd Bischl et al. (*Comp. Statistics, LMU*) <u>https://github.com/mlr-org/mlr</u>
- (Still) Work in progress:
   Paper discussing approach and implementation (Beste, Bethmann & Casalicchio)

## Pass by at our booth



