

Thema	Abstract	Literatur
Updating Intercensal Indicators based on geospatial data	Censuses are fundamental building blocks of most modern-day societies, yet collected every 10 years at best. The thesis assesses different approaches for census updating techniques by incorporating auxiliary information in order to take ongoing subnational population shifts into account. The methods adds satellite imagery as additional data source to derive disaggregated estimates of poverty indicators in Mozambique. The performance of the methods is evaluated using data from two different census periods.	1. T. Koebe, A. Arias-Salazar, N. Rojas-Perilla, and T. Schmid. Intercensal updating using structure-preserving methods and satellite imagery. <i>Journal of the Royal Statistical Society Series A</i> , 185: 170-196, 2022. 2. A.-K. Kreutzmann, S. Pannier, N. Rojas-Perilla, T. Schmid, M. Templ, and N. Tzavidis. The R package emdi for estimating and mapping regionally disaggregated indicators. <i>Journal of Statistical Software</i> 91(7): 1-33, 2019.
Multinomial Mixed Models in Small Area Estimation	The topic deals with small area estimation for the composition of categorical variables. A possible example is, for instance, the estimation of multidimensional poverty in small areas. From a methodological perspective, the thesis deals with small area estimation based on multinomial (mixed) models including the computation of the mean squared error of the estimators. A simulation study could be conducted to investigate the properties of the estimators. Finally, the methods can be applied to poverty data from Mexico.	1. E. Lopez-Vizcaino, M. J. Lombardia, and D. Morales. Multinomial-based small area estimation of labour force indicators. <i>Statistical Modelling</i> , 13(2):153-178, 2013. 2. I. Molina, A. Saei, and M. J. Lombardia. Small area estimates of labour force participation under a multinomial logit mixed model. <i>Journal of the Royal Statistical Society: Series A</i> , 170(4):975-1000, 2007.
Nonparametric regression with application to small area problems	Within the field of small area estimation the best linear unbiased predictor (BLUP) or its empirical version, the empirical BLUP (short EBLUP), is a widely used tool to estimate statistics of small areas, even though it is sensitive to model misspecification. More robust estimators like the penalised spline regression, which avoids functional specification, are brought forward as alternatives. Opsomer et al. (2008) suggest an estimator for unit-level data that treats the penalisation coefficient as a random effect, which allows the use of standard EBLUP theory in further analysis. Their estimator is known as nonparametric EBLUP. The nonparametric EBLUP framework could be compared to the standard EBLUP and to Random forests within a simulation study and an application with real data.	1. J. Opsomer, F. Breidt, G. Claeskens, G. Kauermann, and G. Ranalli. Nonparametric small area estimation using penalized spline regression. <i>Journal of the Royal Statistical Society. Series B (Statistical Methodology)</i> , 70(1):265-286, 2008. 2. D. Pfeiffermann. Small Area Estimation: New Developments and Directions. <i>International Statistical Review/ Revue Internationale de Statistique</i> , 70(1):125-143, 2002.
Poverty Mapping using ELL-Methodology	The World Bank Method, also called ELL (Elbers, Lanjouw and Lanjouw) method, is a small area estimation strategy used to estimate poverty indicators. This estimation technique is implemented in the World Bank software PovMap. In a thesis, a package could be programmed that also makes this estimation technique available to a broader audience via the R programming language. Furthermore, the ELL can be adjusted and analysed based on the comparison to other established methods in the field, such as the empirical best prediction method (EBP).	1. C. Elbers, J. O. Lanjouw, and P. Lanjouw. Micro-level estimation of poverty and inequality. <i>Econometrica</i> 71(1): 355-364, 2003. 2. ANSD. Demographic and Health and Multiple Indicator Cluster Survey (EDS-MICS) 2010-2011. 1. T. Schmid, and R. Munnich. Spatial robust small area estimation. <i>Statistical Papers</i> , 55(3):653-670, 2014. 2. R. B. Papalia, C. Bruch, T. Enderle, S. Falorsi, A. Fasulo, E. Fernandez Vazquez, M. Ferrante, J.-P. Kolb, R. Münnich, S. Pacei, R. Priam, P. Righi, T. Schmid, N. Shlomo, F. Volk, and T. Zimmermann. Best practice recommendations on variance estimation and small area estimation in business surveys. BLUE-ETS. SP1- Cooperation-Collaborative Project - Small or medium-scale focused research project, 2013.
Quantile Regression in Complex Design Samples	The Quantile Regression (QR) is an extension of the (multiple) linear regression that specifies the change of conditional quantiles of a dependent variable, instead of the change the conditional mean. The additional information gathered from this method might uncover new insights in datasets such as the one of the Demographic and Health Survey (DHS) of Senegal. QR could be used to estimate the conditional quantiles of wealth indicators such as the Wealth Index Factor Score for women and men, which uses the total value of an individual's assets to determine wealth, instead of commonly used monetary variables.	1. R. Koenker. <i>Quantile Regression</i> . Cambridge University Press: Cambridge, New York, 2005. 2. ANSD. Demographic and Health and Multiple Indicator Cluster Survey (EDS-MICS) 2010-2011.
Robust Small Area Estimation Methods in the Context of Spatial Correlations	The empirical best linear unbiased predictor (EBLUP), which is commonly used in small area statistics, uses the assumptions of normality and uncorrelated random effects. The robust version of EBLUP (REBLUP) can be used when a violation of normality is assumed. The SREBLUP is another estimation technique that additionally considers spatial correlated area effects in the model. A possible thesis could compare the performance of these estimators for different scenarios according to each estimator's underlying assumptions in a model-based simulation study. Real world data that displays spatial dependency and deviation of normality can be used.	1. T. Schmid, and R. Munnich. Spatial robust small area estimation. <i>Statistical Papers</i> , 55(3):653-670, 2014. 2. R. B. Papalia, C. Bruch, T. Enderle, S. Falorsi, A. Fasulo, E. Fernandez Vazquez, M. Ferrante, J.-P. Kolb, R. Münnich, S. Pacei, R. Priam, P. Righi, T. Schmid, N. Shlomo, F. Volk, and T. Zimmermann. Best practice recommendations on variance estimation and small area estimation in business surveys. BLUE-ETS. SP1- Cooperation-Collaborative Project - Small or medium-scale focused research project, 2013.
Comparing area-level to unit-level models in model-based settings	First, different methods (unit-level models under limited data vs. area-level models (with and without aggregation)) will be compared in 4 model-based settings. If possible, it should then be considered how extreme settings might look like in which one of the methods is most convincing. A little application to real-world data is also possible.	1. Rao, J. N. K. and I. Molina (2015). <i>Small Area Estimation (Second Edition)</i> . Hoboken: John Wiley & Sons. 2. Corral, Paul; Molina, Isabel; Cojocar, Alexandru; Segovia, Sandra. 2022. <i>Guidelines to Small Area Estimation for Poverty Mapping</i> . Washington, DC : World Bank.
Working with saeTrafo: Implementation of an analytical MSE for log-transformed unit-level models	Based on the publication of Molina, I. and N. Martin (2018). Empirical best prediction under a nested error model with log transformation. <i>The Annals of Statistics</i> 46(5), 1961-1993, this analytical MSE is to be implemented for the saeTrafo package. An application to EUSILC data is also intended for this method.	1. Würz, N. (2022). saeTrafo: Transformations for Unit-Level Small Area Models. R package version 1.0.0. 2. Molina, I. and N. Martin (2018). Empirical best prediction under a nested error model with log transformation. <i>The Annals of Statistics</i> 46(5), 1961-1993.
Logistic and multinomial models on the example of gender issues	Topic addresses small-area estimates for the composition of categorical variables. Gender issues in male-dominated occupations can be studied. From a methodological perspective, the thesis deals with estimation based on logistic and multinomial models. Data from NEPS, DZHW or data from Alumnae Tracking can be applied.	1)Christof Wolf, Henning Best (2010): <i>Handbuch der sozialwissenschaftlichen Datenanalyse</i> . Wiesbaden : VS, Verlag für Sozialwissenschaften. 2) Mood, C. (2010). <i>Logistic Regression: Why We Cannot Do What We Think We Can Do, And What We Can Do About It</i> . <i>European Sociological Review</i> , 26, 67-82. <a href="https://doi.org/10.1093/esr/jcp006">https://doi.org/10.1093/esr/jcp006</a> 3) Förtsch, S.; Gärtig-Daug, A.; Buchholz, S.; Schmid, U. (2018): Keep it going Girl! An Empirical Analysis of Gender Differences and Inequalities in Computer Sciences, <i>International Journal of Gender, Science and Technology</i> , 10 (2), pp. 265-286. 4) Förtsch, S., & Gärtig-Daug, A. (2019): Trust yourself: You have the IT-Factor! Career coaching for female computer scientists, <i>International Journal of Gender, Science and Technology</i> , 11 (3), pp. 490-527
Statistical misuse and consequences	Statistical literacy: Selection from typical problem clusters of statistical literacy, such as confounding, lack of measures of significance and/or goodness of fit, correlation versus causality, relative versus absolute risk, various kinds of biases (e.g., self-selection, response bias, attrition, Hawthorne effect, social desirability bias), immortal time bias, etc.	1)Milo Schield. <i>Statistical literacy: Thinking critically about statistics</i> . <i>Of Significance</i> , 11(1):15-20, 1999. 2)Djordje Kadijević & Max Stephens. Modern statistical literacy, data science, dashboards, and automated analytics and its applications. <i>Teaching of Mathematics</i> , 23(1):71-80, 2020. 3) Eric Soweay & Peter Petocz. A panorama of statistics: Perspectives, puzzles and paradoxes in statistics. <i>John Wiley &amp; Sons</i> , 2017.
Convergence in Distribution: Examining the Gelman-Rubin Diagnostic	Convergence in distribution is factually impossible to prove. The GD diagnostic has become a widely-accepted measure to assess convergence in distribution. Investigate the Accuracy of the GD diagnostic using an extensive MC simulation study.	Gelman, A., and D. B. Rubin. 1992. Inference from iterative simulation using multiple sequences. <i>Statistical Science</i> 7: 457-472
Diagnostics for mass imputation: R package	To date none of the imputation packages in R (or Stata, SPSS,...) can visualize missing data in meaningful way if the number of observations/ variables is large.	Meinfelder, F. (2013): Datenfusion: Theoretische Implikationen und praktische Umsetzung. In: <i>Weiterentwicklung der amtlichen Haushaltsstatistiken</i> . (eds. Riede, T.; Ott, N.; Bechthold, S.), pp. 83-100, 1st ed., Berlin, <i>GWl Wissenschaftspolitik Infrastrukturentwicklung</i> .
Inference for the Potential Outcome framework	Different matching strategies impact the variance term of the test statistic. With and without replacement matching as well as weighting affect how independent matched samples are.	Bai, Yuehao; Romano, Joseph P.; Shaikh, Azeem M (2019) Inference in experiments with matched pairs, <i>Econstor Working Paper</i> , <a href="https://www.econstor.eu/bitstream/10419/211112/1/1663702217.pdf">https://www.econstor.eu/bitstream/10419/211112/1/1663702217.pdf</a>
Multiple Imputation of categorical variables using an extension to Predictive Mean Matching	PMM has become a popular MI method and is the default imputation method in 'mice' for 'numeric'-type variables. However, no extension for non-ordered categorical (nominal-scale) variables is available. One problem is the identification of a nearest neighbour over k categories. This can be overcome by using the KL divergence to match via the empirical distribution of donor and recipient predictions. Your task is to implement this method in R and compare it with alternative approaches.	Koller-Meinfelder, F. (2009): <i>Analysis of Incomplete Survey Data – Multiple Imputation via Bayesian Bootstrap Predictive Mean Matching</i> , Dissertation, Otto-Friedrich-Universität Bamberg.
Comparing Generalized Linear Mixed Models and Transformations in the SAE Context	Two main approaches to deal with non normal and a non linear relationship between the independent and the dependent variable are generalized linear mixed models (GLM) and transformations. In a possible thesis the researcher could compare the pros and cons of both approaches in a simulation study.	1) Morales, Domingo, et al. "A course on small area estimation and mixed models." <i>Methods, theory and applications in R</i> (2021). 2) Chandra, Hukum, and Ray Chambers. "Small area estimation under transformation to linearity." (2008).