

Constructing socio-demographic indicators for National Statistical Institutes using mobile phone data: estimating literacy rates in Senegal

Timo Schmid^{*}, Fabian Bruckschen^{*}, Nicola Salvati^{**}, and Till Zbiranski^{*}

^{*}Institute of Statistics and Econometrics, Freie Universität Berlin, Berlin, Germany

^{**}Dipartimento di Economia e Management, Università of Pisa, Pisa, Italy

Abstract

Modern systems of official statistics require the accurate and timely estimation of socio-demographic indicators for disaggregated geographical regions. Traditional data collection methods such as censuses or household surveys impose great financial and organizational burdens for National Statistical Institutes. The rise of new information and communication technologies offers promising sources to mitigate these shortcomings. In this paper we propose a unified approach for National Statistical Institutes in developing countries based on small area estimation that allows for the estimation of socio-demographic indicators by using mobile phone data. In particular, the methodology is applied to mobile phone data from Senegal for deriving sub-national estimates of the share of illiterates disaggregated by gender. The estimates are used to identify hot spots of illiterates with a need for additional infrastructure or policy adjustments. Although the paper focuses on literacy as a particular socio-demographic indicator, the proposed approach is applicable to indicators from national statistics in general.

Keywords: Indicators, Model-based estimation, Official statistics, Small area estimation.

1 Introduction

If you can't measure it, you can't manage it. (Michael Bloomberg, former Mayor of New York City)

A country's budget can hardly be allocated efficiently, if the country does not know where the money is needed the most. Reliable knowledge on the socio-demographic indicators of a country's population is essential for sound evidence-based policymaking. For instance, the geographic distribution of wealth is used to make decisions regarding the allocation of resources. Traditionally, this knowledge is collected via household surveys and is provided by National Statistical Institutes (NSI). The surveys are generally designed to provide reliable estimates for the indicators only for larger domains such as the national or the regional level. One possible way to derive estimates on spatially disaggregated levels, like municipalities or communes, is by using small area methods (Rao and Molina, 2015). During the last decade there has been a substantial growth in the development and application of model-based small area methods for the estimation of indicators. Examples are manifold in literature: Elbers et al. (2003), Molina and Rao (2010) and Pratesi (2016) used small area techniques for the estimation of poverty indicators and,

recently, Lopez-Vizcaino et al. (2015) and Chambers et al. (2016) investigated the estimation of labour force indicators. For a comprehensive review we refer to Pfeffermann (2013) and Rao and Molina (2015). However, the production of precise small area estimates of indicators relies on the availability of predictive auxiliary variables like census or register information. In many countries successive census and national surveys are conducted with long lag times. Both require a well-functioning infrastructure, starting from cars for the interviewers to computers and well-trained personnel for the analysis. With national statistical systems in developing countries often being subject to unstable funding and a lack of human resources, the collection and processing of relevant data imposes a great challenge or often does not exist (Ghosh and Rao, 1994). For instance, in Angola the most recent census before 2014 was conducted in 1970 and the official population grew by more than 400% in that period (Blumenstock et al., 2015).

An alternative to the usage of census information for small area estimation is to investigate different sources of passively collected data like social media sources (e.g. Facebook, Twitter etc.) or mobile phone data. Eagle et al. (2010) used recently social network data to measure economic growth in the UK. Nevertheless, social media data are rare in developing countries whereas mobile phone data are a remarkable exception. The unique subscriber penetration is between 40% – 55% in developing countries with a share of around 40% in Sub-Saharan Africa (GSMA, 2015).

In this paper we investigate how mobile phone data (in combination with survey data) can be used to predict socio-demographic indicators at regionally disaggregated levels when census information is not available. The motivation is that mobile phone data are collected as a by-product and include valuable information on the timing and frequency of communication events and patterns of location and travel choices (Blumenstock et al., 2015). Eagle et al. (2010) and Deville et al. (2014) showed that spatially aggregated measures of mobile phone usage and penetration have a high correlation with spatially aggregated statistics from censuses. At this point we should make clear that the paper does not discuss whether the socio-demographic indicators can be directly estimated using only the mobile data. We are aware of some important recent work by Blumenstock et al. (2015). The authors predict poverty and wealth by using an individual's past history of mobile phone usage in combination with a phone survey. In our paper we had access to the Demographic and Health Survey (DHS) 2011 and mobile phone data covering the year 2013 in Senegal.

The Republic of Senegal is located in West Africa at the Atlantic Ocean between Mauritania to the North and Guinea-Bissau to the South. At the most Western tip lies Dakar, the country's capital and also the largest city. The set-up of administrative areas in Senegal is complex, but can be divided into four different levels: 14 regions, 45 departments, 123 arrondissements and 431 communes. The total population is estimated at about 13.5 million (2013) and consists of several ethnic groups, e.g. the Wolof or the Serer.

From a methodological point of view the present article uses area-level small area models (Fay and Herriot, 1979) in combination with covariates from alternative data sources. The resulting estimates are benchmarked such that the aggregated small area estimates produce the official national estimate for the country. We also apply transformation to restrict the indicator of interest, for instance the literacy rate, to particular intervals when necessary. However, the idea of alternative covariates is not new in literature. Porter et al. (2014) applied functional covariates extracted from Google in spatial Fay-Herriot models (Pratesi and Salvati, 2009). Recently, Marchetti et al. (2015) give a comprehensive overview how alternative data sources can be used in the context of small area estimation. Nevertheless, none of these papers considered in detail the usage of mobile phone data. To the best of our knowledge, this paper

is the first attempt to provide an easily applicable approach for NSIs to model a basket of regionally disaggregated socio-demographic indicators using survey data in combination with mobile phone data. In particular, the paper investigates the usability of mobile phone data, in this case tower-to-tower traffic in Senegal from 2013, for constructing fine granular indicators, like literacy and poverty rates, access to electricity and safe water or religious affiliations. The application here aims at estimating the socio-demographic indicator *literacy rate* for women and men for regionally disaggregated areas because it is a common problem across Africa. From an applied point of view, the paper also discusses the processing, cleaning and handling of the mobile phone data used as additional source of information.

Especially child labour, poverty and poor access to education are common problems across the Africa continent (Ford, 2007). Poverty in developing countries is not only a result of low income, but also of a lack of opportunities to improve the situation (UNESCO, 2015). Literacy is one of the keys to improve people's chances to escape from the lowest poverty levels. Although there are countries with a situation worse than the one of Senegal, the country is only ranked 117th out of 127 countries in the Education for All Development Index (EDI) published by the UNESCO (2012). Especially the literacy rate is quite low compared to other African countries (literacy rate in 2011: 37.8% for women and 60.0% for men, ANSD (2012)). The high number of illiterates can be partially explained by historic reason. Senegal was a former French colony until it gained independence from France in 1960. At that point the school attendance of children in the primary school was at 36%, while the country's average literacy rate was around 34% (Schelle, 2013). The origin for this low share of literacy lies in the little interest of the colonial rulers in educating the indigenous people. Other colonial powers in West Africa like Germany (Togoland) or England (Gold Coast, now called Ghana) had a pupils count which was around four times as high as Senegal's count (Schelle, 2013). Concerning the country's literacy rate from 2011, not much has changed in this regard since the withdrawal of the French power in 1960. Another problem of the educational situation is the slow development of a coherent education system due to opposing education concepts with different traditions. The indigenous African concept coexists next to the Islamic and Western concept. Nowadays, if children visit school, they often visit a public school and additionally a Qur'anic school in Senegal. In 2002 a new system emerged, the so called franco-arabic schools. A hybrid form of a bilingual (French and Arabic) school with a heavy curriculum. Although Villalon and Bodian (2012) predict this franco-arabic schools could be the future and predominant form of public schools, Senegal is after more than 50 years of independence still in the development stage of a coherent education system. The problem is doubtless not only due to a fragmented school system, but also caused by low attendance rates of children at any school. Although primary and secondary education is compulsory and free in Senegal, many parents still do not send their children to school, and drop-out rates are high (Ford, 2007). UNESCO (2012) reported that as the level of education increases especially the enrollment ratios of women in comparison with men strongly decrease. Although Senegal achieved a gender parity in primary education, the disparity for secondary education is even more severe. For every 100 boys attending secondary education in Senegal, only around 79 girls attend (UNESCO, 2012). This is one reason for low literacy rates especially among women. According to UNESCO (2012) more than two million women in Senegal miss skills in basic literacy. Especially in the country's poor regions like Matam and Tambacounda, both located in the East, girls are involved in economic activities and therefore the parents keep the girls out of the school to earn some additional income. Next to economic reasons, gender-based violence, early marriage and pregnancy as well as the traditional role of women in the society are further issues which add to low literacy rates for women (UNESCO, 2012).

The Senegalese government wants to significantly improve the literacy rate, especially for women.

For instance, in the early 2000s, the government built community schools and literacy centers for disadvantaged people, like women who missed a basic school education. However, according to the literacy rates for 2011 there is still a large gender disparity and a persisting need to address this issue in Senegal. Organizations like the UNESCO and UNICEF are constantly working on this educational issue and initiated several projects. Currently the Senegalese government and the UNESCO office in Dakar run a project to improve the literacy rate for women (UNESCO, 2015). In particular, the PAJEF project (Projet d’alphabétisation des jeunes filles et jeunes femmes) provides, for instance, access to organized literacy classes and develops training manuals. The project currently runs in seven regions identified by the National Agency of Statistics and Demography (ANSD - Agence Nationale de Statistique et de la Démographie) in Senegal. Further information is available in UNESCO (2015).

So far Senegal belongs to the most successful countries in advancement of gender equality for the enrollment in primary schools, but the national number of illiterate women remains high. All the efforts mentioned above are experimental and not countrywide because of a lack of spatially disaggregated knowledge where more support is needed. To obtain a higher countrywide literacy rate, areas of high illiteracy have to be identified. In this paper we propose an approach for NSIs based on small area estimation for deriving estimates of the share of literates by gender by using mobile phone data for the 431 communes in Senegal. The estimates are used to identify hot spots of illiterate women for the PAJEF project with a need for additional infrastructure.

The structure of the paper is as follows. In Section 2 we describe the DHS survey and the mobile phone data including the cleaning and preparation. In Section 3 we review small area estimation using Fay-Herriot models. The methodological approach for constructing socio-demographic indicators based on mobile phone data is described and computational details are provided. In Section 4 we present the results of the application for the indicator *literacy rate* in Senegal by using the mobile phone data. The performance of the proposed approach is empirically evaluated in a large-scaled design-based simulation in Section 5. Finally, in Section 6 we conclude the paper with some final remarks and discuss limitations of the proposed approach. Additional results are presented in the supplementary materials.

2 Data sources: survey data and mobile phone data

In this section we describe the data sources used in the analysis. In particular, we had access to the Demographic and Health Survey (DHS) 2011 and mobile phone data covering the year 2013 in Senegal. We present details regarding practical implementation of the time-intensive cleaning and preparation of the mobile phone data and discuss the construction of mobile phone covariates.

2.1 Demographic and Health Survey

The DHS program collects representative data on population, health, HIV and nutrition in over 90 countries. The data that we use are from the DHS survey 2011 carried out by the ANSD in Senegal. The survey includes a section on the production of socio-demographic indicators on household level and another part on assessing the availability of material and human resources. In particular, the DHS survey consists of three questionnaires: (i) a household questionnaire, (ii) a women’s questionnaire and (iii) a men’s questionnaire. The household survey collects information on the usual household members including, for instance, gender, age, education, survival of parents, and child labor. Additional information like household characteristics (source of water, availability of electricity, building material and type of toilet), ownership, use of mosquito nets and several health related questions are collected as well. The

household survey is also used to identify men and women for the individual questionnaires. The questionnaire for women consists of 10 sections covering socio-demographic indicators (like age and date of birth, schooling, literacy, and ethnicity), reproduction, use of contraception, pregnancy, marriage and female genital mutilation. The men's questionnaire is a short version of the questionnaire for women covering socio-demographic characteristics and health related questions. Note as socio-demographic characteristics are only available in the gender-specific questionnaires we focus in the analysis in this paper on the women's and men's questionnaires. For additional information regarding the variables and the questionnaires we refer to ANSD (2012).

The survey aims to cover the complete country and is based on a stratified two-stage cluster sampling design. The 28 strata are defined by a cross-classification of the 14 regions and rural/urban areas in Senegal. The survey is designed to produce reliable results for most indicators for the 14 regions. In the first sampling stage 391 census districts (147 urban and 244 rural) were drawn with probability proportional to size (number of households in the census districts). In the second sampling stage 21 households were selected with equal probability in each of the 391 census districts which were sampled in the DHS survey. Among the 21 households selected for the women's survey, 8 households were drawn for the men's survey. All men (age between 15-59) and women (age between 15-49) in these households were interviewed. The interview was successfully conducted for 15,688 women (response rate of 92.7 percent) and for 4,929 men (response rate of 87 percent) (ANSD, 2012).

Figure 1 presents results based on DHS survey 2011 of the indicator *literacy rate* by gender for the regions in Senegal. In particular, the variable *literacy* is collected by four different categories in the DHS survey. The categories 'able to read only parts of sentence' and 'able to read whole sentence' are grouped as 'literate'. The answers 'blind/ visually impaired', 'cannot read at all' and 'no card with required language' are categorized 'illiterate'. The initial results indicate that the proportion of *literate women* (37.8%) in Senegal is lower than the proportion of literate men (60.0%). The results are consistent with the official published results of the ANSD (2012).

As the ANSD aims to estimate socio-demographic indicators for the 431 communes in Senegal, we allocated the information of the DHS survey to the administrative areas (communes). In particular, we had access to the geographical coordinates of the centroids of the 391 census districts. As the actual coverage of the census districts was not available, we matched the centroids of the census districts with the geographical boundaries of the 431 communes. Six out of the 391 census districts were excluded from the analysis because the coordinates of the centroids were missing. Direct survey estimates are only available for 242 out of the 431 communes given the data from the DHS survey 2011. A summary of the commune specific sample sizes for the women's and men's questionnaires is provided in Table 1. Figure 2 shows direct estimates for the literacy rate by gender on commune level for the capital Dakar (right panel) and for the rest of Senegal (left panel). Communes filled with white color represent areas with zero sample size, so direct estimates based on the DHS survey 2011 are not available. The spatial distribution of literacy on commune level is not clearly visible and the identification of hot spots of illiterates with a need for additional infrastructure might be difficult.

The application of small area methods could significantly improve the interpretation of Figure 2 by providing results for the communes with zero sample size. This requires fitting of an appropriate model to the survey data. The estimated model parameters are then combined with known population information. The reason that we relied on mobile phone data for predicting socio-demographic indicators is twofold: first, the predictive power of the covariates for socio-demographic indicators from the Senegalese census is limited and second, the ANSD is interested in a widely applicable approach based on the DHS survey

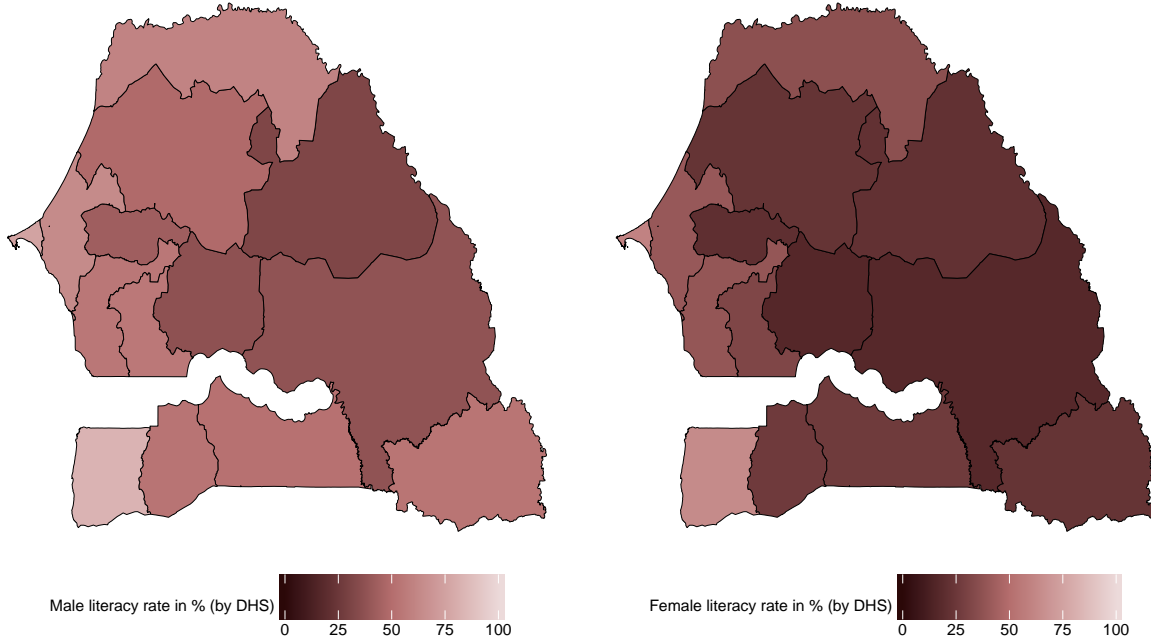


Figure 1: Estimates for the literacy rate by gender on regional level based on DHS survey 2011.

Table 1: Sample sizes over communes.

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA
Women's questionnaire	15	35	44	63.90	61	756	189
Men's questionnaire	2	10	14	19.98	20	160	189

for disaggregated indicators independent of census data.

2.2 Mobile Phone Data

The mobile phone data used in this paper consist of anonymized call detail records (CDR) from the Senegalese telecommunication company Sonatel covering the year 2013. The dataset is based on more than 9 million unique mobile phone numbers and represents a market share of around 60%. In particular, we had access to the tower-to-tower traffic of all 1666 mobile phone towers in Senegal. In the following we discuss the practical implementation of the processing of the mobile phone data and present details regarding the construction of the mobile phone covariates.

2.2.1 Data processing and cleaning

The preprocessing of the mobile phone raw data is essential and accounts for a considerable amount of time in the whole analysis. The dataset is not *dirty* or *noisy* in the sense of an excessive amount of missing values or illogical recorded values. The data is collected automatically by machines and not gathered by human hand. This means errors in the data are more likely a consequence of machine breakdowns than of human failure.

The traffic of all 1666 towers in Senegal for 2013 is about 1.1 Terabyte of data stored in a cloud system. Because of the massive amount of data, the mobile phone records need to be preprocessed directly in the cloud system. In particular, the raw data is organized in a Hadoop cluster with one separate file by hour per day per month. Hadoop is an open-source software for storing and handling massive



Figure 2: Estimates for the literacy rate by gender on commune level based on DHS survey 2011: Senegal (left panel) and Dakar (right panel).

data. Each single row contains an interaction and has several characteristics. For example indicating if it is an incoming or outgoing interaction, if it is a phone call or short message service (SMS), which tower received and sent the interaction, or simply the duration of a call in minutes. To process these data we used Apache Hive (Apache Hive is a data warehouse infrastructure built on top of Hadoop for providing data summarization) and its SQL logic. MapReduce is applied to create daily, monthly and yearly aggregates of the variables of interest on the cluster. The programming model MapReduce is an implementation for processing large datasets with parallel algorithms on a cluster.

For instance, the aggregated dataset for SMS usage includes the number of incoming and outgoing SMS for every tower on an hourly basis for the year 2013. Table 2 shows the head of a preprocessed dataset for the usage of SMS. The first column is the observation indicator which reaches in January

Table 2: Structure of the call detail records for SMS.

	DH	TO	TI	E
1	2013-01-01 00	1	61	1
2	2013-01-01 00	1	340	1
3	2013-01-01 00	1	419	1
4	2013-01-01 00	1	420	1
5	2013-01-01 00	1	447	2
6	2013-01-01 00	1	495	1

2013 alone around 50 million rows. Variable *DH* tracks the day and hour of a sent SMS; *TO* and *TI* are the tower numbers corresponding to outgoing and incoming, respectively; *E* gives the number of events happening, i.e. SMS being sent. So the first row says that on the 1st of January at midnight there was sent 1 SMS from tower 1 to tower 61. We also had access to the exact geo-coordinate (longitude and latitude) of the towers provided by Sonatel.

2.2.2 Construction of mobile phone covariates

Mobile phone data are measured on tower level on an hourly basis with an excessive amount of observations over the year. To construct variables which can be used as covariates for a statistical model for estimating indicators on commune level, the data needs to be aggregated by two dimensions: time and geographic level. First, in order to reduce the amount of data, the aggregation was done up to the whole year 2013 for each tower. Annual aggregates may disregard sub-annual trends, but since most of the socio-demographic indicators, especially the literacy rate, are time insensitive variables, this fact can be neglected. Second, for having the covariates on the same geographical level like the DHS survey, we used the aggregated (by time) covariates on tower level and averaged them for higher geographic levels like communes or regions. Note as the actual coverage of the mobile towers are unknown, we matched the geo-coordinate of the tower with the geographical boundaries of the 431 communes.

In total we constructed around 70 mobile phone covariates on commune level based on the call detail records. The aggregation routine is done in R by using the package `data.table`. The package extends `data.frames` in R based on SQL logic and focuses on fast aggregation of large data (Dowle et al., 2014). For instance, we construct the sum of the number of calls starting from/ending in a specific tower and denote these variables as *outgoing calls* / *incoming calls*, respectively. In addition, we also build the variable *call volume* which sums up the minutes of calls. In the following we label SMS and phone calls together as *events*. For each event we also calculated the *ratios* of the number of outgoing events divided by incoming events. The variable *mean distance* is defined as the average distance in kilometers

for an event. In particular, the distance is computed on the tower level by taking the distance of the outgoing tower to the incoming tower for each event and dividing it by the amount of events between the two towers. The covariate *distance-to-dakar* measures the distance from each tower to a centroid of the region Dakar. In addition we construct the variable *isolation* which quantifies the diversity of interactions by users of a tower. The variable is defined for an outgoing tower t_i by

$$Isolation(t_i) = \sum_{\substack{j=1 \\ j \neq i}}^{1666} \mathbf{I}_{E(t_i, t_j)}, \quad (1)$$

where the indicator function \mathbf{I} is 1 if the condition $E(t_i, t_j)$ is true, i.e. an event happened between the towers t_i and t_j , and 0 otherwise. The variable ranges between 0 and 1666 (total number of towers). We measure the average amount of information an event contains by the variable *Entropy* (Montjoye et al., 2014). The intuition behind Entropy is that the more unlikely an event is to happen, the more information it contains once it happens. Entropy for a tower t_i is defined by

$$Entropy(t_i) = - \sum_{\substack{j=1 \\ j \neq i}}^{1666} p(t_i, t_j) \cdot \log[p(t_i, t_j)], \quad (2)$$

where $p(t_i, t_j)$ is the probability of an event between the towers t_i and t_j . In addition, we calculated the *monthly growth* and the *variation* (i.e. variance) of monthly aggregates for the number and volume of events respectively. Variables *Calls-to-dakar* and *sms-to-dakar* reflect the amount of calls or SMS for each tower that were directed to towers located in the capital Dakar. A complete list and description of the covariates is provided in the supplementary materials.

Additionally to the variables described above and in the supplementary materials, we created behavioral indicators based on the mobile phone data with the open-source python toolkit bandicoot (Montjoye et al., 2013). A list of these variables can be found at <http://bandicoot.mit.edu/docs/reference/index.html>. As the bandicoot indicators are constructed for analyzing individual patterns based on the mobile behavior of each single user, we summarized the information to tower level. In particular, a bandicoot indicator on tower level is calculated as a weighted average of all individuals' indicators where this tower was part of the interaction. The steps are as follows: first, we calculated the bandicoot indicators on a monthly level for all single users. Second, we extracted the number of interactions (calls and SMS) during that month for each user and tower combination from the call detail records. Third, we used the number of interactions as a weight to average the individuals' indicators on tower level for each month. Finally, we averaged the monthly values to obtain a yearly indicator for each tower.

2.2.3 First descriptive statistics

Figure 3 gives a first impression of the spatial distribution of the 1666 mobile phone towers (red points) in Senegal. The towers are spread over the whole country with higher densities in regions with higher population densities. For instance, most of the towers are located in the region of the capital Dakar which itself is located on the Cap-Vert Peninsula on the Atlantic coast in the West. Table 3 shows summary statistics of the number of mobile phone towers over the communes. The mean number of towers per commune is 4.1 with a maximum of 60. Although Figure 3 suggests a good coverage of the country by mobile phone towers, there are 30 communes without mobile phone towers. Most of these communes are quite small and they are mainly covered by towers which are close-by. For instance, the map at

Table 3: Mobile phone towers over communes

	Min.	1st Qu.	Median	Mean	3rd Qu.	90%	Max.	NA
Number of towers	1	1	2	4.11	4	9	60	30

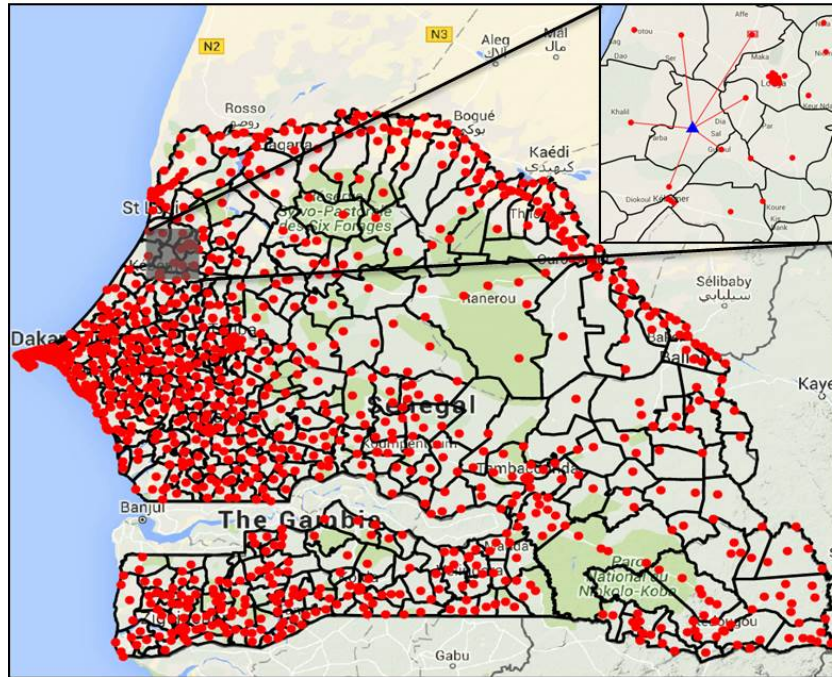


Figure 3: Location of mobile phone towers in Senegal.

the top on the right of Figure 3 shows the area around the commune Badegne Ouolof without tower information. Badegne Ouolof is located in north-western Senegal within the Louga Region on a total of around 300 square kilometers. The centroid of Badegne Ouolof is represented by a blue triangle. In order to apply small area estimation methods for the *out-of-covariate* communes, the covariates are constructed by inverse distance weighting from neighboring mobile towers. In particular, the assigned covariates to *out-of-covariate* communes are calculated by a weighted average of the covariates available at known tower locations. We used the Euclidian distance function and a power parameter of 2 in the weighting.

3 Description of the small area estimation method

In this section we describe the methodological approach for constructing socio-demographic indicators based on mobile phone data. Since our aim is to provide an easy-applicable approach for the production of official statistics, especially for the ANSD in Senegal, we apply relatively simple small area estimation methods and correct for misspecifications by adjustments. The implemented approach should meet three conditions:

1. the method should provide estimates for all 431 communes in Senegal;
2. the estimates should be *close* to the direct estimators for communes with *large* sample sizes;
3. the aggregated estimates for the communes should produce the official national estimate for the country.

Note that the Ministry of Chile recently conducted a small area project for the estimation of poverty in Chile based on similar guidelines (Casas-Cordero et al., 2016). In addition the mobile phone covariates are only available on area-level (communes) and it is not possible to link the individuals in the survey with the mobile phone numbers because of confidentiality constraints. Based on the mentioned conditions and available data we considered a benchmarked transformed Fay-Herriot estimator in this paper.

3.1 Fay-Herriot estimator

We assume that the population U , consisting of N units, is divided into m disjoint small areas. The sample s is selected from the population by using a complex sampling design. The population is separated into n sampled and $N - n$ non-sampled units, indexed by s and r , respectively. We use the subscript i to indicate the restriction to the area i , for instance, n_i and N_i denote the sample size and the population size in area i , respectively. Let y denote a continuous variable of interest and y_{ij} the response value of unit j in area i and ω_{ij} are the corresponding sampling weights. An estimator for the population mean θ_i of the variable of interest y in area i is given by

$$\hat{\theta}_i^{direct} = \frac{\sum_{j=1}^{n_i} \omega_{ij} y_{ij}}{\sum_{j=1}^{n_i} \omega_{ij}}. \quad (3)$$

The area level model proposed by Fay and Herriot (1979) (hereafter FH model) links the direct estimates with area-level covariates. The FH model is based on two stages:

$$\text{Sampling model (first stage): } \hat{\theta}_i^{direct} = \theta_i + \varepsilon_i \quad (4)$$

$$\text{Linking model (second stage): } \theta_i = \mathbf{x}_i^T \boldsymbol{\beta} + u_i, \quad (5)$$

where \mathbf{x}_i^T and $\boldsymbol{\beta}$ denote the $(k \times 1)$ vectors of area-level covariates and regression parameters, respectively. The sampling errors are assumed to be normally distributed and independent with $\varepsilon_i \sim N(0, \sigma_{\varepsilon_i}^2)$. The random effects u_i are assumed to be independently normally distributed with $u_i \sim N(0, \sigma_u^2)$. For additional details we refer to Rao and Molina (2015). The combination of both models leads to an area-level linear mixed model given by

$$\hat{\theta}_i^{direct} = \mathbf{x}_i^T \boldsymbol{\beta} + u_i + \varepsilon_i. \quad (6)$$

Let $\hat{\boldsymbol{\beta}}$ define the empirical best linear unbiased estimator (EBLUE) of $\boldsymbol{\beta}$ and \hat{u}_i the empirical best linear unbiased predictor (EBLUP) of u_i (Henderson, 1950; Searle, 1971), where the variance component σ_u^2 can be estimated by maximum likelihood or residual maximum likelihood (Datta and Lahiri, 2000; Rao and Molina, 2015). The EBLUP of θ_i under the FH model is obtained by

$$\hat{\theta}_i^{FH} = \mathbf{x}_i^T \hat{\boldsymbol{\beta}} + \hat{u}_i \quad (7)$$

$$= \hat{\gamma}_i \hat{\theta}_i^{direct} + (1 - \hat{\gamma}_i) \mathbf{x}_i^T \hat{\boldsymbol{\beta}}, \quad (8)$$

where $\hat{\gamma}_i = \hat{\sigma}_u^2 (\hat{\sigma}_u^2 + \sigma_{\varepsilon_i}^2)^{-1}$ denotes the shrinkage factor for area i and $\hat{u}_i = \hat{\gamma}_i (\hat{\theta}_i^{direct} - \mathbf{x}_i^T \hat{\boldsymbol{\beta}})$. In practice, many of the small areas may have zero sample sizes, so a direct estimator is not available. In this case we rely on synthetic estimation as follows (Rao and Molina, 2015):

$$\hat{\theta}_{i,out}^{FH} = \mathbf{x}_i^T \hat{\boldsymbol{\beta}}. \quad (9)$$

The MSE of the EBLUP in (7) can be obtained by analytic solutions following Prasad and Rao (1990) and Datta et al. (2005). Note that Li and Lahiri (2010) pointed out that standard estimation methods of the variance component in the Fay-Herriot model can produce zero estimates of the strictly positive model variance. Standard methods for the estimation of the variance component considered in the literature are, for instance, the Prasad-Rao method-of-moments estimator (Prasad and Rao, 1990), the Fay-Herriot method-of-moments estimator (Fay and Herriot, 1979), the maximum likelihood estimator or the residual maximum likelihood estimator. As a consequence, the EBLUP estimator (7) can reduce to a regression estimator, which can have an overshrinking problem. Li and Lahiri (2010) propose an adjusted maximum likelihood estimator of the variance component. In particular, an adjusted likelihood of σ_u^2 is defined by

$$L_{adj}(\sigma_u^2) = \sigma_u^2 \times L(\sigma_u^2), \quad (10)$$

where $L(\sigma_u^2)$ can be either the profile likelihood function or the residual likelihood function. Under certain regularity conditions, the adjusted maximum likelihood estimator of σ_u^2 is consistent for a large number of areas m and the shrinkage factors, γ_i , are all strictly greater than 0, even for small m , and are also consistent for large m (Li and Lahiri, 2010). From a Monte-Carlo simulation study carried out by Li and Lahiri (2010) results that in terms of bias and mean squared error, the adjusted maximum profile likelihood method turns out to be better than the adjusted maximum residual likelihood approach. For this reason, we use the adjusted profile likelihood function for estimating the value of σ_u^2 in the paper. Note that Yoshimori and Lahiri (2014) recently proposed an improvement to the adjusted likelihood estimators of Li and Lahiri (2010), showing better performance in a simulation study when σ_u^2 is small relative to the sampling variance $\sigma_{\varepsilon_i}^2$.

3.2 Transformed Fay-Herriot estimator

Some socio-demographic indicators are restricted to a specific range. For instance, the share of literates in an area i should be within the interval $[0, 1]$. However, there is no guarantee that the FH estimates produces estimates in a particular range. In the context of estimating small area proportions Jiang and Lahiri (2001), Liu et al. (2014), Bell and Franco (2015), and others present different modeling options for the linking and sampling distribution for area-level models. In particular, Ha et al. (2014) propose a normal-logistic model (NL) with a logistic distribution for the linking model. In addition, they extend the model by Liu et al. (2014) to general complex survey designs and denote it as a normal-logistic random sampling variance model (NLRS). The NLRS model captures parts of the uncertainty due to the estimation of the small area sampling variance. For additional details we refer to Ha et al. (2014). Following Carter and Rolph (1974), Jiang et al. (2001) and Raghunathan et al. (2007) we use in this paper an arcsine transformation for the modeling. Let y now denote a binary variable of interest and y_{ij} is the 0-1 response value of unit j in area i . The steps of the estimation are as follows:

1. Transform the direct estimator via $\vartheta_i = f(\hat{\theta}_i^{direct}) = \arcsin \sqrt{\hat{\theta}_i^{direct}}$.
2. The sampling variance of ϑ_i is approximated by $\sigma_{\varepsilon_i}^2 = 1/(4\tilde{n}_i)$, where \tilde{n}_i stands for the effective sample size (Jiang et al., 2001). In particular, the effective sample size is the sample size divided by an estimate of the design effect.
3. Estimate $\hat{\theta}_i^{FH} \{ \vartheta_i, 1/(4\tilde{n}_i) \}$ according to (7). $\hat{\theta}_i^{FH}$ is truncated to the interval $[0, \pi/2]$ if necessary .
4. Back-transform the estimator $\hat{\theta}_i^{FH}$ to the original scale via

$$\hat{\theta}_i^{FH,trans} = f^{-1}(\hat{\theta}_i^{FH}) = \sin^2(\hat{\theta}_i^{FH}) \quad \text{for } i = 1, \dots, m, \quad (11)$$

where $\hat{\theta}_i^{FH,trans}$ denotes the transformed FH estimator.

For constructing the confidence intervals for θ_i we use a parametric bootstrap procedure following Casas-Cordero et al. (2016). See also Chatterjee et al. (2008) and Li and Lahiri (2010). The steps are as follows:

1. For given $\hat{\beta}$, $\hat{\sigma}_u^2$ and $\hat{\gamma}_i$ estimated with the transformed direct estimator ϑ_i , sampling variance $1/(4\tilde{n}_i)$ and covariates \mathbf{x}_i , we generate u_i^* from $N(0, \hat{\sigma}_u^2)$ and ε_i^* from $N(0, 1/(4\tilde{n}_i))$.
2. Using u_i^* and ε_i^* to generate the bootstrap sample,

$$\hat{\theta}_i^{*,(b)} = \mathbf{x}_i^T \hat{\beta} + u_i^* + \varepsilon_i^* \quad (12)$$

and the corresponding bootstrap population parameter

$$\theta_i^{*,(b)} = \mathbf{x}_i^T \hat{\beta} + u_i^*. \quad (13)$$

3. Using the bootstrap sample, we estimate the model parameters in (6). Based on the estimated model parameters from the bootstrap sample, we compute the corresponding FH estimator (7) in area i , $\hat{\theta}_i^{FH,(b)}$.
4. Calculate the following pivotal quantity:

$$t_i^{(b)} = \frac{\theta_i^{*,(b)} - \hat{\theta}_i^{FH,(b)}}{\sqrt{\hat{\gamma}_i^{(b)}/(4\tilde{n}_i)}} \quad (14)$$

5. Repeat steps 1-4 B times.
6. For each area i , calculate the $100\alpha/2$ quantile q_{1i} and $100(1 - \alpha/2)$ quantile q_{2i} of $\{t_i^{(b)}, b = 1, \dots, B\}$.
7. An approximate $100(1 - \alpha)$ confidence interval for θ_i is defined as: (lo_i, up_i) , where $lo_i = \hat{\theta}_i^{FH} + q_{1i}\sqrt{\hat{\gamma}_i/(4\tilde{n}_i)}$, and $up_i = \hat{\theta}_i^{FH} + q_{2i}\sqrt{\hat{\gamma}_i/(4\tilde{n}_i)}$. If lo_i is negative, it is truncated to 0 and if up_i is greater than $\pi/2$, it is truncated to $\pi/2$. This truncated confidence interval is defined (lo_i^*, up_i^*) . Back-transform the lower and the upper limits (lo_i^*, up_i^*) for each area to obtain the approximate $100(1 - \alpha)$ confidence interval: $(\sin^2(lo_i^*), \sin^2(up_i^*))$.

Note that the back-transformed confidence interval can be obtained because the function \sin^{-1} and \sin^2 are monotonically increasing functions of the parameters in the ranges of interest (Casas-Cordero et al., 2016). In addition, as the upper and lower bound of the confidence interval depends on the effective sample size \tilde{n}_i in area i , we can only estimate the confidence interval for the in-sample areas. In order to obtain a second-order correct confidence interval for out-of-sample areas, one has to replace equation (14) by the following pivotal quantity:

$$t_{i,out}^{(b)} = \frac{\theta_{i,out}^{*,(b)} - \mathbf{x}_{i,out}^T \hat{\beta}^{(b)}}{\hat{\sigma}_u^{(b)}}, \quad (15)$$

where $\hat{\beta}^{(b)}$ and $\hat{\sigma}_u^{(b)}$ are estimates of β and σ_u based on in-sample b -th parametric bootstrap replicate (Chatterjee et al., 2008). Following Chatterjee et al. (2008), the boundaries of the confidence interval for the out-of-sample areas are given by $lo_{i,out} = \mathbf{x}_{i,out}^T \hat{\beta} + q_{1i,out} \hat{\sigma}_u$, and $up_{i,out} = \mathbf{x}_{i,out}^T \hat{\beta} + q_{2i,out} \hat{\sigma}_u$.

where $q_{1i,out}$ and $q_{2i,out}$ are the $100\alpha/2$ and $100(1 - \alpha/2)$ percent quantiles of $\{t_{i,out}^{(b)}, b = 1, \dots, B\}$, respectively.

An alternative approach is to apply a jackknife method on the transformed scale proposed by Jiang et al. (2001). In particular, Jiang et al. (2001) consider an arcsine transformation and show that the bias of the jackknife MSE estimator is of order $o(m^{-1})$. Then, the authors approximate the MSE in the original scale for $\hat{\theta}_i^{FH,trans}$ (11) by

$$mse[\hat{\theta}_i^{FH,trans}] = f^{-1'}(\hat{\theta}_i^{FH})mse(\hat{\theta}_i^{FH}),$$

where $f^{-1'}$ denotes the derivative of f^{-1} (defined in equation 11) with respect to $\hat{\theta}_i^{FH}$. $mse(\hat{\theta}_i^{FH})$ is an estimate of the MSE obtained by the jackknife method proposed by Jiang et al. (2001). Future work can be devoted to compare the width of the confidence interval and the coverage rate obtained with the parametric bootstrap (Casas-Cordero et al., 2016) and the jackknife procedure (Jiang et al., 2001).

3.3 Benchmarked Fay-Herriot estimators

Although the model-based estimator in (11) provides estimates for all communes (small areas) in Senegal, the aggregated estimates on national level can differ from the corresponding direct estimator. Following Datta et al. (2010) we use a benchmark approach to achieve the internal consistency with the direct estimator on national level.

For the FH model proposed in Section 3.1, we seek for a benchmarked FH estimator $\hat{\theta}_i^{FH,bench}$ such that

$$\sum_{i=1}^m \xi_i \hat{\theta}_i^{FH,bench} = \tau,$$

where

$$\tau = \sum_{i=1}^m \xi_i \hat{\theta}_i^{direct}.$$

We define the weights by $\xi_i = N_i/N$. We define the benchmarked FH estimator (Datta et al., 2010) by

$$\hat{\theta}_i^{FH,bench} = \hat{\theta}_i^{FH} + \left(\sum_{i=1}^m \frac{\xi_i^2}{\phi_i} \right)^{-1} \left(\tau - \sum_{i=1}^m \xi_i \hat{\theta}_i^{FH} \right) \frac{\xi_i}{\phi_i} \quad \text{for } i = 1, \dots, m. \quad (16)$$

There are several ways to define the weight ϕ_i (Datta et al., 2010). For instance, $\phi_i = \xi_i/\hat{\theta}_i^{FH}$ leads to a ratio adjustment of the FH estimator, where small areas with larger estimates will receive a larger adjustment and vice versa. Another option is to define the weights by $\phi_i = \xi_i/\widehat{MSE}(\hat{\theta}_i^{FH})$. That means that small areas with higher variability in terms of MSE will receive a larger adjustment.

For the benchmarked transformed FH estimator proposed in Section 3.2, we use a *naive* approach where the weights are given by $\phi_i = \xi_i$. Thus, the benchmarked transformed FH estimator results as a constant shift from the transformed FH estimator (11) and is given by

$$\hat{\theta}_i^{FH,trans,bench} = \hat{\theta}_i^{FH,trans} + \left(\tau - \sum_{i=1}^m \xi_i \hat{\theta}_i^{FH,trans} \right) \quad \text{for } i = 1, \dots, m. \quad (17)$$

4 Application: estimating literacy rates in Senegal

In this section the benefits of using the presented Fay-Herriot-type estimators in combination with mobile phone covariates for the estimation of socio-demographic indicators are illustrated in an application which uses the data from the DHS survey 2011 and the mobile phone data we described in Section 2. The application aims at estimating the literacy rate by gender on commune level in Senegal. The analysis is carried out by using the variables *literacy women* and *literacy men* from the gender-specific questionnaires introduced in Section 2. The estimates are used to identify hot spots of illiterate women for the PAJEF project with a need for additional infrastructure and financial support from the government.

4.1 Model selection and model checking

Before proceeding with the analysis of literacy in Senegal, we discuss the model selection and present some diagnostic plots. As discussed in Section 2.2 there are various commune specific covariates available from the mobile phone data. To select reasonable covariates in the context of Fay-Herriot models we followed an approach taken by several authors (Jiang et al., 2001; Ha et al., 2014; Casas-Cordero et al., 2016). In particular, we used the Bayesian information criterion (BIC) based on a linear regression model with $\arcsin \sqrt{\hat{\theta}_i^{direct}}$ as dependent variable and considered only data from the 50% largest communes in terms of sample size for the model selection. The reasons for choosing this particular model selection technique are threefold: First, we implicitly assume that for these communes the sampling variability of the direct estimators is reduced, so standard selection techniques are applicable. Second, we apply the BIC to penalize the model complexity more heavily compared to the Akaike information criterion (AIC) in order to enhance the interpretability of the final model. Third, although we are aware of more complex methods for Fay-Herriot model selection discussed in Marhuenda et al. (2014) we use a simple approach which is efficiently implemented by automatic stepwise selection procedures in standard statistical software. The final model on commune level for the variables *literacy women* and *literacy men* include 7 and 8 mobile phone covariates with an adjusted R^2 of 82% and 76% respectively. Note that we report an adjusted R^2 here proposed by Lahiri and Suntornchost (2015) that accounts for the sampling error variability.

Based on the transformed direct estimates from the DHS survey 2011 and the set of selected mobile phone covariates on commune level we fitted area level mixed models (6) by gender. As discussed in Section 3 the sampling variances of the direct estimates are approximated by $1/4\tilde{n}_i$ where \tilde{n}_i denotes the sample size divided by the design effect. Following Casas-Cordero et al. (2016), we used the design effect on regional level as an approximation for the design effect on commune level. The reason here is that the variance estimation of the direct estimator is unstable because of a low number of cluster or even not directly possible because only one cluster is nested in some communes. We refer to Opsomer et al. (2012) for a recent discussion on this issue in the context of forestry data.

Table 4 reports the design effects of the direct estimators by gender on regional level in Senegal. The estimates are consistent with official results published by the ANSD (2012) in Senegal and show a high value of the design effect of the direct estimator using DHS survey 2011.

Figure 4 shows normal probability plots of the standardized residuals (level 1) and the standardized random effects (level 2) obtained from fitting the female model (left panel) and the male model (right panel). The figure indicates some small departures from normality especially in the tails of the distribution. However, the departures are not severe. The Shapiro-Wilk test supports the lack of evidence against the normality assumption for the level 1 standardized residuals (p-values: male model = 0.4453

Table 4: Design effects of the direct estimator in Senegal by region.

Region	Female	Male	Region	Female	Male
Dakar	6.260	2.825	Louga	4.473	3.410
Diourbel	3.186	1.987	Saint Louis	4.584	1.874
Fatick	7.695	2.499	Matam	6.569	3.908
Kaffrine	5.058	2.682	Sedhiou	7.840	3.216
Kaolack	5.153	2.434	Tambacounda	4.281	3.386
Kedougou	2.566	1.962	Thies	5.480	3.227
Kolda	3.434	2.615	Ziguinchor	2.525	2.165

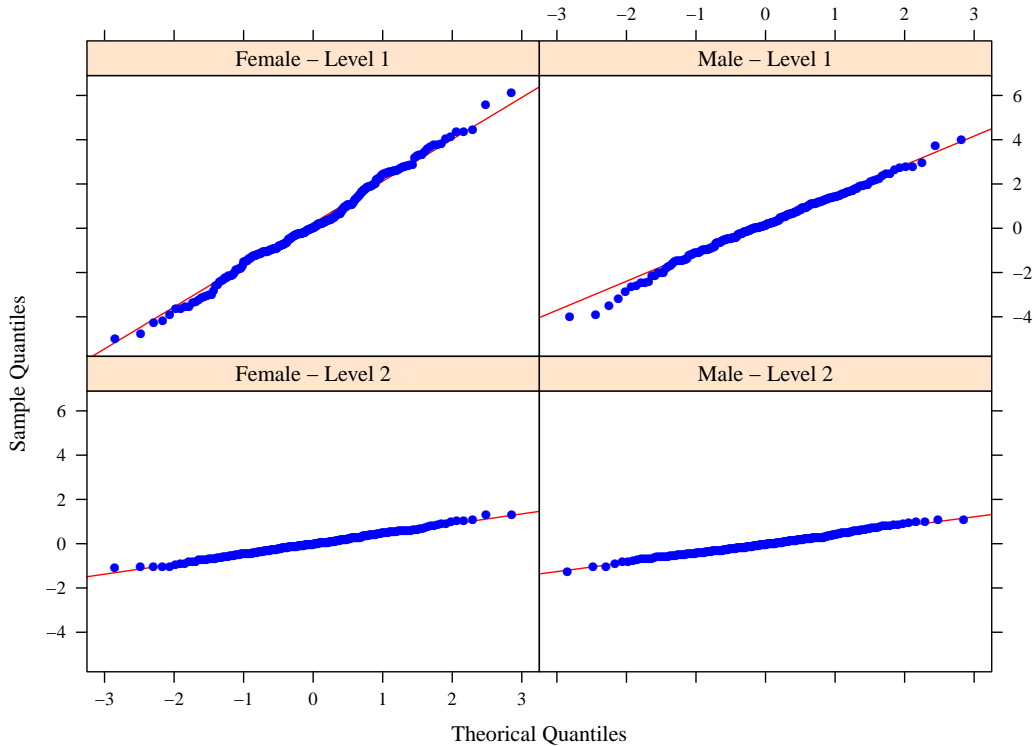


Figure 4: Normal probability plots of standardized residuals (level 1) and the standardized random effects (level 2) for the female model (left panel) and for the male model (right panel).

and female model = 0.4656) and level 2 standardized random effects (p-values: male model = 0.3311 and female model = 0.6059). Using the transformed Fay-Herriot model (11) may be advisable for estimating the literacy of women and men.

4.2 Small area estimates at commune level

Estimates of the literacy rate by gender for each commune are calculated by using the transformed FH estimator (11) (FH Trans) and by the benchmarked transformed FH estimator (17) (FH Bench). For constructing the confidence intervals based on the FH Trans we use the parametric bootstrap approach of Casas-Cordero et al. (2016) discussed in Section 3. We performed $B = 500$ bootstrap replications. We also include the direct estimator to assess the resulting estimates as the model-based estimators should be consistent with the unbiased direct estimators but with a higher precision. Note that direct estimation is not an option for the DHS survey 2011 on commune level because around 45% of the communes are out-of-sample. The estimators are implemented by computationally efficient algorithms using R. The

codes are available from the authors upon request.

Table 5: Distribution of the female and male literacy rates over communes in Senegal.

233 In-sample communes							
Gender	Estimator	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Female	Direct	0.000	0.105	0.234	0.298	0.474	0.839
	FH Trans	0.002	0.151	0.252	0.296	0.434	0.822
	FH Bench	0.002	0.158	0.264	0.310	0.455	0.861
Male	Direct	0.000	0.250	0.533	0.508	0.720	1.000
	FH Trans	0.062	0.368	0.500	0.516	0.662	0.971
	FH Bench	0.064	0.383	0.520	0.537	0.689	1.000
168 Out-of-sample communes							
		Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Female	FH Trans	0.005	0.157	0.219	0.271	0.363	0.732
	FH Bench	0.006	0.165	0.230	0.283	0.381	0.766
Male	FH Trans	0.066	0.309	0.468	0.478	0.633	0.960
	FH Bench	0.069	0.322	0.487	0.497	0.659	0.999
30 Out-of-covariate communes							
		Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Female	FH Trans.	0.130	0.188	0.213	0.227	0.226	0.501
	FH Bench.	0.136	0.198	0.223	0.238	0.237	0.525
Male	FH Trans.	0.346	0.383	0.431	0.444	0.499	0.721
	FH Bench.	0.360	0.398	0.448	0.462	0.519	0.750

Table 5 reports the distribution of estimated literacy rates for women and men in the communes in Senegal, split by in-sample, out-of-sample and out-of-covariate communes. Our first observation is that the estimates for female and male literacy rates are higher for the FH Bench compared to the FH Trans, respectively. The reason is that the aggregated FH Trans estimates (women: 36.1% and men: 57.7%) on national level slightly underestimate the national share of literates (women: 37.8% and men: 60.0%) and, thus, need a small adjustment to meet the national estimate for the country.

In order to judge the quality of the model-based FH Trans, we have a closer look to the Figures 5 and 6. In particular, Figure 5 represents the shrinkage factor for the female (left panel) and the male (right panel) model as well as the corresponding sample sizes (dashed lines). On the x -axis, communes are ordered by their sample size (descending order from left to right). We observe that for communes with larger sample size the direct estimator gets substantial weight for both models. In contrast, for communes with a smaller sample size the FH Trans tends to be highly synthetic. Comparing both models we note that the FH Trans for the female model puts in general more weight on the direct estimator than the male model - mean shrinkage factor: 0.262 (female) vs. 0.206 (male) - as a consequence of the larger sample size in the women's questionnaire (cf. Table 1). In Figure 6 we plot the direct (diamonds) and the FH Trans (dots) estimates of literacy rates against communes ordered by their sample size (descending order from left to right) for the male and female model (top down). For illustration, we show only every fourth commune in the figure. The estimated national literacy rate (based on the DHS survey) is represented by the solid line. The vertical lines show the confidence intervals for each commune. Note that we do not report variance estimates for the direct estimator because there was only one sampling cluster nested in

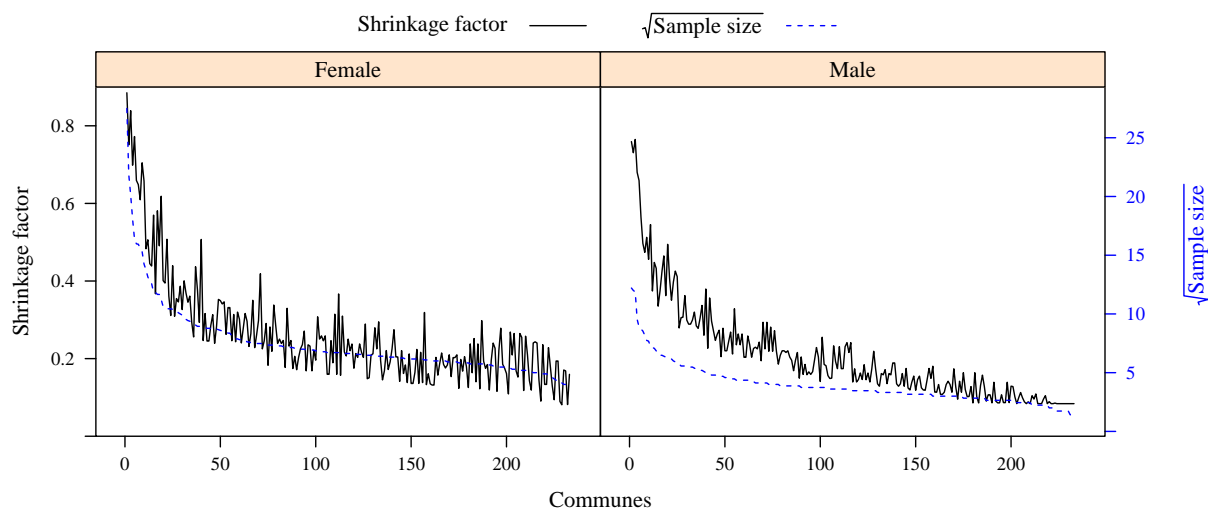


Figure 5: Shrinkage factor for the female model (left panel) and male model (right panel).

most of the communes. We observe that the FH Trans and the direct estimates randomly vary around the national estimate and do not indicate any systematic behaviour to show a possible bias from the modeling. Confirming the findings from Figure 5, the direct estimates are very similar to the model-based estimates for communes with a larger sample size. Most of the direct estimates are contained within the confidence intervals for both models. The length of the confidence interval are larger for the male model than for the female model. Due to the shrinkage the FH Trans estimates tend to be more stable around the national estimate than the direct estimator. The variability of the direct estimates and the length of the confidence intervals increase as we move from left to the right side of the figure.

4.3 Literacy rates by gender in Senegal

Having assessed the results of the estimators from a statistical perspective, we now discuss the results in the context of female and male literacy in Senegal. As the required approach for the ANSD should meet the third guideline which is that the aggregated estimates for the communes should produce the official national estimate for Senegal we focus in the following only on the benchmarked transformed FH. Figure 7 shows the estimates for literacy by gender on commune level for the capital Dakar (right panel) and for the rest of Senegal (left panel). In order to simplify the interpretation of the results, Figure 7 presents geographical maps for Dakar and for Senegal which are extracted from Google Maps. As a first comment, we note that the relative spatial distribution of male and female literacy rates are very similar in the Dakar region and in the rest of Senegal.

Having a closer look to the Dakar region (right panel) we observe that the coastal area, where the city of Dakar and its harbor are located, shows a very high rate of literates for male and female. This trend continues by moving from the peninsula closer to the main land and is only interrupted by a pocket of lower literacy around the district of Pikine (located to the east of the lake in the middle of Dakar). The district was founded in 1952 by the French colonial government for the former residents of the coastal area around the harbor. Since 1967, it is forbidden by law to build houses on this land because of problems with flooding. Today, however, illegal housings of migrant workers and refugees dominate this area, reflected in remarkably low literacy rates. Moving further into the interior of the country, the area gets more rural and the literacy rate shrinks.

We now turn to the estimated literacy rates for the rest of Senegal in Figure 7 (left panel). Next to

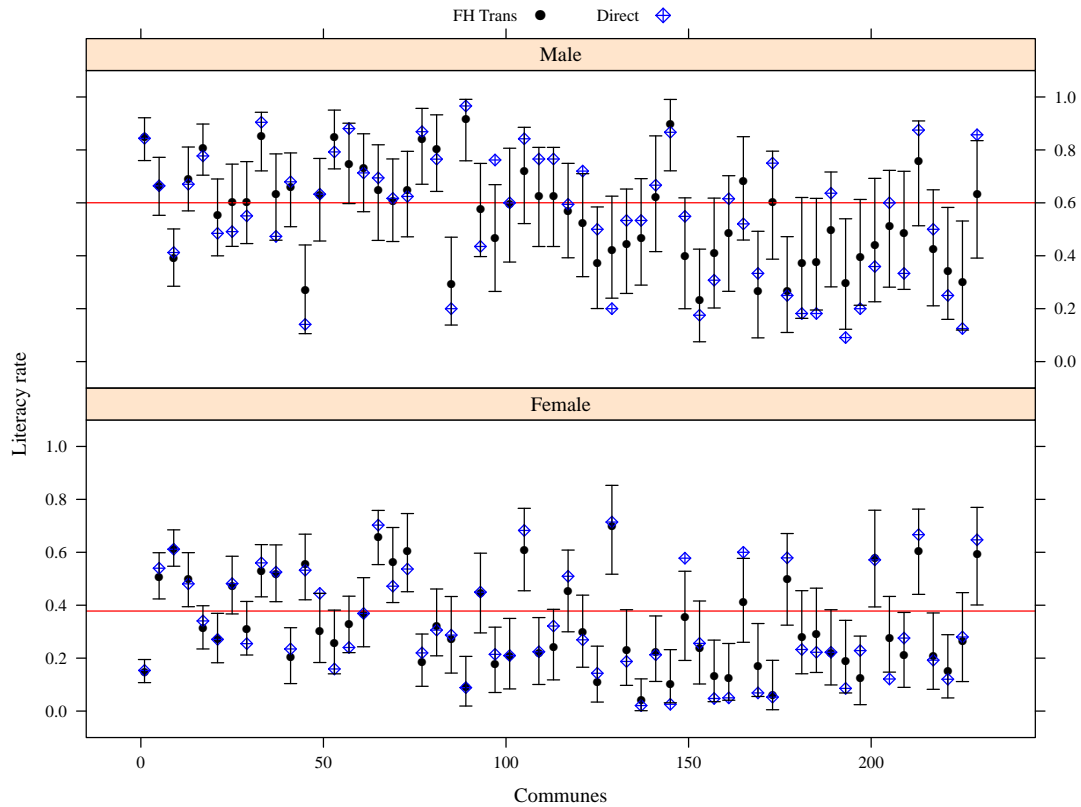


Figure 6: Coverage for the FH Trans for the male and female model (top down).

the Dakar region, the region around Ziguinchor below Gambia reveals a high literacy rate for men and women. The high literacy rates can be explained by the strategic position between the countries Guinea-Bissau and Gambia as well as to its closeness to the Atlantic Ocean. Ziguinchor is Senegal's second largest city and it is also the trade center of the Casamance region (area of Senegal south of Gambia including the Casamance river). Another reason is that the Casamance region is ethnically different from the other parts of Senegal. The region consists mainly of Jola people with a strong influence of Christianity whereas the Islam is the predominant religion in most other parts of the country (Heil, 2014). Another finding is that communes closer to the ocean and to borders in the North to Mauritania and in the South to Guinea-Bissau have higher literacy rates for men and women. In contrast, communes located on the borders to Mali (South-East) and to Gambia tend to have lower ones. As expected, the density of mobile phone towers in Figure 3 is higher in communes with higher literacy rates. Rural communes with a low coverage of mobile phone towers seem to have a lower literacy rates in general. Especially the central part of Senegal in the Matam and Tambacounda region reveals high shares of illiterate men and women.

Although the relative distribution is very similar in Senegal, Figure 7 reveals clear differences in terms of absolute values. The literacy rate for women is around 20% lower compared to men. Reasons are manifold in Senegal: Especially in poor regions of the country like Matam and Tambacounda in the eastern part of the country, girls are involved in economic activities and therefore the parents keep the girls out of the school to earn some additional income. Next to economic reasons, unsafe and long roads to school, gender-based violence, early marriage and pregnancy, the traditional role of women in the society and the low quality of the education system are further issues which add to low literacy rates for women. The PAJEF project, already mentioned in the introduction, aims to boost literacy among

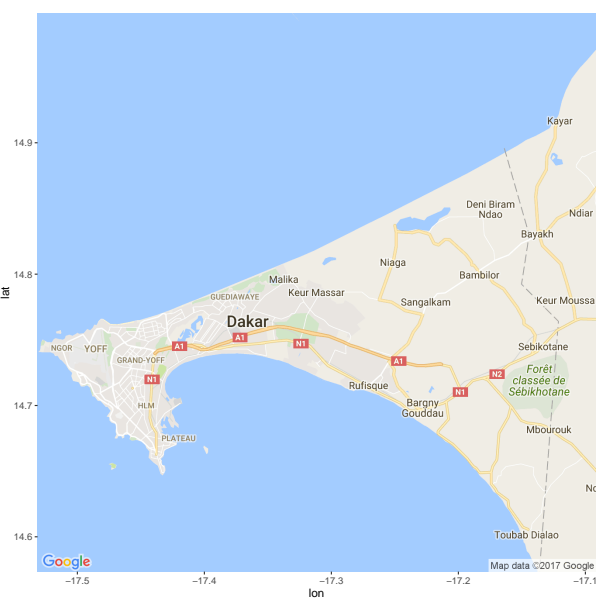
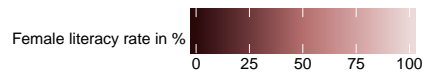
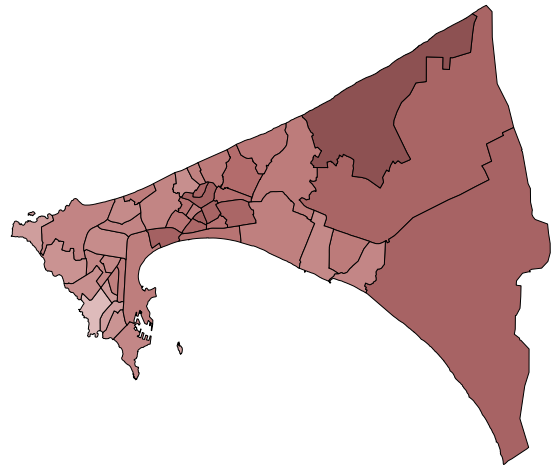
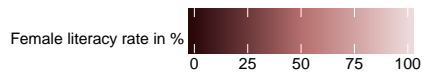
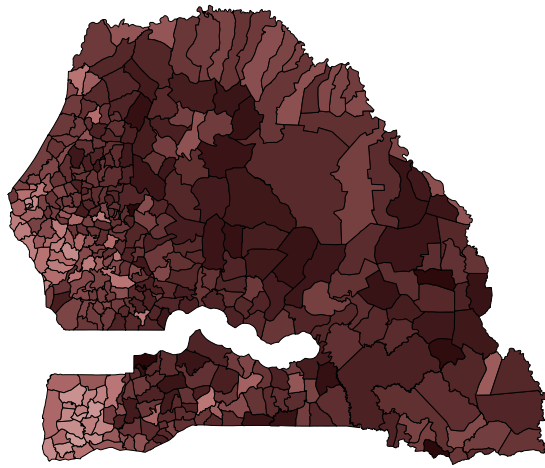
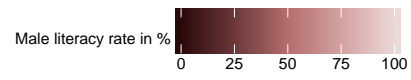
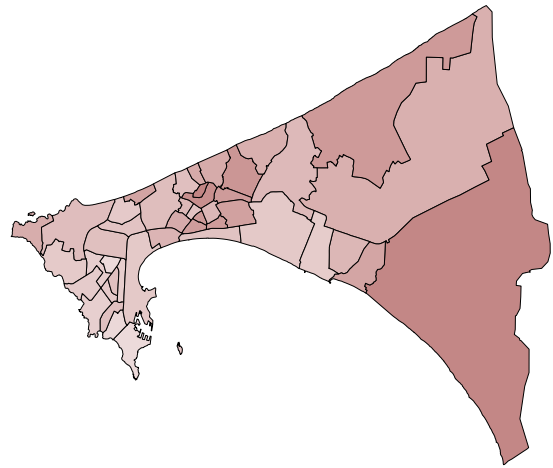
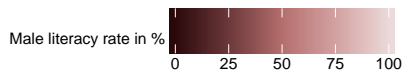
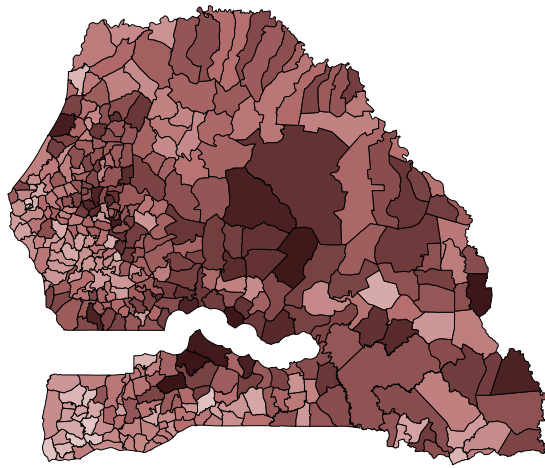


Figure 7: Estimates for the literacy rate by gender on commune level based on a benchmarked FH model: Senegal (left panel) and Dakar (right panel).

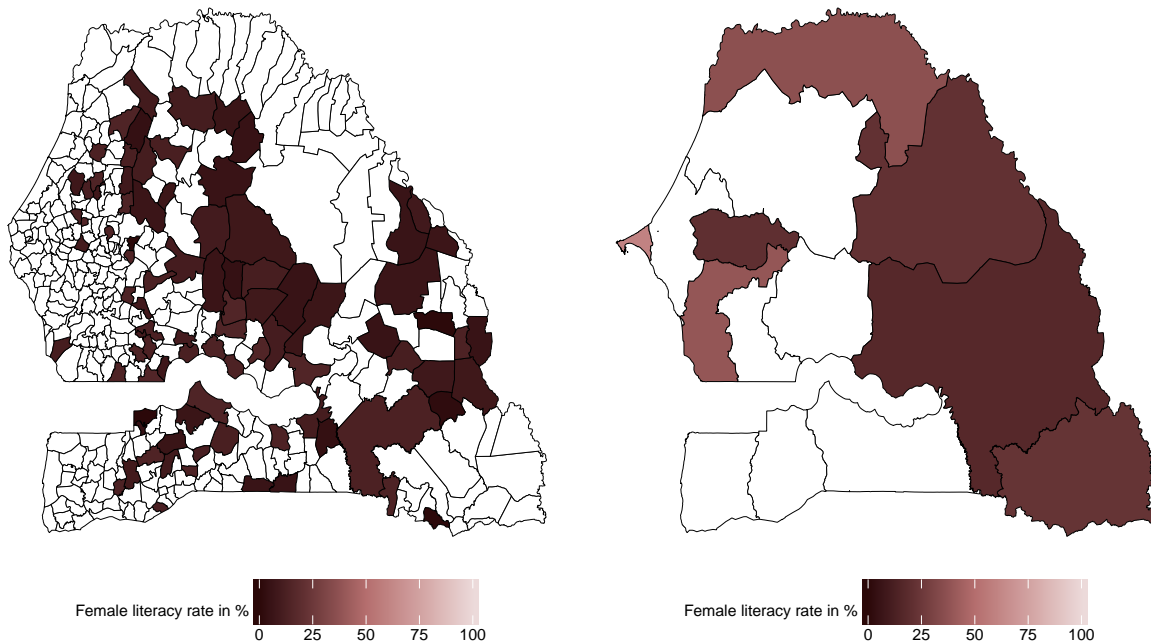


Figure 8: Estimates for the literacy rate for women: 20% of the communes with the lowest literacy rate (left panel) and seven regions in Senegal identified by the ANSD for the PAJEF project (right panel).

women in Senegal is currently conducted by UNESCO Dakar and the government of Senegal (UNESCO, 2015). The project runs in the seven regions (Dakar, Diourbel, Fatick, Kedougou, Matam, Saint-Louis and Tambacounda) with the lowest literacy rate identified by the ANSD based on the DHS survey. The seven regions and the corresponding literacy rates for women are displayed in Figure 8 (right panel). The regions cover around 50% of the country. The left figure shows the literacy rate for women on commune level held by the lowest 20% estimated by using the DHS survey 2011 in combination with mobile phone covariates. There are some hotspots for example in the region around Gambia in the Ziguinchor region or in the Western part of Senegal, with low literacy rates for women but without any financial support. In contrast, the PAJEF project provides financial support to the Saint-Louis region in the north of Senegal or to Dakar where the female literacy rates are above average.

Hence, the use of the proposed approach may enable NSIs and governmental organisations to make sound strategic decisions regarding the best places for investing in creating infrastructure for education. Figures for the indicators *no school education* or *secondary school education or higher* are available from the authors upon request.

5 Design-based simulation for unemployment

The analysis of literacy rates by gender in Section 4 was sample specific which makes conclusions about efficiency and bias difficult. In this section, we present results from a design-based simulation study that was carried out for assessing the performance of the introduced methodology we discussed in Section 3. The aim of the design-based simulation is to investigate the behaviour of the Fay-Herriot type models for estimating socio-demographic indicators based on mobile phone covariates in a controlled environment. In particular, for the evaluation of the approach we had access to the variable *unemployment* from the Senegalese register by collaborating with the staff of the ANSD.

The *pseudo* population in the design-based simulation is based on data collected from a sample of

around 1 million individuals in Senegal. The data was collected by ANSD as part of the census 2013 and is spread across the 431 communes. The *pseudo* population reflects around 10% of the population in Senegal. The variable of interest is defined by 0 = employed and 1 = unemployed. Summaries of the population sizes and unemployment rate over communes are given in Table 6. Given the fixed *pseudo*

Table 6: Summary statistics over communes

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA
Population size	82	717	1303	2257.0	2373	56670	-
Unemployment rate	0.274	0.488	0.550	0.555	0.617	0.898	-
Sample size	3	28	48	79.3	81	1448	235

population we independently drew $T = 500$ samples following a sampling design similar to the one of the DHS survey. The design is a stratified two-stage cluster sampling design, with the 431 communes as primary sampling units (PSUs). Similar to the DHS survey, we used 14 strata corresponding to the 14 regions of Senegal. In the first sampling stage we selected communes within each stratum with a probability proportional to their size. Around 2% of the individuals within each selected commune are drawn using equal probability systematic sampling. This leads to a sample size of around 15,543 individuals with 196 in-sample communes and 235 out-of-sample communes similar to the women’s questionnaire ($n = 15,688$) in the DHS survey (cf. Table 1). The summary statistics of the sample sizes over communes are also provided in Table 6.

We investigate the estimators presented in Section 3 under repeated sampling performance for the unemployment rate on commune level in Senegal using aggregated mobile phone covariates. To do so, we used an area-level linear mixed model (6). The covariates were selected by using the Bayesian information criterion (BIC) and held fixed for the simulation study. In particular, we considered only data from communes with a sample size of more than 30. Like in the application, we implicitly assume that for these communes the sampling variances of the direct estimators are negligible and standard regression model selection tools are applicable. We refer to Jiang et al. (2001) and Ha et al. (2014) for a similar approach for the model selection. The adjusted R^2 by Lahiri and Suntornchost (2015) was on average around 47% depending on the selected sample. We evaluate four estimators for the unemployment rate in the communes in the simulation. These are the direct estimator (3), the transformed FH estimator based on an arcsine transformation (11) (FH Trans) as well as the normal-logistic model (NL) and the normal-logistic random sampling variance model (NLRS) proposed by Ha et al. (2014). The direct estimator and FH Trans are implemented by computationally efficient algorithms using R. The NL and NLRS are implemented by using JAGS with three parallel chains, each with 20000 iterations, a burn-in of 10000 and the samples were thinned by a factor of two (Liu et al., 2014). The codes are available from the authors upon request. Additionally, we also assess the benchmarked transformed FH estimator (17) (FH Bench), the benchmarked NL estimator (NL Bench) and the benchmarked NLRS estimator (NLRS Bench). Note that we also use the *naive* benchmarking approach introduced in Section 3.3 for the NL and NLRS estimators.

The performance of the estimators is assessed by the bias (Bias) and root mean squared errors

(RMSE) given by

$$\text{Bias}(\hat{m}_i) = \frac{1}{T} \sum_{t=1}^T (\hat{m}_{ti} - m_i)$$

$$\text{RMSE}(\hat{m}_i) = \sqrt{\frac{1}{T} \sum_{t=1}^T (\hat{m}_{ti} - m_i)^2},$$

where \hat{m}_i is a generic notation to denote an estimator of the share in commune i and m_i denotes the true population share in commune i .

The results presented in Table 7 are splitted by the 191 in-sample, the 210 out-of-sample and the 30 out-of-covariate communes. The table reports summary statistics of the RMSE and Bias of the estimators (FH Trans, NL, and NLRs) over communes. The results confirm our expectations regarding the performance of the estimators. The direct estimator is almost unbiased but suffers from a higher RMSE compared to the model-based approaches (FH Trans, NL, and NLRs) for the sampled communes. The performance of the FH Trans and NLRs is very comparable regarding Bias and RMSE for the in-sample and out-of-sample communes and outperforms the NL estimator in this particular simulation study. For the out-of-covariate communes, where the covariates are obtained by geographically weighting as described in Section 2, all model-based estimators (FH Trans, NL, and NLRs) reveal on average a small positive bias.

In order to save space, the corresponding results for the benchmarked estimators (FH Bench, NL Bench, and NLRs Bench) are only reported in the supplementary materials. However, the results of the benchmarked estimators (FH Bench, NL Bench, and NLRs Bench) and the non-benchmarked estimators (FH Trans, NL, and NLRs) are close in terms of Bias and RMSE because the average of the commune level estimates required only a small adjustment to meet the national estimate for the country. However, note that the benchmarked and the non-benchmarked results are not directly comparable as the FH Bench, NL Bench, and NLRs Bench fulfil the benchmarking constraint.

The results from the study indicate that combining mobile phone covariates with survey data based on model-based estimators can lead i) to gains in efficiency compared to the direct estimator and ii) to reasonable results for communes with zero sample sizes.

6 Concluding remarks

Modern systems of official statistics require reliable statistics on socio-demographic indicators on regionally disaggregated levels. These statistics are essential for sound evidence-based policymaking. In this paper we have discussed an easy-applicable approach for NSIs for estimating these indicators by small area methods based on survey data and covariates from alternative data sources. The motivation is to reduce the dependence on census or register information for the NSIs. In particular, we used in this paper passively collected mobile phone data in combination with survey data to predict socio-demographic indicators. Although the paper focuses on literacy rates as specific socio-demographic indicator, the proposed approach is applicable to general indicators. For instance, we can provide results for two other indicators for women in Senegal: i) *Body mass index below 18.5* and ii) *Current usage of any contraception method*. One interesting approach for further research would be to predict the indicators purely on the mobile phone data and to further reduce the dependency of NSIs on actively collected data like survey or register data. For instance, Blumenstock et al. (2015) predicted poverty by using an individual's past

Table 7: Performance of predictors over communes in design-based simulations

191 In-sample communes							
Indicator	Estimator	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
RMSE	Direct	0.014	0.053	0.068	0.077	0.091	0.287
	FH Trans	0.016	0.030	0.042	0.053	0.071	0.248
	NL	0.013	0.040	0.049	0.055	0.061	0.260
	NLRS	0.014	0.030	0.041	0.054	0.070	0.251
Bias	Direct	-0.019	-0.002	-0.000	0.000	0.002	0.019
	FH Trans	-0.203	-0.029	0.001	0.001	0.029	0.247
	NL	-0.106	-0.011	0.003	0.009	0.026	0.169
	NLRS	-0.210	-0.029	0.001	0.002	0.029	0.250
210 Out-of-sample communes							
		Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
RMSE	FH Trans	0.012	0.030	0.053	0.073	0.104	0.349
	NL	0.011	0.035	0.062	0.076	0.103	0.325
	NLRS	0.010	0.030	0.057	0.073	0.101	0.349
Bias	FH Trans	-0.349	-0.044	0.011	0.007	0.062	0.245
	NL	-0.324	-0.043	0.012	0.013	0.065	0.281
	NLRS	-0.349	-0.043	0.008	0.008	0.059	0.247
30 Out-of-covariate communes							
		Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
RMSE	FH Trans	0.010	0.033	0.058	0.079	0.098	0.276
	NL	0.011	0.046	0.075	0.085	0.098	0.280
	NLRS	0.010	0.034	0.059	0.080	0.094	0.279
Bias	FH Trans	-0.173	-0.003	0.043	0.036	0.090	0.276
	NL	-0.153	-0.003	0.048	0.044	0.093	0.280
	NLRS	-0.170	-0.002	0.038	0.036	0.083	0.279

history of mobile phone usage in combination with a phone survey. One could extend these results to different indicators. Furthermore, mobile phone data can be used to update the small area estimates in the period between surveys. This would save considerable money, but would require additional assumptions about the model remaining constant between surveys.

For the combination of the survey data and the mobile phone covariates we used an easy-applicable FH small area method for the modeling. Additionally, we have investigated more complex extensions like the spatial FH (Pratesi and Salvati, 2009), the non-parametric FH (Giusti et al., 2012) and the spatial non-stationary FH (Chandra et al., 2015), but the results were comparable. One limitation of our modeling is the approximation of the sampling variance by $1/(4\tilde{n}_i)$ for the transformed FH estimator. As the derivations of the arcsine transformations are based on large-sample theory, we might expect some deficiencies, especially for communes with small sample sizes. Another limitation of our modeling is a potential back-transformation bias in (11) due to the nonlinearity of the arcsine transformation. We noted in Section 4.2 that the aggregation of $\hat{\theta}_i^{FH,trans}$ estimates is slightly lower than the national estimate for the male and female model, leading to an upward adjustment from benchmarking. Slud and Maiti (2006) discuss bias-corrected small area estimation formulas for the Fay-Herriot model in the context of a logarithmically transformed data. In case of an arcsine transformation, rather than back-

transforming the linear model predictions, one could calculate an additive adjustment for the bias as $E(\hat{\theta}_i^{direct}) - E(\hat{\theta}_i^{FH,trans})$, where $E(\cdot)$ is the unconditional expectation under the model. One additional line of research might be to explore the above mentioned bias correction and to extend the MSE estimation to the adjusted predictors. Another line for further work could be to investigate machine learning approaches like random forest for the prediction of socio-demographic indicators and compare them with small area methods.

We have also presented first discussions regarding the time-intensive cleaning, processing and handling of the mobile phone data and available software. However, this can be only a first step in this direction. From a long-run perspective it is necessary to build platforms with open software/ algorithms for NSIs. The aim of such platforms can be twofold: first, NSIs can use code and software to work with large data sources and, second, NSIs can potentially access passively collected data of private companies in a safe environment.

The use of mobile phone covariates has some drawbacks as well. First, additional uncertainty in the mobile phone data arises from the fact that the coverage of the mobile phone tower differs and is unknown. To the best of our knowledge we are not aware of an established way to handle a potential overlap of tower coverage. Second, landlines and the use of internet-based mobile communication services such as Skype, WhatsApp or Viber may cause distortion in communication patterns. However, for Senegal the distortions may be less strongly because of a stagnating landline penetration rate of 2.8% (GSMA, 2015). In addition, the all-time downloads of messaging applications are extremely low in Senegal compared to other countries (e.g. WhatsApp 124,818 and Viber 95,891 on iOS as of December 18th 2014 - extracted from Priori Data). Nevertheless, some types of users may systematically be excluded. Modelling these users is another avenue for further research.

Acknowledgements

First of all, the authors are indebted to the Editor, Associate Editor and referees for comments that significantly improved the paper. In addition, we would like to thank the staff of the Senegalese National Statistics Institute ANSD, especially the Directeur-Général Mr. Aboubacar Sedikh Beye and the Directeur du Management de l'Information Statistique Mr. Mamadou Niang for many fruitful discussions, country insights and data access. Further, we would like to thank Nicolas de Cordes and the Orange Group for the tremendous help with the coordination and logistics of the project as well as Sonatel for providing the call detail records. The project was made possible with the funding of the Bill and Melinda Gates Foundation. The work of Salvati has been developed under the support of the project PRIN-SURWEY <http://www.sp.unipg.it/surwey/> (Grant 2012F42NS8, Italy).

References

- ANSD (2012). Demographic and Health Survey - Multiple Indicator Cluster Survey (2010-2011). https://dhsprogram.com/pubs/pdf/FR258/FR258_English.pdf. Accessed: 2016-02-16.
- Bell, W. R. and C. Franco (2015). Borrowing information over time in binomial/logit normal models for small area estimation. *Statistics in Transition new series* 16(4), 563–584.

- Blumenstock, J., G. Cadamuro, and R. On (2015). Predicting poverty and wealth from mobile phone metadata. *Science* 350, 1073–1076.
- Carter, G. and J. Rolph (1974). Empirical bayes methods applied to estimating fire alarm probabilities. *Journal of the American Statistical Association* 69 (348), 880–885.
- Casas-Cordero, C., J. Encina, and P. Lahiri (2016). *Analysis of Poverty Data by Small Area Estimation*, Chapter Poverty Mapping for the Chilean Comunas, pp. 379–403. Wiley.
- Chambers, R., N. Salvati, and N. Tzavidis (2016). Semiparametric small area estimation for binary outcomes with application to unemployment estimation for local authorities in the uk. *Journal of the Royal Statistical Society: Series A* 179 (2), 453–479.
- Chandra, H., N. Salvati, and R. R. Chambers (2015). A spatially nonstationary fay-herriot model for small area estimation. *Journal of Survey Statistics and Methodology*, 109–135.
- Chatterjee, S., P. Lahiri, and H. Li (2008). Parametric bootstrap approximation to the distribution of eblup and related prediction intervals in linear mixed models. *The Annals of Statistics* 36 (3), 1221–1245.
- Datta, G. S., M. Ghosh, R. Steorts, and J. Maples (2010). Bayesian benchmarking with applications to small area estimation. *TEST* 20 (3), 574–588.
- Datta, G. S. and P. Lahiri (2000). A unified measure of uncertainty of estimated best linear unbiased predictors in small area estimation problems. *Statistica Sinica* 10, 613–627.
- Datta, G. S., J. N. K. Rao, and D. D. Smith (2005). On measuring the variability of small area estimators under a basic area level model. *Biometrika* 92 (1), 183–196.
- Deville, P., C. Linard, S. Martin, M. Gilbert, F. R. Stevens, A. E. Gaughan, V. D. Blondel, and A. J. Tatem (2014). Dynamic population mapping using mobile phone data. *Proceedings of the National Academy of Sciences* 111 (45), 15888–15893.
- Dowle, M., T. Short, S. Lianoglou, and A. Srinivasan (2014). *data.table: Extension of Data.frame*. R package version 1.9.6.
- Eagle, N., M. Macy, and R. Claxton (2010). Network diversity and economic development. *Science* 328, 1029–1031.
- Elbers, C., J. Lanjouw, and P. Lanjouw (2003). Micro-level estimation of poverty and inequality. *Econometrica* 71, 355–364.
- Fay, R. E. and R. A. Herriot (1979). Estimation of income for small places: An application of james-stein procedures to census data. *Journal of the American Statistical Association* 74 (366), 269–277.
- Ford, L. (2007). People treat you as if you are nothing: What good is free education when poverty means children are forced into work? <http://www.theguardian.com/education/2007/aug/21/schoolsworldwide.schools>. Accessed: 2016-04-04.
- Ghosh, M. and J. N. K. Rao (1994). Small area estimation: An appraisal. *Statistical Science* 9 (1), 55–93.
- Giusti, C., S. Marchetti, M. Pratesi, and N. Salvati (2012). Semiparametric fay-herriot model using penalized splines. *Journal of the Indian Society of Agricultural Statistics* 66, 1–14.

- GSMA (2015). The mobile economy 2015. <http://www.gsamobileeconomy.com/GSMA\Global\Mobile\Economy\Report\2015.pdf>. Accessed: 2016-01-30.
- Ha, N. S., P. Lahiri, and V. Parsons (2014). Methods and results for small area estimation using smoking data from the 2008 national health interview survey. *Statistics in Medicine* 33 (22), 3932–3945.
- Heil, T. (2014). Are neighbours alike? Practices of conviviality in Catalonia and Casamance. *European Journal of Cultural Studies* 17 (4), 452–470.
- Henderson, C. R. (1950). Estimation of genetic parameters. *Annals of Mathematical Statistics* 21 (2), 309–310.
- Jiang, J. and P. Lahiri (2001). Empirical best prediction for small area inference with binary data. *Annals of Institute of Statistical Mathematics* 53, 217–243.
- Jiang, J., P. Lahiri, S.-M. Wan, and C.-H. Wu (2001). Jackknifing in the fay–herriot model with an example. In *Proceedings of the Seminar on Funding Opportunity in Survey Research*, Bureau of Labor Statistics, 2 Massachusetts Avenue, NE, Washington, DC, pp. 75–97.
- Lahiri, P. and J. Suntornchost (2015). Variable selection for linear mixed models with applications in small area estimation. *Sankhya B* 77 (2), 312–320.
- Li, H. and P. Lahiri (2010). An adjusted maximum likelihood method for solving small area estimation problems. *Journal of Multivariate Analysis* 101, 882–892.
- Liu, B., P. Lahiri, and G. Kalton (2014). Hierarchical bayes modeling of survey-weighted small area proportions. *Survey Methodology* 40 (1), 1–13.
- Lopez-Vizcaino, E., M. J. Lombardia, and D. Morales (2015). Small area estimation of labour force indicators under a multinomial model with correlated time and area effects. *Journal of the Royal Statistical Society: Series A* 178 (3), 535–565.
- Marchetti, S., C. Giusti, M. Pratesi, N. Salvati, F. Giannotti, D. Pedreschi, S. Rinzivillo, L. Pappalardo, and L. Gabrielli (2015). Small area model-based estimations using big data sources. *Journal of Official Statistics* 31 (2), 263–281.
- Marhuenda, Y., D. Morales, and M. del Carmen Pardo (2014). Information criteria for fay-herriot model selection. *Computational Statistics & Data Analysis* 70, 268 – 280.
- Molina, I. and J. N. K. Rao (2010). Small area estimation of poverty indicators. *The Canadian Journal of Statistics* 38 (3), 369–385.
- Montjoye, Y.-A., J. Quoidbach, F. Robic, and A. S. Pentland (2013). *Social Computing, Behavioral-Cultural Modeling and Prediction: 6th International Conference, SBP 2013, Washington, DC, USA, April 2-5, 2013. Proceedings*, Chapter Predicting Personality Using Novel Mobile Phone-Based Metrics, pp. 48–55. Berlin, Heidelberg: Springer Berlin Heidelberg.
- Montjoye, Y.-A. d., Z. Smoreda, R. Trinquart, C. Ziemlicki, and V. Blondel (2014). D4d-senegal: The second mobile phone data for development challenge. *Working paper arXiv:1407.4885*.
- Opsomer, J. D., M. Francisco-Fernandez, and X. Li (2012). Model-based non-parametric variance estimation for systematic sampling. *Scandinavian Journal of Statistics* 39 (3), 528–542.

- Pfeffermann, D. (2013). New important developments in small area estimation. *Statistical Science* 28 (1), 40–68.
- Porter, A. T., S. H. Holan, C. K. Wikle, and N. Cressie (2014). Spatial fay-herriot models for small area estimation with functional covariates. *Spatial Statistics* 10, 27–42.
- Prasad, N. G. N. and J. N. K. Rao (1990). The estimation of the mean squared error of small area estimators. *Journal of the American Statistical Association* 85 (409), 163–171.
- Pratesi, M. (Ed.) (2016). *Analysis of Poverty Data by Small Area Estimation*. John Wiley & Sons.
- Pratesi, M. and N. Salvati (2009). Small area estimation in the presence of correlated random area effects. *Journal of Official Statistics* 25 (1), 37–53.
- Raghunathan, T. E., D. Xie, N. Schenker, V. L. Parsons, W. W. Davis, K. W. Dodd, and E. J. Feuer (2007). Combining information from two surveys to estimate county-level prevalence rates of cancer risk factors and screening. *Journal of the American Statistical Association* 102 (478), 474–486.
- Rao, J. N. K. and I. Molina (2015). *Small Area Estimation* (2nd Edition ed.). New York: Wiley.
- Schelle, C. (2013). *Schulsysteme, Unterricht und Bildung im mehrsprachigen frankophonen Westen und Norden Afrikas*. Münster, New York, München and Berlin: Waxmann.
- Searle, S. R. (1971). *Linear Models*. New York: Wiley.
- Slud, E. and T. Maiti (2006). Mean-squared error estimation in transformed fay-herriot models. *Journal of the Royal Statistical Society: Series B* 68 (2), 239–257.
- UNESCO (2012). Global partnership for girls' and women's education. http://www.unesco.org/eri/cp/factsheets_ed/SN_EDFactSheet.pdf. Accessed: 2016-03-04.
- UNESCO (2015). Literacy project for girls and women in senegal. <http://www.unesco.org/ui1/litbase/?menu=4&programme=180>. Accessed: 2016-04-04.
- Villalon, L. A. and M. Bodian (2012). Religion, social demand, and educational reforms in senegal. <http://r4d.dfid.gov.uk/PDF/Outputs/APPP/appp-research-report5-avril-2012.pdf>. Accessed: 2016-01-03.
- Yoshimori, M. and P. Lahiri (2014). A new adjusted maximum likelihood method for the fay-herriot small area model. *Journal of Multivariate Analysis* 124, 281–294.